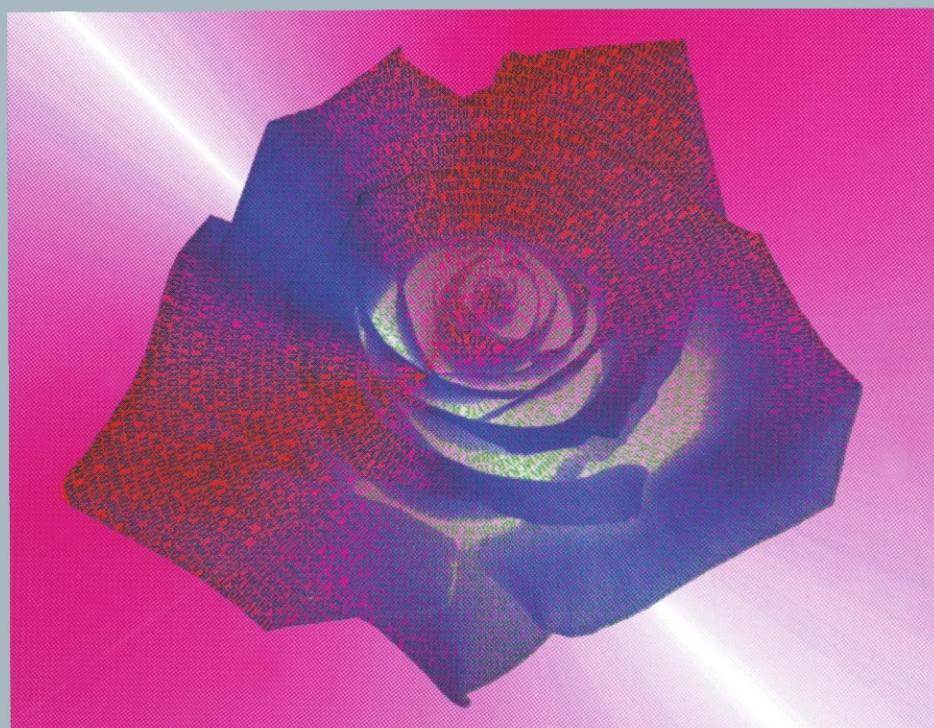


# An Introduction to Kolmogorov Complexity and Its Applications



Ming Li  
Paul Vitányi

SECOND EDITION



Springer Science+Business Media, LLC

# **GRADUATE TEXTS IN COMPUTER SCIENCE**

---

*Editors*

David Gries

Fred B. Schneider

**Springer Science+Business Media, LLC**

## **GRADUATE TEXTS IN COMPUTER SCIENCE**

---

*Fitting*, First-Order Logic and Automated Theorem Proving

*Li and Vitányi*, An Introduction to Kolmogorov Complexity  
and Its Applications

*Nerode and Shore*, Logic for Applications

*Smith*, A Recursive Introduction to the Theory of Computation

*Socher-Ambrosius and Johann*, Deduction Systems



Ming Li Paul Vitányi

**AN INTRODUCTION TO**  
**KOLMOGOROV COMPLEXITY**  
**AND ITS APPLICATIONS**

Second Edition

With 41 Illustrations



Springer

Ming Li  
Department of Computer Science  
University of Waterloo  
Waterloo, Ontario  
N2L 3G1 Canada

Paul Vitányi  
Centrum voor Wiskunde en Informatica  
Kruislaan 413  
1098 SJ Amsterdam  
The Netherlands

*Series Editors*

David Gries  
Fred B. Schneider

Department of Computer Science  
Cornell University  
Upson Hall  
Ithaca, NY 14853-7501, USA

Library of Congress Cataloging-in-Publication Data  
Li, Ming.

An introduction to Kolmogorov complexity and its applications /  
Ming Li, Paul Vitányi. — 2nd ed.  
p. cm. — (Graduate texts in computer science)  
Includes bibliographical references (p. — ) and index.  
ISBN 978-1-4757-2608-4  
1. Kolmogorov complexity. I. Vitányi, P.M.B. II. Title.  
III. Series: Graduate texts in computer science (Springer-Verlag New  
York Inc.)  
QA267.7.L5 1997  
511.3—dc20

96-42357

Printed on acid-free paper.

© 1997, 1993 Springer Science+Business Media New York  
Originally published by Springer-Verlag New York, Inc. in 1997  
Softcover reprint of the hardcover 2nd edition 1997

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher  
(Springer Science+Business Media, LLC), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection  
with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology  
now known or hereafter developed is forbidden.

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

Production managed by Lesley Poliner; manufacturing supervised by Jacqui Ashri.  
Photocomposed copy prepared from the authors' L<sup>A</sup>T<sub>E</sub>X files.

9 8 7 6 5 4 3 2 1

ISBN 978-1-4757-2608-4 ISBN 978-1-4757-2606-0 (eBook)  
DOI 10.1007/978-1-4757-2606-0

---

## Preface to the First Edition

We are to admit no more causes of natural things (as we are told by Newton) than such as are both true and sufficient to explain their appearances. This central theme is basic to the pursuit of science, and goes back to the principle known as Occam's razor: "if presented with a choice between indifferent alternatives, then one ought to select the simplest one." Unconsciously or explicitly, informal applications of this principle in science and mathematics abound.

The conglomerate of different research threads drawing on an objective and absolute form of this approach appears to be part of a single emerging discipline, which will become a major applied science like information theory or probability theory. We aim at providing a unified and comprehensive introduction to the central ideas and applications of this discipline.

Intuitively, the amount of information in a finite string is the size (number of binary digits, or *bits*) of the shortest program that without additional data, computes the string and terminates. A similar definition can be given for infinite strings, but in this case the program produces element after element forever. Thus, a long sequence of 1's such as

$\underbrace{11111\dots1}_{\text{10,000 times}}$

contains little information because a program of size about  $\log 10,000$  bits outputs it:

```
for i := 1 to 10,000
    print 1
```

Likewise, the transcendental number  $\pi = 3.1415\dots$ , an infinite sequence of seemingly "random" decimal digits, contains but a few bits of information. (There is a short program that produces the consecutive digits of  $\pi$  forever.) Such a definition would appear to make the amount of information in a string (or other object) depend on the particular programming language used.

Fortunately, it can be shown that all reasonable choices of programming languages lead to quantification of the amount of "absolute" information in individual objects that is invariant up to an additive constant. We call this quantity the "Kolmogorov complexity" of the object. If an object contains regularities, then it has a shorter description than itself. We call such an object "compressible."

The application of Kolmogorov complexity takes a variety of forms, for example, using the fact that some strings are extremely compressible; using the compressibility of strings as a selection criterion; using the fact that many strings are not compressible at all; and using the fact that

---

some strings may be compressed, but that it takes a lot of effort to do so.

The theory dealing with the quantity of information in individual objects goes by names such as “algorithmic information theory,” “Kolmogorov complexity,” “K-complexity,” “Kolmogorov-Chaitin randomness,” “algorithmic complexity,” “stochastic complexity,” “descriptive complexity,” “minimum description length,” “program-size complexity,” and others. Each such name may represent a variation of the basic underlying idea or a different point of departure. The mathematical formulation in each case tends to reflect the particular traditions of the field that gave birth to it, be it probability theory, information theory, theory of computing, statistics, or artificial intelligence.

This raises the question about the proper name for the area. Although there is a good case to be made for each of the alternatives listed above, and a name like “Solomonoff-Kolmogorov-Chaitin complexity” would give proper credit to the inventors, we regard “Kolmogorov complexity” as well entrenched and commonly understood, and we shall use it hereafter.

The mathematical theory of Kolmogorov complexity contains deep and sophisticated mathematics. Yet one needs to know only a small amount of this mathematics to apply the notions fruitfully in widely divergent areas, from sorting algorithms to combinatorial theory, and from inductive reasoning and machine learning to dissipationless computing.

Formal knowledge of basic principles does not necessarily imply the wherewithal to apply it, perhaps especially so in the case of Kolmogorov complexity. It is our purpose to develop the theory in detail and outline a wide range of illustrative applications. In fact, while the pure theory of the subject will have its appeal to the select few, the surprisingly large field of its applications will, we hope, delight the multitude.

The mathematical theory of Kolmogorov complexity is treated in Chapters 2, 3, and 4; the applications are treated in Chapters 5 through 8. Chapter 1 can be skipped by the reader who wants to proceed immediately to the technicalities. Section 1.1 is meant as a leisurely, informal introduction and peek at the contents of the book. The remainder of Chapter 1 is a compilation of material on diverse notations and disciplines drawn upon.

We define mathematical notions and establish uniform notation to be used throughout. In some cases we choose nonstandard notation since the standard one is homonymous. For instance, the notions “absolute value,” “cardinality of a set,” and “length of a string,” are commonly denoted in the same way as  $|\cdot|$ . We choose distinguishing notations  $| \cdot |$ ,  $d(\cdot)$ , and  $l(\cdot)$ , respectively.

Briefly, we review the basic elements of computability theory and probability theory that are required. Finally, in order to place the subject in the appropriate historical and conceptual context we trace the main roots of Kolmogorov complexity.

This way the stage is set for Chapters 2 and 3, where we introduce the notion of optimal effective descriptions of objects. The length of such a description (or the number of bits of information in it) is its Kolmogorov complexity. We treat all aspects of the elementary mathematical theory of Kolmogorov complexity. This body of knowledge may be called *algorithmic complexity theory*. The theory of Martin-Löf tests for randomness of finite objects and infinite sequences is inextricably intertwined with the theory of Kolmogorov complexity and is completely treated. We also investigate the statistical properties of finite strings with high Kolmogorov complexity. Both of these topics are eminently useful in the applications part of the book. We also investigate the recursion-theoretic properties of Kolmogorov complexity (relations with Gödel's incompleteness result), and the Kolmogorov complexity version of information theory, which we may call “algorithmic information theory” or “absolute information theory.”

The treatment of *algorithmic probability theory* in Chapter 4 presupposes Sections 1.6, 1.11.2, and Chapter 3 (at least Sections 3.1 through 3.4). Just as Chapters 2 and 3 deal with the optimal effective description length of objects, we now turn to optimal (greatest) effective probability of objects. We treat the elementary mathematical theory in detail. Subsequently, we develop the theory of effective randomness tests under arbitrary recursive distributions for both finite and infinite sequences. This leads to several classes of randomness tests, each of which has a universal randomness test. This is the basis for the treatment of a mathematical theory of inductive reasoning in Chapter 5 and the theory of algorithmic entropy in Chapter 8.

Chapter 5 develops a general theory of inductive reasoning and applies the developed notions to particular problems of inductive inference, prediction, mistake bounds, computational learning theory, and minimum description length induction in statistics. This development can be viewed both as a resolution of certain problems in philosophy about the concept and feasibility of induction (and the ambiguous notion of “Occam’s razor”), as well as a mathematical theory underlying computational machine learning and statistical reasoning.

Chapter 6 introduces the incompressibility method. Its utility is demonstrated in a plethora of examples of proving mathematical and computational results. Examples include combinatorial properties, the time complexity of computations, the average-case analysis of algorithms such as Heapsort, language recognition, string matching, “pumping lemmas”

---

in formal language theory, lower bounds in parallel computation, and Turing machine complexity. Chapter 6 assumes only the most basic notions and facts of Sections 2.1, 2.2, 3.1, 3.3.

Some parts of the treatment of resource-bounded Kolmogorov complexity and its many applications in computational complexity theory in Chapter 7 presuppose familiarity with a first-year graduate theory course in computer science or basic understanding of the material in Section 1.7.4. Sections 7.5 and 7.7 on “universal optimal search” and “logical depth” only require material covered in this book. The section on “logical depth” is technical and can be viewed as a mathematical basis with which to study the emergence of life-like phenomena—thus forming a bridge to Chapter 8, which deals with applications of Kolmogorov complexity to relations between physics and computation.

Chapter 8 presupposes parts of Chapters 2, 3, 4, the basics of information theory as given in Section 1.11, and some familiarity with college physics. It treats physical theories like dissipationless reversible computing, information distance and picture similarity, thermodynamics of computation, statistical thermodynamics, entropy, and chaos from a Kolmogorov complexity point of view. At the end of the book there is a comprehensive listing of the literature on theory and applications of Kolmogorov complexity and a detailed index.

## How to Use This Book

The technical content of this book consists of four layers. The main text is the first layer. The second layer consists of examples in the main text. These elaborate the theory developed from the main theorems. The third layer consists of nonindented, smaller-font paragraphs interspersed with the main text. The purpose of such paragraphs is to have an explanatory aside, to raise some technical issues that are important but would distract attention from the main narrative, or to point to alternative or related technical issues. Much of the technical content of the literature on Kolmogorov complexity and related issues appears in the fourth layer, the exercises. When the idea behind a nontrivial exercise is not our own, we have tried to give credit to the person who originated the idea. Corresponding references to the literature are usually given in comments to an exercise or in the historical section of that chapter.

Starred sections are not really required for the understanding of the sequel and should be omitted at first reading. The application sections are not starred. The exercises are grouped together at the end of main sections. Each group relates to the material in between it and the previous group. Each chapter is concluded by an extensive historical section with full references. For convenience, all references in the text to the Kolmogorov complexity literature and other relevant literature are given in full where they occur. The book concludes with a References section intended as a separate exhaustive listing of the literature restricted to the

theory and the direct applications of Kolmogorov complexity. There are reference items that do not occur in the text and text references that do not occur in the References. We added a very detailed index combining the index to notation, the name index, and the concept index. The page number where a notion is defined first is printed in boldface. The initial part of the Index is an index to notation. Names as “J. von Neumann” are indexed European style “Neumann, J. von.”

The exercises are sometimes trivial, sometimes genuine exercises, but more often compilations of entire research papers or even well-known open problems. There are good arguments to include both: the easy and real exercises to let the student exercise his comprehension of the material in the main text; the contents of research papers to have a comprehensive coverage of the field in a small number of pages; and research problems to show where the field is (or could be) heading. To save the reader the problem of having to determine which is which: “I found this simple exercise in number theory that looked like Pythagoras’s Theorem. Seems difficult.” “Oh, that is Fermat’s Last Theorem; it was unsolved for three hundred and fifty years...,” we have adopted the system of *rating numbers* used by D.E. Knuth [*The Art of Computer Programming, Vol. 1: Fundamental Algorithms*, Addison-Wesley, 1973 (2nd Edition), pp. xvii–xix]. The interpretation is as follows:

- 00 A very easy exercise that can be answered immediately, from the top of your head, if the material in the text is understood.
- 10 A simple problem to exercise understanding of the text. Use fifteen minutes to think, and possibly pencil and paper.
- 20 An average problem to test basic understanding of the text and may take one or two hours to answer completely.
- 30 A moderately difficult or complex problem taking perhaps several hours to a day to solve satisfactorily.
- 40 A quite difficult or lengthy problem, suitable for a term project, often a significant result in the research literature. We would expect a very bright student or researcher to be able to solve the problem in a reasonable amount of time, but the solution is not trivial.
- 50 A research problem that, to the authors’ knowledge, is open at the time of writing. If the reader has found a solution, he is urged to write it up for publication; furthermore, the authors of this book would appreciate hearing about the solution as soon as possible (provided it is correct).

This scale is “logarithmic”: a problem of rating 17 is a bit simpler than average. Problems with rating 50, subsequently solved, will appear in

---

a next edition of this book with rating 45. Rates are sometimes based on the use of solutions to earlier problems. The rating of an exercise is based on that of its most difficult item, but not on the number of items. Assigning accurate rating numbers is impossible—one man's meat is another man's poison—and our rating will differ from ratings by others.

An orthogonal rating “M” implies that the problem involves more mathematical concepts and motivation than is necessary for someone who is primarily interested in Kolmogorov complexity and applications. Exercises marked “HM” require the use of calculus or other higher mathematics not developed in this book. Some exercises are marked with “●”; and these are especially instructive or useful. Exercises marked “O” are problems that are, to our knowledge, unsolved at the time of writing. The rating of such exercises is based on our estimate of the difficulty of solving them. Obviously, such an estimate may be totally wrong.

Solutions to exercises, or references to the literature where such solutions can be found, appear in the “Comments” paragraph at the end of each exercise. Nobody is expected to be able to solve all exercises.

The material presented in this book draws on work that until now was available only in the form of advanced research publications, possibly not translated into English, or was unpublished. A large portion of the material is new. The book is appropriate for either a one- or a two-semester introductory course in departments of mathematics, computer science, physics, probability theory and statistics, artificial intelligence, cognitive science, and philosophy. Outlines of possible one-semester courses that can be taught using this book are presented below.

Fortunately, the field of descriptive complexity is fairly young and the basics can still be comprehensively covered. We have tried to the best of our abilities to read, digest, and verify the literature on the topics covered in this book. We have taken pains to establish correctly the history of the main ideas involved. We apologize to those who have been unintentionally slighted in the historical sections. Many people have generously and selflessly contributed to verify and correct drafts of this book. We thank them below and apologize to those we forgot. In a work of this scope and size there are bound to remain factual errors and incorrect attributions. We greatly appreciate notification of errors or any other comments the reader may have, preferably by email to [kolmogorov@cwi.nl](mailto:kolmogorov@cwi.nl), in order that future editions may be corrected.

## Acknowledgments

We thank Greg Chaitin, Péter Gács, Leonid Levin, and Ray Solomonoff for taking the time to tell us about the early history of our subject and for introducing us to many of its applications. Juris Hartmanis and Joel Seiferas initiated us into Kolmogorov complexity in various ways.

Many people gave substantial suggestions for examples and exercises, or pointed out errors in a draft version. Apart from the people already mentioned, these are, in alphabetical order, Eric Allender, Charles Bennett, Piotr Berman, Robert Black, Ron Book, Dany Breslauer, Harry Buhrman, Peter van Emde Boas, William Gasarch, Joe Halpern, Jan Heering, G. Hotz, Tao Jiang, Max Kanovich, Danny Krizanc, Evangelos Kranakis, Michiel van Lambalgen, Luc Longpré, Donald Loveland, Albert Meyer, Lambert Meertens, Ian Munro, Pekka Orponen, Ramamohan Paturi, Jorma Rissanen, Jeff Shallit, A.Kh. Shen', J. Laurie Snell, Th. Tsantilas, John Tromp, Vladimir Uspensky, N.K. Vereshchagin, Osamu Watanabe, and Yaakov Yesha. Apart from them, we thank the many students and colleagues who contributed to this book.

We especially thank Péter Gács for the extraordinary kindness of reading and commenting in detail on the entire manuscript, including the exercises. His expert advice and deep insight saved us from many pitfalls and misunderstandings. Piergiorgio Odifreddi carefully checked and commented on the first three chapters. Parts of the book have been tested in one-semester courses and seminars at the University of Amsterdam in 1988 and 1989, the University of Waterloo in 1989, Dartmouth College in 1990, the Universitat Polytècnica de Catalunya in Barcelona in 1991/1992, the University of California at Santa Barbara, Johns Hopkins University, and Boston University in 1992/1993.

This document has been prepared using the L<sup>A</sup>T<sub>E</sub>X system. We thank Donald Knuth for T<sub>E</sub>X, Leslie Lamport for L<sup>A</sup>T<sub>E</sub>X, and Jan van der Steen at CWI for online help. Some figures were prepared by John Tromp using the xpic program.

The London Mathematical Society kindly gave permission to reproduce a long extract by A.M. Turing. The Indian Statistical Institute, through the editor of *Sankhyā*, kindly gave permission to quote A.N. Kolmogorov.

We gratefully acknowledge the financial support by NSF Grant DCR-8606366, ONR Grant N00014-85-k-0445, ARO Grant DAAL03-86-K-0171, the Natural Sciences and Engineering Research Council of Canada through operating grants OGP-0036747, OGP-046506, and International Scientific Exchange Awards ISE0046203, ISE0125663, and NWO Grant NF 62-376. The book was conceived in late Spring 1986 in the Valley of the Moon in Sonoma County, California. The actual writing lasted on and off from autumn 1987 until summer 1993.

One of us [PV] gives very special thanks to his lovely wife Pauline for insisting from the outset on the significance of this enterprise. The Aiken Computation Laboratory of Harvard University, Cambridge, Massachusetts, USA; the Computer Science Department of York University, Ontario, Canada; the Computer Science Department of the University

---

of Waterloo, Ontario, Canada; and CWI, Amsterdam, the Netherlands provided the working environments in which this book could be written.

## Preface to the Second Edition

---

When this book was conceived ten years ago, few scientists realized the width of scope and the power for applicability of the central ideas. Partially because of the enthusiastic reception of the first edition, open problems have been solved and new applications have been developed. We have added new material on the relation between data compression and minimum description length induction, computational learning, and universal prediction; circuit theory; distributed algorithmics; instance complexity; CD compression; computational complexity; Kolmogorov random graphs; shortest encoding of routing tables in communication networks; computable universal distributions; average case properties; the equality of statistical entropy and expected Kolmogorov complexity; and so on. Apart from being used by researchers and as reference work, the book is now commonly used for graduate courses and seminars. In recognition of this fact, the second edition has been produced in textbook style. We have preserved as much as possible the ordering of the material as it was in the first edition. The many exercises bunched together at the ends of some chapters have been moved to the appropriate sections. The comprehensive bibliography on Kolmogorov complexity at the end of the book has been updated, as have the “History and References” sections of the chapters. Many readers were kind enough to express their appreciation for the first edition and to send notification of typos, errors, and comments. Their number is too large to thank them individually, so we thank them all collectively.

## Outlines of One-Semester Courses

---

We have mapped out a number of one-semester courses on a variety of topics. These topics range from basic courses in theory and applications to special interest courses in learning theory, randomness, or information theory using the Kolmogorov complexity approach.

**PREREQUISITES:** Sections 1.1, 1.2, 1.7 (except Section 1.7.4).

### I. Course on Basic Algorithmic Complexity and Applications

TYPE OF COMPLEXITY	THEORY	APPLICATIONS
plain complexity	2.1, 2.2, 2.3	4.4, Chapter 6
prefix complexity	1.11.2, 3.1 3.3, 3.4	5.1, 5.1.3, 5.2, 5.5 8.2, 8.3 8
resource-bounded complexity	7.1, 7.5, 7.7	7.2, 7.3, 7.6, 7.7

## II. Course on Algorithmic Complexity

TYPE OF COMPLEXITY	BASICS	RANDOMNESS	ALGORITHMIC PROPERTIES
state × symbol	1.12		
plain complexity	2.1, 2.2, 2.3	2.4	2.7
prefix complexity	1.11.2, 3.1 3.3, 3.4	3.5	3.7, 3.8
monotone complexity	4.5 (intro)	4.5.4	

## III. Course on Algorithmic Randomness

RANDOMNESS TESTS ACCORDING TO	COMPLEXITY USED	FINITE STRINGS	INFINITE SEQUENCES
von Mises			1.9
Martin-Löf	2.1, 2.2	2.4	2.5
prefix complexity	1.11.2, 3.1, 3.3, 3.4	3.5	3.6, 4.5.6
general discrete	1.6 (intro), 4.3.1	4.3	
general continuous	1.6 (intro), 4.5 (intro), 4.5.1		4.5

## IV. Course on Algorithmic Information Theory and Applications

TYPE OF COMPLEXITY USED	BASICS	ENTROPY	SYMMETRY OF INFORMATION
classical information theory	1.11	1.11	1.11
plain complexity	2.1, 2.2	2.8	2.8
prefix complexity	3.1, 3.3, 3.4		3.8, 3.9.1
resource-bounded	7.1		Exercises 7.1.11 7.1.12
applications		8.1, 8.4, 8.5	Theorem 7.2.6 Exercise 6.10.15

## V. Course on Algorithmic Probability Theory, Learning, Inference and Prediction

THEORY	BASICS	UNIVERSAL DISTRIBUTION	APPLICATIONS TO INFERENCE
classical probability	1.6, 1.11.2		1.6
algorithmic complexity	2.1, 2.2, 2.3 3.1, 3.3, 3.4		8
algorithmic discrete probability	4.2, 4.1 4.3 (intro)	4.3.1, 4.3.2 4.3.3, 4.3.4, 4.3.6	
algorithmic contin. probability	4.5 (intro)	4.5.1, 4.5.2 4.5.4, 4.5.8	5.2
Solomonoff's inductive inference	5.1, 5.1.3, 5.2 5.4, 8	5.3, 5.4.3, 5.5 5.5.8	5.1.3

**VI. Course on  
the  
Incompressibility  
Method**

Chapter 2 (Sections 2.1, 2.2, 2.4, 2.6, 2.8), Chapter 3 (mainly Sections 3.1, 3.3), Section 4.4, and Chapters 6 and 7. The course covers the basics of the theory with many applications in proving upper and lower bounds on the running time and space use of algorithms.

**VII. Course on  
Randomness,  
Information, and  
Physics**

Course III and Chapter 8. In physics the applications of Kolmogorov complexity include theoretical illuminations of foundational issues. For example, the approximate equality of statistical entropy and expected Kolmogorov complexity, the nature of “entropy,” a fundamental resolution of the “Maxwell’s Demon” paradox. However, also more concrete applications like “information distance” and “thermodynamics of computation” are covered.

---

# Contents

Preface to the First Edition . . . . .	v
How to Use This Book . . . . .	viii
Acknowledgments . . . . .	x
Preface to the Second Edition . . . . .	xii
Outlines of One-Semester Courses . . . . .	xii
List of Figures . . . . .	xix
<b>1 Preliminaries</b>	<b>1</b>
1.1 A Brief Introduction . . . . .	1
1.2 Prerequisites and Notation . . . . .	6
1.3 Numbers and Combinatorics . . . . .	8
1.4 Binary Strings . . . . .	12
1.5 Asymptotic Notation . . . . .	15
1.6 Basics of Probability Theory . . . . .	18
1.7 Basics of Computability Theory . . . . .	24
1.8 The Roots of Kolmogorov Complexity . . . . .	47
1.9 Randomness . . . . .	49
1.10 Prediction and Probability . . . . .	59
1.11 Information Theory and Coding . . . . .	65
1.12 State × Symbol Complexity . . . . .	84
1.13 History and References . . . . .	86

<b>2 Algorithmic Complexity</b>	<b>93</b>
2.1 The Invariance Theorem . . . . .	96
2.2 Incompressibility . . . . .	108
2.3 $C$ as an Integer Function . . . . .	119
2.4 Random Finite Sequences . . . . .	127
2.5 *Random Infinite Sequences . . . . .	136
2.6 Statistical Properties of Finite Sequences . . . . .	158
2.7 Algorithmic Properties of $C$ . . . . .	167
2.8 Algorithmic Information Theory . . . . .	179
2.9 History and References . . . . .	185
<b>3 Algorithmic Prefix Complexity</b>	<b>189</b>
3.1 The Invariance Theorem . . . . .	192
3.2 *Sizes of the Constants . . . . .	197
3.3 Incompressibility . . . . .	202
3.4 $K$ as an Integer Function . . . . .	206
3.5 Random Finite Sequences . . . . .	208
3.6 *Random Infinite Sequences . . . . .	211
3.7 Algorithmic Properties of $K$ . . . . .	224
3.8 *Complexity of Complexity . . . . .	226
3.9 *Symmetry of Algorithmic Information . . . . .	229
3.10 History and References . . . . .	237
<b>4 Algorithmic Probability</b>	<b>239</b>
4.1 Enumerable Functions Revisited . . . . .	240
4.2 Nonclassical Notation of Measures . . . . .	242
4.3 Discrete Sample Space . . . . .	245
4.4 Universal Average-Case Complexity . . . . .	268
4.5 Continuous Sample Space . . . . .	272
4.6 Universal Average-Case Complexity, Continued . . . . .	307
4.7 History and References . . . . .	307

---

<b>5 Inductive Reasoning</b>	<b>315</b>
5.1 Introduction . . . . .	315
5.2 Solomonoff's Theory of Prediction . . . . .	324
5.3 Universal Recursion Induction . . . . .	335
5.4 Simple Pac-Learning . . . . .	339
5.5 Hypothesis Identification by Minimum Description Length	351
5.6 History and References . . . . .	372
<b>6 The Incompressibility Method</b>	<b>379</b>
6.1 Three Examples . . . . .	380
6.2 High-Probability Properties . . . . .	385
6.3 Combinatorics . . . . .	389
6.4 Kolmogorov Random Graphs . . . . .	396
6.5 Compact Routing . . . . .	404
6.6 Average-Case Complexity of Heapsort . . . . .	412
6.7 Longest Common Subsequence . . . . .	417
6.8 Formal Language Theory . . . . .	420
6.9 Online CFL Recognition . . . . .	427
6.10 Turing Machine Time Complexity . . . . .	432
6.11 Parallel Computation . . . . .	445
6.12 Switching Lemma . . . . .	449
6.13 History and References . . . . .	452
<b>7 Resource-Bounded Complexity</b>	<b>459</b>
7.1 Mathematical Theory . . . . .	460
7.2 Language Compression . . . . .	476
7.3 Computational Complexity . . . . .	488
7.4 Instance Complexity . . . . .	495
7.5 $Kt$ Complexity and Universal Optimal Search . . . . .	502
7.6 Time-Limited Universal Distributions . . . . .	506
7.7 Logical Depth . . . . .	510
7.8 History and References . . . . .	516

<b>8 Physics, Information, and Computation</b>	<b>521</b>
8.1 Algorithmic Complexity and Shannon's Entropy . . . . .	522
8.2 Reversible Computation . . . . .	528
8.3 Information Distance . . . . .	537
8.4 Thermodynamics . . . . .	554
8.5 Entropy Revisited . . . . .	565
8.6 Compression in Nature . . . . .	583
8.7 History and References . . . . .	586
<b>References</b>	<b>591</b>
<b>Index</b>	<b>618</b>

---

# List of Figures

1.1	Turing machine	28
1.2	Inferred probability for increasing $n$	60
1.3	Binary tree for $E(1) = 0, E(2) = 10, E(3) = 110,$ $E(4) = 111$	72
1.4	Binary tree for $E(1) = 0, E(2) = 01, E(3) = 011,$ $E(4) = 0111$	73
2.1	The graph of the integer function $C(x)$	121
2.2	The graph of the integer function $C(x l(x))$	123
2.3	Test of Example 2.4.1	128
2.4	Complexity oscillations of initial segments of infinite high-complexity sequences	139
2.5	Three notions of “chaotic” infinite sequences	148
3.1	The 425-bit universal combinator $U'$ in pixels	201
3.2	The graphs of $K(x)$ and $K(x l(x))$	207
3.3	Complexity oscillations of a typical random sequence $\omega$	215
3.4	$K$ -complexity criteria for randomness of infinite sequences	215
3.5	Complexity oscillations of $\Omega$	216
4.1	Graph of $\mathbf{m}(x)$ with lower bound $1/x \cdot \log x \cdot \log \log x \dots$	249

4.2 Relations between five complexities . . . . .	285
5.1 Trivial consistent automaton . . . . .	317
5.2 Smallest consistent automaton . . . . .	317
5.3 Sample data set . . . . .	365
5.4 Imperfect decision tree . . . . .	366
5.5 Perfect decision tree . . . . .	367
6.1 Single-tape Turing machine . . . . .	381
6.2 The two possible nni's on $(u, v)$ : swap $B \leftrightarrow C$ or $B \leftrightarrow D$ . . . . .	416
6.3 The nni distance between (i) and (ii) is 2 . . . . .	416
6.4 Multitape Turing machine . . . . .	428
8.1 Reversible Boolean gates . . . . .	530
8.2 Implementing reversible AND gate and NOT gate . . . . .	531
8.3 Controlling billiard ball movements . . . . .	532
8.4 A billiard ball computer . . . . .	533
8.5 Combining irreversible computations of $y$ from $x$ and $x$ from $y$ to achieve a reversible computation of $y$ from $x$ . . . . .	543
8.6 Reversible execution of concatenated programs for $(y x)$ and $(z y)$ to transform $x$ into $z$ . . . . .	545
8.7 Carnot cycle . . . . .	555
8.8 Heat engine . . . . .	556
8.9 State space . . . . .	559
8.10 Atomic spin in CuO <sub>2</sub> at low temperature . . . . .	562
8.11 Regular "up" and "down" spins . . . . .	563
8.12 Algorithmic entropy: left a random micro state, right a regular micro state . . . . .	567
8.13 Szilard engine . . . . .	568
8.14 Adiabatic demagnetization to achieve low temperature . . . . .	583
8.15 The maze: a binary tree constructed from matches . . . . .	584
8.16 Time required for <i>Formica sanguinea</i> scouts to transmit information about the direction to the syrup to the forager ants . . . . .	585

# Preliminaries

## 1.1 A Brief Introduction

---

Suppose we want to describe a given object by a finite binary string. We do not care whether the object has many descriptions; however, each description should describe but one object. From among all descriptions of an object we can take the length of the shortest description as a measure of the object's complexity. It is natural to call an object "simple" if it has at least one short description, and to call it "complex" if all of its descriptions are long.

But now we are in danger of falling into the trap so eloquently described in the Richard-Berry paradox, where we define a natural number as "the least natural number that cannot be described in less than twenty words." If this number does exist, we have just described it in thirteen words, contradicting its definitional statement. If such a number does not exist, then all natural numbers can be described in fewer than twenty words. We need to look very carefully at the notion of "description."

Assume that each description describes at most one object. That is, there is a specification method  $D$  that associates at most one object  $x$  with a description  $y$ . This means that  $D$  is a function from the set of descriptions, say  $Y$ , into the set of objects, say  $X$ . It seems also reasonable to require that for each object  $x$  in  $X$ , there is a description  $y$  in  $Y$  such that  $D(y) = x$ . (Each object has a description.) To make descriptions useful we like them to be finite. This means that there are only countably many descriptions. Since there is a description for each object, there are also only countably many describable objects. How do we measure the complexity of descriptions?

Taking our cue from the theory of computation, we express descriptions as finite sequences of 0's and 1's. In communication technology, if the specification method  $D$  is known to both a sender and a receiver, then a message  $x$  can be transmitted from sender to receiver by transmitting the sequence of 0's and 1's of a description  $y$  with  $D(y) = x$ . The cost of this transmission is measured by the number of occurrences of 0's and 1's in  $y$ , that is, by the length of  $y$ . The least cost of transmission of  $x$  is given by the length of a shortest  $y$  such that  $D(y) = x$ . We choose this least cost of transmission as the *descriptive complexity* of  $x$  under specification method  $D$ .

Obviously, this descriptive complexity of  $x$  depends crucially on  $D$ . The general principle involved is that the syntactic framework of the description language determines the succinctness of description.

In order to objectively compare descriptive complexities of objects, to be able to say " $x$  is more complex than  $z$ ," the descriptive complexity of  $x$  should depend on  $x$  alone. This complexity can be viewed as related to a universal description method that is *a priori* assumed by all senders and receivers. This complexity is optimal if no other description method assigns a lower complexity to any object.

We are not really interested in optimality with respect to all description methods. For specifications to be useful at all it is necessary that the mapping from  $y$  to  $D(y)$  can be executed in an effective manner. That is, it can at least in principle be performed by humans or machines. This notion has been formalized as that of "partial recursive functions." According to generally accepted mathematical viewpoints it coincides with the intuitive notion of effective computation.

The set of partial recursive functions contains an optimal function that minimizes description length of every other such function. We denote this function by  $D_0$ . Namely, for any other recursive function  $D$ , for all objects  $x$ , there is a description  $y$  of  $x$  under  $D_0$  that is shorter than any description  $z$  of  $x$  under  $D$ . (That is, shorter up to an additive constant that is independent of  $x$ .) Complexity with respect to  $D_0$  minorizes the complexities with respect to all partial recursive functions.

We identify the length of the description of  $x$  with respect to a fixed specification function  $D_0$  with the "algorithmic (descriptive) complexity" of  $x$ . The optimality of  $D_0$  in the sense above means that the complexity of an object  $x$  is invariant (up to an additive constant independent of  $x$ ) under transition from one optimal specification function to another. Its complexity is an objective attribute of the described object alone: it is an intrinsic property of that object, and it does not depend on the description formalism. This complexity can be viewed as "absolute information content": the amount of information that needs to be transmitted between all senders and receivers when they communicate the message in

absence of any other a priori knowledge that restricts the domain of the message.

Broadly speaking, this means that all description syntaxes that are powerful enough to express the partial recursive functions are approximately equally succinct. All algorithms can be expressed in each such programming language equally succinctly, up to a fixed additive constant term. The remarkable usefulness and inherent rightness of the theory of Kolmogorov complexity stems from this independence of the description method.

Thus, we have outlined the program for a general theory of algorithmic complexity. The four major innovations are as follows:

1. In restricting ourselves to formally effective descriptions, our definition covers every form of description that is intuitively acceptable as being effective according to general viewpoints in mathematics and logic.
2. The restriction to effective descriptions entails that there is a universal description method that minorizes the description length or complexity with respect to any other effective description method. This would not be the case if we considered, say, all noneffective description methods. Significantly, this implies Item 3.
3. The description length or complexity of an object is an intrinsic attribute of the object independent of the particular description method or formalizations thereof.
4. The disturbing Richard-Berry paradox above does not disappear, but resurfaces in the form of an alternative approach to proving Kurt Gödel's (1906–1978) famous result that not every true mathematical statement is provable in mathematics.

**Example 1.1.1 (Gödel's incompleteness result)** A formal system (consisting of definitions, axioms, rules of inference) is *consistent* if no statement that can be expressed in the system can be proved to be both true and false in the system. A formal system is *sound* if only true statements can be proved to be true in the system. (Hence, a sound formal system is consistent.)

Let  $x$  be a finite binary string. We write “ $x$  is random” if the shortest binary description of  $x$  with respect to the optimal specification method  $D_0$  has length at least  $x$ . A simple counting argument shows that there are random  $x$ 's of each length.

Fix any sound formal system  $F$  in which we can express statements like “ $x$  is random.” Suppose  $F$  can be described in  $f$  bits—assume, for example, that this is the number of bits used in the exhaustive description of

$F$  in the first chapter of the textbook *Foundations of F*. We claim that for all but finitely many random strings  $x$ , the sentence “ $x$  is random” is not provable in  $F$ . Assume the contrary. Then given  $F$ , we can start to exhaustively search for a proof that some string of length  $n \gg f$  is random, and print it when we find such a string  $x$ . This procedure to print  $x$  of length  $n$  uses only  $\log n + f$  bits of data, which is much less than  $n$ . But  $x$  is random by the proof and the fact that  $F$  is sound. Hence,  $F$  is not consistent, which is a contradiction.  $\diamond$

This shows that although most strings are random, it is impossible to effectively prove them random. In a way, this explains why the incompressibility method in Chapter 6 is so successful. We can argue about a “typical” individual element, which is difficult or impossible by other methods.

**Example 1.1.2 (Lower bounds)** The secret of the successful use of descriptive complexity arguments as a proof technique is due to a simple fact: the overwhelming majority of strings have almost no computable regularities. We have called such a string “random.” There is no shorter description of such a string than the literal description: it is incompressible. Incompressibility is a noneffective property in the sense of Example 1.1.1.

Traditional proofs often involve all instances of a problem in order to conclude that some property holds for at least one instance. The proof would be more simple, if only that one instance could have been used in the first place. Unfortunately, that instance is hard or impossible to find, and the proof has to involve all the instances. In contrast, in a proof by the incompressibility method, we first choose a random (that is, incompressible) individual object that is known to exist (even though we cannot construct it). Then we show that if the assumed property did not hold, then this object could be compressed, and hence it would not be random. Let us give a simple example.

A prime number is a natural number that is not divisible by natural numbers other than itself and 1. We prove that for infinitely many  $n$ , the number of primes less than or equal to  $n$  is at least  $\log n / \log \log n$ . The proof method is as follows. For each  $n$ , we construct a description from which  $n$  can be effectively retrieved. This description will involve the primes less than  $n$ . For some  $n$  this description must be long, which shall give the desired result.

Assume that  $p_1, p_2, \dots, p_m$  is the list of all the primes less than  $n$ . Then,

$$n = p_1^{e_1} p_2^{e_2} \cdots p_m^{e_m}$$

can be reconstructed from the vector of the exponents. Each exponent is at most  $\log n$  and can be represented by  $\log \log n$  bits. The description of  $n$  (given  $\log n$ ) can be given in  $m \log \log n$  bits.

It can be shown that each  $n$  that is random (given  $\log n$ ) cannot be described in fewer than  $\log n$  bits, whence the result follows. Can we do better? This is slightly more complicated. Let  $l(x)$  denote the length of the binary representation of  $x$ . We shall show that for infinitely many  $n$  of the form  $n = m \log^2 m$ , the number of distinct primes less than  $n$  is at least  $m$ .

Firstly, we can describe any given integer  $N$  by  $E(m)N/p_m$ , where  $E(m)$  is a prefix-free encoding (page 71) of  $m$ , and  $p_m$  is the largest prime dividing  $N$ . For random  $N$ , the length of this description,  $l(E(m)) + \log N - \log p_m$ , must exceed  $\log N$ . Therefore,  $\log p_m < l(E(m))$ . It is known (and easy) that we can set  $l(E(m)) \leq \log m + 2 \log \log m$ . Hence,  $p_m < m \log^2 m$ . Setting  $n := m \log^2 m$ , and observing from our previous result that  $p_m$  must grow with  $N$ , we have proven our claim. The claim is equivalent to the statement that for our special sequence of values of  $n$  the number of primes less than  $n$  exceeds  $n/\log^2 n$ . The idea of connecting primality and prefix code-word length is due to P. Berman, and the present proof is due to J. Tromp.

Chapter 6 introduces the incompressibility method. Its utility is demonstrated in a variety of examples of proving mathematical and computational results. These include questions concerning the average case analysis of algorithms (such as Heapsort), sequence analysis, average case complexity in general, formal languages, combinatorics, time and space complexity analysis of various sequential or parallel machine models, language recognition, and string matching. Other topics like the use of resource-bounded Kolmogorov complexity in the analysis of computational complexity classes, the universal optimal search algorithm, and “logical depth” are treated in Chapter 7. ◇

**Example 1.1.3 (Prediction)** We are given an initial segment of an infinite sequence of zeros and ones. Our task is to predict the next element in the sequence: zero or one? The set of possible sequences we are dealing with constitutes the “sample space”; in this case, the set of one-way infinite binary sequences. We assume some probability distribution  $\mu$  over the sample space, where  $\mu(x)$  is the probability of the initial segment of a sequence being  $x$ . Then the probability of the next bit being “0,” after an initial segment  $x$ , is clearly  $\mu(0|x) = \mu(x0)/\mu(x)$ . This problem constitutes, perhaps, the central task of inductive reasoning and artificial intelligence. However, the problem of induction is that in general we do not know the distribution  $\mu$ , preventing us from assessing the actual probability. Hence, we have to use an estimate.

Now assume that  $\mu$  is computable. (This is not very restrictive, since any distribution used in statistics is computable, provided the parameters are computable.) We can use Kolmogorov complexity to give a very good

estimate of  $\mu$ . This involves the so-called “universal distribution”  $\mathbf{M}$ . Roughly speaking,  $\mathbf{M}(x)$  is close to  $2^{-l}$ , where  $l$  is the length in bits of the shortest effective description of  $x$ . Among other things,  $\mathbf{M}$  has the property that it assigns at least as high a probability to  $x$  as any computable  $\mu$  (up to a multiplicative constant factor depending on  $\mu$  but not on  $x$ ). What is particularly important to prediction is the following:

Let  $S_n$  denote the  $\mu$ -expectation of the square of the error we make in estimating the probability of the  $n$ th symbol by  $\mathbf{M}$ . Then it can be shown that the sum  $\sum_n S_n$  is bounded by a constant. In other words,  $S_n$  converges to zero faster than  $1/n$ . Consequently, any actual (computable) distribution can be estimated and predicted with great accuracy using only the single universal distribution.

Chapter 5 develops a general theory of inductive reasoning and applies the notions introduced to particular problems of inductive inference, prediction, mistake bounds, computational learning theory, and minimum description length induction methods in statistics. In particular, it is demonstrated that data compression improves generalization and prediction performance. ◇

The purpose of the remainder of this chapter is to define several concepts we require, if not by way of introduction, then at least to establish notation.

## 1.2 Prerequisites and Notation

We usually deal with nonnegative integers, sets of nonnegative integers, and mappings from nonnegative integers to nonnegative integers.  $A, B, C, \dots$  denote sets.  $\mathcal{N}, \mathcal{Z}, \mathcal{Q}, \mathcal{R}$  denote the sets of nonnegative integers (natural numbers including zero), integers, rational numbers, and real numbers, respectively. For each such set  $A$ , by  $A^+$  we denote the subset of  $A$  consisting of positive numbers.

We use the following set-theoretical notations.  $x \in A$  means that  $x$  is a member of  $A$ . In  $\{x : x \in A\}$ , the symbol “ $:$ ” denotes set formation.  $A \cup B$  is the *union* of  $A$  and  $B$ ,  $A \cap B$  is the *intersection* of  $A$  and  $B$ , and  $\bar{A}$  is the *complement* of  $A$  when the universe  $A \cup \bar{A}$  is understood.  $A \subseteq B$  means  $A$  is a *subset* of  $B$ .  $A = B$  means  $A$  and  $B$  are *identical* as sets (have the same members).

The *cardinality* (or diameter) of a finite set  $A$  is the number of elements it contains and is denoted as  $d(A)$ . If  $A = \{a_1, \dots, a_n\}$ , then  $d(A) = n$ . The *empty* set  $\{\}$ , with no elements in it, is denoted by  $\emptyset$ . In particular,  $d(\emptyset) = 0$ .

Given  $x$  and  $y$ , the ordered pair  $(x, y)$  consists of  $x$  and  $y$  in that order.  $A \times B$  is the *Cartesian product* of  $A$  and  $B$ , the set  $\{(x, y) : x \in A$  and

$y \in B\}$ . The  $n$ -fold Cartesian product of  $A$  with itself is denoted as  $A^n$ . If  $R \subseteq A^2$ , then  $R$  is called a *binary relation*. The same definitions can be given for  $n$ -tuples,  $n > 2$ , and the corresponding relations are *n-ary*. We say that an  $n$ -ary relation  $R$  is *single-valued* if for every  $(x_1, \dots, x_{n-1})$  there is at most one  $y$  such that  $(x_1, \dots, x_{n-1}, y) \in R$ . Consider the *domain*  $\{(x_1, \dots, x_{n-1}) : \text{there is a } y \text{ such that } (x_1, \dots, x_{n-1}, y) \in R\}$  of a single-valued relation  $R$ . Clearly, a single-valued relation  $R \subseteq A^{n-1} \times B$  can be considered as a mapping from its domain into  $B$ . Therefore, we also call a single-valued  $n$ -ary relation a *partial function* of  $n-1$  variables (“partial” because the domain of  $R$  may not comprise all of  $A^{n-1}$ ). We denote functions by  $\phi, \psi, \dots$  or  $f, g, h, \dots$  Functions defined on the  $n$ -fold Cartesian product  $A^n$  are denoted with possibly a superscript denoting the number of variables, like  $\phi^{(n)} = \phi^{(n)}(x_1, \dots, x_n)$ .

We use the notation  $\langle \cdot \rangle$  for some standard one-to-one encoding of  $\mathcal{N}^n$  into  $\mathcal{N}$ . We will use  $\langle \cdot \rangle$  especially as a *pairing function* over  $\mathcal{N}$  to associate a unique natural number  $\langle x, y \rangle$  with each pair  $(x, y)$  of natural numbers. An example is  $\langle x, y \rangle$  defined by  $y + (x + y + 1)(x + y)/2$ . This mapping can be used recursively:  $\langle x, y, z \rangle = \langle x, \langle y, z \rangle \rangle$ .

If  $\phi$  is a *partial function* from  $A$  to  $B$ , then for each  $x \in A$  either  $\phi(x) \in B$  or  $\phi(x)$  is undefined. If  $x$  is a member of the domain of  $\phi$ , then  $\phi(x)$  is called a *value* of  $\phi$ , and we write  $\phi(x) < \infty$  and  $\phi$  is called *convergent* or *defined* at  $x$ ; otherwise we write  $\phi(x) = \infty$  and we call  $\phi$  *divergent* or *undefined* at  $x$ . The set of values of  $\phi$  is called the *range* of  $\phi$ . If  $\phi$  converges at every member of  $A$ , it is a *total function*, otherwise a *strictly partial function*. If each member of a set  $B$  is also a value of  $\phi$ , then  $\phi$  is said to *map onto*  $B$ , otherwise to *map into*  $B$ . If for each pair  $x$  and  $y$ ,  $x \neq y$ , for which  $\phi$  converges  $\phi(x) \neq \phi(y)$  holds, then  $\phi$  is a *one-to-one mapping*, otherwise a *many-to-one mapping*. The function  $f : A \rightarrow \{0, 1\}$  defined by  $f(x) = 1$  if  $\phi(x)$  converges, and  $f(x) = 0$  otherwise, is called the *characteristic function* of the domain of  $\phi$ .

If  $\phi$  and  $\psi$  are two partial functions, then  $\psi\phi$  (equivalently,  $\psi(\phi(x))$ ) denotes their *composition*, the function defined by  $\{(x, y) : \text{there is a } z \text{ such that } \phi(x) = z \text{ and } \psi(z) = y\}$ . The *inverse*  $\phi^{-1}$  of a one-to-one partial function  $\phi$  is defined by  $\phi^{-1}(y) = x$  iff  $\phi(x) = y$ .

A set  $A$  is called *countable* if it is either empty or there is a total one-to-one mapping from  $A$  to the natural numbers  $\mathcal{N}$ . We say  $A$  is *countably infinite* if it is both countable and infinite. By  $2^A$  we denote the *set of all subsets* of  $A$ . The set  $2^{\mathcal{N}}$  has the cardinality of the *continuum* and is therefore uncountably infinite.

For binary relations, we use the terms *reflexive*, *transitive*, *symmetric*, *equivalence*, *partial order*, and *linear* (or *total*) *order* in the usual meaning. Partial orders can be *strict* ( $<$ ) or *nonstrict* ( $\leq$ ).

If we use the logarithm notation  $\log x$  without subscript, then we shall always mean base 2. By  $\ln x$  we mean the *natural logarithm*  $\log_e x$ , where  $e = 2.71\dots$ .

We use the quantifiers  $\exists$  (“there exists”),  $\forall$  (“for all”),  $\exists^\infty$  (“there exist infinitely many”), and the awkward  $\forall^\infty$  (“for all but finitely many”). This way,  $\forall^\infty x[\phi(x)]$  iff  $\neg\exists^\infty x[\neg\phi(x)]$ .

## 1.3 Numbers and Combinatorics

---

The *absolute value* of a real number  $r$  is denoted by  $|r|$  and is defined as  $|r| = -r$  if  $r < 0$  and  $r$  otherwise. The *floor* of a real number  $r$ , denoted by  $\lfloor r \rfloor$ , is the greatest integer  $n$  such that  $n \leq r$ . Analogously, the *ceiling* of a real number  $r$ , denoted by  $\lceil r \rceil$ , is the least integer  $n$  such that  $n \geq r$ .

**Example 1.3.1**  $|-1| = |1| = 1$ .  $\lfloor 0.5 \rfloor = 0$  and  $\lceil 0.5 \rceil = 1$ . Analogously,  $\lfloor -0.5 \rfloor = -1$  and  $\lceil -0.5 \rceil = 0$ . But  $\lfloor 2 \rfloor = \lceil 2 \rceil = 2$  and  $\lfloor -2 \rfloor = \lceil -2 \rceil = -2$ .  $\diamond$

A *permutation* of  $n$  objects is an arrangement of  $n$  distinct objects in an ordered sequence. For example, the six different permutations of objects  $a, b, c$  are

$$abc, acb, bac, bca, cab, cba.$$

The number of permutations of  $n$  objects is found most easily by imagining a sequential process to choose a permutation. There are  $n$  choices of which object to place in the first position; after filling the first position there remain  $n - 1$  objects and therefore  $n - 1$  choices of which object to place in the second position, and so on. Therefore, the number of permutations of  $n$  objects is  $n \times (n - 1) \times \dots \times 2 \times 1$ , denoted by  $n!$  and is referred to as  *$n$  factorial*. In particular,  $0! = 1$ .

A *variation* of  $k$  out of  $n$  objects is an arrangement consisting of the first  $k$  elements of a permutation of  $n$  objects. For example, the twelve variations of two out of four objects  $a, b, c, d$  are

$$ab, ac, ad, ba, bc, bd, ca, cb, cd, da, db, dc.$$

The number of variations of  $k$  out of  $n$  is  $n!/(n - k)!$ , as follows by the previous argument. While there is no accepted standard notation, we denote the number of variations as  $(n)_k$ . In particular,  $(n)_0 = 1$ .

The *combinations* of  $n$  objects taken  $k$  at a time (“ $n$  choose  $k$ ”) are the possible choices of  $k$  different elements from a collection of  $n$  objects. The six different combinations of two out of four objects  $a, b, c, d$  are

$$\{a, b\}, \{a, c\}, \{a, d\}, \{b, c\}, \{b, d\}, \{c, d\}.$$

We can consider a combination as a variation in which the order does not count. We have seen that there are  $n(n - 1) \cdots (n - k + 1)$  ways to choose the first  $k$  elements of a permutation. Every  $k$ -combination appears precisely  $k!$  times in these arrangements, since each combination occurs in all its permutations. Therefore, the number of combinations, denoted by  $\binom{n}{k}$ , is

$$\binom{n}{k} = \frac{n(n - 1) \cdots (n - k + 1)}{k(k - 1) \cdots (1)}.$$

In particular,  $\binom{n}{0} = 1$ . The quantity  $\binom{n}{k}$  is also called a *binomial coefficient*. It has an extraordinary number of applications. Perhaps the foremost relation associated with it is the Binomial Theorem, discovered in 1676 by Isaac Newton

$$(x + y)^n = \sum_k \binom{n}{k} x^k y^{n-k},$$

with  $n$  a positive integer. Note that in the summation  $k$  need not be restricted to  $0 \leq k \leq n$ , but can range over  $-\infty < k < +\infty$ , since for  $k < 0$  or  $k > n$  the terms are all zero.

**Example 1.3.2** An important relation following from the Binomial Theorem is found by substituting  $y = 1$ :

$$(x + 1)^n = \sum_k \binom{n}{k} x^k.$$

Substituting also  $x = 1$  we find

$$2^n = \sum_k \binom{n}{k}.$$

◇

## Exercises

**1.3.1.** [13] A “stock” of bridge cards consists of four suits and thirteen face values, respectively. Each card is defined by its suit and face value.

- (a) How many cards are there?
- (b) Each player gets a “hand” consisting of thirteen cards. How many different hands are there?
- (c) What is the probability of getting a full suit as a hand? Assume this is the probability of obtaining a full suit when drawing thirteen cards, successively without replacement, from a full stock.

*Comments.* (a)  $4 \times 13 = 52$ . (b)  $\binom{52}{13} = 635013559600$ . (c)  $4/\binom{52}{13}$ .

**1.3.2.** [12] Consider a random distribution of  $k$  distinguishable balls in  $n$  cells, that is, each of the  $n^k$  possible arrangements has probability  $n^{-k}$ . Show that the probability  $P_i$  that a specified cell contains exactly  $i$  balls ( $0 \leq i \leq k$ ) is given by  $P_i = \binom{k}{i}(1/n)^i(1 - 1/n)^{k-i}$ .

*Comments.* Source: W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. 1, Wiley, 1968.

**1.3.3.** [08] Show that  $\binom{n}{k} = \frac{\binom{n}{k}}{k!}$  and  $\binom{n}{k} = \binom{n}{n-k}$ .

**1.3.4.** [M34] Prove the following identity, which is very useful in the sequel of this book.

$$\log \binom{n}{k} = k \log \frac{n}{k} + (n - k) \log \frac{n}{n - k} + \frac{1}{2} \log \frac{n}{k(n - k)} + O(1).$$

**1.3.5.** [15] (a) Prove that the number of ways  $n$  distinguishable balls can be placed in  $k$  numbered cells such that the first cell contains  $n_1$  balls, the second cell  $n_2$  balls, up to the  $k$ th cell contains  $n_k$  balls with  $n_1 + \dots + n_k = n$  is

$$\binom{n}{n_1, \dots, n_k} = \frac{n!}{n_1! \cdots n_k!}.$$

This number is called a *multinomial coefficient*. Note that the order of the cells is essential in that the partitions  $(n_1 = 1, n_2 = 2)$  and  $(n_1 = 2, n_2 = 1)$  are different. The order of the elements within a cell is irrelevant.

(b) Show that

$$(x_1 + \dots + x_k)^n = \sum \binom{n}{n_1, \dots, n_k} x_1^{n_1} \cdots x_k^{n_k},$$

with the sum taken for all  $n_1 + \dots + n_k = n$ .

(c) The *number of ordered different partitions* of  $n$  in  $r$  nonnegative integral summands is denoted by  $A_{n,r}$ . Compute  $A_{n,r}$  in the form of a binomial coefficient.

*Comments.*  $(1, 0)$  and  $(0, 1)$  are different partitions, so  $A_{1,2} = 2$ . Source: W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. 1, Wiley, 1968.

**1.3.6.** [14] Define the *occupancy numbers* for  $n$  balls distributed over  $k$  cells as a  $k$ -tuple of integers  $(n_1, n_2, \dots, n_k)$  satisfying  $n_1 + n_2 + \dots + n_k = n$  with  $n_i \geq 0$  ( $1 \leq i \leq k$ ). That is, the first cell contains  $n_1$  balls, the second cell  $n_2$  balls, and so on.

- (a) Show that there are  $\binom{n}{n_1, \dots, n_k}$  placements of  $n$  balls in  $k$  cells resulting in the numbers  $(n_1, \dots, n_k)$ .
- (b) There are  $k^n$  possible placements of  $n$  balls in  $k$  cells altogether. Compute the fraction that results in the given occupancy numbers.
- (c) Assume that all  $k^n$  possible placements of  $n$  balls in  $k$  cells are equally probable. Conclude that the probability of obtaining the given occupancy numbers is

$$\frac{n!}{n_1! \cdots n_k!} k^{-n}.$$

*Comments.* In physics this is known as the *Maxwell-Boltzmann statistics* (here “statistics” is used as a synonym to “distribution”). Source: W. Feller, *Ibid.*

**1.3.7.** [15] We continue with the previous Exercise. In physical situations the assumption of equiprobability of possible placements seems unavoidable, for example, molecules in a volume of gas divided into (hypothetical) cells of equal volume. Numerous attempts have been made to prove that physical particles behave in accordance with the Maxwell-Boltzmann distribution. However, it has been shown conclusively that *no known particles* behave according to this distribution.

- (a) In the *Bose-Einstein* distribution we count only *distinguishable* distributions of  $n$  balls over  $k$  cells without regard for the identities of the balls. We are only interested in the number of solutions of  $n_1 + n_2 + \cdots + n_k = n$ . Show that this number is  $A_{n,k} = \binom{k+n-1}{n} = \binom{k+n-1}{k-1}$ . Conclude that the probability of obtaining each given occupancy number is equally  $1/A_{n,k}$ . (Illustration: the distinguishable distributions of two balls over two cells are |\*\*, \*|\*, and \*\*|. Hence, according to Bose-Einstein statistics there are only three possible outcomes for two coin flips: head-head, head-tail, and tail-tail, and each outcome has equal probability  $\frac{1}{3}$ .)
- (b) In the *Fermi-Dirac* distribution, (1) two or more particles cannot occupy the same cell and (2) all distinguishable arrangements satisfying (1) have the same probability. Note that (1) requires  $n \leq k$ . Prove that in the Fermi-Dirac distribution there are in total  $\binom{k}{n}$  possible arrangements. Conclude that the probability for each possible occupancy number is equally  $1/\binom{k}{n}$ .

*Comments.* According to modern physics, photons, nuclei, and atoms containing an even number of elementary particles behave according to model (a), and electrons, neutrons, and protons behave according to model (b). This shows that nature does not necessarily satisfy our a priori assumptions, however plausible they may be. Source: W. Feller, *Ibid.*

## 1.4 Binary Strings

We are concerned with *strings* over a nonempty set  $\mathcal{B}$  of *basic elements*. Unless otherwise noted, we use  $\mathcal{B} = \{0, 1\}$ . Instead of “string” we also use “word” and “sequence,” synonymously. The way we use it, “strings” and “words” are usually finite, while “sequences” are usually infinite. The set of all finite strings over  $\mathcal{B}$  is denoted by  $\mathcal{B}^*$ , defined as

$$\mathcal{B}^* = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\},$$

with  $\epsilon$  denoting the *empty* string, with no letters. *Concatenation* is a binary operation on the elements of  $\mathcal{B}^*$  that associates  $xy$  with each ordered pair of elements  $(x, y)$  in the Cartesian product  $\mathcal{B}^* \times \mathcal{B}^*$ . Clearly,

1.  $\mathcal{B}^*$  is closed under the operation of concatenation; that is, if  $x$  and  $y$  are elements of  $\mathcal{B}^*$ , then so is  $xy$ ;
2. concatenation is an associative operation on  $\mathcal{B}^*$ ; that is,  $(xy)z = x(yz) = xyz$ ; and
3. concatenation on  $\mathcal{B}^*$  has the *unit* element  $\epsilon$ ; that is,  $\epsilon x = x\epsilon = x$ .

We now consider a *correspondence* of finite binary strings and natural numbers. The standard binary representation has the disadvantage that either some strings do not represent a natural number, or each natural number is represented by more than one string. For example, either “010” does not represent “2,” or both “010” and “10” represent “2.” However, we can map  $\mathcal{B}^*$  one-to-one onto the natural numbers by associating each string with its index in the lexicographical ordering

$$(\epsilon, 0), (0, 1), (1, 2), (00, 3), (01, 4), (10, 5), (11, 6), \dots \quad (1.1)$$

This way we represent  $x = 2^{n+1} - 1 + \sum_{i=0}^n a_i 2^i$  by  $a_n \dots a_1 a_0$ . This is equivalent to  $x = \sum_{i=0}^n b_i 2^i$  with  $b_i \in \{1, 2\}$  and  $b_i = a_i + 1$  for  $0 \leq i \leq n$ .

This way we have a binary representation for the natural numbers that is different from the standard binary representation. It is convenient not to distinguish between the first and second element of the same pair, and call them “string” or “number” arbitrarily. That is, we consider both the string 01 and the natural number 4 as the same object. For example, we may write 01 = 4. We denote these objects in general with lowercase roman letters. A string consisting of  $n$  zeros is denoted by  $0^n$ .

If  $x$  is a string of  $n$  0’s and 1’s, then  $x_i$  denotes the  $i$ th *bit* (binary digit) of  $x$  for all  $i$ ,  $1 \leq i \leq n$ , and  $x_{i:j}$  denotes the  $(j - i + 1)$ -bits segment  $x_i x_{i+1} \dots x_j$ . For  $x = 1010$  we have  $x_1 = x_3 = 1$  and  $x_2 = x_4 = 0$ ; for  $x = x_1 x_2 \dots x_n$  we have  $x_{1:i} = x_1 x_2 \dots x_i$ . The *reverse*,  $x^R$ , of a string  $x = x_1 x_2 \dots x_n$  is  $x_n x_{n-1} \dots x_1$ .

The *length* of a finite binary string  $x$  is the number of bits it contains and is denoted by  $l(x)$ . If  $x = x_1x_2 \dots x_n$ , then  $l(x) = n$ . In particular,  $l(\epsilon) = 0$ .

Thus,  $l(xy) = l(x) + l(y)$ , and  $l(x^R) = l(x)$ . Recall that we use the above pairing of binary strings and natural numbers. Thus,  $l(4) = 2$  and  $01 = 4$ . The *number of elements (cardinality)* in a finite set  $A$  is denoted by  $d(A)$ . Therefore,  $d(\{u : l(u) = n\}) = 2^n$  and  $d(\{u : l(u) \leq n\}) = 2^{n+1} - 1$ .

Let  $D$  be any function  $D : \{0, 1\}^* \rightarrow \mathcal{N}$ . Considering the domain of  $D$  as the set of *code words*, and the range of  $D$  as the set of *source words*,  $D(y) = x$  is interpreted as “ $y$  is a code word for the source word  $x$ , and  $D$  is the *decoding* function.” (In the introduction we called  $D$  a specification method.) The set of all code words for source word  $x$  is the set  $D^{-1}(x) = \{y : D(y) = x\}$ . Hence,  $E = D^{-1}$  can be called the *encoding* substitution ( $E$  is not necessarily a function). Let  $x, y \in \{0, 1\}^*$ . We call  $x$  a *prefix* of  $y$  if there is a  $z$  such that  $y = xz$ . A set  $A \subseteq \{0, 1\}^*$  is *prefix-free*, if no element in  $A$  is the prefix of another element in  $A$ . A function  $D : \{0, 1\}^* \rightarrow \mathcal{N}$  defines a *prefix-code* if its domain is prefix-free. (Coding theory is treated in Section 1.11.1.) A simple prefix-code we use throughout is obtained by reserving one symbol, say 0, as a stop sign and encoding  $x \in \mathcal{N}$  as  $1^x 0$ . We can prefix an object with its length and iterate this idea to obtain ever shorter codes:

$$E_i(x) = \begin{cases} 1^{l(x)} & \text{for } i = 0, \\ E_{i-1}(l(x))x & \text{for } i > 0. \end{cases} \quad (1.2)$$

Thus,  $E_1(x) = 1^{l(x)} 0x$  and has length  $l(E_1(x)) = 2l(x) + 1$ . This encoding is sufficiently important to have a simpler notation:

$$\begin{aligned} \bar{x} &= 1^{l(x)} 0x, \\ l(\bar{x}) &= 2l(x) + 1. \end{aligned}$$

Sometimes we need the shorter prefix-code  $E_2(x)$ ,

$$\begin{aligned} E_2(x) &= \overline{l(x)}x, \\ l(E_2(x)) &= l(x) + 2l(l(x)) + 1. \end{aligned}$$

We call  $\bar{x}$  the *self-delimiting* version of the binary string  $x$ . Now we can effectively recover both  $x$  and  $y$  unambiguously from the binary string  $\bar{x}\bar{y}$ . If  $\bar{x}\bar{y} = 111011011$ , then  $x = 110$  and  $y = 11$ . If  $\bar{x}\bar{y} = 1110110101$  then  $x = 110$  and  $y = 1$ .

**Example 1.4.1** It is convenient to consider also the set of *one-way infinite* sequences  $\mathcal{B}^\infty$ . If  $\omega$  is an element of  $\mathcal{B}^\infty$ , then  $\omega = \omega_1\omega_2\dots$  and  $\omega_{1:n} = \omega_1\omega_2\dots\omega_n$ .

The set of infinite sequences of elements in a finite, nonempty basic set  $\mathcal{B}$  corresponds with the set  $\mathcal{R}$  of real numbers in the following way:

Let  $\mathcal{B} = \{0, 1, \dots, k - 1\}$  with  $k \geq 2$ . If  $r$  is a real number  $0 < r < 1$  then there is a sequence  $\omega_1\omega_2\dots$  of elements  $\omega_n \in \mathcal{B}$  such that

$$r = \sum_n \omega_n/k^n,$$

and that sequence is unique except when  $r$  is of the form  $q/k^n$ , in which case there are exactly two such sequences, one of which has infinitely many 0's. Conversely, if  $\omega_1\omega_2\dots$  is an infinite sequence of integers with  $0 \leq \omega_n < k$ , then the series

$$\sum_n \omega_n/k^n$$

converges to a real number  $r$  with  $0 \leq r \leq 1$ . This sequence is called the *k*-ary expansion of  $r$ . In the following we identify a real number  $r$  with its *k*-ary expansion (if there are two *k*-ary expansions, then we identify  $r$  with the expansion with infinitely many 0's).

Define the set  $S \subseteq \mathcal{B}^\infty$  as the set of sequences that do not end with infinitely many digits " $k - 1$ ." Then,  $S$  is in one-to-one relation with the set of real numbers in the interval  $[0, 1)$ .

Let  $x$  be a finite string over  $\mathcal{B}$ . The set of all one-way infinite sequences starting with  $x$  is called a *cylinder* and is denoted by  $\Gamma_x$  and is defined by  $\Gamma_x = \{\omega : \omega_{1:l(x)} = x\}$ . Geometrically speaking, the cylinder  $\Gamma_x$  can be identified with the *half-open interval*  $[0.x, 0.x + 2^{-l(x)})$  in the real interval  $[0, 1)$ . Observe that the usual geometric length of interval  $\Gamma_x$  equals  $2^{-l(x)}$ . Furthermore,  $\Gamma_y \subseteq \Gamma_x$  iff  $x$  is a prefix of  $y$ . The prefix relation induces a partial order on the cylinders of  $\mathcal{B}^\infty$ .  $\diamond$

## Exercises

---

**1.4.1.** [03] If  $\bar{x}\bar{y}z = 10010111$ , what are  $x, y, z$  in decimal numbers?

*Comments.* 1, 2, 6.

**1.4.2.** [07] Show that for  $x \in \mathcal{N}$  we have  $l(x) = \lfloor \log(x + 1) \rfloor$ .

**1.4.3.** [10] Let  $E : \mathcal{N} \rightarrow \{0, 1\}^*$  be a total one-to-one function whose range is prefix-free.  $E$  defines a prefix-code. Define the mapping  $\langle \cdot \rangle : \mathcal{N} \times \mathcal{N} \rightarrow \mathcal{N}$  by  $\langle x, y \rangle = E(x)y$ .

(a) Show that  $\langle \cdot \rangle$  is total and one-to-one.

(b) Show that we can extend this scheme to  $k$ -tuples  $(n_1, n_2, \dots, n_k)$  of natural numbers to obtain a total one-to-one mapping from  $\mathcal{N} \times \mathcal{N} \times \dots \times \mathcal{N}$  into  $\mathcal{N}$ .

*Comments.* Define the mapping for  $(x, y, z)$  as  $\langle x, \langle y, z \rangle \rangle$ , and iterate this construction.

**1.4.4.** [10] Let  $E$  be as above. Define the mapping  $\langle \cdot \rangle : \mathcal{N} \times \mathcal{N} \rightarrow \mathcal{N}$  by  $\langle x, y \rangle = E(x)E(y)$ .

(a) Show that  $\langle \cdot \rangle$  is a total one-to-one mapping and a prefix-code.

(b) Show that we can extend this scheme to  $k$ -tuples  $(n_1, n_2, \dots, n_k)$  of natural numbers to obtain a total one-to-one mapping from  $\mathcal{N} \times \mathcal{N} \times \dots \times \mathcal{N}$  into  $\mathcal{N}$  that is a prefix-code.

*Comments.* Define the mapping for  $(x, y, z)$  as  $\langle x, \langle y, z \rangle \rangle$  and iterate this construction. Another way is to map  $(x, y, \dots, z)$  to  $E(x)E(y)\dots E(z)$ .

**1.4.5.** [10] (a) Show that  $E(x) = \bar{x}$  is a prefix-code.

(b) Consider a variant of the  $\bar{x}$  code such that  $x = x_1x_2\dots x_n$  is encoded as  $x_11x_21\dots x_{n-1}1x_n0$ . Show that this is a prefix-code for the binary nonempty strings with  $l(\bar{x}) = 2l(x)$ .

(c) Consider  $x = x_1x_2\dots x_n$  encoded as  $x_1x_1x_2x_2\dots x_{n-1}x_{n-1}x_n\bar{x}_n$ . Show that this is a prefix-code for the nonempty binary strings.

(d) Give a prefix-code  $\tilde{x}$  for the set of all binary strings  $x$  including  $\epsilon$ , such that  $l(\tilde{x}) = 2l(x) + 2$ .

## 1.5 Asymptotic Notation

It is often convenient to express approximate equality or inequality of one quantity with another. If  $f$  and  $g$  are functions of a real variable, then it is customary to denote  $\lim_{n \rightarrow \infty} f(n)/g(n) = 1$  by  $f(n) \sim g(n)$ , and we write “ $f$  goes asymptotically to  $g$ .”

P. Bachman introduced a convenient notation for dealing with approximations in his book *Analytische Zahlentheorie* in 1892. This “big- $O$ ” notation allows us to write  $l(x) = \log x + O(1)$  (no subscript on the logarithm means base 2).

We use the notation  $O(f(n))$  whenever we want to denote a quantity that does not exceed  $f(n)$  by more than a fixed multiplicative factor. This is useful in case we want to simplify the expression involving this quantity by suppressing unnecessary detail, but also in case we do not know this quantity explicitly. Bachman’s notation is the first of a family of *order of magnitude* symbols:  $O$ ,  $o$ ,  $\Omega$ , and  $\Theta$ . If  $f$  and  $g$  are functions on the real numbers, then

1.  $f(x) = O(g(x))$  if there are constants  $c, x_0 > 0$  such that  $|f(x)| \leq c|g(x)|$ , for all  $x \geq x_0$ ;
2.  $f(x) = o(g(x))$  if  $\lim_{x \rightarrow \infty} f(x)/g(x) = 0$ ;
3.  $f(x) = \Omega(g(x))$  if  $f(x) \neq o(g(x))$ ; and

4.  $f(x) = \Theta(g(x))$  if both  $f(x) = O(g(x))$  and  $f(x) = \Omega(g(x))$ .

It is straightforward to extend these definitions to functions of more variables. For example,  $f(x, y) = O(g(x, y))$  if there are positive constants  $c, x_0, y_0$  such that  $|f(x, y)| \leq c|g(x, y)|$ , for all  $x \geq x_0, y \geq y_0$ . The definitions are standard, except Item 3. This definition of  $\Omega$  was introduced first by G.H. Hardy and J.E. Littlewood in 1914 and is the one commonly used in mathematics. It has the advantage that  $\Omega$  is the complement of  $o$ . This is not the case with the definition proposed by D.E. Knuth in 1976, which is often referred to in computer science. Namely, Knuth defines  $f(x) = \Omega(g(x))$  if there is a constant  $c > 0$  such that  $|f(x)| \geq c|g(x)|$  from some  $x$  onward. We have defined  $f(x) = \Omega(g(x))$  if there is a constant  $c > 0$  such that  $|f(x)| \geq c|g(x)|$  infinitely often. This use of  $\Omega$  should not be confused with Chaitin's mystery number  $\Omega$ , which we encounter in Section 3.6.2.

- Example 1.5.1** The definition of the “big- $O$ ” notation contains some mysterious absolute value signs. This becomes understandable if we realize that one wants to use a term like  $O(f(x))$  to bound the absolute value of an error term, be it positive or negative. For example, as in

$$x^2 + x \sin x = x^2 + O(x).$$

This avoids the clumsy notation “ $\pm O(f(x))$ ” one would have been forced to use otherwise.  $\diamond$

- Example 1.5.2** If  $f(n) \sim g(n)$ , then  $f(n) = \Theta(g(n))$ , but the converse implication does not hold. For instance, we have  $2x = \Theta(x)$ , but  $2x \sim x$  does not hold. On the other hand,  $-x = \Theta(x)$ .  $\diamond$

- Example 1.5.3** We can use  $O$ -notation to speak generically about  $m$ th degree polynomials, for instance,  $1+2+\dots+n = n(n+1)/2$ . Then  $1+2+\dots+n = O(n^2)$ , but also  $1+2+\dots+n = n^2/2 + O(n)$ . The latter approximation is obviously a stronger statement than the first approximation. Similarly, if  $p(n)$  is any polynomial of degree  $m$ , then  $p(n) = O(n^m)$ ;  $p(n) = \Theta(n^m)$ .  $\diamond$

---

## Exercises

- 1.5.1.** [07] Show that  $x = o(x^2)$ ;  $\sin x = O(1)$ ;  $x^{-1/2} = o(1)$ ;  $x + x^2 \sim x^2$ , and  $\sum_{k=1}^n kn = O(n^3)$ .

- 1.5.2.** [10] Show that  $f(n) = O(f(n))$ ;  $c \cdot O(f(n)) = O(f(n))$  if  $c$  is a constant;  $O(f(n) + O(f(n))) = O(f(n))$ ;  $O(O(f(n))) = O(f(n))$ ;  $O(f(n))O(g(n)) = O(f(n)g(n))$ ;  $O(f(n)g(n)) = f(n)O(g(n))$ .

**1.5.3.** [10] Show that  $f(n) = o(g(n))$  implies  $f(n) = O(g(n))$ , but not vice versa.

**1.5.4.** [M30] It is natural to wonder how large  $100!$  is approximately, without carrying out the multiplications implied by the definition. Prove the approximation  $n! \sim \sqrt{2\pi n}(n/e)^n$ .

*Comments.* This celebrated approximation of the factorial function was found by James Stirling [*Methodus Differentialis* (1730), 137].

**1.5.5.** [15] Denote the Hardy-Littlewood version of  $\Omega$  we gave in the main text by  $\Omega_H$  and the Knuth version by  $\Omega_K$ .

(a) Show that for function  $f$  defined by  $\log \log f(n) = \lfloor \log \log n \rfloor$ , we have  $f(n) = \Omega_H(n)$  but for no  $\epsilon > 0$ ,  $f(n) = \Omega_K(n^{1/2+\epsilon})$  holds.

(b) Show that, nonetheless, while  $f(n) = O(n)$ , for no  $\epsilon > 0$  we have  $f(n) = O(n^{1-\epsilon})$ .

*Comments.* Source: P. Vitányi, L. Meertens, *SIGACT News*, 16:4(1985), 56–59.

**1.5.6.** [20] It is well known that  $n^{1/n} \rightarrow 1$  for  $n \rightarrow \infty$ . Let “ $\ln n$ ” denote the *natural logarithm*, that is,  $\log_e n$  with  $e = 2.7\dots$

(a) Show that  $n^{1/n} = e^{\ln n/n} = 1 + (\ln n/n) + O((\ln n/n)^2)$ .

(b) Use (a) to show that  $\lim_{n \rightarrow \infty} n(n^{1/n} - 1) = \ln n$ .

**1.5.7.** [15] (a) Show that  $f(n) \neq \Omega(g(n))$  iff  $f(n) = o(g(n))$ .

(b) Show that  $f(n) = \Theta(g(n))$  or  $f(n) = o(g(n))$  iff  $f(n) = O(g(n))$ .

(c) Show that  $f(n) \neq O(g(n))$  iff  $f(n) = \Omega(g(n))$  and  $f(n) \neq \Theta(g(n))$ .

**1.5.8.** [HM45] Let  $\pi(n)$  denote the number of primes that do not exceed  $n$ . Show that

(a) a crude approximation is

$$\pi(n) \sim \frac{n}{\ln n};$$

(b) a better approximation is

$$\pi(n) = \frac{n}{\ln n} + \frac{n}{(\ln n)^2} + \frac{2!n}{(\ln n)^3} + \frac{3!n}{(\ln n)^4} + O\left(\frac{n}{(\ln n)^5}\right).$$

*Comments.* The displayed formulas are called *Prime Number Theorems*. Source: R. Graham, D.E. Knuth, O. Patashnik, *Concrete Mathematics*, Addison-Wesley, 1989.

## 1.6 Basics of Probability Theory

---

It is useful to recall briefly the basic notions of probability theory. The calculus of probabilities wants to study mathematical models of situations (experiments, observations) where the *outcome* is not deterministic but is determined by uncertain circumstances. The set of all possible outcomes is called the *sample space*, usually denoted by  $S$ , and an *event*  $E$  is a subset of  $S$ . The sample space  $S$  can be *countable*, which means that it is finite or countably infinite, or *continuous*, which means that it is uncountably infinite.

**Example 1.6.1** The throwing of two dice, one white and one black, gives a sample space  $S$  consisting of all pairs  $(i, j)$  where  $i$  is the number on the top face of the white die and  $j$  is the number on the top face of the black die. If  $A = \{(1, 3), (2, 2), (3, 1)\}$ , then  $A$  is the event that the sum of  $i$  and  $j$  is four. If  $B = \{(1, 1), (1, 2), (2, 1)\}$ , then  $B$  is the event that the sum of  $i$  and  $j$  is less than 4.  $\diamond$

Intuitively, the probability  $p$  of an event  $A$  is the apparent limit of the relative frequency of outcomes in  $A$  in the long run in a sequence of independent repetitions of the experiment. For instance, the probability associated with  $A$  in the example is  $\frac{1}{12}$ .

### 1.6.1 Kolmogorov Axioms

Let  $S$  denote the sample space. Following A.N. Kolmogorov's formalization of 1933, it is customary to use the following axioms.

- (A1) If  $A$  and  $B$  are events, then so is the *intersection*  $A \cap B$ , the *union*  $A \cup B$ , and the *difference*  $A - B$ .
- (A2) The *sample space*  $S$  is an event. We call  $S$  the *certain* event. The *empty set*, denoted by " $\emptyset$ ," is an event. We call  $\emptyset$  the *impossible* event.
- (A3) To each event  $E$  is assigned a nonnegative real number  $P(E)$  that we call the *probability* of event  $E$ .
- (A4)  $P(S) = 1$ .
- (A5) If  $A$  and  $B$  are disjoint, then  $P(A \cup B) = P(A) + P(B)$ .
- (A6) For a decreasing sequence

$$A_1 \supset A_2 \supset \cdots \supset A_n \supset \cdots$$

of events with  $\bigcap_n A_n = \emptyset$  we have  $\lim_{n \rightarrow \infty} P(A_n) = 0$ .

For systems with *finitely* many events Axiom A6 clearly follows from Axioms A1 through A5. For systems with *infinitely* many events, however, it is independent of the first five axioms. Therefore, Axiom A6 is only essential for systems with infinitely many events.

A system  $\mathcal{F}$  of sets  $f \subseteq S$  that is closed under the binary operations union, intersection, and difference, and contains a 1-element (here  $S$ ) and a 0-element (here  $\emptyset$ ) is called a (set) *field*. The set of events  $\mathcal{F}$  together with the associated set function  $P$ , also called a *measure* on  $\mathcal{F}$ , is called a *probability field* and is denoted by  $(\mathcal{F}, P)$ . It is easy to show that the axiom system A1 through A6 is *consistent* (free from contradictions). This is ascertained in the usual way by constructing an example that satisfies the axioms. Let  $S$  consist of a single element, and let the set of events be  $S$  and the empty event  $\emptyset$ , and set  $P(S) = 1$  and  $P(\emptyset) = 0$ . It is easy to verify that the system defined this way satisfies all the axioms above. However, the set of axioms is *incomplete*: for different problems in probability theory we have to construct different probability fields.

**Example 1.6.2** We call  $P$  a *probability distribution* over  $S$ . It follows from the axioms that  $0 \leq P(E) \leq 1$  for every event  $E$ ;  $P(\emptyset) = 0$ ; if  $A \subseteq B$ , then  $P(A) \leq P(B)$ ; if  $\bar{A}$  is the complement  $S - A$  of  $A$ , then  $P(\bar{A}) = 1 - P(A)$ ; and  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .  $\diamond$

### 1.6.2 Conditional Probability

If  $A$  and  $B$  are two events, with  $P(A) > 0$ , then the *conditional probability* that  $B$  occurs given that  $A$  occurs is defined by

$$P(B|A) = \frac{P(A \cap B)}{P(A)}.$$

It follows immediately that  $P(A \cap B) = P(A)P(B|A)$ , and by induction we find the *Multiplication Rule*:

$$P(A \cap B \cap \dots \cap N) = P(A)P(B|A) \dots P(N|A \cap B \cap \dots \cap M).$$

Consider the setup of Example 1.6.1 again. Let  $A$  be the event that at least one out of two dice shows an even number, and let  $B$  be the event that the sum of the numbers shown is even. Then  $P(B|A) = P(A \cap B)/P(A) = \frac{1}{3}$ .

The function  $P(\cdot|A)$  is a probability distribution on  $S$  and is called the *conditional probability distribution* given  $A$ . Clearly,  $P(A|A) = 1$ .

**Example 1.6.3** Rewriting  $P(A \cap B)$  as  $P(B)P(A|B)$  and as  $P(A)P(B|A)$ , substitution yields the formula that incorporates the essence of Bayes's Rule,

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}.$$

We require the *Law of Complete Probabilities*. Let  $A \cup B \cup \dots \cup N = S$  for disjoint events  $A, B, \dots, N$ , and let  $X$  be an arbitrary event. Then

$$P(X) = P(A)P(X|A) + P(B)P(X|B) + \dots + P(N)P(X|N).$$

If  $A, B, \dots, N$  are disjoint and  $X$  is arbitrary, it is easy to derive from the previous two displayed formulas the important *Bayes's Rule*:

$$P(Y|X) = \frac{P(Y)P(X|Y)}{P(A)P(X|A) + P(B)P(X|B) + \dots + P(N)P(X|N)},$$

where  $Y \in \{A, B, \dots, N\}$ . We often call  $A, B, \dots, N$  *hypotheses* and say that Bayes's Rule gives the *posterior* or *inferred* probability  $P(Y|X)$  of hypothesis  $Y$  after occurrence of event  $X$ . It is common to call  $P(Y)$  the *prior* probability (or a priori probability) of hypothesis  $Y$  before the occurrence of event  $X$ . In Sections 1.10 and 5.1.3 we analyze the meaning of Bayes's rule in detail, and in Chapter 5 we apply it extensively in inductive reasoning.  $\diamond$

**Example 1.6.4** The notion of mutual independence of two or more events lies at the heart of probability theory. From a mathematical viewpoint the given axioms specify just a special application of the general theory of additive sets. However, the special nature of this application is to a large part contained in the way we formalize the intuitive notion of mutual independence of events. In the following we assume that the events have positive probabilities. Events  $A$  and  $B$  are *mutually independent* iff  $P(A|B) = P(A)$  and  $P(B|A) = P(B)$ , in other words,  $P(A \cap B) = P(A)P(B)$ . And more generally for  $n > 2$  events, events  $A, B, \dots, N$  are mutually independent iff for all subsets of pairwise distinct  $X, Y, \dots, Z$  in  $A, B, \dots, N$  we have

$$P(X|Y \cap \dots \cap Z) = P(X).$$

In Example 1.6.1, we have  $P(A \cap B) \neq P(A)P(B)$ , and hence the events  $A$  and  $B$  are not independent.

The classical work on probability from Laplace to von Mises is essentially concerned with the investigation of sequences of independent events. For instance, in a sequence of throws of a fair coin the throws are treated as mutually independent events. (We do not consider that after a run of a hundred "heads" the chance on throwing "tails" has increased.) If in newer developments such as so-called Markov processes one often dispenses with complete independence, then still some weaker analogous requirements have to be imposed to obtain meaningful results.  $\diamond$

### 1.6.3 Continuous Sample Spaces

If a field is infinite, and additionally all countable unions  $\bigcup A_n$  of disjoint events  $A_n$  belong to it, then we call it a *Borel field* or  $\sigma$ -algebra in honor of E. Borel (1871–1956). We denote a Borel field by the Greek letter  $\sigma$ . It follows easily that in a Borel field  $\sigma$  all countable unions of not necessarily disjoint events also belong to  $\sigma$ , and the same holds

for countable intersections. The closure of  $\mathcal{F}$ , under the field operations and countable union together, gives the unique smallest Borel field that contains  $\mathcal{F}$ . Suppose  $(\mathcal{F}, P)$  is an infinite probability field. It is a fundamental result of measure theory that an extension of both  $\mathcal{F}$  and  $P$  under countable union and addition, respectively, preserving satisfaction of Axioms A1 through A6, is always possible and, moreover, unique. The result is called the *Borel extension*  $(\sigma, P^*)$  of the probability field. This is best illustrated by the construction of such an extension in a particular case.

**Example 1.6.5** We consider infinite binary sequences. The Borel probability field we aim at has as sample space  $S$  the real numbers in the half open interval  $[0, 1)$ . We start out by identifying a real  $\omega$  with its infinite binary expansion  $0.\omega_1\omega_2\dots$ . In case a real number has two representations, such as  $\frac{1}{2}$  which can be represented by  $0.100\dots$  and  $0.011\dots$ , we choose the representation with infinitely many zeros. A *cylinder*  $\Gamma_x$  consisting of all real numbers that start with  $0.x$ , where  $x$  is a finite binary string, is an event. The probability field  $(\mathcal{F}, P)$  is formed as follows. The set field  $\mathcal{F}$  is the closure of all events under pairwise union, intersection, and difference. It contains the impossible event (by  $\Gamma_0 \cap \Gamma_1 = \emptyset$ ) and the certain event  $\Gamma_\epsilon$ . The *uniform* distribution, or *Lebesgue measure*, usually denoted by  $\lambda$ , associates with each cylinder  $\Gamma_y$  a probability  $\lambda(\Gamma_y) = 2^{-l(y)}$ . By Axioms A1 through A5 all unions and intersections of pairs of cylinders are events, including the empty set, and have associated probabilities. We now consider the closure  $\sigma$  of  $\mathcal{F}$  under countably infinite union and the field operations. Then  $\sigma$  is the smallest Borel field containing  $\mathcal{F}$  as a subfield. Let  $A$  be an arbitrary subset of  $S$ . Define  $P^*(A)$  as the greatest lower bound on  $\sum_n P(A_n)$  for all *coverings*  $A \subseteq \bigcup_n A_n$  of  $A$  by  $A_1, A_2, \dots$  finitely many or countably infinitely elements from  $\mathcal{F}$ . It can be shown that for elements  $A$  in the original field  $\mathcal{F}$  we have  $P^*(A) = P(A)$ . We can also say that the probability distribution  $P(A)$  on the sets in  $\mathcal{F}$  associates the *measure*  $P(A)$  with  $A$ .  $\diamond$

**Example 1.6.6** Sample spaces can be discrete (natural numbers), countable (the rational numbers), or continuous (the real numbers). We will be interested in discrete versus continuous measures. The discrete measures we consider will have as sample space the natural numbers  $\mathcal{N}$  or, equivalently, the set of finite binary sequences  $\{0, 1\}^*$ . The continuous measures will have as sample space the real numbers  $\mathcal{R}$  or, equivalently, the set of one-way infinite binary sequences  $\{0, 1\}^\infty$ . Consider the continuous sample space  $S = \{0, 1\}^\infty$  with measure  $\mu$ . If  $\epsilon$  is the empty word, then  $\mu(\Gamma_\epsilon) = 1$  by Axiom A4, and for all  $x \in \{0, 1\}^*$ ,  $\mu(\Gamma_x) = \mu(\Gamma_{x0}) + \mu(\Gamma_{x1})$  by Axiom A5. For convenience, we will in Chapters 4 and 5 write  $\mu(x)$  instead of  $\mu(\Gamma_x)$  and consider a measure  $\mu$  as a function from the finite binary strings into the positive real numbers satisfying Axioms A4 and A5.  $\diamond$

**Example 1.6.7** A real-valued function on a sample space  $S$  is called a *random variable*. We denote random variables by  $X, Y, Z$ . A random variable maps an element from the sample space to an aspect of it we are interested in (or want to measure). For instance,  $S$  consists of the set of infinite binary sequences, and for each  $\omega = \omega_1\omega_2\dots$  in  $S$  the random variable  $X$  is defined as  $X(\omega) = \omega_2$ . We can also talk about a random variable  $X$  which is a finite vector  $X_1X_2\dots X_n$  or infinite vector  $X_1X_2\dots$  with  $X_i(\omega) = \omega_i$  for all  $i$ . If  $\omega$  is a sequence of outcomes of fair coin tosses, the measure on  $S$  is the Lebesgue measure  $\lambda$ , and  $\lambda\{\omega : X_i(\omega) = 0\} = \frac{1}{2}$  for all  $i$ . Justified by the definitions above, we call the random variables  $X_i$  *independent*. Another example of a random variable is  $Y_i(\omega) = k$ , where  $k$  is the length of the longest uninterrupted subsequence of zeros in  $\omega_{1:i}$ . Clearly, always either  $Y_i(\omega) = Y_{i+1}(\omega)$  or  $Y_i(\omega) < Y_{i+1}(\omega)$ , where both options occur. The random variables  $Y_i$  are *dependent*.

If  $P$  is a measure on  $S$ , then we customarily denote “ $P\{\omega : X(\omega) \leq x\}$ ” by the shorthand “ $P(X \leq x)$ .” The function  $F$  defined by  $F(x) = P(X \leq x)$  is called the *distribution function*. For instance, the random variable  $X$  defined as the outcome of a single throw of a fair die has distribution function  $P(X \leq i) = i/6$ ,  $i = 1, 2, \dots, 6$ . A random variable is *discrete* if the distribution function  $F$  is a step function. The domain of a discrete random variable consists of finitely many or countably infinitely many elements  $\{x_1, x_2, \dots\}$ . The function  $P(X = x_i)$ ,  $i = 1, 2, \dots$ , is the *probability (density) function*. The probability function associated with the outcome of a single throw of a fair die is  $P(X = i) = \frac{1}{6}$ . A random variable is *continuous* if the distribution function  $F$  has a continuous derivative  $f$  (at most discontinuous in finitely many points). This  $f$  is called the *probability density function*. For instance, if the distribution function of a random variable satisfies  $F(x) = 1 - e^{-\lambda x}$  for  $x > 0$ , and  $F(x) = 0$  for  $x \leq 0$ , then the density function is  $f(x) = \lambda e^{-\lambda x}$  for  $x > 0$ , and  $f(x) = 0$  for  $x \leq 0$ .  $\diamond$

---

## Exercises

**1.6.1.** [17] A random sample of size  $k$  is taken from a population of  $n$  elements. We draw the  $k$  elements one after the other and replace each drawn element in the population before drawing the next element. What is the probability of the event that in the sample no element occurs twice, that is, our sample could have been obtained also by sampling without replacement?

*Comments.*  $(n)_k/n^k$ . Source: W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. 1, Wiley, 1968.

**1.6.2.** [15] Consider the population of digits  $\{0, 1, \dots, 9\}$ . Use the formula you have found above to check that the probability that five consecutive random digits are all different is  $p = (10)_5/10^5 = 0.3024$ .

*Comments.* Source: W. Feller, *Ibid.*

**1.6.3.** [15] What is the probability that in a party of 23 people at least 2 people have a common birthday? Assume that birthdays are uniformly distributed over a year of 365 days and that the people at the party constitute a random sample. Use the formula derived above.

*Comments.* Note that contrary to intuition, it is better to bet that there will be shared birthdays than the other way. For a party of 23 the probability is close to 0.5. However, the analogous probability for a party of only 30 people already exceeds 0.7. Source: W. Feller, *Ibid.*

**1.6.4.** [10] Show that the probability of obtaining at least one ace (a 6) in four throws with one die is greater than the probability of obtaining at least one double ace in twenty-four throws with two dice.

*Comments.* This is known as *Chevalier de Méré's paradox*. It was posed by this passionate gambler to Pierre de Fermat, who wrote a solution in a letter to Blaise Pascal in 1654. It was solved earlier by Cardano (1501–1576). Source: [*Amer. Math. Monthly* 67(1960), 409–419.]

**1.6.5.** [08] Which probability is greater: to score at least one ace (a 6) in six throws of a die, or to score at least two aces in twelve throws of a die?

*Comments.* This question was submitted to Isaac Newton by the famous diarist Samuel Pepys in 1693. Newton answered that an easy computation shows that the first event has the greater probability, but failed to convince Pepys. Source: [*The Amer. Statistician*, 14(1960), 27–30.]

**1.6.6.** [M12] The *uniform* distribution over the countable sample space  $S = \mathcal{N}$  can be defined as the probability density function

$$L(x) = 2^{-2l(x)-1} \text{ or alternatively, } \frac{6}{\pi^2(l(x)+1)^2} 2^{-l(x)}.$$

(a) Show that in both cases  $\sum_{x \in S} L(x) = 1$ .

(b) Let  $S_1, S_2, \dots$  be a sequence of sample spaces with  $S_n = \{x : l(x) = n\}$ . Show that the probability density function  $L_n(x) = L(x|l(x) = n)$  assigns probability  $L_n(x) = 1/2^n$  to all  $x$  of length  $n$ , and zero probability to other  $x$ 's, for  $n = 1, 2, \dots$ .

**1.6.7.** [15] The *uniform* distribution  $\lambda$  over the continuous sample space  $S = \{0, 1\}^\infty$ , the set of one-way infinite binary sequences (the half-open interval of real numbers  $[0, 1)$ ), is described in Example 1.6.5. Let the  $L_n$ 's be as above. Show that  $\lambda(\Gamma_x) = L_{l(x)}(x)$ , for all  $x = \{0, 1\}^*$ .

## 1.7

### Basics of Computability Theory

---

While the qualitative ideas immanent in Kolmogorov complexity had been around for a long time, their ultimate applicability in quantitative form became possible only after the rise of computability theory (equivalently, recursive function theory) in the 1930s. In this section we develop and review completely the basic notions of that theory insofar as they are needed in the sequel.

In 1936 Alan M. Turing (1912–1954) exhibited an exceedingly simple type of hypothetical machine and gave a brilliant demonstration that everything that can be reasonably said to be computed by a human computer using a fixed procedure can be computed by such a machine. As Turing claimed, any process that can be naturally called an effective procedure is realized by a Turing machine. This is known as *Turing's Thesis*. Over the years, all serious attempts to give precise yet intuitively satisfactory definitions of a notion of “effective procedure” in the widest possible sense have turned out to be equivalent—to define essentially the same class of processes. (In his original paper, Turing established the equivalence of his notion of “effective procedure” with Alonzo Church's (1903–1995) notion of “effective calculability.”) *Church's Thesis* states that, in this sense, there is an objective notion of effective computability independent of a particular formalization.

While the formal part of Turing's paper is difficult to follow for the contemporary reader, the informal arguments he sets forth are as lucid and convincing now as they were then. To us it seems that it is the best introduction to the subject, and we reproduce this superior piece of expository writing below.

All arguments [for Turing's Thesis] are bound to be, fundamentally, appeals to intuition, and for that reason rather unsatisfactory mathematically. The real question at issue is: “what are the possible processes which can be carried out in computing (a number)?” The arguments which I shall use are of three kinds.

- (a) A direct appeal to intuition.
- (b) A proof of equivalence of two definitions (in case the new definition has a greater intuitive appeal).
- (c) Giving examples of large classes of numbers which are computable.

Once it is granted that computable numbers are all “computable [by Turing machines],” several other propositions of the same character follow. In particular it follows that, if there is a general process for determining whether a formula (of the Hilbert function calculus) is provable, then the determination can be carried out by machine. [...]

Computing is normally done by writing certain symbols on paper. We may suppose this paper to be divided into squares like a child's arithmetic book. In elementary arithmetic the 2-dimensional character of the paper is sometimes used. But such use is always avoidable, and I think it will be agreed that the

two-dimensional character of paper is no essential of computation. I assume then that the computation is carried out on one-dimensional paper, on a tape divided into squares. I also suppose that the number of symbols which may be printed is finite. If we were to allow an infinity of symbols, then there would be symbols differing to an arbitrary small extent.<sup>1</sup>

The effect of this restriction on the number of symbols is not very serious. It is always possible to use sequences of symbols in the place of single symbols. Thus an Arabic numeral such as 17 or 99999999999999 is normally treated as a single symbol. Similarly in any European language words are treated as single symbols (Chinese, however, attempts to have an enumerable infinity of [atomic] symbols). The differences in our point of view between single and compound symbols is that the compound symbols, if they are too lengthy, cannot be observed at one glance. This is in accordance with experience. We cannot tell at a glance whether 9999999999999999 and 9999999999999999 are the same. [...]

The behaviour of the [human] computer at any moment is determined by the symbols he is observing, and his 'state of mind' at that moment. We may suppose that there is a bound  $B$  to the number of symbols or squares which the computer can observe at one moment. If he wishes to observe more, he must use successive observations. We will also suppose that the number of states of mind which need be taken into account is finite. The reasons for this are of the same character as those which restrict the number of symbols. If we admit an infinity of states of mind, some of them will be 'arbitrarily close' and will be confused. Again, the restriction is not one which seriously affects computation, since the use of more complicated states of mind can be avoided by writing more symbols on the tape. [...]

Let us imagine the operations performed by the computer to be split up in 'simple operations' which are so elementary that it is not easy to imagine them further divided. Every such operation consists of some change of the physical system consisting of the computer and his tape. We know the state of the system if we know the sequence of symbols on the tape, which of these are observed by the computer (possibly with a special order), and the state of mind of the computer. We may suppose that in a simple operation not more than one symbol is altered. Any other changes can be split up into simple changes of this kind. The situation in regard to the squares whose symbols may be altered this way is the same as in regard to the observed squares. We may, therefore, without loss of generality, assume that the squares whose symbols are changed are the 'observed' squares. [...]

---

<sup>1</sup>If we regard a symbol as literally printed on a square, we may suppose that the square is  $0 < x < 1$ ,  $0 < y < 1$ . The symbol is defined as the set of points in this square, viz., the set occupied by printer's ink. If these sets are restricted to be measurable, we can define the 'distance' between two symbols as the cost of transforming one symbol into the other if the cost of moving a unit area of printer's ink unit distance is unity, and there is an infinite supply of ink at  $x = 2$ ,  $y = 0$ . With this topology the symbols form a conditionally compact space [Turing's note].

Besides these changes of symbols, the simple operations must include changes of distribution of observed squares. The new observed squares must immediately be recognized by the computer. I think it is reasonable to suppose that they can only be squares whose distance from the closest of the immediately previously observed squares does not exceed a certain fixed amount. Let us say that each of the new observed squares is within  $L$  squares of the immediately previously observed square. [...]

In connection to ‘immediate recognizability,’ it may be thought that there are other kinds of squares which are immediately recognizable. In particular, squares marked by special symbols may be taken as immediately recognizable. Now if these squares are marked only by single symbols there can be only a finite number of them, and we should not upset our theory by adjoining these marked squares to the observed squares. If, on the other hand, they are marked as a sequence of symbols, we cannot regard the process of recognition as a simple process. This is a fundamental point and should be illustrated. In most mathematical papers the equations and theorems are numbered. Normally the numbers go not beyond (say) 1000. It is, therefore, possible to recognize a theorem at a glance by its number. But if the paper was very long, we might reach Theorem 157767733443477; then, further on in the paper, we might find ‘... hence (applying Theorem 157767733443477) we have ...’ In order to make sure which was the relevant theorem we should have to compare the two numbers figure by figure, possibly ticking the figures off in pencil to make sure of their not being counted twice. If in spite of this it is still thought that there are other ‘immediately recognizable’ squares, it does not upset my contention so long as these squares can be found by some process of which my type of machine is capable. [...]

The simple operations must therefore include:

- (a) Changes of the symbol on one of the observed squares.
- (b) Changes of one of the squares observed to another square within  $L$  squares of one of the previously observed squares.

It may be that some of these changes necessarily involve a change of state of mind. The most general single operation must therefore be taken to be one of the following:

- (A) A possible change (a) of symbol together with a possible change of state of mind.
- (B) A possible change (b) of observed squares, together with a possible change of state of mind.

The operation actually performed is determined, as has been suggested [above] by the state of mind of the computer and the observed symbols. In particular, they determine the state of mind of the computer after the operation. [...]

We may now construct a machine to do the work of this computer. To each state of mind of the computer corresponds an ‘ $m$ -configuration’ of the machine. The machine scans  $B$  squares corresponding to the  $B$  squares observed by the computer. In any move the machine can change a symbol on a scanned square or can change any one of the scanned squares to another square distant not more than  $L$  squares from one of the other scanned squares. [Without loss

of generality restrict  $L$  and  $B$  to unity.] The move which is done, and the succeeding configuration, are determined by the scanned symbol and the  $m$ -configuration. The machines just described do not differ very essentially from computing machines as defined (previously) and corresponding to any machine of this type a computing machine can be constructed to compute the same sequence, that is to say, the sequence computed by the computer. [...]

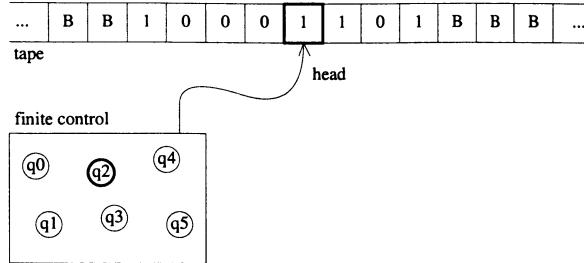
We suppose [above] that the computation is carried out on a tape; but we avoid introducing the “state of mind” by considering a more physical and definitive counterpart of it. It is always possible for the computer to break off from his work, to go away and forget all about it, and later to come back and go on with it. If he does this he must leave a note of instructions (written in some standard form) explaining how the work is to be continued. This note is the counterpart of “the state of mind.” We will suppose that the computer works in such a desultory manner that he never does more than one step and write the next note. Thus the state of progress of the computation at any stage is completely determined by the note of instructions and the symbols on the tape. That is, the state of the system may be described by a single expression (sequence of symbols), consisting of the symbols on the tape followed by a special marker (which we suppose not to appear elsewhere) and then by the note of instructions. This expression may be called the “state formula.” We know that the state formula at any given stage is determined by the state formula before the last step was made, and we assume that the relation of these two formulae is expressible in the functional calculus. In other words, we assume that there is an axiom  $A$  which expresses the rules governing the behaviour of the computer, in terms of the relation of the state formula at any stage to the state formula at the preceding stage. If this is so, we can construct a machine to write down the successive state formulae, and hence to compute the required number.

### 1.7.1 Effective Enumerations and Universal Machines

We formalize Turing’s description as follows: A *Turing machine* consists of a finite program, called the *finite control*, capable of manipulating a linear list of *cells*, called the *tape*, using one access pointer, called the *head* (Figure 1.1). We refer to the two directions on the tape as *right* and *left*. The finite control can be in any one of a finite set of states  $Q$ , and each tape cell can contain a 0, a 1, or a *blank*  $B$ . Time is discrete and the time instants are ordered  $0, 1, 2, \dots$ , with 0 the time at which the machine starts its computation. At any time, the head is positioned over a particular cell, which it is said to *scan*. At time 0 the head is situated on a distinguished cell on the tape called the *start cell*, and the finite control is in a distinguished state  $q_0$ . At time 0 all cells contain  $B$ ’s, except for a contiguous finite sequence of cells, extending from the start cell to the right, which contain 0’s and 1’s. This binary sequence is called the *input*.

The device can perform the following basic operations:

1. it can write an element from  $A = \{0, 1, B\}$  in the cell it scans; and



**FIGURE 1.1.** Turing machine

2. it can shift the head one cell left or right.

When the device is active it executes these operations at the rate of one operation per time unit (a *step*). At the conclusion of each step, the finite control takes on a state out of  $Q$ . The device is constructed so that it behaves according to a finite *list of rules*. These rules determine from the current state of the finite control and the symbol contained in the cell under scan, the operation to be performed next and the state to enter at the end of the next operation execution.

The rules have format  $(p, s, a, q)$ :  $p$  is the current state of the finite control;  $s$  is the symbol under scan;  $a$  is the next operation to be executed of type (1) or (2) designated in the obvious sense by an element from  $S = \{0, 1, B, L, R\}$ ; and  $q$  is the state of the finite control to be entered at the end of this step.

Any two distinct quadruples must differ in the first two elements: the device is *deterministic*. Not every possible combination of the first two elements has to be in the set; in this way we permit the device to perform *no* operation. In this case we say the device *halts*. Hence, we can define a Turing machine by a mapping from a finite subset of  $Q \times A$  into  $S \times Q$ . Given a Turing machine and an input, the Turing machine carries out a uniquely determined succession of operations, which may or may not terminate in a finite number of steps.

We can associate a partial function with each Turing machine in the following way: The input to the Turing machine is presented as an  $n$ -tuple  $(x_1, \dots, x_n)$  of binary strings in the form of a single binary string consisting of self-delimiting versions of the  $x_i$ 's. The integer represented by the maximal binary string (bordered by blanks) of which some bit is scanned, or "0" if a blank is scanned, by the time the machine halts is called the *output* of the computation.

**Definition 1.7.1** Under this convention for inputs and outputs, each Turing machine defines a partial function from  $n$ -tuples of integers onto the integers,  $n \geq 1$ .

We call such a function *partial recursive* or *computable*. If the Turing machine halts for all inputs, then the function computed is defined for all arguments and we call it *total recursive*, or simply *recursive*.

We call a function with range  $\{0, 1\}$  a *predicate*, with the interpretation that the predicate of an  $n$ -tuple of values is “true” if the corresponding function assumes value 1 for that  $n$ -tuple of values for its arguments and is “false” or “undefined” otherwise. Hence, we can talk about *partial (total) recursive predicates*.

**Example 1.7.1** Consider  $x$  as a binary string. It is easy to see that the functions  $l(x)$ ,  $f(x) = \bar{x}$ ,  $g(\bar{x}y) = x$ , and  $h(\bar{x}y) = y$  are partial recursive. Functions  $g$  and  $h$  are not total since the value for input “1111” is not defined. The function  $g'(\bar{x}y)$  defined as 1 if  $x = y$  and as 0 if  $x \neq y$  is a recursive predicate. Consider  $x$  as an integer. The following functions are basic  $n$ -place total recursive functions: the *successor* function  $\gamma^{(1)}(x) = x + 1$ , the *zero* function  $\zeta^{(n)}(x_1, \dots, x_n) = 0$ , and the *projection* function  $\pi_m^{(n)}(x_1, \dots, x_n) = x_m$  ( $1 \leq m \leq n$ ). ◇

**Example 1.7.2** The function  $\langle x, y \rangle = \bar{x}y$  is a total recursive one-to-one mapping from  $\mathcal{N} \times \mathcal{N}$  into  $\mathcal{N}$ . We can easily extend this scheme to obtain a total recursive one-to-one mapping from  $k$ -tuples of integers into the integers, for each fixed  $k$ . Define  $\langle n_1, n_2, \dots, n_k \rangle = \langle n_1, \langle n_2, \dots, n_k \rangle \rangle$ .

Another total recursive one-to-one mapping from  $k$ -tuples of integers into the integers is  $\langle n_1, n_2, \dots, n_k \rangle = \bar{n}_1 \dots \bar{n}_{k-1} \bar{n}_k$ . ◇

**Church’s Thesis.** *The class of algorithmically computable numerical functions (in the intuitive sense) coincides with the class of partial recursive functions.*

Originally intended as a *proposal* to henceforth supply intuitive terms like “computable” and “effective procedure” with a precise meaning as “recursive” and “recursive function,” it has come into use as shorthand for a claim that from a given description of a procedure in terms of an informal set of instructions we can derive a formal one in terms of Turing machines.

It is possible to give an effective (computable) one-to-one pairing between natural numbers and Turing machines. This is called an *effective enumeration*. One way to do this is to encode the table of rules of each Turing machine in binary, in a canonical way.

**Example 1.7.3** The only thing we have to do for every Turing machine, is to encode the defining mapping  $T$  from  $Q \times A$  into  $S \times Q$ . Giving each element of

$Q \cup S \cup A$  a unique binary code requires  $s$  bits for each such element, with  $s \leq \lceil \log(d(Q) + 5) \rceil$ . Denote the encoding function by  $e$ . Then the quadruple  $(p, 0, B, q)$  is encoded as  $e(p)e(0)e(B)e(q)$ . If the number of rules is  $r$ , then  $r \leq 3d(Q)$ . We agree to consider the state of the first rule as the start state. The entire list of quadruples,

$$T = (p_1, t_1, s_1, q_1), (p_2, t_2, s_2, q_2), \dots, (p_r, t_r, s_r, q_r),$$

is encoded as

$$E(T) = \bar{s}\bar{r}e(p_1)e(t_1)e(s_1)e(q_1)\dots e(p_r)e(t_r)e(s_r)e(q_r).$$

Note that  $l(E(T)) \leq 4rs + 2\log rs + 4$ . (Moreover,  $E$  is self-delimiting, which is convenient in situations where we want to recognize the substring  $E(T)$  as prefix of a larger string; see Section 1.4.)

We order the resulting binary strings lexicographically (according to increasing length). We assign an *index*, or Gödel number,  $n(T)$  to each Turing machine  $T$  by defining  $n(T) = i$  if  $E(T)$  is the  $i$ th element in the lexicographic order of Turing machine codes. This yields a sequence of Turing machines  $T_1, T_2, \dots$  that constitutes the effective enumeration. One can construct a Turing machine to decide whether a given binary string  $x$  encodes a Turing machine, by checking whether it can be decoded according to the scheme above, followed by a check whether any two different rules start with the same two elements. This observation enables us to construct “universal” Turing machines.  $\diamond$

A *universal* Turing machine  $U$  is a Turing machine that can imitate the behavior of any other Turing machine  $T$ . It is a fundamental result that such machines exist and can be constructed effectively. Only a suitable description of  $T$ 's finite program and input needs to be entered on  $U$ 's tape initially. To execute the consecutive actions that  $T$  would perform on its *own* tape,  $U$  uses  $T$ 's description to simulate  $T$ 's actions on a representation of  $T$ 's tape contents. Such a machine  $U$  is also called “computation universal.” In fact, there are infinitely many such  $U$ 's.

- Example 1.7.4** We focus on a universal Turing machine  $U$  that uses the encoding above. It is not difficult, but tedious, to define a Turing machine in quadruple format that expects inputs of the format  $1^i0p$  and acts as follows: On the tape left of the input,  $U$  starts to generate the successive strings in  $\{0, 1\}^*$  in lexicographic order according to increasing length and checks for each such string whether or not it encodes a Turing machine. For each Turing machine it finds, it replaces one 1 in  $1^i$  in the input by a  $B$ . The binary string  $E(T)$  that causes the delimiter 0 to be read encodes  $T$  with  $n(T) = i$ . Subsequently,  $U$  starts to execute the successive operations of  $T$  using  $p$  as input and the description  $E(T)$  of  $T$  it has found. We omit the explicit construction of  $U$ .  $\diamond$

For the contemporary reader there should be nothing mysterious in the concept of a general purpose computer which can perform any computation when supplied with an appropriate program. The surprising thing is that a general-purpose computer can be *very simple*: it has been shown that 4 tape symbols and 7 states suffice easily in the above scheme [M. Minsky, *Computation: Finite and Infinite Machines*, Prentice-Hall, 1967]. This machine can be changed to, in the sense of being simulated by, our format using tape symbols  $\{0, 1, B\}$  at the cost of an increase of states; see Section 1.12.

The effective enumeration of Turing machines  $T_1, T_2, \dots$  determines an effective enumeration of partial recursive functions  $\phi_1, \phi_2, \dots$  such that  $\phi_i$  is the function computed by  $T_i$ , for all  $i$ . It is important to distinguish between a *function*  $\psi$  and a *name* for  $\psi$ . A name for  $\psi$  can be an *algorithm* that computes  $\psi$ , in the form of a Turing machine  $T$ . It can also be a natural number  $i$ , such that  $\psi$  equals  $\phi_i$  in the above list. We call  $i$  an *index* for  $\psi$ . Thus, each partial recursive  $\psi$  occurs many times in the given effective enumeration, that is, it has many indices.

**Definition 1.7.2** The partial recursive function  $\nu^{(2)}(i, x)$  computed by the universal Turing machine  $U$  is called the *universal partial recursive function*.

The generalization to  $n$ -place functions is straightforward. A partial recursive function  $\nu^{(n+1)}(i, x_1, \dots, x_n)$  is *universal* for all  $n$ -place partial recursive functions, if for each partial recursive function  $\phi^{(n)}(x_1, \dots, x_n)$  there exists an  $i$  such that the mapping  $\nu^{(n+1)}$  with the first argument fixed to  $i$  is identical to the mapping  $\phi^{(n)}$ . Here  $i$  is an index of  $\phi^{(n)}$  with respect to  $\nu^{(n+1)}$ . For each  $n$ , we fix a partial recursive  $(n + 1)$ -place function that is universal for all  $n$ -place partial recursive functions. The following lemma is usually called the *Enumeration Theorem* for  $n$ -place partial recursive functions. Here  $z$  is the index of the universal function.

**Lemma 1.7.1** *For each  $n$  there exists a  $z$  such that for all  $i$  and  $x_1, \dots, x_n$ , if  $\phi_i^{(n)}(x_1, \dots, x_n)$  is defined, then  $\phi_z^{(n+1)}(i, x_1, \dots, x_n) := \phi_i^{(n)}(x_1, \dots, x_n)$ , and  $\phi_z^{(n+1)}(i, x_1, \dots, x_n)$  is undefined otherwise. ( $\phi_z^{(n+1)}$  is a universal partial recursive function that enumerates the partial recursive functions of  $n$  variables.)*

**Proof.** The machine of Example 1.7.4 is universal for all  $n$ -place partially recursive functions, for all  $n$ .  $\square$

We say that a partial recursive function  $\phi^{(n)}$  on arguments  $x_1, \dots, x_n$  is computed in *time*  $t$  ( $t$  steps or  $t$  operations) by Turing machine  $T$  if  $T$  computes  $\phi^{(n)}$  and  $T$  halts within  $t$  steps when it is started on the input corresponding to these arguments.

**Lemma 1.7.2** *For all partial recursive functions  $\phi^{(n)}$  there exists a recursive function  $\psi^{(n+1)}$  such that  $\psi^{(n+1)}(t, x_1, \dots, x_n) = 1$  if  $\phi^{(n)}(x_1, \dots, x_n)$  can be computed in not more than  $t$  steps, and  $\psi^{(n+1)}(t, x_1, \dots, x_n) = 0$  otherwise.*

**Proof.** Let  $T$  be a Turing machine that computes  $\phi$ . Modify  $T$  to a Turing machine  $T'$  that computes  $\psi$  from an input with one extra variable, the *clock*, containing the natural number  $t$ . Machine  $T'$  works exactly like  $T$ , except that it decrements the clock by one for every simulated step of  $T$ . If simulator  $T'$  enters a halting configuration of  $T$  with a positive clock value, then  $T'$  outputs 1. If  $T'$  decrements its clock to zero, then it halts and outputs 0.  $\square$

A set  $A$  is *recursively enumerable* if it is empty or the range of some total recursive function  $f$ . We say that  $f$  *enumerates*  $A$ . The intuition behind this definition is that there is a Turing machine for listing the elements of  $A$  in some arbitrary order with repetitions allowed. An equivalent definition is,  $A$  is recursively enumerable iff it is accepted by a Turing machine. That is, for each element in  $A$  the Turing machine halts in a distinguished accepting state, and for each element not in  $A$  the machine either halts in a nonaccepting state or computes forever.

A set  $A$  is *recursive* if it possesses a recursive characteristic function. That is,  $A$  is recursive iff there exists a recursive function  $f$  such that for all  $x$ , if  $x \in A$ , then  $f(x) = 1$ , and if  $x \in \bar{A}$ , then  $f(x) = 0$  ( $\bar{A}$  is the complement of  $A$ ). An equivalent definition is:  $A$  is recursive iff  $A$  is accepted by a Turing machine that always halts. Obviously, all recursive sets are recursively enumerable.

**Example 1.7.5** The following sets are recursive: (i) the set of odd integers; (ii) the set of natural numbers and the empty set; (iii) the set of primes; (iv) any finite set; (v) any set with a finite complement. The following sets are recursively enumerable: (i) any recursive set; (ii) the set of indices  $i$  such that the range of  $\phi_i$  is nonempty; (iii) if  $\pi = 3.1415\dots$ , the set  $\{x : \text{a run of at least } x \text{ consecutive } 0\text{'s occurs in } \pi\}$ ; (iv) the set of  $x$  such that there are positive integer 4-tuples  $(n, x, y, z)$  that satisfy  $x^n + y^n = z^n$ ,  $n > 2$ .  $\diamond$

**Lemma 1.7.3** (i) *A set  $A$  is recursive iff both  $A$  and its complement  $\bar{A}$  are recursively enumerable.*

(ii) *An infinite set  $A$  is recursive iff it is recursively enumerable in increasing order. (Here we have postulated a total order on the elements of  $A$ . For instance, if  $A \subseteq \mathbb{N}$  with the usual order, then  $\phi$  enumerates  $A$  in increasing order if  $\phi(i) < \phi(i + 1)$ , for all  $i$ .)*

**Proof.** Let  $A \cup \bar{A} = \mathcal{N}$  and neither  $A$  or  $\bar{A}$  is the empty set. (Otherwise the lemma trivially holds.)

(i) (**ONLY IF**) If  $A$  is recursive, then it is recursively enumerable. Let  $f$  be the characteristic function of  $A$ . By assumption,  $f$  is total recursive with range  $\{0, 1\}$ . Let the list  $a_1, a_2, \dots$  be the ordered sequence of values such that  $f(a_i) = 0$ . That is,  $A = \{a : a \text{ is not in the list}\}$ . Define a function  $g$  such that  $g(i) = a_i$  for all  $i$ . The function  $g$  is clearly recursive. Hence, the complement  $\bar{A}$  of  $A$  is recursively enumerable (in fact, recursive).

(**IF**) By assumption,  $A$  is the range of  $f$  and  $\bar{A}$  is the range of  $g$ , for two recursive functions  $f, g$ . Therefore, we can generate the list  $f(0), g(0), f(1), g(1), f(2), \dots$ , and examine each element in turn. To decide whether a given  $x$  in  $\mathcal{N}$  belongs to  $A$  we just compare it with each element in the above list. If  $x = f(i)$  for some  $i$ , then  $x \in A$ . If  $x = g(i)$  for some  $i$ , then  $x \in \bar{A}$ . Element  $x$  must occur in this list, since  $A \cup \bar{A} = \mathcal{N}$ .

(ii) (**ONLY IF**) Trivial.

(**IF**) Assume  $A$  is recursively enumerated in increasing order  $a_1, a_2, \dots$ . Then this yields a procedure to decide for each  $x$  whether or not  $x$  belongs to  $A$  by examining the first  $i$  elements of the list, for  $i$  determined by  $a_{i-1} < x \leq a_i$ .  $\square$

**Lemma 1.7.4** *Every infinite recursively enumerable set contains an infinite recursive subset.*

**Proof.** Let  $A$  be infinite and recursively enumerable. Let  $f$  be a recursive function with range  $A$ . Define a new recursive function  $g$  with  $g(0) = f(0)$  and  $g(x+1) = f(y)$ , where  $y$  is the least value such that  $f(y) > g(x)$ . Let  $B$  be the range of  $g$ . Since  $A$  is infinite, the definition of  $B$  implies that  $B$  is infinite. Clearly  $g$  enumerates  $B$  in increasing order, so  $B$  is recursive by Lemma 1.7.3, Item (ii).  $\square$

The equivalent lemmas hold for recursive and recursively enumerable sets of  $n$ -tuples.

### 1.7.2 Undecidability of the Halting Problem

Turing's paper, and more so Kurt Gödel's paper, where such a result first appeared, are celebrated for showing that certain well-defined questions in the mathematical domain cannot be settled by *any* effective procedure for answering questions. One such question is "which machine computations eventually terminate with a definite result, and which machine computations go on forever without a definite conclusion?" This is sometimes called the *halting problem*. Since all machines can be simulated by the universal Turing machine  $U$ , this question cannot be decided in the case of the single machine  $U$ , or more generally for *any* individual universal machine. The following lemma, due to Turing in 1936, formalizes

this discussion. Let  $\phi_1, \phi_2, \dots$  be the standard enumeration of partial recursive functions.

**Lemma 1.7.5** *There is no recursive function  $g$  such that for all  $x, y$ , we have  $g(x, y) = 1$  if  $\phi_x(y)$  is defined, and  $g(x, y) = 0$  otherwise.*

**Proof.** Suppose the contrary, and define a partial recursive function  $\psi$  by  $\psi(x) = 1$  if  $g(x, x) = 0$ , and  $\psi(x)$  is undefined otherwise. (The definition of  $\psi$  gives by the assumption of total recursiveness of  $g$  an algorithm, and by Church's Thesis or by explicit construction we can find a Turing machine to compute  $\psi$ .) Let  $\psi$  have an index  $y$  in the fixed enumeration of partial recursive functions,  $\psi = \phi_y$ . Then,  $\phi_y(y)$  is defined iff  $g(y, y) = 0$ , according to  $\psi$ 's definition. But this contradicts the assumption of existence of  $g$  as defined in the statement of the lemma.  $\square$

The trick used in this proof is called *diagonalization*; see Exercise 1.7.3, page 40. Define the *halting set*  $K_0 = \{\langle x, y \rangle : \phi_x(y) < \infty\}$ . Then Lemma 1.7.5 can be rephrased as, “The halting set  $K_0$  is not recursive.” It is easy to see that  $K_0$  is recursively enumerable. The halting set is so ubiquitous that it merits the standard notation “ $K_0$ .” We shall also use the *diagonal halting set*  $K = \{x : \phi_x(x) < \infty\}$ . Just like  $K_0$ , the diagonal halting set is recursively enumerable; and the proof of Lemma 1.7.5 shows that  $K$  is not a recursive set.

Lemma 1.7.5 was preceded by the famous (first) incompleteness theorem of Kurt Gödel in 1931. Recall that a formal theory  $T$  consists of a set of well-formed formulas, *formulas* for short. For convenience these formulas are taken to be finite binary strings. Invariably, the formulas are specified in such a way that an effective procedure exists that decides which strings are formulas and which strings are not.

The formulas are the objects of interest of the theory and constitute the “meaningful” statements. With each theory we associate a set of *true* formulas and a set of *provable* formulas. The set of true formulas is “true” according to some (often nonconstructive) criterion of truth. The set of provable formulas is “provable” according to some (usually effective) syntactic notion of proof.

A *theory*  $T$  is simply any set of formulas. A theory is *axiomatizable* if it can be effectively enumerated. For instance, its axioms (initial formulas) can be effectively enumerated and there is an effective procedure that enumerates all proofs for formulas in  $T$  from the axioms. A theory is *decidable* if it is a recursive set. A theory  $T$  is *consistent* if not both formula  $x$  and its negation  $\neg x$  are in  $T$ . A theory  $T$  is *sound* if each formula  $x$  in  $T$  is true (with respect to the standard model of the natural numbers).

Hence, soundness implies consistency. A particularly important example of an axiomatizable theory is *Peano Arithmetic*, which axiomatizes the standard elementary theory of the natural numbers.

**Lemma 1.7.6** *There is a recursively enumerable set  $K_0$  such that for every axiomatizable theory  $T$  that is sound and extends Peano Arithmetic, there is a number  $n$  such that the formula “ $n \notin K_0$ ” is true but not provable in  $T$ .*

**Proof.** Let  $\phi_1, \phi_2, \dots$  be the standard enumeration of partial recursive functions above, and let  $K_0 = \{\langle x, y \rangle : \phi_x(y) < \infty\}$ . That is,  $K_0$  is the domain of a universal partial recursive function, Lemma 1.7.1. It is easy to see that the set  $K_0$  is recursively enumerable. Moreover, all true statements of the form “ $n \in K_0$ ” belong to Peano Arithmetic.

Assume by way of contradiction that all true statements of the form “ $n \notin K_0$ ” are provable in  $T$ . Then the complement of  $K_0$  is recursively enumerable by enumerating the set of all provable statements in  $T$ . By Lemma 1.7.3, Item (i), if both  $K_0$  and its complement are recursively enumerable, then  $K_0$  is recursive. However, this contradicts the fact that  $K_0$  is nonrecursive (Lemma 1.7.5).  $\square$

In his original proof Gödel uses diagonalization to prove the incompleteness of any sufficiently rich logical theory  $T$  with a recursively enumerable axiom system, such as Peano Arithmetic. By his technique he exhibits for such a theory an explicit construction of an *undecidable statement*  $y$  that says of itself “I am unprovable in  $T$ .” The formulation in terms of recursive function theory, Lemma 1.7.6, is due to A. Church and S.C. Kleene. In the proof diagonalization is needed to show that  $K_0$  is not recursive. And elaboration of this proof would yield a similar explicit construction to the above one. Using Kolmogorov complexity we will be able to derive a new proof of Lemma 1.7.6, with essentially different examples of undecidable statements.

### 1.7.3 Enumerable Functions

We discuss recursive real-valued functions and enumerable functions. These notions will play a prominent role throughout this book. Recall that  $\mathcal{N}$ ,  $\mathcal{Q}$ , and  $\mathcal{R}$ , denote the nonnegative integers, the rational numbers, and the real numbers, respectively. We consider recursive functions  $g(\langle x, k \rangle) = \langle p, q \rangle$  and write  $g(x, k) = p/q$ . The interpretation is that  $g$  is a *rational-valued* function of two nonnegative integer arguments.

**Definition 1.7.3** A real function  $f$  is *enumerable* if there exists a recursive function  $g(x, k)$ , nondecreasing in  $k$ , with  $f(x) = \lim_{k \rightarrow \infty} g(x, k)$ . A function  $f$  is *co-enumerable*, if  $-f$  is enumerable. A real function  $f$  is *recursive* iff there is a recursive function  $g(x, k)$  such that  $|f(x) - g(x, k)| < 1/k$ .

This way we have extended the notion of integer-valued recursive function to real-valued recursive function and real-valued enumerable function. The idea is that enumerable functions can be approximated from one side by a recursive function over the natural numbers, but we may never know how close we are to the real value. Recursive functions can be approximated to any degree of precision by a recursive function over the natural numbers.

**Example 1.7.6** The following properties are easily proven: A function  $f : \mathcal{N} \rightarrow \mathcal{R}$  is enumerable if the set  $\{(x, r) : r \leq f(x), r \in \mathbb{Q}\}$  is recursively enumerable. Therefore, an enumerable function is “recursively enumerable from below,” and a co-enumerable function is “recursively enumerable from above.” In the literature we also meet “semicomputable from below” for “enumerable,” “semicomputable from above” for “co-enumerable,” and “computable” for “recursive.” A function is recursive iff it is both enumerable and co-enumerable. Not all enumerable or co-enumerable functions are recursive.  $\diamond$

**Example 1.7.7** We give an example of an enumerable function that is not recursive. Let  $K = \{x : \phi_x(x) < \infty\}$  be the diagonal halting set. Define  $f(x) = 1$  if  $x \in K$ , and  $f(x) = 0$  otherwise. We first show that  $f(x)$  is enumerable. Define  $g(x, k) = 1$  if the Turing machine computing  $\phi_x$  halts in at most  $k$  steps on input  $x$ , and  $g(x, k) = 0$  otherwise. Obviously,  $g$  is a rational-valued recursive function. Moreover, for all  $x$  and  $k$  we have  $g(x, k+1) \geq g(x, k)$ , and  $\lim_{k \rightarrow \infty} g(x, k) = f(x)$ . Hence,  $f$  is enumerable. However, if  $f(x)$  were recursive, then the set  $\{x : f(x) = 1\}$ , that is the diagonal halting set  $K$ , would be recursive. But we have shown that it is not.  $\diamond$

**Example 1.7.8** In Section 1.6 we have defined the notion of measure functions as  $\mu$  functions which map subsets of the real interval  $[0, 1)$  to  $\mathcal{R}$ . Considering  $[0, 1)$  as the isomorphic  $S = \{0, 1\}^\infty$ , such functions are defined by the values  $\mu(\Gamma_x)$  where  $\Gamma_x = \{x\omega : x \in \{0, 1\}^n, \omega \in \{0, 1\}^\infty\}$  are the so-called cylinder sets. We can extend the notions of recursiveness to set functions. A *measure*  $\mu$  is *recursive* (*enumerable*) iff the function  $f : \mathcal{N} \rightarrow \mathcal{R}$  defined by  $f(x) = \mu(\Gamma_x)$  is recursive (*enumerable*).  $\diamond$

This is enough on enumerable functions to tide us over Chapters 2 and 3, and we go into more detail in Chapter 4.

### 1.7.4 Feasible Computations

We give a brief introduction to computational complexity theory in order to provide some basic ideas and establish terminology required for applications of Kolmogorov complexity in Chapters 6, 7, and 8.

Theoretically, any recursive function is computable by an IBM PC or by a Turing machine as shown in Figure 1.1. But a computation that takes

$2^n$  steps on an input of length  $n$  would not be regarded as *practical* or *feasible*. No computer would ever finish such a computation in the lifetime of the universe even with  $n$  merely 1000. Computational complexity theory tries to identify problems that are feasibly computable.

If we have  $10^9$  processors taking  $10^9$  steps/second, then we can execute  $3.1 \times 10^{25} < 2^{100}$  steps/year.

In computational complexity theory, we are often concerned with languages. A *language* over a finite alphabet  $\Sigma$  is simply a subset of  $\Sigma^*$ . We say a Turing machine *accepts* a language  $L$  if it outputs 1 when the input is a member of  $L$  and outputs 0 otherwise. That is, the Turing machine computes a predicate.

**Definition 1.7.4 (Computational complexity)** Let  $T$  be a Turing machine. For each input of length  $n$ , if  $T$  makes at most  $t(n)$  moves before it stops, then we say that  $T$  runs in time  $t(n)$ , or has *time complexity*  $t(n)$ . If  $T$  uses at most  $s(n)$  tape cells in above computation, then we say that  $T$  uses  $s(n)$  space, or has *space complexity*  $s(n)$ .

For convenience, we often give the Turing machine in Figure 1.1 a few more work tapes and designate one tape as read-only input tape. Thus, each transition rule will be of form  $(p, \bar{s}, a, q)$ , where  $\bar{s}$  contains the scanned symbols on all the tapes, and  $p, a, q$  as in Section 1.7.1, except that an operation now involves moving maybe more than one head.

We sometimes also make a Turing machine *nondeterministic* by allowing two distinct transition rules to have identical first two components. That is, a nondeterministic Turing machine may have different alternative moves at each step. Several other versions of Turing machines will be discussed in later chapters. Turing machines are deterministic unless it is explicitly stated otherwise.

It is a fundamental and easy result that any  $k$ -tape Turing machine running in  $t(n)$  time can be simulated by a Turing machine with just one work tape running in  $t^2(n)$  time. Any Turing machine using  $s(n)$  space can be simulated by a Turing machine with just one work tape using  $s(n)$  space. For each  $k$ , if a language is accepted by a  $k$ -tape Turing machine running in time  $t(n)$  (space  $s(n)$ ), then it can also be accepted by another  $k$ -tape Turing machine running in time  $ct(n)$  (space  $cs(n)$ ), for any constant  $c > 0$ . This leads to the following definitions:

**Definition 1.7.5 (Complexity Classes)** DTIME[ $t(n)$ ] is the set of languages accepted by multitape deterministic Turing machines in time  $O(t(n))$ ;  
NTIME[ $t(n)$ ] is the set of languages accepted by multitape nondeterministic Turing machines in time  $O(t(n))$ ;

$\text{DSPACE}[s(n)]$  is the set of languages accepted by multitape deterministic Turing machines in  $O(s(n))$  space;

$\text{NSPACE}[s(n)]$  is the set of languages accepted by multitape nondeterministic Turing machines in  $O(s(n))$  space;

$\text{P}$  is the complexity class  $\bigcup_{c \in \mathcal{N}} \text{DTIME}[n^c]$ ;

$\text{NP}$  is the complexity class  $\bigcup_{c \in \mathcal{N}} \text{NTIME}[n^c]$ ;

$\text{PSPACE}$  is the complexity class  $\bigcup_{c \in \mathcal{N}} \text{DSPACE}[n^c]$ .

We will define more complexity classes in Chapter 7. Languages in  $\text{P}$ , that is, languages acceptable in *polynomial* time, are considered as *feasibly* computable. The nondeterministic version for  $\text{PSPACE}$  turns out to be identical to  $\text{PSPACE}$ . The following relationships hold trivially,

$$\text{P} \subseteq \text{NP} \subseteq \text{PSPACE}.$$

It is one of the most fundamental open questions in computer science and mathematics to prove whether or not either of the above inclusions is proper. Research in computational complexity theory centers around these questions. In order to solve these problems, one can identify the hardest problems in  $\text{NP}$  or  $\text{PSPACE}$ .

#### Definition 1.7.6

**(Oracle machine)** A Turing machine  $T$  with an *oracle*  $A$ , where  $A$  is a language over  $T$ 's work tape alphabet, is denoted as  $T^A$ . Such a machine operates as a normal Turing machine with the exception that after it has computed a finite string  $x$  it can enter a special “oracular” state and ask if  $x \in A$ . The machine  $T^A$  gets the correct yes/no answer in one step. An oracle machine can use this feature one or more times during each computation.

In Exercise 1.7.16 on page 43 we define the recursion-theoretic notion of reducing one language to another. Here, we scale this notion down to feasible size in computational complexity, by limiting the computational power used in the reduction to polynomial time.

A language  $A$  is called *polynomial time Turing-reducible* to a language  $B$ , denoted as  $A \leq_T^P B$ , if given  $B$  as an *oracle*, there is a deterministic Turing machine that accepts  $A$  in polynomial time. That is, we can accept  $A$  in polynomial time given answers to membership of  $B$  for free.

A language  $A$  is called *polynomial time many-to-one reducible* to a language  $B$ , denoted as  $A \leq_m^P B$ , if there is a function  $r$  that is polynomial time computable, and for every  $a$ ,  $a \in A$  iff  $r(a) \in B$ . In both cases, if  $B \in \text{P}$ , then so is  $A$ .

**Definition 1.7.7 (NP-completeness)** A language  $A$  is NP-hard if all languages in NP are Turing polynomial time (equivalently in this case, many-to-one polynomial time) reducible to  $A$ . Consequently, if any NP-hard language is in P, then P = NP. If  $A$  is NP-hard and  $A \in \text{NP}$ , then we say  $A$  is NP-complete.

NP is the set of problems where it is easy to show (give a certificate) that the answer is “yes,” and P is the set of “yes–no” problems where it is easy to find the answer. The technical sense of “easy” is “doable by a deterministic Turing machine in polynomial time.” The “P versus NP” question can be understood as whether problems for which it is easy to certify the answer are the same problems for which it is easy to find the answer. The relevance is this:

Normally, we do not ask questions unless we can recognize easily in a certain sense when we are handed the correct answer. We are not normally interested in questions where it would take a lifetime of work to check whether or not you got the answer you want. NP is about those questions that we are likely to want answers to.

This excellent explanation was given by one of the inventors of the notions P and NP, J.E. Edmonds [Interview, FAUW Forum, University of Waterloo, January 1993].

**Example 1.7.9** A Boolean formula is in conjunctive normal form if it is a conjunction of disjunctions. For example,

$$f(x_1, x_2, x_3) = (x_1 + \bar{x}_2 + x_3)(\bar{x}_2 + x_3)(x_1 + x_3)$$

is in conjunctive normal form, and  $x_1x_2 + x_2\bar{x}_3$  is not in conjunctive normal form. A Boolean formula  $f(x_1, \dots, x_n)$  is *satisfiable* if there is a Boolean-valued truth assignment  $a_1, \dots, a_n$  such that  $f(a_1, \dots, a_n) = 1$ .

**Definition 1.7.8 (SAT)** Let SAT be the set of satisfiable Boolean formulas in conjunctive normal form. The *SAT problem* is to decide whether or not a given Boolean formula is in SAT.

This problem was the first natural problem shown to be NP-complete. Many practical issues seem to depend on fast solutions to this problem. Given a Boolean formula, a nondeterministic Turing machine can “guess” a correct truth assignment, and verify it. This takes only linear time. However, if we have to deterministically search for a satisfying truth assignment, there are  $2^n$  Boolean vectors to test.

Intuitively, and as far as is known now, a deterministic Turing machine cannot do much better than simply searching through these Boolean vectors one by one, using an exponential amount of time. ◇

For each class, say  $P$ , if a language  $L$  is accepted in deterministic polynomial time using an oracle  $A$ , then we write  $L \in P^A$ . If  $A$  is an NP-complete set, we also write  $L \in P^{NP}$ .

**Definition 1.7.9** (**Polynomial hierarchy**) The so-called *polynomial hierarchy* consists of the following hierarchy of language classes:  $\Sigma_1^P = NP$ ;  $\Delta_1^P = P$ ;  $\Sigma_{i+1}^P = NP^{\Sigma_i^P}$ ;  $\Delta_{i+1}^P = P^{\Sigma_i^P}$ ; and  $\Pi_i^P = \{\bar{L} : L \in \Sigma_i^P\}$ .

## Exercises

---

**1.7.1.** [10] We can rework the effective enumeration of Turing machines above and the definition of universal Turing machines using only tape alphabet  $A = \{0, 1\}$  (as opposed to  $\{0, 1, B\}$ ) such that the  $d(A)d(Q)$  product increases at most by a fixed constant under the change of the original Turing machine to its simulator. Namely, we simply replace each program  $p$  in  $\{0, 1\}^*$  by  $\bar{p}$ , and in simulated computation use only the alternate (say odd) cells that contain the original program  $p$ . To skip over the even “administrative” cells containing 0’s and to detect even cells containing delimiter 1 takes only a few extra states. In particular, it requires exactly the same number of extra states in *each* Turing machine modification. Prove the assertion concerning the  $d(A)d(Q)$  product. What is the “new”  $l(c(T))$ ? Since clearly  $\bar{c}$  suffices, we find  $l(\bar{c}) \leq 2l(c(T)) + 1$ , for all  $T$ .

**1.7.2.** [08] Show that there are exactly countably infinitely many partial recursive functions, and that there are exactly countably infinitely many (total) recursive functions.

**1.7.3.** [15] Georg Cantor proved that the total functions are not countable by introducing his famous *diagonalization* argument. Suppose the contrary, and count the functions in order  $f_1, f_2, \dots, f_i, \dots$ . Define a new function  $g$ , which we shall prove to be not in this list, by  $g(i) = f_i(i) + 1$ , for all natural numbers  $i$ . By contradictory assumption,  $g$  occurs in the list, say,  $g = f_i$ . But by definition  $g(i) \neq f_i(i)$ , which gives the required contradiction. Use a similar argument to prove that there are functions that are not partial recursive.

**1.7.4.** [11] It is important to distinguish between computable *functions* and the *algorithms* that compute those functions. To each function there correspond many different algorithms that compute it. Show that in the effective enumeration of Turing machines, as we have treated in this section, each partial recursive function is computed by countably infinitely many different Turing machines. If  $T_1, T_2, \dots, T_i, \dots$  is an effective enumeration of Turing machines and  $T_i$  computes partial recursive function  $\phi_i$ , for all  $i$ , then each partial recursive function  $f$  occurs infinitely many times in the list  $\phi_1, \phi_2, \dots, \phi_i, \dots$ . This result is not an accident of our

formalism and effective enumeration, but holds in general for all effective enumerations of algorithms that compute all partial recursive functions.

**1.7.5.** [25] Show that for every  $m, n \geq 1$ , there exists a recursive function  $\psi = s_n^{(m+1)}$  of  $m + 1$  variables such that for all parameters  $x, y_1, \dots, y_n$ ,

$$\phi_x^{(m+n)}(y_1, \dots, y_m, z_1, \dots, z_n) = \phi_{\psi(x, y_1, \dots, y_m)}^{(n)}(z_1, \dots, z_n),$$

for all variables  $z_1, \dots, z_n$ . (Hint: prove the case  $m = n = 1$ . The proof is analogous for the other cases.)

*Comments.* This important result, due to Stephen C. Kleene, is usually called the *s-m-n Theorem*. Source: H. Rogers, Jr., *Theory of Recursive Functions and Effective Computability*, McGraw-Hill, 1967; P. Odifreddi, *Classical Recursion Theory*, North-Holland, 1989.

**1.7.6.** [35] If  $P$  is the class of all partial recursive functions (for convenience restricted to one variable), then any map  $\pi$  from  $\mathbb{N}$  (the natural numbers) *onto*  $P$  is called a *numbering*. The standard indexing of the Turing machines provides such a numbering, say  $\pi_0$ , and the indexing of the recursive function definitions another, say  $\pi_1$ . A numbering  $\pi$  is *acceptable*, or a *Gödel numbering*, if we can go back and forth effectively between  $\pi$  and  $\pi_0$ , that is,

- (i) There is a recursive function  $f$  (not necessarily one-to-one) such that  $f\pi_0 = \pi$ .
- (ii) There is a recursive function  $g$  (not necessarily one-to-one) such that  $g\pi = \pi_0$ .
  - (a) Show that  $\pi_1$  is acceptable.
  - (b) Show that (i) is a necessary and sufficient condition that any  $\pi$  have a universal partial recursive function (satisfies an appropriate version of the enumeration theorem).
  - (c) Show that (ii) is a necessary and sufficient condition that any  $\pi$  have an appropriate version of the *s-m-n Theorem* (Exercise 1.7.5).
  - (d) Show that (ii) implies that  $\pi^{-1}(\phi)$  is infinite for every partial recursive  $\phi$  of one variable.
  - (e) Show that we can replace (i) and (ii) by the requirement that there be a recursive isomorphism between  $\pi$  and  $\pi_0$ .

*Comments.* These results, due to H. Rogers, Jr., in 1958, give an abstract formulation of the basic work of Church, Kleene, Post, Turing, Markov, and others, that their basic formalizations of the notion of partial recursive functions are effectively isomorphic. It gives invariant significance to the the notion of acceptable numbering in that major properties like the

Enumeration Theorem or the *s-m-n* Theorem hold for any acceptable numbering. Note that (i) may be viewed as requiring that the numbering is “algorithmic” in that each number yields an algorithm; and (ii) requires that the numbering is “complete” in that it includes all algorithms. Source: H. Rogers, Jr., *Ibid*.

**1.7.7.** [20] Show that there is no recursive function  $f$  such that  $f(x) = 1$  if  $\phi_x(x)$  is defined, and  $f(x) = 0$  otherwise.

*Comments.* This fact is known as the *recursive unsolvability of the halting problem*.

**1.7.8.** [20] Prove that the predicate  $f$  defined as  $f(x) = 1$  if  $\phi_x$  is total, and  $f(x) = 0$  otherwise, is not total recursive.

**1.7.9.** [15] Prove that a set  $A$  is recursively enumerable iff  $A$  is the domain of a partial recursive function.

*Comments.* This is often called the *Basic Theorem of recursively enumerable sets*.

**1.7.10.** [16] Prove that a set  $A$  is recursively enumerable iff  $A$  is the range of some partial recursive function iff  $A$  is the range of a total recursive function or  $\emptyset$ .

**1.7.11.** [34] (a) Show that it is possible to effectively enumerate the partial recursive functions without repetition.

(b) Let  $A = \{x : \phi_x \text{ is a total function}\}$ . Prove that  $A$  is not recursively enumerable.

*Comments.* Items (a) and (b) are not contradictory. Hint: in Item (a) dovetail the computations of all partial recursive functions on all arguments. Attributed to R.A. Friedberg, 1958. Source: H. Rogers, Jr., *Theory of Recursive Functions and Effective Computability*, McGraw-Hill, 1967; P. Odifreddi, *Classical Recursion Theory*, North-Holland, 1989, pp. 230–232.

**1.7.12.** [20] Let  $K = \{x : \phi_x(x) < \infty\}$ . Prove that  $K$  is a recursively enumerable set that is not recursive.

**1.7.13.** [15] (a) Show that the function  $\tau(x, y) = (x^2 + 2xy + y^2 + 3x + y)/2$  is a recursive one-to-one mapping from  $\mathcal{N}^2$  onto  $\mathcal{N}$ . Show that this is not a prefix-code.

(b) Show that  $\tau^{(k)}$  defined by  $\tau^{(2)} = \tau$  and  $\tau^{(k)}(x_1, x_2, \dots, x_k) = \tau(\tau^{(k-1)}(x_1, x_2, \dots, x_{k-1}), x_k)$  is a recursive one-to-one mapping from  $\mathcal{N}^k$  onto  $\mathcal{N}$ . Show that this is not a prefix-code.

(c) Let  $E : \mathcal{N} \rightarrow \mathcal{N}$  be an effective prefix-code with  $E(x)$  the code word for source word  $x$ . Show that  $E(\tau(x, y))$ , and also  $E(\tau^{(k)}(x_1, \dots, x_k))$ , for  $k > 2$ , are effective prefix-codes.

*Comments.* How good are these prefix-codes of  $k$ -tuples of integers in terms of density of the range of the code in  $\mathcal{N}$ ? Clearly, Item (c) is the best possible (in terms of fixed  $E$ ).

**1.7.14.** [20] A set  $A$  is recursively enumerable *without repetitions* if  $A$  equals the range of  $f$ , for some  $f$  that is recursive and one-to-one. Prove that  $A$  is infinite and recursively enumerable iff  $A$  is recursively enumerable without repetitions.

**1.7.15.** [25] (a) Show that there exists an infinite set having no infinite recursively enumerable subset. Such sets have been called *immune* by J.C.E. Dekker.

(b) If a set with this property has a recursively enumerable complement, then this complement was called *simple* by E.L. Post in 1944. Show that there exists a simple set.

*Comments.* By definition a simple set is recursively enumerable. A simple set is not recursive, since its complement is infinite but not recursively enumerable (Lemma 1.7.3, Item (i)), Source: H. Rogers, Jr., *Ibid*.

**1.7.16.** [35] In order of increasing generality we define some notions of reducibilities among sets:

A set  $A$  is *one-to-one reducible* to a set  $B$  ( $A \leq_1 B$ ) if there exists a one-to-one recursive function  $f$  such that for all  $x$  we have  $x \in A$  iff  $f(x) \in B$ .

A set  $A$  is *many-to-one reducible* to a set  $B$  ( $A \leq_m B$ ) if there exists a many-to-one recursive function  $f$  such that for all  $x$  we have  $x \in A$  iff  $f(x) \in B$ .

A set  $A$  is *Turing reducible* to a set  $B$  ( $A \leq_T B$ ) if we can solve membership in  $A$  by a Turing machine that gets the solution to membership in  $B$  for free. (This is commonly accepted as formalizing the most general intuitive notion of reducibility. It is also considered the most significant and useful such notion.)

Intuitively speaking, for  $r$  equal 1,  $m$ , or  $T$ , reducibility  $A \leq_r B$  means that to solve membership in  $B$  is at least as “hard” as to solve membership in  $A$ , up to the reducibility notion  $r$ .

(a) Consider  $\{x : \phi_x(y) < \infty \text{ for infinitely many } y\}$  and  $\{x : \phi_x \text{ is total}\}$ . Show that each of these sets is reducible to the other in all three senses discussed.

(b) Show that  $\leq_r$  is reflexive and transitive for all discussed reducibilities  $r$ . Hence,  $\equiv_r$  (both  $\leq_r$  and  $\geq_r$ ) is an equivalence relation, and  $\leq_r$  induces a partial order of the equivalence classes of  $\equiv_r$ . One equivalence class is *below* another one if members of the first are reducible to members of the

second, but not vice versa. We say that a set  $A$  on a lower  $r$ -equivalence class has a lower *degree of unsolvability* with respect to  $\leq_r$ .

- (c) Consider the diagonal halting set  $K = \{x : \phi_x(x) < \infty\}$ , and let  $A = \{x : \phi_x(y) < \infty \text{ for finitely many } y\}$ . Show that  $K$  is  $r$ -reducible (all  $r = 1, m, T$ ) to  $A$  (easy). Show that contrary to intuition,  $A$  is not  $r$ -reducible (any  $r = 1, m, T$ ) to  $K$  (hard). (Hint: show that  $A$  is not recursively enumerable.) Therefore,  $A$  is of a higher  $r$ -degree of unsolvability than  $K$ .
- (d) Show that the halting set  $K_0$  and  $K$  are of the same  $r$ -degree of unsolvability (all  $r = 1, m, T$ ).
- (e) Show that all recursive sets are of lower  $r$ -degree of unsolvability than  $K_0$  (all  $r = 1, m, T$ ).
- (f) The following notion is due to E.L. Post. A set  $A$  such that each recursively enumerable set is  $r$ -reducible to it is called  *$r$ -hard*. If  $A$  is both recursively enumerable and  $r$ -hard, then  $A$  is called  *$r$ -complete*. Show that the halting set  $K_0$  in the proof of Lemma 1.7.6 is an example of an  $r$ -complete set (all  $r = 1, m, T$ ). Show that  $K$  in Item (c) is another example.

*Comments.* Source: H. Rogers, Jr., *Ibid.*

**1.7.17.** [35] Use the definitions above. The following is known as *Post's Problem* (1944). The recursive sets have lower  $r$ -degree of unsolvability than  $K_0$ , which has the highest  $r$ -degree of unsolvability in the recursively enumerable sets. Are there other  $r$ -degrees of unsolvability (all  $r = 1, m, T$ )? For  $r = 1, m$ , the first examples of such sets were the simple sets of the above exercise; see Item (c) below. (For  $r = T$  the question is much harder, and the affirmative answer was only provided (independently) by R.A. Friedberg and A.A. Muchnik in 1956.)

- (a) Show that for all  $A$ ,  $A$  is  $m$ -complete iff  $A$  is 1-complete.
- (b) Show that  $\{x : \phi_x \text{ is total}\}$  and  $\{x : \phi_x \text{ is not total}\}$  are incomparable under  $\leq_m$ . (Hint: see Exercise 7-11 in H. Rogers, Jr., *Ibid.*)
- (c) Show that a simple set is neither recursive nor  $m$ -complete.
- (d) Show that  $\equiv_1$  and  $\equiv_m$  do not coincide on the nonrecursive recursively enumerable sets, and hence  $\leq_1$  and  $\leq_m$  do not coincide on these sets either.
- (e) (J.C.E. Dekker) Show that the  $m$ -degree of a simple set includes an infinite collection of distinct 1-degrees consisting entirely of simple sets.

*Comments.* From Item (c) it follows that there exist nonrecursive recursively enumerable sets (simple sets) that are not  $m$ -complete (and so by Item (a) not 1-complete). Source: H. Rogers, Jr., *Ibid.* The term

“Post’s Problem” is now used among recursion theorists only for the Turing reduction version.

**1.7.18.** [35] Consider the *generalized exponential* function  $f$  informally having the following property:  $f(0, x, y) = y + x$ ,  $f(1, x, y) = y \times x$ ,  $f(2, x, y) = y^x$ , . . . . A more formal definition of  $f$  is given by  $f(0, 0, y) = y$ ,  $f(0, x + 1, y) = f(0, x, y) + 1$ ,  $f(1, 0, y) = 0$ ,  $f(z + 2, 0, y) = 1$ ,  $f(z + 1, x + 1, y) = f(z, f(z + 1, x, y), y)$ .

- (a) Show that this function is recursive but not primitive recursive.
- (b) Show that the function  $A(x) = f(x, x, x)$ , called the *Ackermann* generalized exponential function, is recursive but not primitive recursive. (It rises faster than any primitive recursive function.)

*Comments.* The function  $f$  was given by Wilhelm Ackermann (1926) as a first example of a recursive function that is not primitive recursive [*Math. Ann.* 99(1928), 118–133]. See also D. Hilbert [*Math. Ann.*, 95(1926), 161–190], who credits the proof to Ackermann. There are many variants of definitions of the Ackermann function. A common recursive definition is  $A'(0, n) = n + 1$ ;  $A'(i, 0) = A'(i - 1, 1)$  for  $i \geq 0$ ; and  $A'(i, n) = A'(i - 1, A'(i, n - 1))$  for  $i, n > 0$ . Then  $A(x) = A'(x, x)$ . This definition is apparently due to R.M. Robinson [*Bull. Amer. Math. Soc.* 54(1948), 987–993], and an earlier variant is due to Rózsa Péter [*Math. Ann.* 111(1935), 42–60]. An inherently iterative algorithm to compute  $A'$  is given by J.W. Grossman and R.Z. Zeitman [*Theoret. Comput. Sci.*, 37(1988), 327–330]. Another definition for the Ackermann function (source: P. Odifreddi) is  $h_\omega$  defined as follows:  $h_0(x) = x + 1$ ;  $h_{n+1}(x) = h_n^{(x)}(x)$ ;  $h_\omega(x) = h_x(x)$  where  $h_n^{(0)}(x) = x$ ;  $h_n^{(z+1)}(x) = h_n(h_n^{(z)}(x))$ . One of the advantages of this version is that it can be translated into the transfinite. Note:  $h_1(x) = 2x$ ;  $h_2(x) = x2^x$ .

**1.7.19.** The function  $BB$  is defined in terms of Turing machines with a purely binary tape alphabet (no blanks) in quintuple format (rather than quadruple format as above). For each  $n$  define the set  $A_n = \{i : T_i \text{ has } n \text{ states and } \phi_i(0) < \infty\}$ . That is,  $T_i$  with  $i$  in  $A_n$  halts when it is started with  $p = \epsilon$ . We define  $BB(n)$  as the maximal number of 1’s in the output of any Turing machine in  $A_n$  when it is started on input  $\epsilon$ . It is easy to see that  $BB(1) = 1$ , but more difficult to see that  $BB(2) = 4$ . The  $BB$  function is known as the *Busy Beaver* function. It is due to T. Rado [*Bell System Tech. J.* (1962), 877–884]. This was one of the first “well-defined” noncomputable or nonrecursive total functions. The definition is natural and uses no overt diagonalization.

- (a) [12] Notice that the  $BB$  function is well-defined, since  $d(A_n) \leq f(n)$  for some easily computable function  $f$ . Compute  $f(n)$ . (Hint: it is doubly exponential.)

- (b) [15] Show that the  $BB$  function is noncomputable. (Hint: it grows faster than any recursive function.)
- (c) [35] Show that  $BB(3) = 6$ .
- (d) [40] Show that  $BB(4) = 13$ .
- (e) [30] Can you give a lower bound like the Ackermann generalized exponential function on  $BB$ ?
- (f) [O45] Is  $BB(5) > 4,098$ ?

*Comments.* See A.K. Dewdney, *Scientific American*, 251 (August 1984), 19–23; *Ibid*, 252 (March 1985), 23; A.H. Brady, pp. 259–278 in: *The Universal Turing Machine: A Half-Century Survey*, R. Herken, Ed., Oxford Univ. Press, 1988; H. Marzen and J. Buntrock, *EATCS Bull.*, 40(1990), 247–251. For the quadruple format see [A. Oberschelp, K. Schmidt-Goettsch, G. Todt, Castor Quadruplorum, *Archive for Math. Logic*, 27(1988), 35–44].

- 1.7.20.** [39] Let  $f$  be any recursive function. Prove that there exists an  $n$  such that  $\phi_n = \phi_{f(n)}$ .

*Comments.* This result, and its elaborations to more complicated versions, is usually called the *Second Recursion Theorem* or the *Fixed-Point Theorem for Recursion Theory*. The  $n$  is called a *fixed-point* value for  $f$ . Standard applications include: There exists an  $e$  such that the only element in the domain of  $\phi_e$  is  $e$  itself, and more generally, it allows us to write programs that know their own index. Source: H. Rogers, Jr., *Theory of Recursive Functions and Effective Computability*, McGraw-Hill, 1967.

- 1.7.21.** [37] In Exercise 1.7.16 we have studied the notions of reducibility and degree. We now look at a coarser classification. An  $n$ -ary relation  $R$  is in the *arithmetic hierarchy* iff it is either recursive or, for some  $m$ , can be expressed as

$$\{\langle x_1, \dots, x_n \rangle : (Q_1 y_1) \dots (Q_m y_m) S(x_1, \dots, x_n, y_1, \dots, y_m)\}, \quad (1.3)$$

where each  $Q_i$  denotes the existential quantifier “there exists” ( $\exists$ ) or the universal quantifier “for all” ( $\forall$ ), and  $S$  is an  $(n+m)$ -ary recursive relation. The *number of alternations* in the prefix sequence of quantifiers is the number of pairs of adjacent but unlike quantifiers. For each  $n > 0$ , Formula 1.3 is a  $\Sigma_n^0$ -form if the first quantifier is  $\exists$  and the number of alternations is  $n - 1$ . A  $\Sigma_0^0$ -form has no quantifiers. For each  $n > 0$ , Formula 1.3 is a  $\Pi_n^0$ -form if the first quantifier is  $\forall$  and the number of alternations is  $n - 1$ . A  $\Pi_0^0$ -form has no quantifiers. A relation  $R$  is in  $\Sigma_n^0$  if it can be expressed by a  $\Sigma_n^0$ -form. A relation  $R$  is in  $\Pi_n^0$  if it can be expressed by a  $\Pi_n^0$ -form. For each  $n \geq 0$ ,

$$\Delta_n^0 = \Sigma_n^0 \bigcap \Pi_n^0.$$

- (a) Show that  $\Sigma_0^0 = \Pi_0^0$  = the class of recursive sets.
- (b) Show that  $\Sigma_1^0$  is the class of recursively enumerable sets.
- (c) Show that  $R$  is in  $\Sigma_n^0$  iff the complement of  $R$  is in  $\Pi_n^0$ .
- (d) Show that  $\Sigma_n^0 \cup \Pi_n^0 \subseteq \Sigma_{n+1}^0 \cap \Pi_{n+1}^0$ .
- (e) Show that for each  $n > 0$ , we have  $\Sigma_n^0 - \Pi_n^0 \neq \emptyset$ , and hence, by (c), that for each  $n > 0$  we have  $\Pi_n^0 - \Sigma_n^0 \neq \emptyset$ .

*Comments.* The result mentioned as Item (e) is called the *Hierarchy Theorem*. In conjunction with Item (d) it shows that the classes  $\Sigma_0^0, \Sigma_1^0, \dots$  form a strictly increasing sequence. Source: H. Rogers, Jr., *Ibid.*

**1.7.22.** [13] A real number  $r \in [0, 1]$  is called a *recursive real number* if there exists a total recursive function  $\phi$  such that  $r = 0.\omega$  with  $\phi(i) = \omega_i$ , for all  $i$ . We call  $\omega$  a *recursive sequence*.

- (a) Show that not all real numbers are recursive.
- (b) Show that there are only countably many recursive numbers.
- (c) Show that there is a recursive sequence of recursive reals that converges to a real, but not to a recursive one.

## 1.8 The Roots of Kolmogorov Complexity

The notion of Kolmogorov complexity has its roots in probability theory, information theory, and philosophical notions of randomness, and came to fruition using the recent development of the theory of algorithms. The idea is intimately related to problems in both probability theory and information theory. These problems as outlined below can be interpreted as saying that the related disciplines are not “tight” enough; they leave things unspecified that our intuition tells us should be dealt with.

### 1.8.1 A Lacuna of Classical Probability Theory

An adversary claims to have a true random coin and invites us to bet on the outcome. The coin produces a hundred heads in a row. We say that the coin cannot be fair. The adversary, however, appeals to probability theory which says that each sequence of outcomes of a hundred coin flips is equally likely,  $1/2^{100}$ , and one sequence had to come up.

Probability theory gives us no basis to challenge an outcome *after* it has happened. We could only exclude unfairness in advance by putting a penalty side-bet on an outcome of 100 heads. But what about 1010...? What about an initial segment of the binary expansion of  $\pi$ ?

**Regular sequence**  $\Pr(00000000000000000000000000) = 1/2^{26}$ ,

**Regular sequence**  $\Pr(0100011011000010100111001) = 1/2^{26}$ ,

**Random sequence**  $\Pr(10010011011000111011010000) = 1/2^{26}$ .

The first sequence is regular, but what is the distinction of the second sequence and the third? The third sequence was generated by flipping a quarter. The second sequence is very regular: 0, 1, 00, 01, .... The third sequence will pass (pseudo-)randomness tests.

In fact, classical probability theory cannot express the notion of *randomness of an individual sequence*. It can only express expectations of properties of outcomes of random processes, that is, the expectations of properties of the total set of sequences under some distribution.

Only relatively recently, this problem has found a satisfactory resolution by combining notions of computability and statistics to express the complexity of a finite object. This complexity is the length of the shortest binary program from which the object can be effectively reconstructed. It may be called the *algorithmic information content* of the object. This quantity turns out to be an attribute of the object alone, and absolute (in the technical sense of being recursively invariant). It is the *Kolmogorov complexity* of the object.

### 1.8.2 A Lacuna of Information Theory

Shannon's classical information theory assigns a quantity of information to an ensemble of possible messages. All messages in the ensemble being equally probable, this quantity is the number of bits needed to count all possibilities. This expresses the fact that each message in the ensemble can be communicated using this number of bits. However, it does not say anything about the number of bits needed to convey any individual message in the ensemble. To illustrate this, consider the ensemble consisting of all binary strings of length 9999999999999999.

By Shannon's measure, we require 9999999999999999 bits on the average to encode a string in such an ensemble. However, the string consisting of 9999999999999999 1's can be encoded in about 55 bits by expressing 9999999999999999 in binary and adding the repeated pattern "1." A requirement for this to work is that we have agreed on an algorithm that decodes the encoded string. We can compress the string still further when we note that 9999999999999999 equals  $3^2 \times 11111111111111$ , and that 1111111111111111 consists of  $2^4$  1's.

Thus, we have discovered an interesting phenomenon: the description of some strings can be compressed considerably, provided they exhibit enough regularity. This observation, of course, is the basis of all systems to express very large numbers and was exploited early on by Archimedes in his treatise *The Sand Reckoner*, in which he proposes a system to name very large numbers:

"There are some, King Golon, who think that the number of sand is infinite in multitude [...] or] that no number has been named which is great enough

to exceed its multitude.[...] But I will try to show you, by geometrical proofs, which you will be able to follow, that, of the numbers named by me [...] some exceed not only the mass of sand equal in magnitude to the earth filled up in the way described, but also that of a mass equal in magnitude to the universe."

However, if regularity is lacking, it becomes more cumbersome to express large numbers. For instance, it seems easier to compress the number "one billion," than the number "one billion seven hundred thirty-five million two hundred sixty-eight thousand and three hundred ninety-four," even though they are of the same order of magnitude.

## 1.9 Randomness

This brings us to the main root of Kolmogorov complexity, the notion of randomness. There is a certain inevitability in the development that led A.N. Kolmogorov (1903–1987) to use the recently developed theory of effective computability to resolve the problems attending the proper definition of a random sequence. Indeed, the main idea involved had already been formulated with unerring intuition by P.S. Laplace, but could not be properly quantified at the time (see Chapter 4).

In the context of the above discussion, random sequences are sequences that cannot be compressed. Now let us compare this with the common notions of mathematical randomness. To measure randomness, criteria have been developed that certify this quality. Yet, in recognition that they do not measure "true" randomness, we call these criteria "pseudo" randomness tests. For instance, statistical surveys of initial sequences of decimal digits of  $\pi$  have failed to disclose any significant deviations from randomness. But clearly, this sequence is so regular that it can be described by a simple program to compute it, and this program can be expressed in a few bits.

"Any one who considers arithmetical methods of producing random digits is, of course, in a state of sin. For, as has been pointed out several times, there is no such thing as a random number—there are only methods to produce random numbers, and a strict arithmetical procedure is of course not such a method. (It is true that a problem we suspect of being solvable by random methods may be solvable by some rigorously defined sequence, but this is a deeper mathematical question than we can go into now.)" [von Neumann]

This fact prompts more sophisticated definitions of randomness. Notably, Richard von Mises (1883–1953) proposed notions that approach the very essence of true randomness. This is related to the construction of a formal mathematical theory of probability, to form a basis for real applications, in the early part of this century. While von Mises's objective was to justify the applications to real phenomena, Kolmogorov's classic 1933 treatment constructs a purely axiomatic theory of probability on the basis of set-theoretic axioms.

"This theory was so successful, that the problem of finding the basis of real applications of the results of the mathematical theory of probability became rather secondary to many investigators. [...] however] the basis for the applicability of the results of the mathematical theory of probability to real "random phenomena" must depend in some form on the *frequency concept of probability*, the unavoidable nature of which has been established by von Mises in a spirited manner." [Kolmogorov]

The point made is that the axioms of probability theory are designed so that abstract probabilities can be computed, but nothing is said about what probability really means, or how the concept can be applied meaningfully to the actual world. Von Mises analyzed this issue in detail, and suggested that a proper definition of probability depends on obtaining a proper definition of a random sequence. This makes him a "frequentist"—a supporter of the frequency theory.

The frequency theory to interpret probability says, roughly, that if we perform an experiment many times, then the ratio of favorable outcomes to the total number  $n$  of experiments will, *with certainty*, tend to a limit,  $p$  say, as  $n \rightarrow \infty$ . This tells us something about the *meaning* of probability, namely, the measure of the positive outcomes is  $p$ . But suppose we throw a coin 1000 times and wish to know what to expect. Is 1000 enough for convergence to happen? The statement above does not say. So we have to add something about the rate of convergence. But we cannot assert a *certainty* about a particular number of  $n$  throws, such as "the proportion of heads will be  $p \pm \epsilon$  for large enough  $n$  (with  $\epsilon$  depending on  $n$ )."  
We can at best say "the proportion will lie between  $p \pm \epsilon$  with at least such and such probability (depending on  $\epsilon$  and  $n_0$ ) whenever  $n > n_0$ ."  
But now we have defined probability in an obviously circular fashion.

In 1919 von Mises proposed to eliminate the problem by simply dividing all infinite sequences into special random sequences (called *collectives*), having relative frequency limits, which are the proper subject of the calculus of probabilities and other sequences. He postulates the existence of random sequences as certified by abundant empirical evidence, in the manner of physical laws, and derives mathematical laws of probability as a consequence. In his view a naturally occurring sequence can be nonrandom or unlawful in the sense that it is not a proper collective.

Von Mises views the theory of probabilities insofar as they are numerically representable as a physical theory of definitely observable phenomena, repetitive or mass events, for instance, as found in games of chance, population statistics, Brownian motion. "Probability" is a primitive notion of the theory comparable to those of "energy" or "mass" in other physical theories.

Whereas energy or mass exist in fields or material objects, probabilities exist only in the similarly mathematical idealization of collectives (random sequences). All problems of the theory of probability consist of deriving, according to certain rules, new collectives from given ones and calculating the

distributions of these new collectives. The exact formulation of the properties of the collectives is secondary and must be based on empirical evidence. These properties are the existence of a limiting relative frequency and randomness. The property of randomness is a generalization of the abundant experience in gambling houses, namely, the impossibility of a successful gambling system. Including this principle in the foundation of probability, von Mises argues, we proceed in the same way as the physicists did in the case of the energy principle. Here too, the experience of hunters of fortune is complemented by solid experience of insurance companies, and so forth.

A fundamentally different approach is to justify *a posteriori* the application of a purely mathematically constructed theory of probability, such as the theory resulting from the Kolmogorov axioms. Suppose we can show that the appropriately defined random sequences form a set of measure one, and without exception satisfy all laws of a given axiomatic theory of probability. Then it appears practically justifiable to assume that as a result of an (infinite) experiment only random sequences appear.

Von Mises's notion of infinite random sequences of 0's and 1's (collective) essentially appeals to the idea that no gambler, making a fixed number of wagers of "heads," at fixed odds [say  $p$  versus  $1 - p$ ] and in fixed amounts, on the flips of a coin [with bias  $p$  versus  $1 - p$ ], can have profit in the long run from betting according to a system instead of betting at random. Says Church: "this definition [below] ... while clear as to general intent, is too inexact in form to serve satisfactorily as the basis of a mathematical theory."

**Definition 1.9.1** An infinite sequence  $a_1, a_2, \dots$  of 0's and 1's is a random sequence in the special meaning of *collective* if the following two conditions are satisfied:

1. Let  $f_n$  be the number of 1's among the first  $n$  terms of the sequence. Then

$$\lim_{n \rightarrow \infty} \frac{f_n}{n} = p, \text{ for some } p, 0 < p < 1.$$

2. A *place-selection rule* is a partial function  $\phi$  from the finite binary sequences to the values 0 and 1 with the purpose of selecting one after another those indices  $n$  for which  $\phi(a_1 a_2 \dots a_{n-1}) = 1$ . We require (1), with the same limit  $p$ , also for every infinite subsequence

$$a_{n_1} a_{n_2} \dots$$

obtained from the sequence by some *admissible* place-selection rule. (We have not yet formally stated which place-selection rules are admissible.)

The existence of a relative frequency limit is a strong assumption. Empirical evidence from long runs of dice throws in gambling houses or with death statistics in insurance mathematics suggests that the relative frequencies are *apparently convergent*. But clearly, no empirical evidence can be given for the existence of a definite limit for the relative frequency. However long the test run, in practice it will always be finite, and whatever the apparent behavior in the observed initial segment of the run, it is always possible that the relative frequencies keep oscillating forever if we continue.

The second condition ensures that no strategy using an admissible place-selection rule can select a subsequence that allows different odds for gambling than a subsequence that is selected by flipping a fair coin. For example, let a casino use a coin with probability  $p = \frac{1}{4}$  of coming up heads and a payoff for heads equal to three times the pay-off for tails. This “Law of Excluded Gambling Strategy” says that a gambler betting in fixed amounts cannot make more profit in the long run betting according to a system than from betting at random.

“In everyday language we call random those phenomena where we cannot find a regularity allowing us to predict precisely their results. Generally speaking, there is no ground to believe that random phenomena should possess any definite probability. Therefore, we should distinguish between randomness proper (as absence of any regularity) and stochastic randomness (which is the subject of probability theory). There emerges the problem of finding reasons for the applicability of the mathematical theory of probability to the real world.” [Kolmogorov]

Intuitively, we can distinguish between sequences that are irregular and do not satisfy the regularity implicit in stochastic randomness, and sequences that are irregular but do satisfy the regularities associated with stochastic randomness. Formally, we will distinguish the second type from the first type by whether or not a certain complexity measure of the initial segments goes to a definite limit. The complexity measure referred to is the length of the shortest description of the prefix (in the precise sense of Kolmogorov complexity) divided by its length. It will turn out that almost all infinite strings are irregular of the second type and satisfy all regularities of stochastic randomness.

“In applying probability theory we do not confine ourselves to negating regularity, but from the hypothesis of randomness of the observed phenomena we draw definite positive conclusions.” [Kolmogorov]

Considering the sequence as fair coin tosses with  $p = \frac{1}{2}$ , the second condition in Definition 1.9.1 says that there is no *strategy*  $\phi$  (*principle of excluded gambling system*) that assures a player betting at fixed odds and in fixed amounts on the tosses of the coin to make infinite gain. That is, no advantage is gained in the long run by following some system, such as betting “heads” after each run of seven consecutive tails, or (more

plausibly) by placing the  $n$ th bet “heads” after the appearance of  $n + 7$  tails in succession. According to von Mises, the above conditions are sufficiently familiar and an uncontested empirical generalization to serve as the basis of an applicable calculus of probabilities.

**Example 1.9.1** It turns out that the naive mathematical approach to a concrete formulation, admitting simply *all* partial functions, comes to grief as follows: Let  $a = a_1a_2\dots$  be any collective. Define  $\phi_1$  as  $\phi_1(a_1\dots a_{i-1}) = 1$  if  $a_i = 1$ , and undefined otherwise. But then  $p = 1$ . Defining  $\phi_0$  by  $\phi_0(a_1\dots a_{i-1}) = b_i$ , with  $b_i$  the complement of  $a_i$ , for all  $i$ , we obtain by the second condition of Definition 1.9.1 that  $p = 0$ . Consequently, if we allow functions like  $\phi_1$  and  $\phi_0$  as strategy, then von Mises’s definition cannot be satisfied at all. ◇

In the thirties, Abraham Wald proposed to restrict the a priori admissible  $\phi$  to any fixed countable set of functions. Then collectives do exist. But which countable set? In 1940, Alonzo Church proposed to choose a set of functions representing “computable” strategies. According to Church’s Thesis, Section 1.7, this is precisely the set of *recursive functions*. With recursive  $\phi$ , not only is the definition completely rigorous, and random infinite sequences do exist, but moreover they are abundant since the infinite random sequences with  $p = \frac{1}{2}$  form a set of measure one. From the existence of random sequences with probability  $\frac{1}{2}$ , the existence of random sequences associated with other probabilities can be derived. Let us call sequences satisfying Definition 1.9.1 with recursive  $\phi$  *Mises-Wald-Church random*. That is, the involved *Mises-Wald-Church place-selection rules* consist of the partial recursive functions.

Appeal to a theorem of Wald yields as a corollary that the set of Mises-Wald-Church random sequences associated with any fixed probability has the cardinality of the continuum. Moreover, each Mises-Wald-Church random sequence qualifies as a normal number. (A number is *normal* if each digit of the base, and each block of digits of any length, occurs with equal asymptotic frequency.) Note, however, that not every normal number is Mises-Wald-Church random. This follows, for instance, from Champernowne’s sequence (or number),

$$0.1234567891011121314151617181920\dots,$$

due to D.G. Champernowne, which is normal in the scale of 10 and where the  $i$ th digit is easily calculated from  $i$ . The definition of a Mises-Wald-Church random sequence implies that its consecutive digits cannot be effectively computed. Thus, an existence proof for Mises-Wald-Church random sequences is necessarily nonconstructive.

Unfortunately, the von Mises-Wald-Church definition is not yet good enough, as was shown by J. Ville in 1939. There exist sequences that

satisfy the Mises-Wald-Church definition of randomness, with limiting relative frequency of ones of  $\frac{1}{2}$ , but nonetheless have the property that

$$\frac{f_n}{n} \geq \frac{1}{2} \text{ for all } n.$$

The probability of such a sequence of outcomes in random flips of a fair coin is zero. Intuition: if you bet “1” all the time against such a sequence of outcomes, then your accumulated gain is always positive! Similarly, other properties of randomness in probability theory such as the Law of the Iterated Logarithm do not follow from the Mises-Wald-Church definition.

For a better understanding of the problem revealed by Ville, and its subsequent solution by P. Martin-Löf in 1966, we look at some aspects of the methodology of probability theory. Consider the sample space of all one-way infinite binary sequences generated by fair coin tosses. We call a sequence “random” iff it is “typical.” It is not “typical,” say “special,” if it has a particular distinguishing property. An example of such a property is that an infinite sequence contains only finitely many ones. There are infinitely many such sequences. But the probability that such a sequence occurs as the outcome of fair coin tosses is zero. “Typical” infinite sequences will have the converse property, namely, they contain infinitely many ones.

In fact, one would like to say that “typical” infinite sequences will have *all converse properties* of the properties that can be enjoyed by “special” infinite sequences. This is formalized as follows: If a particular property, such as containing infinitely many occurrences of ones (or zeros), the Law of Large Numbers, or the Law of the Iterated Logarithm, has been shown to have probability one, then one calls this a *Law of Randomness*. A sequence is “typical,” or “random,” if it satisfies all Laws of Randomness.

But now we are in trouble. Since *all* complements of singleton sets in the sample space have probability one, it follows that the intersection of all sets of probability one is empty. Thus, there are no random infinite sequences!

Martin-Löf, using ideas related to Kolmogorov complexity, succeeded in defining random infinite sequences in a manner that is free of such difficulties. He observed that all laws that are proven in probability theory to hold with probability one are effective (as defined in Section 1.7). That is, we can effectively test whether a particular infinite sequence does not satisfy a particular Law of Randomness by effectively testing whether the law is violated on increasingly long initial segments.

The natural formalization is to identify the effective test with a partial recursive function. This suggests that one ought to consider not the

intersection of all sets of measure one, but only the intersection of all sets of measure one with recursively enumerable complements. (Such a complement set is expressed as the union of a recursively enumerable set of cylinders). It turns out that this intersection has again measure one. Hence, almost all infinite sequences satisfy all effective Laws of Randomness with probability one. This notion of infinite random sequences is related to infinite sequences of which all finite initial segments have high Kolmogorov complexity.

The notion of randomness satisfied by both the Mises-Wald-Church collectives and the Martin-Löf random infinite sequences is roughly that *effective tests* cannot detect regularity. This does not mean that a sequence may not exhibit regularities that cannot be effectively tested. Collectives generated by Nature, as postulated by von Mises, may very well always satisfy stricter criteria of randomness. Why should collectives generated by quantum mechanical phenomena care about mathematical notions of computability? Again, satisfaction of all effectively testable prerequisites for randomness is some form of regularity. Maybe nature is more lawless than adhering strictly to regularities imposed by the statistics of randomness.

Until now the discussion has centered on infinite random sequences where the randomness is defined in terms of limits of relative frequencies. However,

“The frequency concept based on the notion of *limiting frequency* as the number of trials increases to infinity does not contribute anything to substantiate the application of the results of probability theory to real practical problems where we always have to deal with a finite number of trials.” [Kolmogorov]

The practical objection against both the relevance of considering infinite sequences of trials and the existence of a relative frequency limit is concisely put in J.M. Keynes’s famous phrase, “In the long run we shall all be dead.” It seems more appealing to try to define randomness for finite strings first, and only then define random infinite strings in terms of randomness of initial segments.

The approach of von Mises to define randomness of infinite sequences in terms of *unpredictability* of continuations of finite initial sequences under certain laws (like recursive functions) did not lead to satisfying results. Although certainly inspired by the random sequence debate, the introduction of Kolmogorov complexity marks a definite shift of point of departure. Namely, to define randomness of sequences by the fact that no program from which an initial segment of the sequence can be computed is significantly shorter than the initial segment itself, rather than that no program can predict the next elements of the sequence.

Finite sequences that cannot be effectively described in a significant shorter description than their literal representation are called random. Our aim is to characterize random infinite sequences as sequences of

which all initial finite segments are random in this sense (Section 3.6). A related approach characterizes random infinite sequences as sequences of which all initial finite segments pass all effective randomness tests (Section 2.5).

Initially, before the idea of complexity, Kolmogorov proposed a close analogy to von Mises's notions in the finite domain. Consider a generalization of place-selection rules insofar as the selection of  $a_i$  can depend on  $a_j$ , with  $j > i$  [A.N. Kolmogorov, *Sankhyā*, Series A, 25(1963), 369–376]. Let  $\Phi$  be a finite set of such generalized place-selection rules. Kolmogorov suggested that an arbitrary finite binary sequence  $a$  of length  $n \geq m$  can be called  $(m, \epsilon)$ -random with respect to  $\Phi$ , if there exists some  $p$  such that the relative frequency of the 1's in the subsequences  $a_{i_1} \dots a_{i_r}$ , with  $r \geq m$ , selected by applying some  $\phi$  in  $\Phi$  to  $a$ , all lie within  $\epsilon$  of  $p$ . (We discard  $\phi$  that yield subsequences shorter than  $m$ .) Stated differently, the relative frequency in this finite subsequence is approximately (to within  $\epsilon$ ) invariant under any of the methods of subsequence selection that yield subsequences of length at least  $m$ . Kolmogorov has shown that if the cardinality of  $\Phi$  satisfies

$$d(\Phi) \leq \frac{1}{2} e^{2m\epsilon^2(1-\epsilon)},$$

then for any  $p$  and any  $n \geq m$  there is some sequence  $a$  of length  $n$  that is  $(m, \epsilon)$ -random with respect to  $\Phi$ .

## Exercises

---

**1.9.1.** [08] Consider the sequence 101001000100001000001..., that is, increasing subsequences of 0's separated by single 1's. What is the limiting relative frequency of 1's and 0's? Is this sequence a collective?

**1.9.2.** [15] Suppose we are given a coin with an unknown bias, the probability of heads is  $p$  and of tails is  $(1 - p)$ ,  $p$  some unknown real number  $0 < p < 1$ . Can we use this coin to simulate a perfectly fair coin?

*Comments.* The question asks for an effective construction of a collective of 0's and 1's with limiting frequency  $\frac{1}{2}$  from a collective of 0's and 1's with unknown limiting frequency  $p$ . Source: J. von Neumann, Various techniques used in connection with random digits, in: *Collected Works*, Vol. V, A.H. Traub, Ed., Macmillan, 1963; W. Hoeffding and G. Simons, *Ann. Math. Statist.*, 41(1970), 341–352; T.S. Ferguson, *Ann. Math. Statist.*, 41(1970), 352–362; P. Elias, *Ann. Math. Statist.*, 43(1972), 865–870.

**1.9.3.** [15] Suppose the sequence  $a_1 a_2 \dots a_i \dots$  has all the properties of a collective with  $p = \frac{1}{2}$ .

(a) The limiting relative frequency of subsequence “01” equals the limiting frequency of subsequence “11.” Compute these limiting relative frequencies.

- (b) Form a new sequence  $b_1 b_2 \dots b_i \dots$  defined by  $b_i = a_i + a_{i+1}$  for all  $i$ . Then the  $b_i$ 's are 0, 1, or 2. Show that the limiting relative frequencies of 0, 1, and 2 are  $\frac{1}{4}$ ,  $\frac{1}{2}$ , and  $\frac{1}{4}$ , respectively.
- (c) The new sequence constructed in Item (b) satisfies requirement (1) of a collective: the limiting relative frequencies constituent elements of the sample space  $\{0, 1, 2\}$  exist. Prove that it does *not* satisfy requirement (2) of a collective.
- (d) Show that the relative frequency of the subsequence “21” will generally differ from the relative frequency of the subsequence “11” or “10.”

*Comments.* Hint for Item (c): show that the subsequences “02” and “20” do not occur. Use this to select effectively a subsequence of  $b_1 b_2 \dots b_i \dots$  with limiting relative frequency of 2's equal to zero. This phenomenon was first observed by M. von Smoluchowski [*Sitzungsber. Wien. Akad. Wiss., math.-naturw. Kl., Abt. IIa*, 123(1914) 2381–2405; 124(1915) 339–368] in connection with Brownian motion, and called probability “after-effect.” Also: R. von Mises, *Probability, Statistics, and Truth*, Macmillan, 1933.

- 1.9.4.** [10] Given an infinite binary string that is a collective with limiting frequency of 1's equal to  $p = \frac{1}{2}$ . Show how to construct a collective over symbols 0, 1, and 2 (a ternary collective) with equal limiting frequencies of  $\frac{1}{3}$  for the number of occurrences of 0's, 1's, and 2's.

*Comments.* Hint: use Exercise 1.9.2.

- 1.9.5.** [M40] Let  $\omega_1 \omega_2 \dots$  be an infinite binary sequence, and let  $f_n = \omega_1 + \omega_2 + \dots + \omega_n$ .

(a) Sequence  $\omega$  is said to satisfy the *Infinite Recurrence Law* if  $f_n = n/2$  infinitely often. It can be shown that the set of infinite binary sequences having the infinite recurrence property has measure one in the set of all infinite binary sequences with respect to the usual binary measure. Show that there are infinite binary sequences that are Mises-Wald-Church random and  $f_n > n/2$  for all  $n$ .

(b) Show that there are infinite binary sequences that are Mises-Wald-Church random and  $\limsup_{n \rightarrow \infty} (\sum_{i=1}^n \omega_i - n/2) / \sqrt{n \log \log n} > 1/\sqrt{2}$  (that is, they do not satisfy the Law of the Iterated Logarithm).

*Comments.* Since Items (a) and (b) are satisfied with probability zero in the set of all infinite binary strings, we can conclude that Mises-Wald-Church random strings do not satisfy all laws of probability that hold with probability one (the Laws of Randomness). Source: J. Ville, *Étude Critique de la Concept de Collectif*, Gauthier-Villars, 1939.

- 1.9.6.** [32] What happens if we restrict our set of admissible selection functions to those computable by finite-state machines instead of Turing

machines? First we need some definitions. Let  $\omega = \omega_1\omega_2\dots$  be an infinite binary string. For each  $m = 0, 1, 00, 01, \dots$ , let  $f_n$  be the number of occurrences of  $m$  in  $\omega_{1:n}$ . We say that  $\omega$  is  $k$ -distributed if  $\lim_{n \rightarrow \infty} f_n/n = 1/2^{l(m)}$  for all  $m$  with  $l(m) \leq k$ . We say that  $\omega$  is  $\infty$ -distributed, or *normal*, if it is  $k$ -distributed for all integers  $k$ .

A *finite-state place-selection rule*  $\phi$  is a function computed by a finite-state machine with value 0 or 1 for each finite binary sequence. Given an infinite sequence  $\omega$ ,  $\phi$  determines a subsequence  $\omega_{i_1}\omega_{i_2}\dots$  by selecting one after another the indices  $n$  for which  $\phi(\omega_1\omega_2\dots\omega_{n-1}) = 1$ . (Formally, in terms of the definition of a Turing machine in Section 1.7, we can think of  $\phi$  as being computed by a Turing machine  $T : Q \times A \rightarrow S \times Q$  with  $T(\cdot, d) = (R, \cdot)$  for all  $d = 0, 1$  and  $T(\cdot, B) = (a, \cdot)$  for some  $a = 0, 1$ .) Let  $\omega_{i_1}\omega_{i_2}\dots$  be the subsequence of  $\omega$  selected by  $\phi$ , and let  $c_n$  be the number of ones in the first  $n$  bits. Assuming that  $\phi$  has infinitely many values 1 on prefixes of  $\omega$ , define  $\limsup_{n \rightarrow \infty} c_n/n$  as the *prediction ratio* of  $\phi$  for  $\omega$ . Finally, a finite-state prediction function  $\phi$  is a *predictor* for  $\omega$  if its prediction ratio is greater than  $\frac{1}{2}$ .

- (a) Show that every  $\omega$  that is  $k$ -distributed but not  $(k+1)$ -distributed has a  $(k+1)$ -state predictor.
- (b) Show that there are no predictors for  $\infty$ -distributed  $\omega$ .
- (c) Show that the subsequence selected from an  $\infty$ -distributed  $\omega$  by a finite-state selection function is again  $\infty$ -distributed.
- (d) Show that there are  $\infty$ -distributed sequences that can be predicted by stronger models of computation such as Turing machines. Conclude that there are  $\infty$ -distributed sequences that are not random in the sense of Mises-Wald-Church. (Hint: use Champernowne's sequence in the main text.)

*Comments.* Since predictors are machines that in the long run, have some success in making correct predictions on  $\omega$ , we can say that  $\omega$  appears *random* to  $\phi$  if  $\phi$  is not a predictor of  $\omega$ . Then,  $\infty$ -distributed sequences are precisely the sequences that appear random to all finite-state selection functions. In terms of gambling: let a gambler pay \$1.00 for each prediction he makes and receives \$2.00 for each correct prediction. If the sequence supplied by the house is  $\infty$ -distributed, and the gambler makes unboundedly many predictions, then no matter what finite-state selection function he uses, the limit superior of the ratio (paid \$)/(received \$) goes to one. Source: V.N. Agafonoff, *Soviet Math. Dokl.*, 9(1968), 324–325 (English transl.); C.P. Schnorr and H. Stimm, *Acta Informatica*, 1(1972), 345–359; and, apparently independently, M.G. O'Connor, *J. Comput. System Sci.* 37(1988), 324–336. See also: T. Kamae, *Israel J. Math.*, 16(1973), 121–149; T. Kamae and B. Weiss, *Israel J. Math.*, 21(1975), 101–111.

**1.9.7.** [O35] Investigate the related problems as in Exercise 1.9.6 by replacing finite state machines by, for instance, push-down automata, time- or space-bounded classes, and primitive recursive functions.

## 1.10 Prediction and Probability

The question of quantitative probability based on complexity was first raised and treated by R.J. Solomonoff, in an attempt to obtain a completely general theory of inductive reasoning. Let us look at some predecessors in this line of thought.

The so-called *Weak Law of Large Numbers*, formulated by J. Bernoulli in 1713, states that if an experiment with probability of success  $p$  is repeated  $n$  times, then the proportion of successful outcomes will approach  $p$  for large  $n$ . Such a repetitive experiment is called a sequence of *Bernoulli trials* generated by a  $(p, 1 - p)$  *Bernoulli process*, and the generated sequence of outcomes is called a *Bernoulli sequence*.

Thomas Bayes, in “An essay towards solving a problem in the doctrine of chances” [*Philos. Trans.*, London, 53(1763), 376–398, and 54(1764), 298–310], suggested the “inverse of Bernoulli’s problem.” The resulting method, sometimes referred to as *inverse probability*, was further analyzed by P.S. Laplace, who also attached Bayes’s name to it. In Bayesian approaches it is assumed that there is some true, or *a priori* (prior), distribution of probabilities over objects. Then an object with unknown probability  $p$  is drawn. Provided with a (nonempirical) prior probability, together with empirical data and the probabilistic model of these data, Bayes’s Rule supplies a way to calculate a *posterior* or *inferred* probability distribution. We then can give a numerical estimate for  $p$ , for example, by choosing the maximum posterior probability.

The formal statement of Bayes’s Rule was given in Section 1.6. The procedure is most easily explained by example. Suppose we have an urn containing a large number of dice with the faces numbered 1 through 6. Each die has a (possibly different) unknown probability  $p$  of casting “6,” which may be different from the  $\frac{1}{6}$  that it is for a “true” die. A die is drawn from the urn and cast  $n$  times in total and produces the result “6” in  $m$  of those casts. Let  $P(X = p)$  be the probability of drawing a die with attribute  $p$  from the urn. This  $P(X = p)$  is the prior distribution. In von Mises’s interpretation, if we repeatedly draw a die from the urn, with replacement, then the relative frequency with which a die with given value of  $p$  appears in these drawings has the limiting value  $P(X = p)$ . The probability of obtaining  $m$  outcomes “6” in  $n$  throws of a die with attribute  $p$  is

$$P(Y = m|n, p) = \binom{n}{m} p^m (1 - p)^{n-m},$$

the number of ways to select  $m$  items from  $n$  items, multiplied by the probability of  $m$  successes and  $(n-m)$  failures. Hence, the probability of drawing a die with attribute  $p$  and subsequently throwing  $m$  outcomes “6” in  $n$  throws with it is the product  $P(X = p) P(Y = m|n, p)$ .

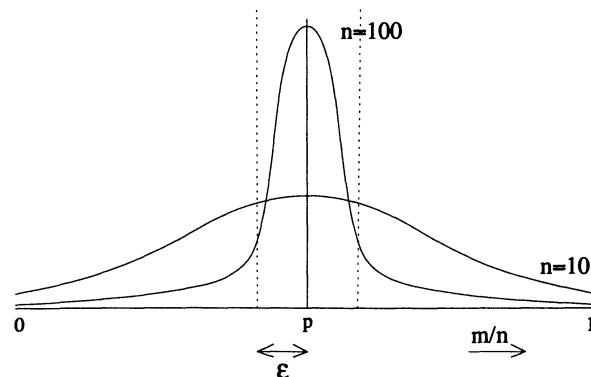
For the case under discussion, Bayes’s problem consists in determining the probability of  $m$  outcomes “6” in  $n$  casts being due to a die with a certain given value of  $p$ . The answer is given by the posterior, or inferred, probability distribution

$$P(X = p|n, m) = \frac{P(X = p)P(Y = m|n, p)}{\sum_p P(X = p)P(Y = m|n, p)},$$

the sum taken over all values of attribute  $p$ . If we repeat this experiment many times, then the limiting value of the probability of drawing a die with attribute value  $p$ , given that we throw  $m$  6’s out of  $n$ , is  $P(X = p|n, m)$ .

The interesting feature of this approach is that it quantifies the intuition that if the number of trials  $n$  is small, then the inferred distribution  $P(X = p|n, m)$  depends heavily on the prior distribution  $P(X = p)$ . However, if  $n$  is large, then *irrespective* of the prior  $P(X = p)$ , the inferred probability  $P(X = p|n, m)$  condenses more and more around  $m/n = p$ . To analyze this, we consider  $P(X = p|n, m)$  as a continuous distribution with fixed  $n$ . Clearly,  $P(X = p|n, m) = 0$  for  $m < 0$  and  $m > n$ . Let  $\epsilon > 0$  be some constant.

Consider the area under the “tails” of  $P(Y = m|n, p)$  for  $m \leq (p - \epsilon)n$  and  $m \geq (p + \epsilon)n$  (the area such that  $|p - m/n| \geq \epsilon$ ). Whatever  $\epsilon$  we choose, for each  $\delta > 0$  we can find an  $n_0$  such that this area is smaller than  $\delta$  for all  $n \geq n_0$ . This can be shown in several ways. We show this by appealing to a result that will be used several times later on.



**FIGURE 1.2.** Inferred probability for increasing  $n$

The probability of  $m$  successes out of  $n$  independent trials with probability  $p$  of success is given by the *binomial distribution*

$$P(Y = m|n, p) = \binom{n}{m} p^m (1-p)^{n-m}. \quad (1.4)$$

The deviation  $\epsilon$  (where  $0 \leq \epsilon$ ) from the average number of successes  $np$  in  $n$  experiments is analyzed by estimating the combined tail probability

$$P(|m - np| > \epsilon pn) = \sum_{|m-np|>\epsilon pn} \binom{n}{m} p^m (1-p)^{n-m}$$

of the binomial distribution, Figure 1.2. We give a variant of the classical estimate and omit the proof.

**Lemma 1.10.1 (Chernoff bounds)** *Assume the notation above. For  $0 \leq \epsilon \leq 1$ ,*

$$P(|m - pn| > \epsilon pn) \leq 2e^{-\epsilon^2 pn/3}. \quad (1.5)$$

*Each tail separately can be bounded by half of the right-hand side.*

This shows that for every  $\epsilon > 0$  (and  $\epsilon \leq 1$ )

$$\lim_{n \rightarrow \infty} \int_{(p-\epsilon)n}^{(p+\epsilon)n} P(X = p|n, m) dm = 1. \quad (1.6)$$

**Example 1.10.1** Let us give a numerical example. Let  $p$  take values  $0.1, 0.2, \dots, 0.9$  with equal probability  $P(X = 0.1) = P(X = 0.2) = \dots = P(X = 0.9) = \frac{1}{9}$ . Let  $n = 5$  and  $m = 3$ . Then the inferred probabilities are  $P(X = p|5, 3) = 0.005$  for  $p = 0.1$ ,  $0.031$  for  $p = 0.2$ , up to  $0.21$  for  $p = 0.6$ , and down again to  $0.005$  for  $p = 0.9$ , the combined probabilities summing up to 1. If we pick a die and do no experiments, then the probability that it is from any particular category is  $\frac{1}{9} \approx 0.11$ . If, however, we know already that it has had three throws of “6” out of five throws, then the probability that it belongs to category  $p = 0.1$  becomes smaller than 0.11, namely, 0.005, and the probability that it belongs to category  $p = 0.6$  increases to 0.21. In fact, the inferred probability that  $0.5 \leq p \leq 0.7$  is 0.59, while the inferred probability for the other six  $p$  values is only 0.41.

Consider the same prior distribution  $P(X = p)$  but set  $n = 500$  and  $m = 300$ . Then, the inferred probability for  $p = 0.6$  becomes  $P(X = 0.6|500, 300) = 0.99995$ . This means that it is now almost certain that a die which throws 60 percent 6’s has  $p = 0.6$ .

We have seen that the probability of inference depends on (a) the prior probability  $P(X = p)$  and (b) the observed results from which the inference is drawn. We have varied (b), but what happens if we start from

a different  $P(X = p)$ ? Let the new prior distribution be  $P(Y = p)$  be  $P(Y = 0.i) = i/45$ ,  $i = 1, \dots, 9$ . Then the corresponding inferred probability for  $n = 5$  and  $m = 3$  is  $P(Y = 0.1|5, 3) = 0.001$ ,  $P(Y = 0.2|5, 3) = 0.011$ , up to  $P(Y = 0.6|5, 3) = 0.21$ , and finally  $P(Y = 0.9|5, 3) = 0.07$ . For this small sequence, these values are markedly different from the previous  $X$  values, although the highest values are still around  $p = 0.6$ . But if we now increase the number of observations to  $n = 500$  with the same relative frequency of success  $m = 300$ , then the resulting  $Y$  values correspond to the  $X$  values but for negligible differences. ◇

Equation 1.6 shows that as the number of trials increases indefinitely, the limiting value of the observed relative frequency of success in the trials approaches the true probability of success, with probability one. This holds no matter what prior distribution the die was selected from. In case the initial probabilities of the events are unknown, Bayes's Rule is a correct tool to make inferences about the probability of events from frequencies based on many observations. For small sequences of observations, however, we need to know the initial probability to make justified inferences.

Solomonoff addresses precisely this issue. Suppose we are faced with a problem we have to solve in which there has been much experience. Then either we know outright how to solve it, or we know the frequency of success for different possible methods. However, if the problem has never occurred before, or only a small number of times, and the prior distribution is unknown, as it usually is, the inference method above is undefined or of poor accuracy. To solve this quandary Solomonoff proposes a *universal* prior probability. The idea is that this universal prior probability serves in a well-defined sense as well as the *true* prior probability, provided this true prior probability is computable in the sense of Section 1.7.

Solomonoff argues that all inference problems can be cast in the form of extrapolation from an ordered sequence of binary symbols. A principle to enable us to extrapolate from an initial segment of a sequence to its continuation will either require some hypothesis about the source of the sequence or a definition of what we mean by extrapolation. Two popular and useful metaphysical principles for extrapolation are those of simplicity (Occam's razor, attributed to the thirteenth century scholastic philosopher William of Ockham, but emphasized about twenty years before Ockham by John Duns Scotus), and indifference. The Principle of Simplicity asserts that the "simplest" explanation is the most reliable. The Principle of Indifference asserts that in the absence of grounds enabling us to choose between explanations we should treat them as equally reliable.

We do not supply any details here, because we shall extensively return to this matter in Chapter 4 and Chapter 5. Roughly, the idea is to define the *universal probability*,  $M(x)$ , as the probability that a fixed reference universal Turing machine outputs a sequence starting with  $x$  when its input is supplied by tosses of a fair coin. Using this as a sort of “universal prior probability” we then can formally do the extrapolation by Bayes’s Rule. The probability that  $x$  will be followed by a 1 rather than by a 0 turns out to be

$$\frac{M(x1)}{M(x0) + M(x1)}.$$

## Exercises

---

**1.10.1.** [12] Assume that a uniform prior probability  $P(X = x) = c$  is a fixed constant independent of  $x$ , for all individual outcomes  $x$ . Prove that after  $r$  successful outcomes  $x = a$  in  $n$  independent trials, the inferred posterior probability  $P(X = a|n, r)$  (according to Bayes’s Rule) is given by  $(r + 1)/(n + 2)$ .

*Comments.* This is P.S. Laplace’s so-called *Law of Succession*. Source: W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. 1, Wiley, 1968.

**1.10.2.** [25] Let an experiment where the outcomes are 0 or 1 with fixed probability  $p$  for outcome 1 and  $1 - p$  for outcome 0 be repeated  $n$  times. Such an experiment consists of a sequence of *Bernoulli trials* generated by a  $(p, 1 - p)$  Bernoulli process, see page 59.

Show that for each  $\epsilon > 0$  the probability that the number  $S_n$  of outcomes 1 in the first  $n$  trials of a single sequence of trials satisfies  $n(p - \epsilon) < S_n < n(p + \epsilon)$  goes to 1 as  $n$  goes to infinity.

*Comments.* This is J. Bernoulli’s Law of Large Numbers [*Ars Conjectandi*, Basel, 1713, Part IV, Ch. 5, p. 236], the so-called *Weak Law of Large Numbers*. This law shows that with great confidence in a series of  $n$  trials the proportion of successful outcomes will approximate  $p$  as  $n$  grows larger. The following interpretation of the weak law is false: “if Alice and Buck toss a perfect coin  $n$  times, then we can expect Alice to be in the lead roughly half of the time, regardless of who wins.” It can be shown that if Buck wins, then it is likely that he has been in the lead for practically the whole game. Thus, contrary to common belief, the time average of  $S_n$  ( $1 \leq n \leq m$ ) over an individual game of length  $m$  has nothing to do with the so-called “ensemble” average of the different  $S_n$ ’s associated with all possible games (the ensemble consisting of  $2^n$  games) at a given moment  $n$ , which is the subject of the weak law. Source: W. Feller, *Ibid.*

**1.10.3.** [25] In an infinite sequence generated by a  $(p, 1-p)$  Bernoulli process, let  $A_n$  denote the event that a run of  $n$  consecutive 1's occurs in between the  $2^n$ th and the  $2^{n+1}$ th trials. Prove that if  $p \geq \frac{1}{2}$ , then with probability one there occur infinitely many  $A_n$ ; if  $p < \frac{1}{2}$ , then with probability one there occur only finitely many  $A_n$ .

*Comments.* Hint: use elementary combinatorics.

**1.10.4.** [30] In an infinite sequence of outcomes generated by a  $(p, 1-p)$  Bernoulli process let  $A_1, A_2, \dots$  be an infinite sequence of events each of which depends only on a finite number of trials in the sequence. Denote the probability of  $A_k$  occurring by  $P_k$ . ( $A_k$  may be the event that  $k$  consecutive 1's occur between the  $2^k$ th trial and the  $2^{k+1}$ th trial. Then  $P_k \leq (2p)^k$ .)

- (a) Prove that if  $\sum P_k$  converges, then with probability one only finitely many  $A_k$  occur.
- (b) Prove that if the events  $A_k$  are mutually independent, and if  $\sum P_k$  diverges, then with probability one infinitely many  $A_k$  occur.

*Comments.* These two assertions are known as the *Borel-Cantelli Lemmas*. Source: W. Feller, *Ibid.*

**1.10.5.** [M30] Prove the limit in Equation 1.6 associated with the condensation of the posterior probability.

*Comments.* This may be called the *Inverse Weak Law of Large Numbers*, since it shows that we can infer with great certainty the original probability (drawn from an unknown distribution) by performing a single sequence of a large number of trials. Note that this is a different statement from the Weak Law.

**1.10.6.** [M37] Consider one-way infinite binary sequences generated by a  $(p, 1-p)$  Bernoulli process. Let  $S_n$  be as in Exercise 1.10.2.

- (a) Show that, for every  $\epsilon > 0$ , we have probability one that  $|pn - S_n| < \epsilon n$  for all but finitely many  $n$ .
- (b) Define the *reduced number of successes*  $S_n^* = (S_n - pn)/\sqrt{pn(1-p)}$ . Show the much stronger statement than Item (a), that with probability one  $|S_n^*| < \sqrt{2a \log n}$  (where  $a > 1$ ) holds for all but finitely many  $n$ .

*Comments.* Item (a) is a form of the *Strong Law of Large Numbers* due to F.P. Cantelli (1917) and G. Pólya (1921). Note that this statement is stronger than the Weak Law of Large Numbers. The latter says that  $S_n/n$  is likely to be near  $p$ , but does not say that  $S_n/n$  is bound to stay near  $p$  as  $n$  increases. The weak law allows that for infinitely many  $n$ , there is a  $k$  with  $n < k < 2n$ , such that  $S_k/k < p - \epsilon$ . In contrast, the Strong Law asserts that with probability one  $p - S_n/n$  becomes small

and remains small. Item (b) is due to A.N. Kolmogorov [*Math. Ann.*, 101(1929), 126–135]. Source: W. Feller, *Ibid.*

**1.10.7.** [M42] Consider an infinite sequence generated by a  $(p, 1-p)$  Bernoulli process. Show that  $\limsup_{n \rightarrow \infty} S_n^*/\sqrt{2 \log \log n} = 1$  with probability one, and symmetrically,  $\liminf_{n \rightarrow \infty} S_n^*/\sqrt{2 \log \log n} = -1$ .

*Comments.* This remarkable statement, known as the *Law of the Iterated Logarithm* is due to A.I. Khintchin [*Fundamenta Mathematicae* 6(1924), 9–20] and was generalized by A.N. Kolmogorov [*Math. Ann.*, 101(1929), 126–135]. For an explanation of its profundity, implications, and applications see also W. Feller, *Ibid.*

## 1.11 Information Theory and Coding

---

It seldom happens that a detailed mathematical theory springs forth in essentially final form from a single publication. Such was the case with information theory, which properly started only with the appearance of C.E. Shannon's paper “The mathematical theory of communication” [*Bell System Technical J.*, 27(1948), 379–423, 623–656]. In this theory we ignore the *meaning* of a message, we are *only* interested in the problem of *communicating* a message between a *sender* and a *receiver* under the assumption that the universe of possible messages is known to both the sender and the receiver.

This notion of information is a measure of one's freedom of choice when one selects a message. Given the choice of transmitting a message consisting of the contents of this entire book, and the message “let's get a beer,” the information concerned is precisely one bit. Obviously this does not capture the information content of the individual object itself. Kolmogorov's intention for introducing algorithmic complexity is as a measure for the information content of individual objects.

We develop the basic ideas in a purely combinatorial manner. This is easier and more fundamental, suffices for our purpose, and does not need extra probabilistic assumptions. The set of possible messages from which the selection takes place is often called an *ensemble*. Information, according to Shannon, is an ensemble notion. For our purpose it is sufficient to consider only *countable* ensembles.

The *entropy* of a random variable  $X$  with outcomes in an ensemble  $S$  is the quantity  $H(X) = \log d(S)$ . This is a measure of the uncertainty in choice before we have selected a particular value for  $X$ , and of the information produced from the set if we assign a specific value to  $X$ . By choosing a particular message  $a$  from  $S$ , we remove the entropy from  $X$  by the assignment  $X := a$  and produce or transmit *information*  $I = \log d(S)$  by our selection of  $a$ . Since the information is usually measured in the number of bits  $I'$  needed to be transmitted from sender to receiver,  $I' = \lceil \log d(S) \rceil$ .

**Example 1.11.1** The number of different binary strings  $\bar{u}$  with  $l(\bar{u}) = 2n + 1$  is  $2^n$ . This gives an information content in each such message of  $I = n$ , and encoding in a purely binary system requires  $I' = n$  bits.  $\diamond$

Note that while a random variable  $X$  usually ranges over a finite set of alternatives, say  $a, b, \dots, c$ , the derived theory is so general that it also holds if we let  $X$  range over a set of sequences composed from these alternatives, which may even be infinite.

If we have  $k$  independent random variables  $X_i$ , each of which can take  $n_i$  values, respectively, for  $i = 1, 2, \dots, k$ , then the number of combinations possible is  $n = n_1 n_2 \dots n_k$  and the entropy is given by

$$\begin{aligned} H(X_1, X_2, \dots, X_k) &= \log n_1 + \log n_2 + \dots + \log n_k \\ &= \log n. \end{aligned} \quad (1.7)$$

Let us look at the efficiency with which an individual message consisting of a *sequence*  $x_1 x_2 \dots x_k$  of symbols, each  $x_i$  being a selection of a random variable  $X$  drawn from the same ensemble of  $s$  alternatives, can be transmitted. We use the derivation to motivate the formal definition of the entropy of a random variable. The *Morse code* used in telegraphy suggests the general idea. In, say, English, the frequency of use of the letter “e” is 0.12, while the frequency of the letter “w” is only 0.02. Hence, a considerable savings on average encoded message length can result by encoding “e” by a shorter binary string than “w.”

Assume that the random variable  $X$  can take on the alternative values  $\{a_1, a_2, \dots, a_s\}$  and that  $a_i$  occurs  $k_i$  times in the message  $x = x_1 x_2 \dots x_k$ , with  $k_1 + k_2 + \dots + k_s = k$ . Under these constraints there is altogether an ensemble of

$$\binom{k}{k_1, k_2, \dots, k_s} = \frac{k!}{k_1! k_2! \dots k_s!}$$

possible messages of length  $k$ , one of which is  $x$ . In the combinatorial approach we define the entropy of an ensemble as the efficiency with which any message from this ensemble can be transmitted. To determine the actual message  $x_1 x_2 \dots x_k$ , we must at least give its ordinal in the ensemble. To reconstruct the message it suffices to give first the ordinal  $(k_1, k_2, \dots, k_s)$  of the ensemble in  $s \log k$  bits, and then give the ordinal of the message in the ensemble. Therefore, we can transmit the message in  $H(x)$  bits, with

$$\log \frac{k!}{k_1! k_2! \dots k_s!} \leq H(x) \leq s \log k + \log \frac{k!}{k_1! k_2! \dots k_s!}.$$

The *frequency* of each symbol  $a_i$  is defined as  $p_i = k_i/k$ . Recall the approximation  $k \log k$  for  $\log(k!)$  from Stirling’s formula, Exercise 1.5.4,

page 17. For fixed frequencies  $p_1, p_2, \dots, p_s$  and large  $k$  we find

$$H(x) \sim k \left( - \sum p_i \log p_i \right),$$

the sum taken in the obvious way. In information-theoretic terminology it is customary to say that the messages are produced by a “stochastic source” that “emits symbols”  $a_i$  with given “probabilities”  $p_i$ . With abuse of terminology and notions, henceforth we use “*probability*” for “frequency.” (Under certain conditions on the stochastic nature of the source this transition can be rigorously justified.)

**Definition 1.11.1** We define the *entropy* of random variable  $X$  with the given probabilities  $P(X = a_i) = p_i$  by

$$H(X) = - \sum p_i \log p_i, \quad (1.8)$$

such that

$$H(x) \sim kH(X). \quad (1.9)$$

Is there a coding method that actually achieves the economy in average message length implied by Equation 1.9? Clearly, we have to encode symbols with high probabilities as short binary strings and symbols with low probabilities as long binary strings.

**Example 1.11.2** We explain the *Shannon-Fano code*. Suppose we want to map messages over a fixed alphabet to binary strings. Let there be  $n$  symbols (also called basic messages or source words). Order these symbols according to decreasing probability, say  $N = \{1, 2, \dots, n\}$  with probabilities  $p_1, p_2, \dots, p_n$ . Let  $P_r = \sum_{i=1}^r p_i$ , for  $r = 1, \dots, n$ . The binary code  $E : N \rightarrow \{0, 1\}^*$  is obtained by coding  $r$  as a binary number  $E(r)$ , obtained by truncating the binary expansion of  $P_r$  at length  $l(E(r))$  such that

$$-\log p_r \leq l(E(r)) < 1 - \log p_r.$$

This code is the *Shannon-Fano code*. It has the property that highly probable symbols are mapped to short code words and symbols with low probability are mapped to longer code words. Moreover,

$$2^{-l(E(r))} \leq p_r < 2^{-l(E(r))+1}.$$

Note that the code for symbol  $r$  differs from all codes of symbols  $r+1$  through  $n$  in one or more bit positions, since for all  $i$  with  $r+1 \leq i \leq n$ ,

$$P_i \geq P_r + 2^{-l(E(r))}.$$

Therefore the binary expansions of  $P_r$  and  $P_i$  differ in the first  $l(E(r))$  positions. This means that  $E$  is one-to-one, and it has an inverse: the decoding mapping  $E^{-1}$ . Even better, since no value of  $E$  is a prefix of any other value of  $E$ , the set of code words is a prefix-code. This means we can recover the source message from the code message by scanning it from left to right without look-ahead.

If  $H_1$  is the average number of bits used per symbol of an original message, then  $H_1 = \sum_r p_{rl}(E(r))$ . Combining this with the previous inequality we obtain

$$-\sum_r p_r \log p_r \leq H_1 < \sum_r (1 - \log p_r) p_r = 1 - \sum_r p_r \log p_r.$$

From this it follows that  $H_1 \sim H(x)$  for large  $n$ , with  $H(X)$  the entropy per symbol of the source.  $\diamond$

**Example 1.11.3** How much information can a random variable  $X$  convey about a random variable  $Y$ ? Taking again a purely combinatorial approach, this notion is captured as follows: If  $X$  ranges over  $S_X$  and  $Y$  ranges over  $S_Y$ , then we look at the set  $U$  of possible events  $(X, Y)$  consisting of joint occurrences of event  $X$  and event  $Y$ . If  $U$  does not equal the Cartesian product  $S_X \times S_Y$ , then this means there is some dependency between  $X$  and  $Y$ . Considering the set  $U_a = \{(a, y) : (a, y) \in U, a \in S_X\}$ ; it is natural to define the *conditional entropy* of  $Y$  given  $X = a$  as  $H(Y|a) = \log d(U_a)$ . This suggests immediately that the information given by  $X$  about  $Y$  is

$$I(X : Y) = H(Y) - H(Y|X).$$

For example, if  $U = \{(1, 1), (1, 2), (2, 3)\}$ ,  $U \subseteq S_X \times S_Y$  with  $S_X = \{1, 2\}$  and  $S_Y = \{1, 2, 3, 4\}$ , then  $I(X = 1 : Y) = 1$  and  $I(X = 2 : Y) = 2$ .

In this formulation it is obvious that  $H(X|X) = 0$ , and that  $I(X : X) = H(X)$ . This approach amounts to the assumption of *uniform distribution* of the probabilities concerned.  $\diamond$

We can develop the generalization of Example 1.11.3, taking into account the frequencies or probabilities of the occurrences of the different values  $X$  and  $Y$  can assume. Let the *joint probability*  $p(X, Y)$  be defined by, “ $p(a, b)$  is the probability of the joint occurrence of event  $X = a$  and event  $Y = b$ .” This leads to the self-evident formulas for joint variables  $X, Y$ :

$$H(X, Y) = -\sum_{X,Y} p(X, Y) \log p(X, Y),$$

$$H(X) = -\sum_{X,Y} p(X, Y) \log \sum_Y p(X, Y),$$

$$H(Y) = -\sum_{X,Y} p(X, Y) \log \sum_X p(X, Y).$$

We can show easily that

$$H(X, Y) \leq H(X) + H(Y), \quad (1.10)$$

with equality only in the case that  $X$  and  $Y$  are independent. In all of these equations the entropy quantity on the left-hand side increases if we choose the probabilities on the right-hand side more equally.

The conditional probability  $p(Y|X)$  of outcome  $Y = b$  given outcome  $X = a$  for random variables  $X$  and  $Y$  (not necessarily independent) is defined by

$$p(Y|X) = \frac{p(X, Y)}{\sum_Y p(X, Y)},$$

Section 1.6.2. This leads to the following analysis of the information in  $X$  about  $Y$  by first considering the *conditional entropy* of  $Y$  given  $X$  as the average of the entropy for  $Y$  for each value of  $X$  weighted by the probability of getting that particular value:

$$\begin{aligned} H(Y|X) &= - \sum_{X,Y} p(X, Y) \log p(Y|X) \\ &= - \sum_Y p(Y|X) \log p(Y|X). \end{aligned}$$

The quantity on the left-hand side tells us how uncertain we are of  $Y$  on average when we know  $X$ . With

$$\begin{aligned} H(Y) &= - \sum_Y p(Y) \log p(Y) \\ &= - \sum_{X,Y} p(X, Y) \log \sum_x p(X, Y)), \end{aligned}$$

and substituting the formula for  $p(Y|X)$ , we find  $H(Y|X) = H(X, Y) - H(X)$ . Rewrite this expression as

$$H(X, Y) = H(X) + H(Y|X). \quad (1.11)$$

This can be interpreted as, “the uncertainty of the joint event  $(X, Y)$  is the uncertainty of  $X$  plus the uncertainty of  $Y$  given  $X$ .” Combining Equations 1.10, 1.11 gives  $H(Y) \geq H(Y|X)$ , which can be taken to imply that knowledge of  $X$  can never increase uncertainty of  $Y$ . In fact, uncertainty in  $Y$  will be decreased unless  $X$  and  $Y$  are independent. Finally, the *information* in  $X$  about  $Y$  is defined as

$$I(X : Y) = H(Y) - H(Y|X). \quad (1.12)$$

Here the quantities  $H(Y)$  and  $H(Y|X)$  on the right-hand side of the equations are always equal to or less than the corresponding quantities under the uniform distribution we analyzed first. The values of the quantities  $I(X : Y)$  under the assumption of uniform distribution of  $Y$  and  $Y|X$  versus any other distribution are not related by inequality in a particular direction. The equalities  $H(X|X) = 0$  and  $I(X : X) = H(X)$  hold under any distribution of the variables. Since  $I(X : Y)$  is a function of  $X$ , while  $I(Y : X)$  is a function of  $Y$ , we do not compare them directly. However, forming the expectation defined as

$$\mathbf{E}(I(X : Y)) = \sum_X p(X)I(X : Y),$$

$$\mathbf{E}(I(Y : X)) = \sum_Y p(Y)I(Y : X),$$

and combining Equations 1.11, 1.12, we see that the resulting quantities are equal. Denoting this quantity by  $I(X; Y)$  and calling it the *mutual information* in  $X$  and  $Y$ , we see that this information is *symmetric*:

$$I(X; Y) = \mathbf{E}(I(X : Y)) = \mathbf{E}(I(Y : X)). \quad (1.13)$$

The quantity  $I(X; Y)$  symmetrically characterizes to what extent random variables  $X$  and  $Y$  are correlated. An inherent problem with probabilistic definitions (which is avoided by the combinatorial approach) is that although  $\mathbf{E}(I(X : Y))$  is always positive, for some probability distributions  $I(X : Y)$  can turn out to be negative—which definitely contradicts our naive notion of information content.

The development of this theory immediately gave rise to at least two different questions. The first observation is that the concept of information as used in the theory of communication is a probabilistic notion, which is natural for information transmission over communication channels. Nonetheless, as we have seen from the discussion, we tend to identify *probabilities* of messages with *frequencies* of messages in a sufficiently long sequence, which under some conditions on the stochastic source can be rigorously justified. For instance, Morse code transmissions of English telegrams over a communication channel can be validly treated by probabilistic methods even if we (as is usual) use empirical frequencies for probabilities. The great probabilist, Kolmogorov, remarks, “If something goes wrong here, the problem lies in the vagueness of our ideas of the relation between mathematical probability theory and real random events in general.” (See also the discussion about the foundations of probability in Section 1.9.)

The second observation is more important and is exemplified in Shannon’s statement, “Messages have *meaning* [...] however [...] the semantic aspects of communication are irrelevant to the engineering problem.” In

other words, can we answer a question like “what is the information in this book” by viewing it as an element of a set of possible books with a probability distribution on it? Or that the individual sections in this book form a random sequence with stochastic relations that damp out rapidly over a distance of several pages? And how to measure the quantity of hereditary information in biological organisms, as encoded in DNA? Again there is the possibility of seeing a particular form of animal as one of a set of possible forms with a probability distribution on it. This seems to be contradicted by the fact that the calculation of all possible lifeforms in existence at any one time on earth would give a ridiculously low figure like  $2^{100}$ .

We are interested in a measure of information content of an *individual finite object*, and in the information conveyed about an individual finite object by another individual finite object. Here, we want the information content of an object  $x$  to be an attribute of  $x$  *alone*, and not to depend on, for instance, the means chosen to describe this information content. Making the natural restriction that the description method should be effective, the information content of an object should be a recursively invariant property (Section 1.7) among the different description systems. Pursuing this thought leads straightforwardly to Kolmogorov complexity.

### 1.11.1 Prefix-Codes

The main issues we treat here, apart from the basic definitions of prefix-codes, are Kraft’s Inequality, the Noiseless Coding Theorem, and optimal and universal codes for infinite source word sets.

We repeat some definitions of Section 1.4. Let  $D$  be any function  $D : V^* \rightarrow \mathcal{N}$ , where  $V$  is a finite code word alphabet. It is common to use  $V = \{0, 1\}$ . The domain of  $D$  is the set of code words and the range of  $D$  is the set of source words. If  $D(y) = x$ , then  $y$  is a code word for source word  $x$ , and  $D$  is the decoding function. The set of all code words for source word  $x$  is the set  $D^{-1}(x) = \{y : D(y) = x\}$  and  $E = D^{-1}$  is the encoding relation (or encoding function if  $D^{-1}$  happens to be a function).

We may identify the source words (natural numbers) with their corresponding finite binary strings according to the enumeration in Equation 1.1. We often identify a code with its code word set (the domain of  $D$ ).

We consider the natural extension of  $E$  to a relation  $E' \subseteq \mathcal{N}^* \times V^*$  defined by

1.  $E'(\epsilon) = \epsilon$ ; and
2. if  $x$  and  $y$  are in  $\mathcal{N}$ , then  $E'(xy) = E(x)E(y)$ .

We ignore the difference between  $E$  and  $E'$  and denote both by  $E$ .

**Example 1.11.4** It is immediately clear that in general we cannot uniquely recover  $x$  and  $y$  from  $E(xy)$ . Let  $E$  be the identity mapping. Then we have  $E(00)E(00) = 0000 = E(0)E(000)$ .  $\diamond$

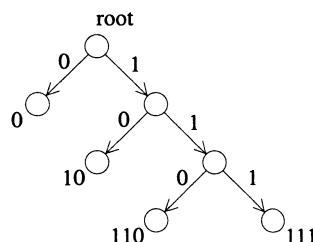
If we want to encode a sequence  $x_1x_2\dots x_n$  with  $x_i \in \mathcal{N}$  ( $i = 1, 2, \dots$ ), then we call  $x_1x_2\dots x_n$  the *source sequence* and  $y_1y_2\dots y_n$  with  $y_i = E(x_i)$  ( $i = 1, 2, \dots$ ) the *code sequence*. A code is *uniquely decodable* if for each source sequence of finite length, the code sequence corresponding to that source sequence is different from the code sequence corresponding to any other source sequence.

**Example 1.11.5** In coding theory attention is often restricted to the case where the source word set is finite, say the range of  $D$  equals  $\{1, 2, \dots, n\}$ . If there is a constant  $l_0$  such that  $l(y) = l_0$  for all code words  $y$ , then we call  $D$  a *fixed-length* code. It is easy to see that  $l_0 \geq \log n$ . For instance, in teletype transmissions the source has an alphabet of  $n = 32$  letters, consisting of the 26 letters in the Latin alphabet plus 6 special characters. Hence, we need  $l_0 = 5$  binary digits per source letter. In electronic computers we often use the fixed-length ASCII code with  $l_0 = 8$ .  $\diamond$

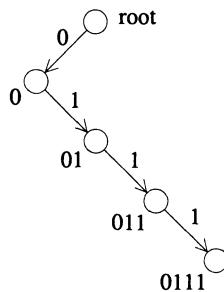
When the set of source words is infinite, say  $\mathcal{N}$ , we have to use *variable-length codes*. If no code word is the prefix of another code word, then each code sequence is uniquely decodable. This explains our interest in so-called prefix-codes.

**Definition 1.11.2** A code is a *prefix-code* or *instantaneous code* if the set of code words is prefix-free (no code word is a prefix of another code word, Section 1.4).

In order to decode a code sequence of a prefix-code, we simply start at the beginning and decode one code word at a time. When we come to the end of a code word, we know it is the end, since no code word is the prefix



**FIGURE 1.3.** Binary tree for  $E(1) = 0$ ,  $E(2) = 10$ ,  $E(3) = 110$ ,  $E(4) = 111$



**FIGURE 1.4.** Binary tree for  $E(1) = 0$ ,  $E(2) = 01$ ,  $E(3) = 011$ ,  $E(4) = 0111$

of any other code word in a prefix-code. Every prefix-code is a uniquely decodable code. For example, if  $E(1) = 0$ ,  $E(2) = 10$ ,  $E(3) = 110$ ,  $E(4) = 111$  as in Figure 1.3, then 1421 is encoded as 0111100, which can be easily decoded from left to right in a unique way.

Not every uniquely decodable code satisfies the prefix condition. For example, if  $E(1) = 0$ ,  $E(2) = 01$ ,  $E(3) = 011$ ,  $E(4) = 0111$ , then every code word is a prefix of every longer code word as in Figure 1.4. But unique decoding is trivial, since the beginning of a new code word is always indicated by a zero. Prefix-codes are distinguished from other uniquely decodable codes by the property that the end of a code word is always recognizable as such. This means that decoding can be accomplished without the delay of observing subsequent code words, which is why prefix-codes are also called instantaneous codes.

**Example 1.11.6** A convenient graphical way to consider codes is by representing each code word as a node of a directed binary tree. If a node has two outgoing arcs, one of them is labeled with zero and the other with one. If a node has one outgoing arc, it is labeled either by zero or by one. The tree may be finite or infinite. If it is finite, then there are nodes without outgoing arcs. Nodes without outgoing arcs are called *external nodes*, and the nodes with outgoing arcs are called *internal nodes*. Each code word is represented by a node such that the consecutive zeros and ones on the branch from the root to that node form that code word. Clearly, for each code there is a tree with a node representing each code word, and such that there is only one node corresponding with each code word. We simplify each tree representation for a code such that it contains only nodes corresponding to code words, together with the intermediate nodes on a branch between root and code word nodes. For  $E$  to be a prefix-code, it is a necessary and sufficient condition that in the simplified tree representation the nodes corresponding to code words are precisely the nodes without outgoing arcs.  $\diamond$

### 1.11.2 The Kraft Inequality

It requires little reflection to realize that prefix-codes waste potential code words since the internal nodes of the representation tree cannot be used, and in fact neither are the potential descendants of the external nodes used. Hence, we can expect that the code-word length exceeds the (binary) source-word length in prefix-codes. Quantification of this intuition leads to a precise constraint on code-word lengths for codes satisfying the prefix condition. This important relation is known as the *Kraft Inequality* and is due to L.G. Kraft.

**Theorem 1.11.1** *Let  $l_1, l_2, \dots$  be a finite or infinite sequence of natural numbers. There is a prefix-code with this sequence as lengths of its binary code words iff*

$$\sum_n 2^{-l_n} \leq 1.$$

**Proof. (ONLY IF)** Recall the standard one-to-one correspondence between a finite binary string  $x$  and the interval  $\Gamma_x = [0.x, 0.x + 2^{-l(x)})$  on the real line, Sections 1.4, 1.6, 2.5. Observe that the length of the interval corresponding to  $x$  is  $2^{-l(x)}$ . A prefix-code corresponds to a set of disjoint such intervals in  $[0, 1)$ , which proves that the inequality holds for prefix-codes.

**(IF)** Suppose  $l_1, l_2, \dots$  are given such that the inequality holds. We can also assume that the sequence is nondecreasing. Choose disjoint adjacent intervals  $I_1, I_2, \dots$  of lengths  $2^{-l_1}, 2^{-l_2}, \dots$  from the left end of the interval  $[0, 1)$ . This way, for each  $n \geq 1$  the right end of  $I_n$  is  $\sum_{i=1}^n 2^{-l_i}$ . Note that the right end of  $I_n$  is the left end of  $I_{n+1}$ . Since the sequence of  $l_i$ 's is nondecreasing, each interval  $I_n$  equals  $\Gamma_x$  for some binary string  $x$  of length  $l(x) = l_n$ . Take the binary string  $x$  corresponding to  $I_n$  as the  $n$ th code word.  $\square$

**Example 1.11.7** Not every code of which the code-word lengths satisfy the inequality is a prefix-code. For instance, the code words 0, 00, and 11 satisfy the inequality, but 0 is a prefix of 00.  $\diamond$

**Example 1.11.8** There is good reason for our emphasis on prefix-codes. Namely, Theorem 1.11.1 stays valid if we replace “prefix-code” by “uniquely decodable code.” This follows directly from the observation (proof omitted) that if a code has code-word lengths  $l_1, l_2, \dots$  and it is uniquely decodable, then the Kraft Inequality must be satisfied.

This important fact means that every uniquely decodable code can be replaced by a prefix-code without changing the set of code-word lengths. Hence, all propositions concerning code-word lengths apply to uniquely decodable codes and to the subclass of prefix-codes. Accordingly, in looking for uniquely decodable codes with minimal average code-word length we can restrict ourselves to prefix-codes.  $\diamond$

### 1.11.3 Optimal Codes

A uniquely decodable code is *complete* if the addition of any new code word to its code word set results in a nonuniquely decodable code. It is easy to see that a code is complete iff the associated Kraft Inequality is satisfied with equality. Does completeness imply optimality in any reasonable sense? Given a source that produces source words from  $\mathcal{N}$  according to probability distribution  $P$ , it is possible to assign code words to source words in such a way that any code word sequence is uniquely decodable, and moreover the average code-word length is minimal.

**Definition 1.11.3** Let  $D : \{0, 1\}^* \rightarrow \mathcal{N}$  be a prefix-code with one code word per source word. Let  $P(x)$  be the probability of source word  $x$ , and let the length of the code word for  $x$  be  $l_x$ . We want to minimize the number of bits we have to transmit. In order to do so, we must minimize the *average code-word length*  $L_{D,P} = \sum_x P(x)l_x$ . We define the *minimal average code word length* as  $L = \min\{L_{D,P} : D \text{ is a prefix-code}\}$ . A prefix-code  $D$  such that  $L_{D,P} = L$  is called an *optimal prefix-code* with respect to prior probability  $P$  of the source words.

The (minimal) average code length of an (optimal) code does not depend on the details of the set of code words, but only on the set of code-word lengths. It is just the expected code-word length with respect to the given distribution. C.E. Shannon discovered that the minimal average code word length is about equal to the entropy of the source word set. This is known as the *Noiseless Coding Theorem*. The adjective “noiseless” emphasizes that we ignore the possibility of errors.

**Theorem 1.11.2** Let  $L$  and  $P$  be as above. If  $H(P) = -\sum_x P(x) \log P(x)$  is the entropy, then

$$H(P) \leq L \leq H(P) + 1.$$

**Proof.** First prove the upper bound  $L \leq H(P) + 1$ . Let  $l_x = \lceil -\log P(x) \rceil$  for  $x = 1, 2, \dots$ . Therefore,  $1 \geq \sum_x P(x) \geq \sum_x 2^{-l_x}$ . By Kraft’s Inequality, Theorem 1.11.1, there exists a prefix-code with code-word lengths  $l_1, l_2, \dots$ . Hence,

$$L \leq \sum_x P(x)l_x \leq \sum_x P(x)(-\log P(x) + 1) = H(P) + 1,$$

which finishes the proof of the second inequality  $L \leq H(P) + 1$ .

We now prove the lower bound  $H(P) \leq L$ . Let  $L = \sum_x P(x)l_x$ . Since  $\sum_x P(x) = 1$  and  $\sum_x (2^{-l_x} / \sum_x 2^{-l_x}) = 1$ , by concavity of the logarithm function (see Equations 5.5 and 5.6 on page 329), we have

$$-\sum_x P(x) \log P(x) \leq -\sum_x P(x) \log \frac{2^{-l_x}}{\sum_x 2^{-l_x}}. \quad (1.14)$$

Equation 1.14 implies

$$-\sum_x P(x) \log P(x) \leq \sum_x P(x)l_x + \left( \sum_x P(x) \right) \log \sum_x 2^{-l_x}. \quad (1.15)$$

Since  $\sum_x P(x) = 1$ ,  $L = \sum_x P(x)l_x$ , and  $H(P) = -\sum_x P(x) \log P(x)$ , Equation 1.15 can be rewritten as

$$H(P) \leq L + \log \sum_x 2^{-l_x}. \quad (1.16)$$

As  $D$  is a prefix-code, it follows from Kraft's Inequality, Theorem 1.11.1, that  $\sum_x 2^{-l_x} \leq 1$ . Thus,  $\log \sum_x 2^{-l_x} \leq 0$ . Hence by Equation 1.16,  $H(P) \leq L$ .  $\square$

**Example 1.11.9** We can now settle the question of whether complete code word sets are necessarily optimal for all prior distributions in the negative. Let  $E$  be a prefix-code with source alphabet  $\{0, \dots, k+1\}$  defined by  $E(x) = 1^{x-1}0$  for  $x = 1, 2, \dots, k$  and  $E(k+1) = 1^k$ . Then  $E$  is obviously complete. It has an average code-word length  $L_{E,P} = \sum_x P(x)x$  with respect to probability distribution  $P$ . If  $P(x) = 2^{-x}$  for  $x = 1, 2, \dots, k$  and  $P(k+1) = 2^{-k}$ , then  $L_{E,P} = -\sum_x P(x) \log P(x)$  so that the expected code-word length is exactly equal to the entropy and hence to the minimal code-word length  $L$  by the Noiseless Coding Theorem 1.11.2. But for the uniform distribution  $P(x) = 1/(k+1)$  for  $x = 1, 2, \dots, k+1$  we find  $L_{E,P} \gg L$  so that the code is not optimal with respect to this distribution.  $\diamond$

It is obviously important to find optimal prefix-codes. We have seen that completeness has not much to do with it. For optimality of finite codes we must choose code-word lengths corresponding to the probabilities of the encoded source words. This idea is implemented in the Shannon-Fano code, Example 1.11.2 on page 67 and Lemma 4.3.3 on page 253.

**Example 1.11.10** Another issue is the effectiveness of the decoding process. The decoder needs to match the code word patterns to code word sequence in order to retrieve the source word sequence. For a finite code word set, the code can be stored in a finite table. For infinite code word sets we must recognize the code words.

**Definition 1.11.4** A prefix-code is called *self-delimiting* if there is a Turing machine that decides whether a given word is a code word, never reading beyond the word itself, and moreover, computes the decoding function. (With respect to the Turing machine each code word has an implicit end marker.)

As an example, let us define the minimal description length for elements in  $\mathcal{N}$  with respect to the class of Turing machines  $\mathcal{T}$ . Fix a self-delimiting code  $E : \mathcal{T} \rightarrow \mathcal{N}$ . Denote the code word for Turing machine  $T$  by  $E(T)$ . Then the minimum description length of  $x \in \mathcal{N}$  with respect to  $E$  is defined as  $\min\{l(E(T)y) : T \text{ on input } y \text{ halts with output } x\}$ . This minimum description length of  $x$  is actually an alternative definition of the Kolmogorov complexity  $C(x)$ .  $\diamond$

#### 1.11.4 Universal Codes

For finite codes the optimality is governed by how closely the set of code-word lengths corresponds to the probability distribution on the set of source words. In the proof of the Noiseless Coding Theorem 1.11.2, we chose a code that corresponds to the probability distribution of the source words. But the actual probability distribution may be unknown, nonrecursive, or it may be unclear how to determine the characteristics of the source. For example, consider the encoding of the printed English language, which emanates from many different sources with (it is to be assumed) different characteristics. Can we find a code that is optimal for *any* probability distribution, rather than for a *particular* one?

##### Example 1.11.11

How does one transmit any sufficiently long source word sequence in an optimal-length code word sequence without knowing the characteristics of the source, and in particular, without knowing in advance the relative frequencies of source words for the entire sequence? The solution is as follows:

Split the source word sequence  $x = x_1x_2\dots x_n$  of length  $n$  into  $m \leq n$  blocks  $b_1b_2\dots b_m$ . Then we encode each  $b_i$  by, for instance, first giving the numbers of occurrences of the source words in  $b_i$  and second the index of  $b_i$  in the lexicographically ordered ensemble of source words determined by these numbers. For instance, let  $b_i$  have length  $n_i$ , and let there be  $q$  source words where the numbers of occurrences of the different source words in  $b_i$  are, respectively,  $k_1, k_2, \dots, k_q$ ,  $\sum_j k_j = n_i$ . Encode  $b_i$  by the frequency vector  $(k_1, k_2, \dots, k_q)$ , together with the index  $h$  of  $b_i$  in the ensemble of all possible sequences of length  $n_i$  in which the source words occur with these frequencies,

$$h \leq \binom{n_i}{k_1, k_2, \dots, k_q}.$$

The encoding of  $b_i$  in standard decodable format with all items except the last one self-delimiting as, say,  $\overline{k_1} \dots \overline{k_q} h$ , takes at most

$$O(q \log n_i) + \log \frac{n_i!}{k_1! k_2! \dots k_q!}$$

bits. Defining the frequency of each source word  $a_j$  in block  $b_i$  as  $p_j = k_j/n_i$  we find that the length of this encoding of  $b_i$  for large  $n_i$  and

fixed  $p_j$ 's approaches  $n_i(-\sum_j p_j \log p_j)$ . By the Noiseless Coding Theorem 1.11.2, the minimal average code-word length for a source word sequence  $b_i$  produced by a stochastic source that emits source word  $a_j$  with probability  $P^{(i)}(a_j) = p_j$  is given by  $-\sum_j p_j \log p_j = H(P^{(i)})$ .

That is, we can separately encode each block  $b_i$  asymptotically optimal, without knowing anything about the overall relative frequencies. As a result, the overall message  $x$  is encoded asymptotically in length

$$n_1 H(P^{(1)}) + n_2 H(P^{(2)}) + \cdots + n_m H(P^{(m)}).$$

It turns out that this implies that  $x$  is optimally encoded as well, since calculation shows that

$$nH(P) \geq n_1 H(P^{(1)}) + n_2 H(P^{(2)}) + \cdots + n_m H(P^{(m)}),$$

with  $H(P)$  the entropy based on the overall source word sequence  $x$  and  $H(P^{(i)})$  the entropy of the individual blocks  $b_i$ .  $\diamond$

Suppose we use variable-length binary blocks  $b_i$  as in Example 1.11.11 as code words for a countably infinite set of source words such as  $\mathcal{N}$ . Then a universal code is a code that optimizes the average code-word length, independent of the distribution of the source words, in the following sense:

**Definition 1.11.5** Let  $C = \{w_1, w_2, \dots\} \subseteq \{0, 1\}^*$  be a set of code words with  $l(w_1) < l(w_2) < \dots$ , and let  $\mathcal{N}$  be a set of source words with probability distribution  $P$ . Code  $C$  is *universal* if there is a constant  $c$ , independent of  $P$ , such that

$$\frac{\sum_x P(x)l(w_x)}{\max\{H(P), 1\}} \leq c,$$

for all  $P$  with  $0 < H(P) < \infty$ . A universal code  $C$  is *asymptotically optimal* if there is an  $f$  such that

$$\frac{\sum_x P(x)l(w_x)}{\max\{H(P), 1\}} \leq f(H(P)) \leq c,$$

with  $\lim_{H(P) \rightarrow \infty} f(H(P)) = 1$ .

**Example 1.11.12** For a nonempty finite binary string  $x = x_1 x_2 \dots x_n$  we defined  $\bar{x} = 1^n 0 x_1 x_2 \dots x_n$ . For example, 01011 is coded as 11111001011. Let  $C$  be defined as the set of binary strings resulting from this construction:  $C = \{\bar{x} : x \in \{0, 1\}^*\}$ . The proof that this is a universal set of code words, but not an asymptotically optimal one, is omitted.

However, a relatively minor improvement yields an asymptotically optimal code word set. This time we encode  $x$  not by  $\bar{x}$ , but by  $E(x) = l(x)x$ , that is, by encoding first the length of  $x$  in prefix-free form, followed by the literal representation of  $x$ . For example, 01011 is now coded as 1101001011, encoding  $l(x)$  as 10 according to Equation 1.1. Code  $E$  is prefix-free, since if we know the length of  $x$  as well as the start of its literal representation, then we also know where it ends. The length set of this code is given by  $l(E(x)) = l(x) + 2l(l(x)) + 1$  for  $x \in \mathcal{N}$ . The proof that this code is asymptotically optimal universal is omitted.  $\diamond$

This shows that there are asymptotically optimal prefix-codes. We now inquire how far we can push the idea involved.

- Example 1.11.13** It is straightforward to improve on Example 1.11.12 by iterating the same idea, as in Equation 1.2, page 13. That is, we precede  $x$  by its length  $l(x)$ , and in turn precede  $l(x)$  by its own length  $l(l(x))$ , and so on. That is, the prefix-code  $E_i$ , for all finite strings  $x$ , is defined by  $E_i(x) = E_{i-1}(l(x))x$  and  $E_1(x) = \bar{x}$ . For example, with  $i = 3$ , the string 01011 is coded as 1011001011, using the correspondence of Formula 1.1. This is a kind of “ladder” code, where the value of  $i$  is supposed to be fixed and known to coder and decoder. The length of  $E_i(x)$  is given by

$$l(E_i(x)) = \begin{cases} 2l(x) + 1 & \text{if } i = 1, \\ l(x) + l(E_{i-1}(l(x))) & \text{if } i > 1. \end{cases}$$

For each fixed  $i > 1$  the prefix-code  $E_i$  is self-delimiting, universal, and asymptotically optimal.

Can we improve this? Denote  $l^k(x) = l(l^{k-1}(x))$ , the  $k$ -fold iteration of taking the length, and  $l^1(x) = l(x)$ . Define a coding

$$E^*(x) = l^k(x)0l^{k-1}(x)0\ldots0x1,$$

with  $k$  the number of steps necessary to get  $l^k(x) = 1$ . For instance, the string 01011 is coded as 10100010111, using the correspondence of Equation 1.1. Again, this code is self-delimiting, universal, and asymptotically optimal. The code-word length is given by

$$l(E^*(x)) = 1 + \sum_{i=1}^k (l^i(x) + 1).$$

This is within an  $O(k)$  additive term of  $l^*(x)$ , defined as

$$l^*(x) = \log x + \log \log x + \log \log \log x + \dots, \quad (1.17)$$

where the sum involves only the positive terms. The number of terms is denoted by  $\log^* x$ . It can be proved that  $\sum_x 2^{-l^*(x)} = c$  is finite, with  $c =$

2.865064.... If we put  $l_x = l^*(x) + \log c$ , then the Kraft Inequality is satisfied with equality.

However, a lower bound on the code length is set by  $\ell^*(x)$ , defined as

$$\ell^*(x) = \begin{cases} l(x) + \ell^*(l(x)) & \text{if } l(x) > 1, \\ l(x) & \text{if } l(x) = 0, 1. \end{cases} \quad (1.18)$$

It can be shown that

$$\begin{aligned} l^*(x) - \ell^*(x) &\leq \log^* x, \\ \sum_x 2^{-\ell^*(x)} &= \infty; \quad \sum_x 2^{-l^*(x)} < 3. \end{aligned} \quad (1.19)$$

Hence, although  $\ell^*$  is pretty close to  $l^*$ , by the divergence of the  $\ell^*$  series in Equation 1.19 it follows by the Kraft Inequality, Theorem 1.11.1, that there is no prefix-code with length set  $\{\ell^*(x) : x \in \mathcal{N}\}$ .

◇

## Exercises

**1.11.1.** [10] To use an extra symbol like 2 is costly when expressed in bits. Show that the coding of strings consisting of  $k$  zeros and ones and one 2 requires messages of about  $k + \log k$  bits.

*Comments* Hint: there are  $2^k(k+1)$  such strings.

**1.11.2.** [13] Suppose we have a random variable  $X$  that can assume values  $a, b, c, d$  with probabilities  $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}$ , and  $\frac{1}{8}$ , respectively, with no dependency between the consecutive occurrences.

(a) Show that the entropy  $H(X) = \frac{7}{4}$  (bits per symbol).

(b) Show that the code  $E$  with  $E(a) = 0, E(b) = 10, E(c) = 110$ , and  $E(d) = 111$  achieves this limiting value.

**1.11.3.** [10] Let  $M$  be the set of possible source messages, and let  $E : M \rightarrow \{0, 1\}^*$  be a prefix-code.

(a) Let  $M$  be a set of messages using symbols in an alphabet  $A$ . Show that if  $E$  is a prefix-code on  $A$ , then the homomorphism induced by  $E$  is a prefix-code on  $M$ .

(b) Show that the Shannon-Fano code in the main text is a prefix-code.

**1.11.4.** [26] Prove that the entropy function  $H$  has the following four properties:

(a) For given  $n$  and  $\sum_{i=1}^n p_i = 1$ , the function  $H(p_1, p_2, \dots, p_n)$  takes its largest value for  $p_i = 1/n$  ( $i = 1, 2, \dots, n$ ). That is, the scheme with the most uncertainty is the one with equally likely outcomes.

- (b)  $H(X, Y) = H(X) + H(Y|X)$ . That is, the uncertainty in the product of schemes  $x$  and  $y$  equals the uncertainty in scheme  $x$  plus the uncertainty of  $y$  given that  $x$  occurs.
- (c)  $H(p_1, p_2, \dots, p_n) = H(p_1, p_2, \dots, p_n, 0)$ . That is, adding the impossible event or any number of impossible events to the scheme does not change its entropy.
- (d)  $H(p_1, p_2, \dots, p_n) = 0$  iff one of the numbers  $p_1, p_2, \dots, p_n$  is one and all the others are zero. That is, if the result of the experiment can be predicted beforehand with complete certainty, then there is no uncertainty as to its outcome. In all other cases the entropy is positive.

*Comments.* Source: C.E. Shannon, *Bell System Tech. J.*, 27(1948), 379–423, 623–656.

**1.11.5.** [32] Prove the following theorem. Let  $H(p_1, p_2, \dots, p_n)$  be a function defined for any integer  $n$  and for all values  $p_1, p_2, \dots, p_n$  such that  $p_i \geq 0$  ( $i = 1, 2, \dots, n$ ),  $\sum_{i=1}^n p_i = 1$ . If for any  $n$  this function is continuous with respect to all its arguments, and if it has the properties in Items (a), (b), and (c) of the previous exercise, then  $H(p_1, p_2, \dots, p_n) = -\lambda \sum_{i=1}^n p_i \log p_i$ , where  $\lambda$  is a positive constant.

*Comments.* This is called the *Entropy Uniqueness Theorem*. It shows that our choice of expression for the entropy for a finite probability scheme is the only one possible, if we want to have certain general properties that seem necessary in view of the intended meaning of the notion of entropy as a measure of uncertainty or as amount of information. Source: A.I. Khintchin, *Mathematical Foundations of Information Theory*, Dover, 1957.

**1.11.6.** [08] Define a code  $c$  by  $E(1) = 00$ ,  $E(2) = 01$ ,  $E(3) = 10$ ,  $E(4) = 11$ ,  $E(5) = 100$ .

- (a) Show that the code sequence 00011011100 is uniquely decodable.
- (b) Show that the code sequence 100100 can be decoded into two different source sequences.

**1.11.7.** [09] (a) Define code  $E$  by  $E(x) = 01^x$  for  $x = 1, 2, \dots$ . Show that  $E$  is uniquely decodable, but it is not a prefix-code.

(b) Define a code  $E$  by  $E(x) = 1^x 0$  for  $x = 1, 2, \dots$ . Show that this is a prefix-code and (hence) uniquely decodable.

**1.11.8.** [19] Let  $K = \{x : \phi_x(x) < \infty\}$  be the diagonal halting set (Section 1.7). Let  $k_1, k_2, \dots$  be the list of elements of  $K$  in increasing order. Define code  $E$  by  $E(x) = 1^{k_x} 0$ .

- (a) Show that  $E$  is a prefix-code and uniquely decodable.
- (b) Show that  $E$  is not recursive.

- (c) Show that the decoding function  $E^{-1}$  is not recursive.

*Comments.* Hint: the set  $K$  is recursively enumerable but not recursive. Codes with nonrecursive code word sets are abundant since there are uncountably many codes (prefix-codes), while there can be only countably many recursive code word sets.

- 1.11.9.** [22] Let a code have the set of code-word lengths  $l_1, l_2, \dots$ . Show that if the code is uniquely decodable, then the Kraft Inequality, Theorem 1.11.1, page 74, must be satisfied.

*Comments.* Thus, Theorem 1.11.1 holds for the wider class of uniquely decodable codes. This is called the *McMillan-Kraft Theorem*. Source: R.G. Gallager, *Information Theory and Reliable Communication*, Wiley, 1968. Attributed to B. McMillan.

- 1.11.10.** [26] (a) Show that the set of code words  $\{\bar{x} : x \in \mathcal{N}\}$  is a universal code word set. Show that it is *not* asymptotically optimal.

- (b) Show that the set of code words  $\{\overline{l(x)}x : x \in \mathcal{N}\}$  is a universal code word set. Show that it is also asymptotically optimal.

*Comments.* Source: P. Elias, *IEEE Trans. Inform. Theory*, IT-21(1975), 194–203.

- 1.11.11.** [37] The prefix-code  $E^*$  of Example 1.11.13 is an asymptotically optimal universal code because already  $E_2$  is one.

- (a) Show that  $\sum_n 2^{-l^*(n)} = c$  with  $c = 2.865064\dots$ . Show that the Kraft Inequality, Theorem 1.11.1, is satisfied with equality for the length set  $l_n = l^*(n) + \log c$ ,  $n = 1, 2, \dots$ .

- (b) Let the code  $E^+$  be defined by

$$E^+(x) = \bar{k}l^k(x)l^{k-1}(x)\dots x,$$

where the length function is iterated until the value  $l^k(x) = 1$ . Show that  $E^+$  is prefix-free. Show that this representation of the integers is even more compact than  $E^*$ .

*Comments.* Source for Item (a): J. Rissanen, *Ann. Statist.*, 11(1982), 416–431.

- 1.11.12.** [29] The function  $\log^* n$  denotes the number of times we can iterate taking the binary logarithm with a positive result, starting from  $n$ . This function grows extremely slowly. It is related to the Ackermann function of Exercise 1.7.18. In that notation it is a sort of inverse of  $f(3, x, 2)$ .

- (a) Let  $l_1, l_2, \dots$  be any infinite integer sequence that satisfies the Kraft Inequality, Theorem 1.11.1. Show that  $l_n > l^*(n) - 2\log^* n$  for infinitely many  $n$ .

- (b) Show that  $\log^* n$  is unbounded and primitive recursive. In particular, show that although  $\log^* n$  grows very slowly, it does not grow more slowly than any unbounded primitive recursive function.

*Comments.* Hint: use exercises in Section 1.7. Because  $\log^* n$  grows very slowly, we conclude that  $l^*(n)$  is not far from the least asymptotic upper bound on the code-word length set for all probability sequences on the positive integers. In this sense it plays a similar role for binary prefix-codes as our one-to-one pairing of natural numbers and binary strings in Equation 1.1 plays with respect to arbitrary binary codes. Source: J. Rissanen, *Ibid.*

**1.11.13.** [M30] We consider convergence of the series  $\sum_n 2^{-f(k,\alpha;n)}$  for  $f(k,\alpha;n) = \log n + \log \log n + \dots + \alpha \log^{(k)} n$ , with  $\log^{(k)}$  the  $k$ -fold iteration of the logarithmic function defined by  $\log^{(1)} n = \log n$  and  $\log^{(k)} n = \log \log^{(k-1)} n$  for  $k > 1$ .

- (a) Show that for each  $k \geq 2$  the series above diverges if  $\alpha \leq 1$  and that the series converges if  $\alpha > 1$ .
- (b) Show that the series  $\sum_n n^{-\alpha}$  diverges for  $\alpha \leq 1$  and converges for  $\alpha > 1$ . (This is the case  $k = 1$ .)

*Comments.* This gives an exact borderline between convergence and divergence for the series in Kraft's Inequality, Theorem 1.11.1. Hint: use Cauchy's condensation test for convergence of series. Source: K. Knopp, *Infinite Sequences and Series*, Dover, 1956. Attributed to N.H. Abel.

**1.11.14.** [22] We derive a lower bound on the efficiency of prefix-codes. Consider the standard correspondence between binary strings and integers as in Equation 1.1. Define  $f(n) = l(n) + l(l(n)) + \dots + 2$ . Show that  $\sum_n 2^{-f(n)} = \infty$ .

*Comments.* Hint: let the number of terms in  $f(n)$ , apart from the 2, be  $h(n) = 0$  for  $n = 2$  and  $1 + h(h(n))$  for  $n > 2$  (this defines function  $h$ ). Define  $s_i = \sum_{h(n)=i} 2^{-f(n)}$  and  $\sum_{n=3}^{\infty} 2^{-f(n)} = \sum_{i=0}^{\infty} s_i$ . We show all  $s_i$  equal to 1, using induction on  $i$ . Clearly  $s_0 = 2^{-0} = 1$ . Also

$$\begin{aligned} s_{i+1} &= \sum_{h(n)=i+1} 2^{-f(n)} = \sum_{h(l(n))=i} 2^{-f(n)} \\ &= \sum_{h(m)=i} \sum_{l(n)=m} 2^{-(m+f(m))} = \sum_{h(m)=i} 2^m 2^{-(m+f(m))} = s_i. \end{aligned}$$

**1.11.15.** [24] The *Prime Number Theorem* states that the number  $\pi(n)$  of primes less than  $n$  is given by  $\pi(n) \sim n / \ln n$ , where “ $\ln$ ” is the natural logarithm. Let  $E : \mathcal{N} \rightarrow \{0, 1\}^*$  be a partial recursive prefix-code that assigns code word  $E(n)$  to source word  $n$ .

(a) Show that  $E(k) = \Omega(p_k)$  (interpret  $E(k)$  as the corresponding integer).

(b) Show that Item (a) implies  $\pi(n) = \Omega(n / \log^2 n)$ .

*Comments.* The estimate in Item (b) can be improved using efficient prefix-free codes. This way we find  $\pi(n) = \Omega(n / (\log n \log \log n \dots))$ , all positive factors. Hint: see page 5. Source: P. Berman, personal communication, November 1987. Current proof by J. Tromp, personal communication, July 1990.

## 1.12

### State × Symbol Complexity

---

A fourth historical root of Kolmogorov complexity seems to be another issue (other than information theory) raised by Shannon. In his paper “A universal Turing machine with two internal states” [pp. 129–153 in: *Automata Studies*, C.E. Shannon and J. McCarthy (Eds.), Princeton Univ. Press, 1956] he showed that there is a simple way of changing each Turing machine using  $m > 2$  states to a Turing machine using only two states and computing essentially the same function. This requires expanding the number of tape symbols, since the state of the original machine needs to be stored and updated on the tape of the corresponding two-state machine. The main result of this exercise is the construction of a two-state universal Turing machine. It turns out that a one-state universal Turing machine does not exist. Trying to minimize the other main resource, the number of tape symbols, is resolved even more simply: two tape symbols suffice, namely, by the expedient of encoding the  $m$  different tape symbols in  $O(\log m)$  length binary strings. The string of all 0’s is reserved to encode the distinguished blank symbol  $B$ . It is not very difficult, but tedious, now to adapt the finite control to make the whole arrangement work. Again this is the minimum, since it is impossible to construct a universal Turing machine with only one tape symbol.

It turns out that in both cases (reducing the number of states to two, or the number of tape symbols to two) the *product* of the number of states and the number of tape symbols increases at most eightfold. This suggests that this product is a relatively stable measure of the complexity of description of algorithms in the syntax of our Turing machine formalism. Following this idea, Shannon proposed the state-symbol product as a measure of complexity of description of algorithms. In particular, we can classify computable functions by the smallest state-symbol product of Turing machines that compute it. Here we assume of course a fixed formalism to express the Turing machines. It is a straightforward insight that this product is closely related to the number of bits to syntactically specify the Turing machine in the usual notation.

Gregory Chaitin in papers appearing in 1966 and 1969 analyzed this question in great detail. In his 1969 paper, observing that each Turing

machine can be coded as a program for a fixed-reference universal machine, he formulated as a variant of Shannon's approach the issue of the lengths of the shortest programs of the reference Turing machine to calculate particular finite binary strings.

## Exercises

---

**1.12.1.** [12] It is usual to allow Turing machines with arbitrarily large tape alphabets  $A$  (with the distinguished blank symbol  $B$  serving the analogous role as before). Use the quadruple formalism for Turing machines as defined earlier. How many Turing machines with  $m$  states and  $n$  tape symbols are there? (Count the blank tape symbol  $B$  as one of the  $n$  tape symbols.)

**1.12.2.** [20] Define Turing machines in quadruple format with arbitrarily large tape alphabets  $A$ , and state sets  $Q$ ,  $d(A), d(Q) < \infty$ .

Show that each such Turing machine with state set  $Q$  and tape alphabet  $A$  can be simulated by a Turing machine with tape alphabet  $A'$ ,  $d(A') = 2$ , and state set  $Q'$  such that  $d(A')d(Q') \leq cd(A)d(Q)$ , for some small constant  $c$ . Determine  $c$ . Show that the analogous simulation with  $d(A') = 1$  is impossible ( $d(Q') = \infty$ ).

*Comments.* See C.E. Shannon, pp. 129–153 in: *Automata Studies*, C.E. Shannon and J. McCarthy (Eds.), Princeton University Press, 1956. This is also the source for the next exercise.

**1.12.3.** [25] Show that each such Turing machine with state set  $Q$  and tape alphabet  $A$  can be simulated by a Turing machine with state set  $Q'$ ,  $d(Q') = 2$ , and tape alphabet  $A'$  such that  $d(A')d(Q') \leq cd(A)d(Q)$ , for some small positive constant  $c$ . Determine  $c$ . Show that the analogous simulation with  $d(Q') = 1$  is impossible (equivalently, implies  $d(A') = \infty$ ).

**1.12.4.** [37] Construct a universal Turing machine with  $d(A)d(Q) \leq 35$ .

*Comments.* Such a construction was first found by M. Minsky [*Ann. Math.*, 74(1961), 437–455].

## 1.13 History and References

---

In notations concerning binary strings and so forth we follow A.K. Zvonkin and L.A. Levin [*Russ. Math. Surveys*, 25:6(1970), 83–124]. (According to the Cyrillic alphabet of the original version of this important survey of Kolmogorov complexity “Zvonkin” precedes “Levin.” This author order was maintained in the English translation.) The big- $O$  notation is discussed in [D.E. Knuth, *SIGACT News*, 4:2(1976), 18–24; P.M.B. Vitányi and L. Meertens, *SIGACT News*, 16:4(1985), 56–59; D.E. Knuth, *Fundamental Algorithms*, Addison-Wesley, 1973].

The basics of combinatorics can be found in many textbooks, for instance [D.E. Knuth, *Fundamental Algorithms*, Addison-Wesley, 1973; W. Feller, *An Introduction to Probability Theory and Its Applications*, Wiley, 1968]. Turing machines were introduced by A.M. Turing in an important paper [*Proc. London Math. Soc.*, 42(1936), 230–265; Correction, *Ibid.*], where also the quoted material can be found. Related notions were introduced also in [E.L. Post, *J. Symb. Logic*, 1936]. The standard textbook references for basic computability theory are [H. Rogers, Jr., *Theory of Recursive Functions and Effective Computability*, McGraw-Hill, 1967; P. Odifreddi, *Classical Recursion Theory*, North-Holland, 1989]. The notion of enumerable functions was, perhaps, first used in [A.K. Zvonkin and L.A. Levin, *Ibid.*]. There they are called “semicomputable functions.” The material on computable (recursive) and enumerable real numbers partly arises from this context. For related work see [M.B. Pour-El and J.I. Richards, *Computability in Analysis and Physics*, Springer-Verlag, 1989].

In Section 1.7.4, we have introduced the basic terminologies of computational complexity theory that we will need in Chapter 7. For computational complexity theory see [J. Hartmanis, *Feasible computations and provable complexity properties*, SIAM, 1978; J.L. Balcázar, J. Diaz, and J. Gabarró, *Structural complexity*, Springer-Verlag, 1988; and M.R. Garey and D.S. Johnson, *Computers and Intractability*, Freeman, 1979].

A.N. Kolmogorov’s classic treatment of the set-theoretic axioms of the calculus of probabilities is his slim book *Grundbegriffe der Wahrscheinlichkeitsrechnung*, Springer-Verlag, 1933; English translation published by Chelsea, New York, 1956. A standard textbook in probability theory is [W. Feller, *An Introduction to Probability Theory and Its Applications*, Wiley, 1968].

General statistical tests for (pseudo-)randomness of sequences are extensively treated in [D.E. Knuth, *Seminumerical Algorithms*, Addison-Wesley, 1981]. Statistical testing for randomness of the decimal representations of  $\pi$ ,  $e$ , and  $\sqrt{2}$  specifically is found in [D.E. Knuth, *Ibid.*, 144–145; S.S. Skiena, *Complex Systems*, 1(1987), 361–366]. The quoted remark of J. von Neumann (1903–1957) appears in [Various techniques used in connection with random digits, *Collected Works*, Vol.V, Macmil-

lan, 1963]. R. von Mises's foundation of probability theory based on frequencies is set forth in [*Mathemat. Zeitsch.*, 5(1919), 52–99; Correction, *Ibid.*, 6(1920); *Probability, Statistics and Truth*, Macmillan, 1939]. In the 1920 correction to this paper von Mises writes, “In der Arbeit von Prof. Helm aus Dresden in Bd. 1 der Naturphilosophie 1902, 364– ‘Wahrscheinlichkeitstheorie des Kollektivbegriffes’ dort wird auch die Umkehrung der Wahrscheinlichkeitsdefinition, nämlich ihre Zurückführung auf Kollektive statt auf Bereiches gleich möglicher Fälle, in der Hauptsache bereits vorgebildet.” (“In the paper of Prof. Helm from Dresden in volume 1 of Naturphilosophie, 1902, 364– ‘Wahrscheinlichkeitstheorie des Kollektivbegriffes’ the key elements of the inverse definition of probability, namely its foundation on collectives instead of its foundation on domains of equally probable events, are also already exhibited.”) A critical comparison of different (foundational) theories of probabilities is [T.L. Fine, *Theories of Probability*, Academic Press, 1973]. The formulation of the apparent circularity in the frequency foundation of probability is from [J.E. Littlewood, *Littlewood's Miscellany*, Cambridge Univ. Press, 1986, 71–73]. The quotations of Kolmogorov, as well as the finitary version of von Mises's collectives, are from [*Sankhyā*, Series A, 25(1963), 369–376]. The work in improving von Mises's notion of admissible place selections is due to A. Wald [*Ergebnisse eines Mathematischen Kolloquiums*, Vol. 8, 1937, 38–72], who proved existence of collectives under the restriction of the number of admissible place-selection rules to countably infinite; to A. Church [*Bull. Amer. Math. Soc.*, 46(1940), 130–135], who further restricted admissible place selections to computable functions; and to J. Ville [*Etude Critique de la Notion de Collectif*, Gauthier-Villars, 1939], who showed that the von Mises-Wald-Church definitions still fail to satisfy randomness properties like the Law of the Iterated Logarithm. Champernowne's number was found by D.G. Champernowne [*J. London Math. Soc.*, 8(1933), 254–260]. The entire issue of randomness of individual finite and infinite sequences is thoroughly reviewed by D.E. Knuth [*Ibid.*, 142–169; summary, history, and references: 164–166]. A recent survey on the foundational issues of randomness was given by M. van Lambalgen [*J. Symb. Logic*, 52(1987), 725–755; *Random Sequences*, Ph.D. Thesis, Universiteit van Amsterdam, 1987].

Our treatment of Bayes's Rule follows R. von Mises [*Probability, Statistics and Truth*, Macmillan, 1939]. The idea of a universal a priori probability is due to R.J. Solomonoff. Lemma 1.10.1, the Chernoff bound, is due to H. Chernoff [*Ann. Math. Stat.*, 23(1952), 493–509]. The form we use is from L.G. Valiant and D. Angluin [*J. Comput. System Sci.*, 18:2(1979), 155–193] who base themselves in turn on [P. Erdős and J. Spencer, *Probabilistic methods in combinatorics*, Academic Press, 1974, p. 18].

Information theory was introduced as an essentially complete new mathematical discipline in C.E. Shannon's classic paper [*Bell System Tech. J.*, 27(1948), 379–423, 623–656]. Standard textbooks are [R.G. Gallager, *Information Theory and Reliable Communication*, Wiley & Sons, 1968; T.M. Cover and J.A. Thomas, *Elements of Information Theory*, Wiley & Sons, New York, 1991]. Our discussion of information theory used the original papers of Shannon. The Shannon-Fano code is due to C.E. Shannon [*Bell System Tech. J.*, 27(1948), 379–423, 623–656]. It is also attributed to R.M. Fano. The discussion is nonstandard insofar as we followed Kolmogorov's idea that the combinatorial approach to information theory should precede the probabilistic approach, in order to emphasize the logical independence of the development of information theory from probabilistic assumptions [*Problems Inform. Transmission*, 1(1965), 1–7; *Russian Math. Surveys*, 38:4(1983), 29–40]. The latter papers mention the following justification. In linguistic analysis it is natural to take such a purely combinatorial approach to the notion of the “entropy of a language.” This entropy is an estimate of the flexibility of a language, a number that measures the diversity of possibilities for developing grammatically correct sentences from a given dictionary and grammar. Using S.I. Ozhegov's Russian dictionary, M. Ratner and N.D. Svetlova obtained the estimate  $h = (\log N)/n = 1.9 \pm 0.1$  with  $N$  the number of Russian texts of length  $n$  (number of letters including spaces). This turns out to be much larger than the upper estimate for entropy of literary texts that can be obtained by various methods of “guessing continuations.” Naturally so, since literary texts must meet many requirements other than grammatical correctness. For this discussion and further remarks about change of entropy under translations, and the entropy cost of adhering to a given meter and rhyme scheme, see A.N. Kolmogorov [*Sankhyā*, Series A, 25(1963), 369–376]. For instance, Kolmogorov reports that classic rhyming iambic tetrameter requires a freedom in handling verbal material characterized by a residual entropy of ca. 0.4. “The broader problem of measuring the information connected with creative human endeavor is of the utmost significance.”

The general theory of coding and prefix-codes as in Section 1.11.1 is treated in [R.G. Gallager, *Information Theory and Reliable Communication*, Wiley, 1968]. The ubiquitous Kraft Inequality for prefix-codes, Theorem 1.11.1, is due to L.G. Kraft [*A Device for Quantizing, Grouping, and Coding Amplitude Modulated Pulses*, M. Sc. Thesis, Dept. Electr. Eng., MIT, Cambridge, Mass., 1949]. C.E. Shannon's Noiseless Coding Theorem 1.11.2 establishing the minimal average code-word length is from [*Bell System Tech. J.*, 27(1948), 379–423, 623–656]. Universal codes for infinite sets and unknown probability distributions were first proposed by Kolmogorov [*Problems Inform. Transmission*, 1:1(1965), 1–7]. Further developments were given in [P. Elias, *IEEE Trans. Inform. The-*

*ory*, IT-21(1975), 1094–203; S.K. Leung-Yan-Cheong and T.M. Cover *IEEE Trans. Inform. Theory*, IT-24(1978), 331–339; J. Rissanen, *Ann. Stat.*, 11(1982), 416–431; *Stochastic Complexity in Statistical Inquiry*, World Scientific, 1989].

The state-symbol complexity measure for Turing machines (and recursive functions) was apparently first discussed by C.E. Shannon [*Automata Studies*, C.E. Shannon and J. McCarthy (Eds.), Princeton Univ. Press, 1956, 129–153].

Kolmogorov complexity originated with the discovery of universal descriptions, and a recursively invariant approach to the concepts of complexity of description, randomness, and a priori probability. Historically, it is firmly rooted in R. von Mises's notion of random infinite sequences as discussed above. The first work in this direction is possibly K. Gödel's *On the length of proofs* of 1936, in which he proves that adding axioms to undecidable systems shortens the proofs of many theorems (thus using length as a measure of the complexity of proofs). With the advent of electronic computers in the 1950s, a new emphasis on computer algorithms and a maturing general recursive function theory, ideas tantamount to Kolmogorov complexity, came to many people's minds, because “when the time is ripe for certain things, these things appear in different places in the manner of violets coming to light in early spring” (Wolfgang Bolyai to his son Johann in urging him to claim the invention of non-Euclidean geometry without delay). Thus, with Kolmogorov complexity one can associate three inventors, in chronological order: R.J. Solomonoff, of Cambridge, Massachusetts, USA; A.N. Kolmogorov, of Moscow, Russia; and G.J. Chaitin, of New York City, USA.

Already in November 1960 R.J. Solomonoff published a *Zator Company* technical report [A preliminary report on a general theory of inductive inference, Tech. Rept. ZTB-138, Zator Company, Cambridge, Mass., November 1960] that presented the basic ideas of the theory of algorithmic complexity as a means to overcome serious problems associated with the application of Bayes's Rule in statistics. Ray Solomonoff was born on July 25, 1926, in Cleveland, Ohio, in the United States. He studied physics during about 1946–1950 at the University of Chicago (he recalls the lectures of E. Fermi and R. Carnap) and obtained an M.Sc. from that university. From about 1951–1956 he worked in the electronics industry doing math and physics and designing analog computers, working half-time. His scientific autobiography is published as [pp. 1–18 in: *Proc. 2nd European Conf. Comput. Learning Theory*, Lect. Notes Artific. Intell., Vol. 904, Springer-Verlag, 1995].

Solomonoff's objective was to formulate a completely general theory of inductive reasoning that would overcome shortcomings in Carnap's [*Logical Foundations of Probability*, Univ. Chicago Press, 1950]. Follow-

ing some more technical reports, in a long journal paper [*Inform. Contr.*, 7(1964), 1–22, 224–254], he introduced “Kolmogorov” complexity as an auxiliary concept to obtain a universal a priori probability and proved the Invariance Theorem 2.1.1. The mathematical setting of these ideas is described in some detail in Section 1.10. Solomonoff’s work has led to a novel approach in statistics [T.L. Fine, *Ibid.*], leading to applicable inference procedures such as the “Minimal Description Length” principle; see Chapter 5.

This makes Solomonoff the first inventor and raises the question whether we ought to talk about “Solomonoff complexity.” However, the name “Kolmogorov complexity” for shortest effective description length has become well entrenched and commonly understood. Solomonoff was primarily interested in universal a priori probability, while Kolmogorov later, independently, discovered and investigated the associated complexity for its own sake. We will associate Solomonoff’s name with the universal distribution and Kolmogorov’s name with the descriptive complexity. It has become customary to designate the entire area dealing with descriptive complexity, algorithmic information, and algorithmic probability loosely by the name of “Kolmogorov complexity” or “algorithmic information theory.” (Associating Kolmogorov’s name with the area may be viewed as an example in the sociology of science of the “Matthew Effect,” first noted in the Gospel according to Matthew, 25: 29–30, “For to every one who has more will be given, and he will have in abundance; but from him who has not, even what he has will be taken away.”) A comprehensive account of Solomonoff’s ideas and their genesis is presented in the History Section of Chapter 4 and in Chapter 5.

Solomonoff’s publications apparently received little attention until Kolmogorov started to refer to them from 1968 onward. These papers contain in veiled form suggestions about randomness of finite strings, non-computability of Kolmogorov complexity, computability of approximations to the Kolmogorov complexity, and resource-bounded Kolmogorov complexity. M. Minsky referred to Solomonoff’s work [*Proc. I.R.E.*, January 1961, 8–30; p. 43 in *Proc. Symp. Appl. Math. XIV*, Amer. Math. Soc., 1962]. To our knowledge, these are the earliest documents outlining an algorithmic theory of descriptions.

The great Russian mathematician Andrei N. Kolmogorov, born 25 April 1903 in Tambov, Russia, died 20 October 1987 in Moscow. Many biographical details can be found in the Soviet Union’s foremost mathematics journal, *Uspenki Mat. Nauk*, translated into English as *Russian Math. Surveys* [B.V. Gnedenko, 28:5(1973), 5–16, P.S. Aleksandrov, 38:4(1983), 5–7; N.N. Bogolyubov, B.V. Gnedenko, and S.L. Sobolev, 38:4(1983), 9–27; A.N. Kolmogorov, 41:6(1986), 225–246; and the entire memorial issue 43:6(1988), especially 1–39 by V.M. Tikhomirov]. Three volumes of Kolmogorov’s (mathematical) *Selected Works* have been published

by Nauka, Moscow (in Russian) in 1985 through 1987; they are to be translated in English and published by D. Reidel. The writings on algorithmic complexity are collected in volume 3 of the *Selected Works*. In the western literature, see the memorial issue [*The Annals of Probability*, 17:3(1989)], especially the scientific biography on pp. 866–944 by A.N. Shiryaev, and an evaluation of Kolmogorov's contributions to information theory and to algorithmic complexity, on pp. 840–865 by T.M. Cover, P. Gács, and R.M. Gray. See also the obituary in [*Bull. London Math. Soc.*, 22:1(1990), 31–100] and [V.A. Uspensky, *J. Symb. Logic*, 57:2(1992), 385–412].

In 1933 A.N. Kolmogorov supplied probability theory with a powerful mathematical foundation in his [*Grundbegriffe der Wahrscheinlichkeitsrechnung*, Springer-Verlag, 1933]. Following a four-decade-long controversy on von Mises's conception of randomness, in which Kolmogorov played little part, the content of which is set forth in some detail in Section 1.9, Kolmogorov finally introduced complexity of description of finite individual objects as a measure of individual information content and randomness, and proved the Invariance Theorem 2.1.1 in his paper of spring 1965 [*Problems Inform. Transmission*, 1(1965), 1–7]. His objective was primarily, apart from resolving the question of randomness of objects, to revise information theory by an algorithmic approach to the information content of individual objects, in contrast to the traditional way discussed in Section 1.11. Kolmogorov described the essence and background to the algorithmic approach in his report to the Probability Section of the Moscow Mathematical Society on April 24, 1963: “One often has to deal with very long sequences of symbols. Some of them, for example, the sequences of symbols in a 5-digit logarithm table, permit a simple logical definition and therefore might be obtained by the computations (though clumsy at times) of a simple pattern... Others seem not to admit any sufficiently simple ‘legitimate’ way to construct them. It is supposed that such is the case for a rather long segment in a table of random numbers... There arises the question of constructing a rigorous mathematical theory to account for these differences in behavior” [A.N. Shiryaev, *Ibid.*, 921]. Says Kolmogorov: “I came to similar conclusions [as Solomonoff], before becoming aware of Solomonoff's work, in 1963–1964” [*IEEE Trans. Inform. Theory*, IT 14:5(1968), 662–664]. And again “The basic discovery, which I have accomplished independently from and simultaneously with R. Solomonoff lies in the fact that the theory of algorithms enables us to eliminate this arbitrariness [of interpretive mechanisms for descriptions leading to different lengths of shortest descriptions for the same object, and hence to different ‘complexities’ with respect to different interpretive mechanisms] by the determination of a ‘complexity’ which is almost invariant (the replacement of one method by another leads only to the supplement of the bounded

term)" [A.N. Shiryaev, *Ibid.*, 921]. In the case of the other two inventors the subject we are concerned with appears, as it were, out of the blue. But Kolmogorov's involvement strikes one as the inevitable confluence of interests of this great scientist: his lifelong fascination with the foundations of probability theory and randomness, his immediate appreciation of information theory upon its formulation by Shannon, and his vested interest in the theory of algorithms witnessed by, for instance, A.N. Kolmogorov and V.A. Uspensky, *Uspekhi Mat. Nauk*, 13:4(1958), 3–28 [in Russian; translated *Amer. Math. Soc. Transl. (2)*, 29(1963), 217–245]. The new ideas were vigorously investigated by his associates. These included the Swedish mathematician P. Martin-Löf, visiting Kolmogorov in Moscow during 1964–1965, who investigated the complexity oscillations of infinite sequences and proposed a definition of infinite random sequences that is based on constructive measure theory [*Inform. Contr.*, 9(1966), 602–619; *Z. Wahrsch. Verw. Geb.*, 19(1971), 225–230]. There is a related development in set theory and recursion theory, namely, the notion of "generic object" in the context of "forcing." For example, being a member of the arithmetical generic set is analogous (but not precisely) to being a member of the intersection of all arithmetical sets of measure 1. There is a notion, called "1-genericity," that in restricted version calls for the intersection of all recursively enumerable sets of measure 1. This is obviously related to the approach of Martin-Löf. Forcing was introduced by P. Cohen in 1963 to show the independence of the continuum hypothesis, and using sets of positive measure as forcing conditions is due to R.M. Solovay soon afterwards.

G.J. Chaitin had finished the Bronx High School of Science, and was an eighteen-year-old undergraduate student at the City College of the City University of New York, when he submitted two papers [*J. ACM*, 13(1966), 547–569; *J. ACM*, 16(1969), 145–159] for publication, in October and November 1965, respectively. In the 1966 paper he addresses the "state × symbol" complexity of algorithms following Shannon's coding concepts, as described in Section 1.12, but does not introduce an invariant notion of complexity. Continuing this work in the 1969 paper, in the final part, Chaitin puts forward the notion of Kolmogorov complexity, proves the Invariance Theorem 2.1.1, and studies infinite random sequences (in the sense of having maximally random finite initial segments) and their complexity oscillations. As Chaitin [*Scientific American*, 232:5(1975), 47–52] formulates it: "this definition [of Kolmogorov complexity] was independently proposed about 1965 by A.N. Kolmogorov and me ... Both Kolmogorov and I were then unaware of related proposals made in 1960 by Ray Solomonoff." A short autobiography appears in [G.J. Chaitin, *Information-Theoretic Incompleteness*, World Scientific, Singapore, 1992].

## Algorithmic Complexity

The most natural approach to defining the quantity of information is clearly to define it in relation to the individual object (be it Homer's *Odyssey* or a particular type of dodo) rather than in relation to a set of objects from which the individual object may be selected. To do so, one could define the quantity of information in an object in terms of the number of bits required to describe it. A description of an object is evidently only useful if we can reconstruct the object from this description.

We aim at something different from C.E. Shannon's theory of communication, which deals with the specific technological problem of data transmission, that is, with the information that needs to be transmitted in order to select an object from a previously agreed-upon set of alternatives, Section 1.11. Our task is to widen the limited set of alternatives until it is universal. We aim at a notion of "absolute" information of individual objects, that is the information that by itself describes the object completely.

Intuition tells us that some objects are complicated and some objects are simple. For instance, a number like  $2^{1000}$  is certainly very simple (we have just expressed it in a few bits); yet evidently there are numbers of a thousand bits for which it is hard to see how we can find a description requiring much less than a thousand bits. Such hard to describe numbers would be their own shortest descriptions.

We require both an agreed-upon universal description method and an agreed-upon mechanism to produce the object from its alleged description. This would appear to make the information content of an object depend on whether it is particularly favored by the description method

we have selected. By “favor” we mean to produce short descriptions in terms of bits.

For instance, it is well-known that certain programming languages favor symbolic computations while other programming languages favor arithmetic computations, even though all of them are universal. The notion of information content of individual objects can only be useful if the quantity of information is an attribute of the object *alone* and is independent of the means of description. It is *a priori* by no means obvious that this is possible. Relatively recent advances resulting in the great ideas of computability theory from the 1930s onwards have made it possible to design a universal description method that appears to meet our goals.

Denote the set of objects by  $S$ , and assume some standard enumeration of objects  $x$  by natural numbers  $n(x)$ . We are interested in the fact that  $n(x)$  may not be the most economical way to specify  $x$ . To compare methods of specification, we view such a method as a partial function over the nonnegative integers defined by  $n = f(p)$ . We do not yet assume that  $f$  is recursive, but maintain full generality to show to what extent such a theory can also be developed with noneffective notions, and at which point effectiveness is required. With each natural number  $p$  associate the length of the finite binary string identified with  $p$  as in Equation 1.1. Denote this length by  $l(p)$ .

For each object  $x$  in  $S$ , the complexity of object  $x$  with respect to the specifying method  $f$  is defined as

$$C_f(x) = \min\{l(p) : f(p) = n(x)\},$$

and  $C_f(x) = \infty$  if there are no such  $p$ . In computer science terminology we would say that  $p$  is a program and  $f$  a computer, so that  $C_f(x)$  is the minimal length of a program for  $f$  (without additional input) to compute output  $x$ .

Considering distinct methods  $f_1, f_2, \dots, f_r$  of specifying the objects of  $S$ , it is easy to construct a new method  $f$  that assigns to each object  $x$  in  $S$  a complexity  $C_f(x)$  that exceeds only by  $c$  (less than about  $\log r$ ) the minimum of  $C_{f_1}(x), C_{f_2}(x), \dots, C_{f_r}(x)$ . The only thing we have to do is to reserve the first  $\log r$  bits of  $p$  to identify the method  $f_i$  that should be followed, using as a program the remaining bits of  $p$ .

We say that a method  $f$  *minorizes* a method  $g$  (additively) if there is a constant  $c$  such that for all  $x$

$$C_f(x) \leq C_g(x) + c.$$

Above we have shown how to construct a method  $f$  that minorizes each of the methods  $f_1, \dots, f_r$  with constant  $c \approx \log r$ . Two methods  $f$  and  $g$  are called *equivalent* if each of them minorizes the other.

Consider the hierarchy of equivalence classes of methods with respect to minorization. Kolmogorov has remarked that the idea of “description length” would be useless if the constructed hierarchy did not have certain niceness properties. In particular, we would like such a hierarchy to have a unique minimal element: the equivalence class of description methods that minorize all other description methods. Some sets of description methods do have a unique minimal element while other sets of description methods don’t.

**Definition 2.0.1** Let  $\mathcal{C}$  be a subclass of the partial functions over the nonnegative integers. A function  $f$  is *universal* (or *additively optimal*) for  $\mathcal{C}$  if it belongs to  $\mathcal{C}$  and if for every function  $g \in \mathcal{C}$  there is a constant  $c_{f,g}$  such that  $C_f(x) \leq C_g(x) + c_{f,g}$ , for all  $x$ . (Here  $c_{f,g}$  depends on  $f$  and  $g$ , but not on  $x$ .) Replacing  $x$  by  $\langle x, y \rangle$  with  $\langle \cdot \rangle$  the standard recursive bijective pairing function yields the definition for a class of two-variable functions.

Clearly, all additively optimal methods  $f, g$  of specifying objects in  $S$  are equivalent in the following way:

$$|C_f(x) - C_g(x)| \leq c_{f,g},$$

for all  $x$ , where  $c_{f,g}$  is a constant depending only on  $f$  and  $g$ . Thus, from an asymptotic point of view, the complexity  $C(x)$  of an object  $x$ , when we restrict ourselves to optimal methods of specification, does not depend on accidental peculiarities of the chosen optimal method.

**Example 2.0.1** Consider the class of description methods consisting of *all* partial functions over the nonnegative integers. Assume, by way of contradiction, that this class includes a universal element  $f$ . Take an infinite sequence  $p_1 < p_2 < \dots$  such that  $f(p_i) = x_i$  with  $p_i$  minimal, for all  $i \geq 1$ . Select a subsequence  $q_1, q_2, \dots$  from  $p_1, p_2, \dots$  such that  $\log p_i < (\log q_i)/2$ . Define another function  $g$  by  $g(p_i) = f(q_i)$  for all  $i \geq 1$  and that coincides with  $f$  otherwise. Then, for infinitely many  $x$  we have  $C_g(x) \leq C_f(x)/2$ , which is the required contradiction. Hence, there is no universal partial function and the hierarchy of complexities with respect to the partial functions does not have *any* minimal element.

The development of the theory of Kolmogorov complexity is made possible by the remarkable fact that the class of partial *recursive* functions (defined in Section 1.7) possesses a universal element. Under this relatively natural restriction on the class of description methods (that is, partial functions) we obtain a well-behaved hierarchy of complexities.  $\diamond$

We begin by worrying about notation. There are several variants of Kolmogorov complexity, with notations that are not used consistently among different authors or even by the same author at different times.

In the main text of this book we shall concentrate on two major variants of Kolmogorov complexity. It seems educationally the right approach to first study Kolmogorov complexity as originally defined by Solomonoff, Kolmogorov, and Chaitin because it is intuitively clearer.

Some mathematical technicalities will naturally lead up to, and justify, the less intuitive version of Kolmogorov complexity. The first type we call *plain* Kolmogorov complexity, and the second type we call *prefix* Kolmogorov complexity. We use “ $C$ ” to denote the plain Kolmogorov complexity. We reserve “ $K$ ” for the prefix type. Fortunately, the majority of theorems we derive for plain Kolmogorov complexity carry over unchanged and with the same proofs to the prefix version. The difference is that the prefix version is tweaked to have just the right quantitative properties for some desired usages and applications.

## 2.1 The Invariance Theorem

Identify an object  $x$  from a countably infinite sample space  $S$  with its index  $n(x)$ . Consider the class of description methods

$$\{\phi : \phi \text{ is a partial recursive function}\}.$$

Consider the particular problem of describing objects consisting of natural numbers in terms of programs consisting of finite strings of 0's and 1's. Just as in information theory, Section 1.11, where the entropy and information of a message over any size alphabet are expressed in the normalized format of bits, the restriction of the programs to a binary alphabet does not imply any loss of generality. In both cases changing alphabet size leaves all statements invariant up to an appropriate logarithmic multiplicative factor related to the alphabet sizes involved, see Exercise 2.1.9.

The following *Invariance Theorem* is the cornerstone for the subsequent development of the theory. In fact, for many later applications it embodies the *entire* theoretical foundation. Recall Definition 2.0.1 of a function that is universal (or additively optimal) for a class of functions. We give the unconditional version as a preliminary lemma.

**Lemma 2.1.1** *There is a universal partial recursive function.*

**Proof.** Let  $\phi_0$  be the function computed by a universal Turing machine  $U$ . Machine  $U$  expects inputs of the format

$$\langle n, p \rangle = \underbrace{11\dots1}_{l(n) \text{ times}} 0 np .$$

The interpretation is that the total program  $\langle n, p \rangle$  is a two-part code of which the first part consists of a self-delimiting encoding of  $T_n$  and the second part is the literally rendered program  $p$ . This way  $U$  can first

parse the binary input into the  $T_n$ -part and the  $p$ -part, and subsequently simulate the computation of  $T_n$  started with program  $p$  as its input (Section 1.7). That is,  $\phi_0(\langle n, p \rangle) = \phi_n(p)$ . What happens if  $U$  gets the program “ $0pU = T_0$  and therefore  $U(0p) = U(p)$ . Altogether, if  $T_n$  computes partial recursive function  $\phi_n$ , then

$$C_{\phi_0}(x) \leq C_{\phi_n}(x) + c_{\phi_n},$$

where  $c_{\phi_n}$  can be set to  $2l(n) + 1$ .  $\square$

This is the simple formulation of the Invariance Theorem. However, for many applications we require a generalization to a “conditional” version, as follows. The difficulty of specifying an object can be facilitated when another object is already specified. We define the complexity of an object  $x$ , given an object  $y$ . Fix an effective enumeration of Turing machines  $T_1, T_2, \dots$  as in Section 1.7. The Turing machines use a tape alphabet  $\{0, 1, B\}$ , and the input to a Turing machine is a program consisting of a contiguous string of 0’s and 1’s, delimited by blanks  $B$  on both sides. This way a Turing machine can detect the end of its program. The effective enumeration of Turing machines induces an effective enumeration of partial recursive functions  $\phi_1, \phi_2, \dots$  such that  $T_i$  computes  $\phi_i$  for all  $i$ . As above,  $\langle \cdot \rangle : \mathcal{N} \times \mathcal{N} \rightarrow \mathcal{N}$  is a standard recursive bijective pairing function mapping the pair  $(x, y)$  to the singleton  $\langle x, y \rangle$ . We can iterate this as  $(x, y, z) = \langle x, \langle y, z \rangle \rangle$ .

**Definition 2.1.1** Let  $x, y, p$  be natural numbers. Any partial recursive function  $\phi$ , together with  $p$  and  $y$ , such that  $\phi(\langle y, p \rangle) = x$ , is a *description* of  $x$ . The complexity  $C_\phi$  of  $x$  conditional to  $y$  is defined by

$$C_\phi(x|y) = \min\{l(p) : \phi(\langle y, p \rangle) = x\},$$

and  $C_\phi(x|y) = \infty$  if there are no such  $p$ . We call  $p$  a *program* to compute  $x$  by  $\phi$ , given  $y$ .

**Theorem 2.1.1** *There is a universal partial recursive function  $\phi_0$  for the class of partial recursive functions to compute  $x$  given  $y$ . Formally this says that  $C_{\phi_0}(x|y) \leq C_\phi(x|y) + c_\phi$  for all partial recursive functions  $\phi$  and all  $x$  and  $y$ , where  $c_\phi$  is a constant depending on  $\phi$  but not on  $x$  or  $y$ .*

**Proof.** Let  $\phi_0$  be the function computed by a universal Turing machine  $U$  such that  $U$  started on input  $\langle y, \langle n, p \rangle \rangle$  simulates  $T_n$  on input  $\langle y, p \rangle$  (Section 1.7). That is, if  $T_n$  computes the partial recursive function  $\phi_n$ , then  $\phi_0(\langle y, \langle n, p \rangle \rangle) = \phi_n(\langle y, p \rangle)$ . Hence, for all  $n$ ,

$$C_{\phi_0}(x|y) \leq C_{\phi_n}(x|y) + c_{\phi_n},$$

where  $c_{\phi_n} = 2l(n) + 1$ .  $\square$

The key point is not that the universal description method does necessarily give the shortest description in each case, but that no other description method can improve on it infinitely often. Note also that the optimal complexity  $C_{\phi_0}(x|y)$  is defined for all  $x$  and  $y$ . Namely, for each  $x$  and  $y$  we can find a Turing machine that computes output  $x$ , given  $y$ , for some input  $p$  (such as the Turing machine that outputs  $x$  for all inputs).

For *any* pair  $\psi, \psi'$  of additively optimal functions, there is a fixed constant  $c_{\psi, \psi'}$ , depending only on  $\psi$  and  $\psi'$ , such that for all  $x, y$  we have

$$|C_\psi(x|y) - C_{\psi'}(x|y)| \leq c_{\psi, \psi'}.$$

To see this, first substitute  $\phi_0 = \psi$  and  $\phi = \psi'$  in Theorem 2.1.1, then substitute  $\phi = \psi$  and  $\phi_0 = \psi'$  in Theorem 2.1.1, and combine the two resulting inequalities. While the complexities according to  $\psi$  and  $\psi'$  are not exactly equal, they are *equal up to a fixed constant* for all  $x$  and  $y$ .

**Definition 2.1.2** Fix a universal  $\phi_0$  and dispense with the subscript by defining the *conditional Kolmogorov complexity*  $C(\cdot| \cdot)$  by

$$C(x|y) = C_{\phi_0}(x|y).$$

This particular  $\phi_0$  is called the *reference function* for  $C$ . We also fix a particular Turing machine  $U$  that computes  $\phi_0$  and call  $U$  the *reference machine*. The *unconditional Kolmogorov complexity*  $C(\cdot)$  is defined by

$$C(x) = C(x|\epsilon).$$

**Example 2.1.1** Programmers are generally aware that programs for symbolic manipulation tend to be shorter when they are expressed in the LISP programming language than if they are expressed in FORTRAN, while for numerical calculations the opposite is the case. Or is it? The Invariance Theorem in fact shows that to express an algorithm succinctly in a program, it does not matter which programming language we use (up to a fixed additive constant that depends only on the two programming languages compared).

To see this, as an example consider the lexicographical enumeration of all syntactically correct LISP programs  $\lambda_1, \lambda_2, \dots$  and the lexicographical enumeration of all syntactically correct FORTRAN programs  $\pi_1, \pi_2, \dots$ . With proper definitions we can view the programs in both enumerations as computing partial recursive functions from their inputs to their outputs. Choosing reference machines in both enumerations we can define complexities  $C_{\text{LISP}}(x)$  and  $C_{\text{FORTRAN}}(x)$  completely analogous to  $C(x)$ . All of these measures of the descriptional complexity of  $x$  coincide up

to a fixed additive constant. Let us show this directly for  $C_{\text{LISP}}(x)$  and  $C_{\text{FORTRAN}}(x)$ .

It is well-known and also easy to see that each enumeration contains a universal program; the LISP enumeration contains a LISP interpreter program that interprets any LISP program. But there is also a LISP program  $\lambda_P$  that is a FORTRAN interpreter in the sense that it interprets any FORTRAN program. Consequently,  $C_{\text{LISP}}(x) \leq C_{\text{FORTRAN}}(x) + 2l(\lambda_P)$ . Similarly, there is a FORTRAN program  $\pi_L$  that is a LISP interpreter, which yields  $C_{\text{FORTRAN}}(x) \leq C_{\text{LISP}}(x) + 2l(\pi_L)$ . Consequently,  $|C_{\text{LISP}}(x) - C_{\text{FORTRAN}}(x)| \leq 2l(\lambda_P) + 2l(\pi_L)$  for all  $x$ .  $\diamond$

## 2.1.1 Two-Part Codes

It is a deep and useful fact that the shortest effective description of an object  $x$  can be expressed in terms of a *two-part code*: the first part describing an appropriate Turing machine and the second part describing the program that interpreted by the Turing machine reconstructs  $x$ . The essence of the Invariance Theorem is as follows: For the fixed reference universal Turing machine  $U$  the length of the shortest program to compute  $x$  is  $\min\{l(p) : U(p) = x\}$ . Looking back at the proof of Lemma 2.1.1 we notice that  $U(0p) = U(p)$ . From the definitions it therefore follows that

$$C(x) = \min\{l(T) + l(p) : T(p) = x\} \pm 1,$$

where  $l(T)$  is the length of a self-delimiting encoding for a Turing machine  $T$ , which provides an alternative definition of Kolmogorov complexity (similarly, for conditional Kolmogorov complexity). The above expression for Kolmogorov complexity can be rewritten as

$$C(x) = \min\{l(T) + C(x|T) : T \in \{T_0, T_1, \dots\}\} + O(1), \quad (2.1)$$

that emphasizes the two-part code nature of Kolmogorov complexity: using the regular aspects of  $x$  to maximally compress. In the example

we can encode  $x$  by a small Turing machine which computes  $x$  from the program “13.” Intuitively, the Turing machine part of the code squeezes out the *regularities* in  $x$ . What is left are irregularities, or *random aspects*, of  $x$  relative to that Turing machine. The minimal-length two-part code squeezes out regularity only insofar as the reduction in the length of the description of random aspects is greater than the increase in the regularity description.

This interpretation of  $C(x)$  as the shortest length of a two-part code for  $x$ , one part describing a Turing machine, or *model*, for the *regular* aspects of  $x$  and the second part describing the *irregular* aspects of  $x$  in

the form of a program to be interpreted by  $T$ , has profound applications. We can interpret this as that the regular, or “valuable,” information in  $x$  is constituted by the bits in the “model” while the random or “useless” information of  $x$  constitutes the remainder.

The “right model” is a Turing machine  $T$  among the ones that reach the minimum description length

$$\min_T \{l(T) + C(x|T) : T \in \{T_0, T_1, \dots\}\}.$$

This  $T$  embodies the amount of useful information contained in  $x$ . The main remaining question is which such  $T$  to select among the ones that satisfy the requirement. The problem is how to separate a shortest program  $x^*$  for  $x$  in parts  $x^*q = pr$  such that  $p$  represents an appropriate  $T$ . This idea has spawned the “minimum description length” principle in statistics and inductive reasoning, Section 5.5, and the notion of algorithmic entropy in Section 8.5.

M. Koppel [*Complex Systems*, 1(1987), 1087–1091; *The Universal Turing Machine: A Half-Century Survey*, R. Herken (Ed.), Oxford Univ. Press, 1988, pp. 435–452; M. Koppel and H. Atlan, Manuscript] defined a concept called *sophistication*. We roughly define sophistication for finite strings. The universal Turing machine  $U$  has four tapes: program tape, data tape, output tape, and work tape. The first three tapes are one-way. A pair  $(p, D)$  is a *description* of a finite string  $x$  if  $U(p, D)$  prints  $x$  and  $p$  is a *total* self-delimiting program. A total program  $p$  is a program that computes an output and halts for each input  $D$ . Define  $CS(x) = \min\{l(p) + l(D) : (p, D) \text{ is a description of } x\}$ . The  $c$ -*sophistication* of  $x$  is

$$\min\{l(p) : \exists D \text{ } (p, D) \text{ is a description of } x, l(p) + l(D) \leq CS(x) + c\}.$$

Note that since we required  $p$  to be total, one cannot shift  $p$  to  $D$  and write  $(q, pD)$ , where  $q$  is a program that considers the first self-delimiting prefix of the data as a program to be simulated. Such a program  $q$  cannot be total. The key to this definition is to separate program from data. Compare also  $(k, \delta)$ -stochasticity in Section 2.2.2.

### 2.1.2 Upper Bounds

Theorem 2.1.1 has a wider importance than just showing that the hierarchy of  $C_\phi$  complexity measures contains an additively optimal one. It is also our principal tool in finding upper bounds on  $C(x)$ . Such upper bounds depend on the choice of reference function, and hence are proved only to within an additive constant.

Intuitively, the Kolmogorov complexity of a binary string cannot exceed its own length, because the string is obviously a (literal) description of itself.

**Theorem 2.1.2** *There is a constant  $c$  such that for all  $x$  and  $y$ ,*

$$C(x) \leq l(x) + c \text{ and } C(x|y) \leq C(x) + c.$$

**Proof.** The first inequality is supremely obvious: define a Turing machine  $T$  that copies the input to the output. Then for all  $x$ , we have  $C_T(x) = l(x)$ . By Theorem 2.1.1 on page 97 the result follows.

To prove the second inequality, construct a Turing machine  $T$  that for all  $y, z$  computes output  $x$  on input  $\langle z, y \rangle$  iff the universal reference machine  $U$  computes output  $x$  for input  $\langle z, \epsilon \rangle$ . Then  $C_T(x|y) = C(x)$ . By Theorem 2.1.1, there is a constant  $c$  such that  $C(x|y) \leq C_T(x|y) + c = C(x) + c$ .  $\square$

Note that the additive constants in these inequalities are fudge terms related to the reference machine  $U$ . For example, we need to indicate to the reference machine that a given description is the object itself, and this information adds a number of bits to the literal description. In Section 3.2 we will calculate the constants explicitly as 8 and 2, respectively. Let us look at some more examples in order to develop our intuition about the notion of complexity of description.

**Example 2.1.2** For each finite binary string  $x$  we have  $C(xx) \leq C(x) + O(1)$ . Construct a Turing machine  $V$  such that  $V(p) = U(p)U(p)$ , for all programs  $p$ , where  $U$  is the reference machine in the proof of Theorem 2.1.1. In particular, if  $U(p) = x$ , then  $V(p) = xx$ . Let  $V = T_m$  in the standard enumeration of Turing machines  $T_1, T_2, \dots$ . With  $\bar{m}$  denoting the self-delimiting description  $1^{l(m)}0m$  of  $m$ , we have  $U(\bar{m}p) = T_m(p) = xx$  and  $l(\bar{m}p) = l(p) + 2l(m) + 1$ . Hence,  $C(xx) \leq C(x) + 2l(m) + 1$ . From now on we leave the more obvious details of this type of argument for the reader to fill in.  $\diamond$

**Example 2.1.3** Recall that  $x^R$  denotes the reverse of  $x$ . Clearly, the complexities of  $x$  and  $x^R$  can differ by at most a fixed constant  $c$  independent of  $x$ . That is,  $|C(x) - C(x^R)| < c$  holds for all  $x$ . We can generalize this example as follows: For every total recursive function  $\phi$  that is one-to-one there is (another) constant  $c$  such that  $|C(\phi(x)) - C(x)| < c$  for all  $x$ .

In fact, if  $\phi$  is computed by Turing machine  $T_n$  and  $U(p) = x$ , then there is a Turing machine  $V$  such that  $V(\bar{n}p) = \phi(x)$ . If  $V = T_m$ , then  $U(\bar{m}\bar{n}p) = \phi(x)$ , and therefore  $|C(\phi(x)) - C(x)| < 2l(m) + 2l(n) + 2$ . Similar relations hold for the conditional complexity  $C(x|y)$ .  $\diamond$

**Example 2.1.4** Can the complexity of a pair of strings exceed the sum of the complexities of the individual strings? In other words, is  $C$  *subadditive*? Let  $\langle \cdot \rangle : \mathcal{N} \times \mathcal{N} \rightarrow \mathcal{N}$  be the standard recursive bijection over the natural numbers that encodes  $x$  and  $y$  as  $\langle x, y \rangle$ . Define  $C(x, y) = C(\langle x, y \rangle)$ . That is, up to a fixed constant,  $C(x, y)$  is the length of the shortest program such that  $U$  computes both  $x$  and  $y$  and a way to tell them apart. It is seductive to conjecture  $C(x, y) \leq C(x) + C(y) + O(1)$ , the obvious (but false)

argument running as follows: Suppose we have a shortest program  $p$  to produce  $x$ , and a shortest program  $q$  to produce  $y$ . Then with  $O(1)$  extra bits to account for some Turing machine  $T$  that schedules the two programs, we have a program to produce  $x$  followed by  $y$ . However, any such  $T$  will have to know where to divide its input to identify  $p$  and  $q$ . One way to do this is by using input  $\bar{l}(p)pq$  or input  $\bar{l}(q)qp$ . This way we can show that for all  $x, y$ , we have

$$C(x, y) \leq C(x) + C(y) + 2 \log(\min(C(x), C(y))). \quad (2.2)$$

We cannot eliminate the logarithmic error term for the general case. Namely, in Example 2.2.3 on page 111 we show that there is a constant  $c$  such that for all  $n$  there are  $x$  and  $y$  of length at most  $n$  such that

$$C(x, y) \geq C(x) + C(y) + \log n - c.$$

We can eliminate the logarithmic error term at the cost of entering the length of one of the programs in the conditional,

$$C(x, y|C(x)) \leq C(x) + C(y) + O(1).$$

Analogous considerations hold for the complexity  $C(xy)$  of the unmarked concatenation  $xy$ .  $\diamond$

**Example 2.1.5** If we know  $C(x)$  and  $x$ , then we can run all programs of length  $C(x)$  in parallel on the reference machine  $U$  in dovetail fashion (in stage  $k$  of the overall computation execute the  $i$ th computation step of program  $k-i$ ). By definition of  $C(\cdot)$ , there must be a program of length  $C(x)$  that halts with output  $x$ . The first such program is the *first shortest program* for  $x$  in enumeration order, and is denoted by  $x^*$ .

Therefore, a program to compute  $C(x)$ , given  $x$ , can be converted to a program to compute  $x^*$ , given  $x$ , at the cost of a constant number of extra bits. If we have computed  $x^*$ , then  $C(x)$  is simply its length, so the converse is trivial. Furthermore, to describe  $C(x)$  from scratch takes at least as many bits as to describe  $C(x)$  using  $x$ . Altogether we have, up to additional constant terms,

$$C(x^*|x) = C(C(x)|x) \leq C(C(x)) \leq \log l(x).$$

$\diamond$

The upper bound on  $C(x^*|x)$  cannot be improved to  $O(1)$ . If it could, then one could show that  $C(x)$  is a recursive function. However, in Theorem 2.3.2 on page 121 we shall show that  $C(x)$  is not partial recursive. It is a curious fact that for some  $x$ , knowledge of  $x$  does not help much in computing  $x^*$ . In fact, the upper bound is nearly optimal. In Theorem 3.8.1 on page 227 we shall show that for some  $x$  of each length  $n$  the quantity  $C(C(x)|x)$ , and hence also  $C(x^*|x)$ , is almost  $\log n$ .

Clearly, the information that an element belongs to a particular set can severely curtail the complexity of that element. The following simple observation due to Kolmogorov turns out to be very useful. We show that for every easily describable set the conditional complexity of every one of its elements is at most equal to the logarithm of the cardinality of that set. (We will observe later, in Theorem 2.2.1 on page 109, that the conditional complexities of the majority of elements in a finite set cannot be significantly less than the logarithm of the cardinality of that set: we will say that they are “incompressible” and have a small “randomness deficiency.”)

**Theorem 2.1.3** *Let  $A \subseteq \mathcal{N} \times \mathcal{N}$  be recursively enumerable, and  $y \in \mathcal{N}$ . Suppose  $Y = \{x : (x, y) \in A\}$  is finite. Then, for some constant  $c$  depending only on  $A$ , for all  $x$  in  $Y$ , we have  $C(x|y) \leq l(d(Y)) + c$ .*

**Proof.** Let  $A$  be enumerated without repetition as  $(x_1, y_1), (x_2, y_2), \dots$  by a Turing machine  $T$ . Let  $(x_{i_1}, y_{i_1}), \dots, (x_{i_k}, y_{i_k})$  be the subsequence in which the elements of  $Y$  are enumerated,  $k = d(Y)$ . Using the fixed  $y$ , modify  $T$  to  $T_y$  such that  $T_y$ , on input  $1 \leq p \leq d(Y)$ , outputs  $x_{i_p}$ ,  $T_y(p) = x_{i_p}$ . Therefore, we have by the Invariance Theorem 2.1.1 that  $C(x|y) \leq C_{T_y}(x) + c \leq l(d(Y)) + c$ , with  $c$  depending only on  $A$ .  $\square$

Let us illustrate the use of this theorem. Let  $A$  be recursively enumerable and  $d(A^{\leq n}) \leq p(n)$ , with  $p$  a polynomial. Then, for all  $x \in A$  of length at most  $n$  we have  $C(x|n) \leq l(p(n)) + O(1)$ , by Theorem 2.1.3. For all  $x$  of length at most  $n$  we have  $C(x) \leq C(x|n) + 2l(n) + O(1)$ . Therefore, for  $x \in A^{\leq n}$  we find  $C(x) = O(\log n)$ .

### 2.1.3 Invariance of Kolmogorov Complexity

The complexity  $C(x)$  is only invariant up to a constant depending on the reference function  $\phi_0$ . Thus, one may object, for *every* string  $x$  there is an additively optimal recursive function  $\psi_0$  such that  $C_{\psi_0}(x) = 0$ . So how can one claim that  $C(x)$  is an objective notion?

A mathematically clean solution to this problem is as follows: Call two complexities  $C_\phi$  and  $C_\psi$  *equivalent*,  $C_\phi \equiv C_\psi$ , if there is a constant  $c$  such that for all  $x$ ,

$$|C_\phi(x) - C_\psi(x)| \leq c.$$

Then the equivalence relation  $\equiv$  induces equivalence classes

$$[C_\phi] = \{C_\psi : C_\psi \equiv C_\phi\}.$$

We order the equivalence classes by  $[C_\phi] \leq [C_\psi]$  iff for all  $x$  we have  $C_\phi(x) \leq C_\psi(x)$ . Then the ordering of the equivalence classes according to  $\leq$  has a single minimal element, namely,  $[C_{\phi_0}]$ , such that for all  $C_\psi$ ,

$$[C_{\phi_0}] \leq [C_\psi].$$

We have somewhat glibly overlooked the fact that our definition of Kolmogorov complexity is relative to the particular effective enumeration of Turing machines as used in the proof of the Invariance Theorem 2.1.1, page 97. We have claimed that the quantity of information in an object depends on itself alone. That is, it should be independent of the particular enumeration of Turing machines.

Consider two different enumerations of all partial recursive functions, say  $\phi_1, \phi_2, \dots$  and  $\psi_1, \psi_2, \dots$ . Assume that the  $\phi$  enumeration is the enumeration corresponding to our effective enumeration of Turing machines as used in the proof of the Invariance Theorem.

Let the standard enumeration  $\phi_1, \phi_2, \dots$  and the other enumeration  $\psi_1, \psi_2, \dots$  be related by  $\psi_i = \phi_{f(i)}$  and  $\phi_i = \psi_{g(i)}$ ,  $i = 1, 2, \dots$ . If both  $f$  and  $g$  are partial recursive, then the enumerations are called *recursively isomorphic* and are both *acceptable numberings* (Section 1.7, Exercise 1.7.6 on page 41). Let  $C(x)$  be the complexity with respect to the reference function in the  $\phi$  enumeration, and let  $C'(x)$  be the complexity with respect to the reference function in the  $\psi$  enumeration. It is an easy exercise to show that there is a constant  $c$  such that  $|C(x) - C'(x)| < c$  for all  $x$ . (Hint: use the indexes of  $f$  and  $g$  in the enumerations.)

Therefore, not only additively optimal functions in the *same* acceptable numberings yield complexities that are equal up to a fixed constant, but additively optimal functions in two *different* acceptable numberings do so as well. Hence, Kolmogorov complexity is *recursively invariant* between acceptable numberings, even though we have chosen to define it using the specific enumeration of Turing machines of Section 1.7. Using an analogy due to Hartley Rogers, Jr., the fixed choice of effective enumeration of Turing machines can be compared with using a particular coordinate system to establish coordinate-free results in geometry.

A contradiction is possible only if there is no recursive isomorphism between the  $\phi$  enumeration and the  $\psi$  enumeration. We give an example of an enumeration of all partial recursive functions for which an additively optimal function yields a complexity  $C'(x)$  such that  $|C(x) - C'(x)|$  is unbounded. Let  $C(x)$  be defined with respect to the  $\phi$  enumeration as in Theorem 2.1.1. Define the  $\psi$  enumeration as follows: The even functions  $\psi_{2i}$  are defined by  $\psi_{2i}(1) := y_i$  for some  $y_i$  with  $C(y) \geq i^2$  and  $\psi_{2i}(x) := \phi_i(x)$  for all  $x > 1$ . The odd functions  $\psi_{2i+1}$  are given by  $\psi_{2i+1} := \phi_i$ .

Clearly, the  $\psi$  enumeration contains all partial recursive functions. By way of contradiction, assume that  $C'(\cdot)$  is the Kolmogorov complexity in the  $\psi$ -enumeration defined as in Theorem 2.1.1. Then,  $C'(y_i) \leq C'_{\psi_{2i}}(y_i) + c_{\psi_{2i}}$ . By construction,  $C'_{\psi_{2i}}(y_i) = 1$  and  $c_{\psi_{2i}} \leq 2 \log 2i + O(1)$ . On the other hand,  $C(y_i) > i^2$  by construction. Hence,  $|C'(y_i) - C(y_i)|$  rises unboundedly with  $i$ .

### 2.1.4 Concrete Kolmogorov Complexity

It is possible to eliminate the indeterminacy of “equality up to a constant” everywhere by using a fixed domain of objects, a fixed effective enumeration of Turing machines, and a fixed choice of additively optimal function (rather the universal Turing machine that computes it). Start from the enumeration of Turing machines in Section 1.7. Fix any small universal machine, say  $U$ , with state-symbol product less than 30. It is known that there exists at least one  $7 \times 4$  universal Turing machine as mentioned in the comment on page 31. In Section 3.2 we exhibit a universal reference machine  $U$  to fix a concrete Kolmogorov complexity with  $C(x|y) \leq l(x) + 2$  and  $C(x) \leq l(x) + 8$ .

For every  $x$  it is of course possible to choose a universal Turing machine  $U'$  such that  $C_{U'}(x) = 0$  (in this notation identifying  $U'$  with the function it computes). For every such universal Turing machine  $U'$ , we have for all  $x$  that

$$C(x) \leq C_{U'}(x) + C(U').$$

Here  $C(U')$  is at least the length of the shortest program  $p$  such that for all programs  $q$  we have  $U(pq) = U'(q)$ . This means that if  $C_{U'}(x) = 0$ , then  $C(U') \geq C(x)$ . That is,  $C_{U'}(x) = 0$  unavoidably means that the description of  $U'$  contains a description of  $x$ . Therefore, in order to assign low complexity to a large and complicated object a universal machine has to be large and complicated as well.

## Exercises

- 2.1.1.** [15] (a) Show that  $C(0^n|n) \leq c$ , where  $c$  is a constant independent of  $n$ .  
 (b) Show that  $C(\pi_{1:n}|n) \leq c$ , where  $\pi = 3.1415\dots$  and  $c$  is some constant independent of  $n$ .  
 (c) Show that  $C(a_{1:n}|n) \approx n/4$ , where  $a_i$  is the  $i$ th bit in Shakespeare’s *Romeo and Juliet*.  
 (d) What is  $C(a_{1:n}|n)$ , where  $a_i$  is the  $i$ th bit in the expansion of the fine structure constant  $a = e^2/\hbar c$ , in physics.

*Comments.* Hint: for Item (c) use known facts concerning the letter frequencies (entropy) in written English. Source: T.M. Cover, *The Impact of Processing Technique on Communications*, J.K. Skwirzynski, Ed., Martinus Nijhof, 1985, pp. 23–33.

- 2.1.2.** [10] Let  $x$  be a finite binary string with  $C(x) = q$ . What is the complexity  $C(x^q)$ , where  $x^q$  denotes the concatenation of  $q$  copies of  $x$ ?

- 2.1.3.** [14] Show that there are infinite binary sequences  $\omega$  such that the length of the shortest program for reference Turing machine  $U$  to compute the consecutive digits of  $\omega$  one after another can be significantly

shorter than the length of the shortest program to compute an initial  $n$ -length segment  $\omega_{1:n}$  of  $\omega$ , for any large enough  $n$ .

*Comments Hint:* choose  $\omega$  a recursive sequence with shortest program of length  $O(1)$ . Then  $C(\omega_{1:n}) = C(n) + O(1)$ , which goes to  $\infty$  with  $n$ .

**2.1.4.** [12] Prove that for every  $x$ , there is an additively optimal function  $\phi_0$  (as in Theorem 2.1.1) such that  $C_{\phi_0}(x) = 0$ . Prove the analogous statement for  $x$  under condition  $y$ .

**2.1.5.** [07] Below,  $x$ ,  $y$ , and  $z$  are arbitrary elements of  $\mathcal{N}$ . Prove the following:

- (a)  $C(x|y) \leq C(x) + O(1)$ .
- (b)  $C(x|y) \leq C(x, z|y) + O(1)$ .
- (c)  $C(x|y, z) \leq C(x|y) + O(1)$ .
- (d)  $C(x, x) = C(x) + O(1)$ .
- (e)  $C(x, y|z) = C(y, x|z) + O(1)$ .
- (f)  $C(x|y, z) = C(x|z, y) + O(1)$ .
- (g)  $C(x, y|x, z) = C(y|x, z) + O(1)$ .
- (h)  $C(x|x, z) = C(x|x) + O(1) = O(1)$ .

**2.1.6.** [14] Let  $\phi_k$  be any partial recursive function in the effective enumeration  $\phi_1, \phi_2, \dots$ . Let  $x, y, z$  be arbitrary elements of  $\mathcal{N}$ . Show the following:

- (a)  $C(\phi_k(x)|y) \leq C(x|y) + 2l(k) + O(1)$ .

Assume  $\phi_k$  is also one-to-one. Show

- (b)  $C(y|\phi_k(x)) \geq C(y|x) - 2l(k) + O(1)$ .
- (c)  $|C(x) - C(\phi_k(x))| \leq 2l(k) + O(1)$ .
- (d)  $C(x|y, z) \leq C(x|\phi_k(y), z) + 2l(k) + O(1)$ .

**2.1.7.** [12] Let  $x, y, z$ , and  $\phi_k$  be as before. Show the following.

- (a)  $C(x, y) \leq C(x) + 2l(C(x)) + C(y|x) + O(1)$ .
- (b)  $C(\phi_k(x, y)) \leq C(x) + 2l(C(x)) + C(y|x) + 2l(k) + O(1) \leq C(x) + 2l(C(x)) + C(y) + 2l(k) + O(1)$ .

**2.1.8.** [12] Show that if  $\phi$  is a fixed one-to-one and onto recursive function  $\phi : \{0, 1\}^* \rightarrow \{0, 1\}^*$ , then for every  $x \in \{0, 1\}^*$ ,

$$C(x) - C(x|\phi(x)) = C(x) + O(1) = C(\phi(x)) + O(1).$$

**2.1.9.** • [19] We investigate the invariance of  $C$  under change of program representations from 2-ary to  $r$ -ary representations. Let  $A_r = \{0, 1, \dots, r-1\}^*$ ,  $r \geq 2$ , and  $A = \mathcal{N}^*$  with  $\mathcal{N}$  the set of natural numbers. A function  $\phi : A_r \times A \rightarrow A$  is called an  $r$ -ary *decoder*. In order not to hide too much information in the decoder we want it to be a “simple” function, a partial recursive one. Analogous to the definitions in the main text, for any binary decoder  $\phi$  and  $x, y$  in  $A$ ,

$$C_\phi(x|y) = \min\{l(p) : \phi(p, y) = x\},$$

or  $\infty$  if such  $p$  does not exist.

- (a) Prove the Invariance Theorem 2.1.1 under this definition of  $C$ .
- (b) Define for each pair of natural numbers  $r, s \geq 2$  a standard encoding  $E$  of strings  $x$  in base  $r$  to strings  $E(x)$  in base  $s$  such that  $l(E(x)) \leq l(x) \log r / \log s + 1$ .
- (c) Prove the Invariance Theorem for  $r$ -ary decoders  $\phi$ . First, let us define  $C_\phi(x|y) = \min\{l(p) \log r : \phi(p, y) = x\}$  and  $C_\phi(x|y) = \infty$  if such  $p$  does not exist. Then prove that there exists an additively optimal (universal)  $r$ -ary decoder  $\phi_0$  such that for all  $s$ , for all  $s$ -ary decoders  $\phi$ , there exists a constant  $c_\phi$  such that for all  $x, y \in A$  we have

$$C_{\phi_0}(x|y) \leq C_\phi(x|y) + c_\phi.$$

- (d) Show that for any  $x \in A_r$  of length  $n$ , we have  $C(x) \leq n \log r + 2 \log r + c$  for some fixed  $c$ , independent of  $x$  and  $r$ .
- (e) Fix natural numbers  $r, s \geq 2$  and choose an additively optimal  $r$ -ary decoder and an additively optimal  $s$ -ary decoder. Call the associated canonical  $C$  measures respectively  $C_r$  and  $C_s$ . Show that there exists a constant  $c$  such that for all  $x, y$  in  $A$  we have

$$|C_r(x|y) - C_s(x|y)| \leq c,$$

where  $c$  is independent of  $x$  and  $y$ . Conclude that  $C_2$ , the  $C$  measure treated in the main text, is universal in the sense that neither the restriction to binary objects to be described nor the restriction to binary descriptions (programs) results in any loss of generality.

*Comments.* In general, if we denote by  $C_r(x)$  the analogous complexity of  $x$  in terms of programs over alphabets of  $r$  letters ( $C_2(x) = C(x)$ ) but for  $r > 2$  without the  $\log r$  normalizing multiplicative factor as in Item (c)), then by the same analysis as of Item (c) we find  $C_r(x) \sim C(x) / \log r$ . Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

- 2.1.10.** [12] (a) Show that  $C(x + C(x)) \leq C(x) + O(1)$ .

- (b) Show that if  $m \leq n$ , then  $m + C(m) \leq n + C(n) + O(1)$ .

*Comments.* Hint for Item (a): if  $U(p) = x$  with  $l(p) = C(x)$ , then  $p$  also suffices to reconstruct  $x + l(p)$ . Hint for Item (b): use Item (a). Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987; result is attributed to C.P. Schnorr.

**2.1.11.** [13] Let  $\phi_1, \phi_2, \dots$  be the standard enumeration of the partial recursive functions, and let  $a$  be a fixed natural number such that the set  $A = \{x : \phi_k(y) = \langle a, x \rangle \text{ for some } y \in \mathcal{N}\}$  is finite. Show that for each  $x$  in  $A$  we have  $C(x|a) \leq l(d(A)) + 2l(k) + O(1)$ .

**2.1.12.** [18] Define the *function complexity* of a function  $f : \mathcal{N} \rightarrow \mathcal{N}$ , restricted to a finite domain  $D$ , as

$$C(f|D) = \min\{l(p) : \forall_{x \in D}[U(p, x) = f(x)]\}.$$

(a) Show that for all recursive functions  $f$ , there exists a constant  $c_f$  such that for all finite  $D \subseteq \mathcal{N}$ , we have  $C(f|D) \leq c_f$ .

(b) Show that for all enumerable functions, for all  $D = \{i : i \leq n\}$ , we have  $C(f|D) \leq \log n + c_f$ , where  $c_f$  depends on  $f$  but not on  $D$ .

*Comments.* Compare Theorem 2.7.2. Source: J.M. Barzdins, *Soviet Math. Dokl.*, 9(1968), 1251–1254.

**2.1.13.** [26] Show that  $2C(a, b, c) \leq C(a, b) + C(b, c) + C(c, a) + O(\log n)$ .

*Comments.* For an application relating the 3-dimensional volume of a geometric object in Euclidean space with the 2-dimensional volumes of its projections see page 457. Source: D. Hammer and A.Kh. Shen', 'A strange application of Kolmogorov complexity,' *Math. Systems Theory*, to appear.

## 2.2 Incompressibility

---

It is easy to see that there are strings that can be described by programs much shorter than themselves. For instance, the function defined by  $f(1) = 2$  and  $f(i) = 2^{f(i-1)}$  for  $i > 1$  grows very fast,  $f(k)$  is a “stack” of  $k$  twos. Yet for each  $k$  it is clear that the string  $x = 1^{f(k)}$ , or the integer  $y = 2^{f(k)}$ , have at most complexity  $C(k) + c$  for some constant  $c$  independent of  $k$ .

Trivially, this simple argument can be generalized to show the following fact: for every recursive function  $\phi$ , no matter how fast it grows, there is a constant  $c$  such that for each value of  $n$  there is a string  $x$  such that  $l(x) = \phi(n)$  but  $C(x) \leq n + c$ . That is, for an appropriate sequence of strings, the ratio of string length to description length can increase as fast as any recursive function—some strings are very *compressible*.

What about incompressibility? By a simple counting argument one can show that whereas some strings can be very far compressed, the majority of strings cannot be compressed at all.

For each  $n$  there are  $2^n$  binary strings of length  $n$ , but only  $\sum_{i=0}^n 2^i = 2^n - 1$  possible shorter descriptions. Therefore, there is at least one binary string  $x$  of length  $n$  such that  $C(x) \geq n$ . We call such strings *incompressible*. It also follows that for any length  $n$  and any binary string  $y$ , there is a binary string  $x$  of length  $n$  such that  $C(x|y) \geq n$ .

**Definition 2.2.1** For each constant  $c$  we say a string  $x$  is *c-incompressible* if  $C(x) \geq l(x) - c$ .

Strings that are incompressible (say,  $c$ -incompressible with small  $c$ ) are patternless, since a pattern could be used to reduce the description length. Intuitively, we think of such patternless sequences as being random, and we use “random sequence” synonymously with “incompressible sequence.” Later we give a formalization of the intuitive notion of a random sequence as a sequence that passes all effective tests for randomness.

How many strings of length  $n$  are  $c$ -incompressible? By the same counting argument we find that the number of strings of length  $n$  that are  $c$ -incompressible is at least  $2^n - 2^{n-c} + 1$ . Hence there is at least one 0-incompressible string of length  $n$ , at least one-half of all strings of length  $n$  are 1-incompressible, at least three-fourths of all strings of length  $n$  are 2-incompressible, . . . , and at least the  $(1 - 1/2^c)$ th part of all  $2^n$  strings of length  $n$  are  $c$ -incompressible. This means that for each constant  $c > 1$  the majority of all strings of length  $n$  with  $n > c$  is  $c$ -incompressible. We generalize this to the following simple but extremely useful *Incompressibility Theorem*.

**Theorem 2.2.1** Let  $c$  be a positive integer. For each fixed  $y$ , every finite set  $A$  of cardinality  $m$  has at least  $m(1 - 2^{-c}) + 1$  elements  $x$  with  $C(x|y) \geq \log m - c$ .

**Proof.** The number of programs of length less than  $\log m - c$  is

$$\sum_{i=0}^{\log m - c - 1} 2^i = 2^{\log m - c} - 1.$$

Hence, there are at least  $m - m2^{-c} + 1$  elements in  $A$  that have no program of length less than  $\log m - c$ .  $\square$

As an example, set  $A = \{x : l(x) = n\}$ . Then the cardinality of  $A$  is  $m = 2^n$ . Since Theorem 2.1.2 asserts that  $C(x) \leq n + c$  for some fixed  $c$  and all  $x$  in  $A$ , Theorem 2.2.1 demonstrates that this trivial estimate is

quite sharp. The deeper reason is that since there are few short programs, there can be only few objects of low complexity.

It is important to realize that Theorem 2.1.1 and Theorem 2.2.1, together with the trivial upper bound of Theorem 2.1.2, give us already all we need for most applications.

**Example 2.2.1** Are all substrings of incompressible strings also incompressible? A string  $x = uvw$  of length  $n$  can be specified by a short program  $p$  for  $v$  and the string  $uw$  itself. Additionally, we need information on how to tell these items apart. For instance,  $q = \overline{l(p)pl(u)}uw$  is a program for  $x$ . There exists a machine  $T$ , that starting on the left end of  $q$ , first determines  $l(p)$ , then uses  $l(p)$  to delimit  $p$ , and computes  $v$  from  $p$ . Continuing on its input,  $T$  determines  $l(u)$  and uses this to delimit  $u$  on the remainder of its input. Subsequently,  $T$  reassembles  $x$  from the three pieces  $v$ ,  $u$ , and  $w$  it has determined. It follows that  $C(x) \leq C_T(x) + O(1) \leq l(q) + O(1)$ , since  $l(q) \leq C(v) + 2l(C(v)) + 2l(n) + n - l(v) + 2$ . Therefore,

$$C(x) \leq C(v) + n - l(v) + 4 \log n + O(1).$$

Hence, for  $c$ -incompressible strings  $x$  with  $C(x) \geq n - c$  we obtain

$$C(v) \geq l(v) - O(\log n).$$

Thus, we have shown that  $v$  is incompressible up to a additive term logarithmic in  $n$ .

Can we hope to prove  $C(v) \geq l(v) - O(1)$  for all  $x$  and  $v$ ? If this were true, then  $x$  could not contain long regular subsequences, for instance, a subsequence of  $k$  zeros has complexity  $O(\log k)$  and not  $k - O(1)$ . However, the very restriction on  $x$  of not having long regular subsequences imposes some regularity on  $x$  by making it a member of a relatively small set. Namely, we can describe  $x$  by stating that it does not contain certain subsequences, followed by  $x$ 's index in the set that is determined by these constraints. But the set of which  $x$  is a member is so small that  $C(x)$  drops below  $n - c$ , and  $x$  is compressible. Hence, the very incompressibility of  $x$  requires that it has compressible substrings. This corresponds to a fact we know from probability theory: a random sequence must contain long runs of zeros.  $\diamond$

**Example 2.2.2** If  $p$  is a shortest program for  $x$ , so that  $C(x) = l(p)$ , then we would like to assert that  $p$  is incompressible. This time our intuition corresponds with the truth. There is a constant  $c > 0$  such that for all strings  $x$  we have  $C(p) \geq l(p) - c$ . For suppose the contrary, and for every constant  $c$  there is an  $x$  and a shortest program  $q$  that generates  $p$  with  $l(q) < l(p) - c$ . Define a universal machine  $V$  that works just like the reference machine  $U$ , except that  $V$  first simulates  $U$  on its input to obtain the output,

and then uses this output as input on which to simulate  $U$  once more. Let  $V = T_i$ , the  $i$ th Turing machine in the standard enumeration. Then,  $U$  with input  $1^i 0q$  computes  $x$ , and therefore  $C(x) < l(p) - c + i + 1$ . But this contradicts  $l(p) = C(x)$  for  $c \geq i + 1$ .  $\diamond$

**Example 2.2.3** We continue Example 2.1.4 on page 101 that  $C(x, y)$  is not *subadditive* since the logarithmic term in Equation 2.2 cannot be eliminated. Namely, there are  $(n+1)2^n$  pairs  $(x, y)$  of binary strings whose sum of lengths is  $n$ . By Theorem 2.2.1 there is a pair  $(x, y)$  with  $l(x) + l(y) = n$  such that  $C(x, y) \geq n + \log n - 1$ . But Theorem 2.1.2 on page 100 states that  $C(x) + C(y) \leq l(x) + l(y) + c$  for some constant  $c$  independent of  $x$  and  $y$ . Hence, for all  $n$  there are  $x$  and  $y$  of length at most  $n$  such that

$$C(x, y) > C(x) + C(y) + \log n - c,$$

where  $c$  is a constant independent of  $x$  and  $y$ . For the unmarked concatenation  $xy$  with  $l(xy) = n$ , if  $C(xy) \geq n$ , then  $xy$  contains a block of 0's or 1's of length about  $\log n$  (as in Example 2.2.1). We can split up  $xy$  in  $x$  ending with  $\log n$  0's or 1's. This means that  $C(x) \leq l(x) - \log n + 2 \log \log n$ . Then,  $C(xy) \geq C(x) + C(y) + \log n - 2 \log \log n - c$ .  $\diamond$

There is a particular use we had in mind when defining conditional Kolmogorov complexity. Namely, we often want to speak about the complexity of  $x$  given its length  $n$ . This is because a string  $x$  of length  $n$  carries in a sense *two* quantities of information, one associated with the irregularity of the pattern of 0's and 1's in  $x$ , and one associated with the length  $n$  of  $x$ .

**Example 2.2.4** One effect of the information quantity associated with the length of strings is that  $C(x)$  is *nonmonotonic on prefixes*. This can be due to the information contained in the length of  $x$ . That is, for  $m < n$  we can still have  $C(m) > C(n)$ . But then  $C(y) > C(x)$  for  $x = 1^n$  and  $y = 1^m$ , notwithstanding that  $y$  is a proper prefix of  $x$ . For example, if  $n = 2^k$ , then  $C(1^n) \leq \log \log n + O(1)$ , while Theorem 2.2.1 shows that there are  $m < n$  for which  $C(1^m) \geq \log n - O(1)$ . Therefore, the complexity of a part can turn out to be bigger than the complexity of the whole. In an initial attempt to solve this problem we may try to eliminate the effect of the length of the string on the complexity measure by treating the length as given.  $\diamond$

**Definition 2.2.2** The *length-conditional* Kolmogorov complexity of  $x$  is  $C(x|l(x))$ .

Roughly speaking this means the length of the shortest program for  $x$  may save up to  $\log l(x)$  bits in comparison with the shortest program in the unconditional case. Clearly, there is a constant  $c$  such that for all  $x$ ,

$$C(x|l(x)) \leq C(x) + c.$$

While on the face of it the measure  $C(x|l(x))$  gives a pure estimate of the quantity of information in solely the pattern of 0's and 1's of  $x$ , this is not always true. Namely, sometimes the information contained in  $l(x)$  can be used to determine the pattern of zeros and ones of  $x$ . This effect is noticeable especially in the low-complexity region.

**Example 2.2.5** For each integer  $n$ , the  $n$ -string is defined by  $n0^{n-l(n)}$  (using the binary string  $n$ ). There is a constant  $c$  such that for all  $n$ , if  $x$  is the  $n$ -string, then  $C(x|n) \leq c$ . Namely, given  $n$  we can find the  $n$ th binary string according to Equation 1.1 and pad the string with zeros up to overall length  $n$ . We use  $n$ -strings to show that unfortunately, like the original  $C(x)$ , the complexity measure  $C(x|l(x))$  is *not monotonic over the prefixes*. Namely, if we choose  $n$  such that its pattern of 0's and 1's is very irregular,  $C(n) \geq l(n)$ , then for  $x = n0^{n-l(n)}$ , we still find  $C(x|l(x)) \leq c$ . But clearly  $C(n|l(n)) \geq C(n) - C(l(n)) \geq \log n - 2 \log \log n$ .  $\diamond$

**Example 2.2.6** Consider the complexity of a string  $x$ , with  $x$  an element of a given set  $A$ . Clearly, the information that an element belongs to a particular set severely curtails the complexity of that element if that set is small or sparse. The following is a simple application of the very useful Theorem 2.1.3 on page 103. Let  $A$  be a subset of  $\mathcal{N}$ . Define  $A^{\leq n} = \{x \in A : l(x) \leq n\}$ . We call  $A$  *meager* if  $\lim d(A^{\leq n})/2^n = 0$  for  $n \rightarrow \infty$ . For example, the set of all finite strings that have twice as many 0's as 1's is meager. We show that meagerness may imply that almost all strings in the meager set have short descriptions.

**Claim 2.2.1** If  $A$  is recursive and meager, then for each constant  $c$  there are only finitely many  $x$  in  $A$  that are  $c$ -incompressible ( $C(x) \geq l(x) - c$ ).

**Proof.** Consider the lexicographic enumeration of all elements of  $A$ . Because  $A$  is recursive, there is a total recursive function  $\phi_i$  that enumerates  $A$  in increasing order. Hence, for the  $j$ th element  $x$  of  $A$  we have  $C(x) \leq C(j) + 2l(i) + 1$ . If  $x$  has length  $n$ , then the meagerness of  $A$  implies that for each constant  $c'$ , no matter how large,  $n - C(j) > c'$  from some  $n$  onwards. Hence,  $C(x) < n - c' + 2l(i)$ . The proof is completed by setting  $c' = c + 2l(i)$ .  $\square$   $\diamond$

### 2.2.1 Randomness Deficiency

If we know that  $x$  belongs to a subset  $A$  of the natural numbers, then we can consider its complexity  $C(x|A)$ . For instance,  $C(x) = C(x|\mathcal{N})$ , because it is understood that  $x$  is a natural number. If  $x$  is an element of a finite set  $A$ , then Theorem 2.1.3 asserts that  $C(x|A) \leq l(d(A)) + c$  for some  $c$  independent of  $x$  but possibly dependent on  $A$ . For instance, the infinite meager sets of Example 2.2.6 contain only finitely many incompressible strings.

**Definition 2.2.3** The *randomness deficiency* of  $x$  relative to  $A$  is defined as  $\delta(x|A) = l(d(A)) - C(x|A)$ . It follows that  $\delta(x|A) \geq -c$  for some fixed constant  $c$  independent of  $x$ .

If  $\delta(x|A)$  is large, then this means that there is a description of  $x$  with the help of  $A$  that is considerably shorter than just giving  $x$ 's serial number in  $A$ . There are comparatively few objects with large randomness deficiency—this is the substance of Martin-Löf's notion of a statistical test in Section 2.4. Quantitatively this is expressed as follows:

**Theorem 2.2.2** Assume the discussion above. Then,  $d(\{x : \delta(x|A) \geq k\}) \leq d(A)/2^{k-1}$ .

**Proof.** There are less than  $2^{l+1}$  descriptions of length at most  $l$ .  $\square$

By Theorem 2.1.3 on page 103, the complexity of a string  $x$  in a given finite section of a recursively enumerable set is bounded above by the logarithm of the cardinality of that finite section. Let  $\langle \cdot \rangle : \mathcal{N}^2 \rightarrow \mathcal{N}$  be the standard recursive bijective pairing function. Let  $R = \{(x, y) : \phi(i) = \langle x, y \rangle, i \geq 1\}$  with  $\phi$  a partial recursive function, say  $\phi = \phi_r$  in the standard enumeration of partial recursive functions. Then  $R$  is recursively enumerable. Let the set  $A = \{x : (x, y) \in R\}$  be finite. We can assume that  $A$  is enumerated without repetition, and that  $j \leq d(A)$  is the position of  $x$  in this enumeration. Clearly,

$$C(x|y) \leq \log d(A) + \log r + 2 \log \log r + O(1).$$

As above, define  $C(x|A) = C(x|y)$  with the obvious interpretation. The randomness deficiency of  $x$  relative to  $y$  is

$$\delta(x|y) = \log d(A) - C(x|y).$$

The randomness deficiency measures the difference between the maximal complexity of a string in  $A$  and the complexity of  $x$  in  $A$ . Now, the defect of randomness is positive up to a fixed constant independent of  $x$  and  $A$  (but dependent on  $r$ ). We may consider  $x$  to be random in the set  $A$  iff  $\delta(x|y) = O(1)$ . If  $A$  is the set of binary strings of length  $n$ , or equivalently,  $R$  is the set  $\{(x, n) : l(x) = n\}$  and  $A = \{x : l(x) = n\}$ , then we note that

$$\delta(x|n) = n - C(x|n) + O(1).$$

That is,  $x$  is a random finite string in our informal sense iff  $\delta(x|n) = O(1)$ . It will turn out that this coincides with Martin-Löf's notion of randomness in Section 2.4.

2.2.2  
**Kolmogorov Minimal Sufficient Statistic**

Let  $k$  and  $\delta$  be natural numbers. A string  $x$  is called  $(k, \delta)$ -stochastic if there is a finite set  $A \subseteq \mathcal{N}$  and  $x \in A$  such that

$$x \in A, \quad C(A) \leq k, \quad C(x|A) \geq \log d(A) - \delta.$$

The first inequality (with  $k$  not too large) means that  $A$  is sufficiently simple. The second inequality (with the randomness deficiency  $\delta$  not too large) means that  $x$  is an undistinguished (typical) element of  $A$ . Indeed, if  $x$  had properties defining a very small subset  $B$  of  $A$ , then these could be used to obtain a simple description of  $x$  by determining its ordinal number in  $B$ , which would require  $\log d(B)$  bits, which is much less than  $\log d(A)$ .

This notion of  $(k, \delta)$ -stochasticity is connected with the very foundations of mathematical statistics. Suppose we carry out some probabilistic experiment of which the outcome can be a priori every natural number. Suppose this number is  $x$ . Knowing  $x$ , we want to recover the probability distribution  $P$  on the set of natural numbers  $\mathcal{N}$ . It seems reasonable to require that first,  $P$  has a simple description, and second, that  $x$  would be a “typical” outcome of an experiment with probability distribution  $P$ —that is,  $x$  is maximally random with respect to  $P$ . The analysis above addresses a simplified form of this issue where  $A$  plays the part of  $P$ —for example,  $A$  is a finite set of high-probability elements. In  $n$  tosses of a coin with probability  $p > 0$  of coming up “heads,” the set  $A$  of outcomes consisting of binary strings of length  $n$  with  $n/2$  1’s constitutes a set of cardinality  $\binom{n}{n/2} = \Theta(2^n/\sqrt{n})$ . To describe an element  $x \in A$  requires  $n - \frac{1}{2} \log n + O(1)$  bits. To describe  $A \subseteq \{0, 1\}^n$  given  $n$  requires  $O(1)$  bits (small  $k$ ). Conditioning everything on the length  $n$ , we have

$$C(x|n) \leq C(x|A, n) + C(A|n) + O(1) \leq n - \frac{1}{2} \log n + O(1),$$

and for the overwhelming majority of the  $x$ ’s in  $A$ ,

$$C(x|n) \geq n - \frac{1}{2} \log n - O(1).$$

These latter  $x$ ’s are  $(O(1), O(1))$ -stochastic. The *Kolmogorov structure function*  $C_k(x|n)$  of  $x \in \{0, 1\}^n$  is defined by

$$C_k(x|n) = \min\{\log d(A) : x \in A, \quad C(A|n) \leq k\}.$$

For a given small constant  $c$ , let  $k_0$  be the least  $k$  such that

$$C_k(x|n) + k \leq C(x|n) + c.$$

Let  $A_0$  be the corresponding set, and with some abuse of notation let  $A_0^*$  be its shortest program. Then  $A_0^*$  is the *Kolmogorov minimal sufficient statistic* for  $x$ .

This  $k_0$  with  $C(A_0|n) \leq k_0$  is the least  $k$  for which the two-part description of  $x$  is as parsimonious as the best single part description of  $x$ . If  $x$  is maximally complex in  $A_0$ , then  $A_0$  captures all structure in  $x$ . The  $C_{k_0}(x|n)$ -stage of the description just provides an index for  $x$  in  $A_0$ —essentially the description of the randomness within the string. The set  $A_0$ —rather the shortest program  $A_0^*$  that prints out the characteristic sequence of  $A_0 \in \{0,1\}^n$ —is called the *Kolmogorov minimal sufficient statistic* for  $x$ , given  $n$ . Note that all programs describing sets  $A$  with  $C(A|n) \leq k$  such that  $C_k(x|n) + k \approx C(x|n)$  are sufficient statistics in the sense that the complexity  $C(x|A, n) \approx \log d(A)$  is maximal and the randomness deficiency  $\delta(x|A) = \log d(A) - C(x|A, n)$  is minimal. But the “minimal” sufficient statistic is induced by the set  $A$  having the shortest description among them. This is closely related to the hypothesis identification method satisfying the Fundamental Inequality Equation 5.23 using minimum description length in Section 5.5.

The behavior of the Kolmogorov structure function is as follows. Consider  $x \in \{0,1\}^n$ . If  $k = 0$  then the smallest set  $A$  containing  $x$  that can be described by  $k$  bits (and given  $n$ ) is  $A = \{0,1\}^n$ . Therefore,  $C_0(x|n) = n$ . For increasing  $k$ , the size of the set  $A$  one can describe decreases rapidly until  $C_k(x|n) + k \approx C(x|n)$ . Further increase of  $k$  halves the set  $A$  for each additional bit of  $k$  until  $k = C(x|n)$ . For  $k \geq C(x|n)$  obviously the smallest set  $A$  containing  $x$  that one can describe using  $k$  bits (given  $n$ ) is the singleton set  $A = \{x\}$ .

Let us look again at the coin toss example. If the probability  $p$  of tossing “1” is unknown, then the best two-part description of a string  $x$  representing the sequence of  $n$  outcomes is to give the number  $k$  of 1’s in  $x$  first, followed by the index  $j \leq d(A)$  of  $x$  in the set  $A$  of strings with  $k$  1’s. If  $k$  is incompressible with  $C(k|n) \approx \log n$  and  $C(j|k, n) \approx \log \binom{n}{k}$  then the Kolmogorov minimal sufficient statistic is given by the characteristic function of  $A$  (described by  $k$  and  $n$ ). However if  $p$  is a simple value like  $\frac{1}{2}$  or  $1/\pi$ , then with overwhelming probability the Kolmogorov minimal sufficient characteristic is given by a description of  $p$  and  $C(k|n) = O(1)$ .

Surprisingly enough, there are outcomes  $x$  that are atypical for *any simple* distribution  $P$  (Exercise 2.2.13, page 119). To obtain such a pathological case we have to assume that our only existing information about  $P$  is the value of the outcome  $x$ . But in practice the specific nature of the problem often suggests beforehand a possible form of  $P$ , and it only remains to select some parameters. We will pick up the thread of the relations between Kolmogorov complexity and inductive reasoning in Chapter 5.

## Exercises

---

**2.2.1.** [08] Show the following continuity property of  $C(x)$ . For all natural numbers  $x, y$  we have  $|C(x + y) - C(x)| \leq 2l(y) + O(1)$ .

**2.2.2.** [15] Let  $x$  satisfy  $C(x) \geq n - O(1)$ , where  $n = l(x)$ .

(a) Show that  $C(y), C(z) \geq n/2 - O(1)$  for  $x = yz$  and  $l(y) = l(z)$ .

(b) Show that  $C(y) \geq n/3 - O(1)$  and  $C(z) \geq 2n/3 - O(1)$  for  $x = yz$  and  $l(z) = 2l(y)$ .

(c) Let  $x = x_1 \dots x_{\log n}$  with  $l(x_i) = n/\log n$  for all  $1 \leq i \leq \log n$ . Show that  $C(x_i) \geq n/\log n - O(1)$  for all  $1 \leq i \leq \log n$ .

**2.2.3.** [21] Let  $x$  satisfy  $C(x) \geq n - O(1)$ , where  $n = l(x)$ . Show that for all divisions  $x = yz$  we have  $n - \log n - 2\log \log n \leq C(y) + C(z)$  and for some divisions we have  $C(y) + C(z) \leq n - \log n + \log \log n$ .

**2.2.4.** [23] Assume that the elements of  $\{1, \dots, n\}$  are uniformly distributed with probability  $1/n$ . Compute the expected value of  $C(x)$  for  $1 \leq x \leq n$ .

*Comments.* Hint:  $\sum_{x=1}^n \frac{C(x)}{n} \geq \sum_{i=1}^{\log n} 2^{-i} \left(1 - \frac{i}{\log n}\right) = \log n + O(1)$ .

**2.2.5.** [14]  $x$  is an  $n$ -string if  $x$  has length  $n$  and  $x = n00\dots 0$ .

(a) Show that there is a constant  $c$  such that for all  $n$ -strings  $x$  we have  $C(x|n) \leq c$ . (Of course,  $c$  depends on the reference Turing machine  $U$  used to define  $C$ .)

(b) Show there is a constant  $c$  such that  $C(x|n) \leq c$  for all  $x$  in the form of the  $n$ -length prefix of  $nn\dots n$ .

(c) Let  $c$  be as in Item (a). Consider any string  $x$  of length  $n$  with  $C(x|n) \gg c$ . Prove that the extension of  $x$  to a string  $y = x00\dots 0$  of length  $x$  has complexity  $C(y|x) \leq c$ . Conclude that there is a constant  $c$  such that each string  $x$ , no matter how high its  $C(x|l(x))$  complexity, can be extended to a string  $y$  with  $C(y|l(y)) < c$ .

*Comments.* The  $C(x)$  measure contains the information about the *pattern* of 0's and 1's in  $x$  and information about the *length*  $n$  of  $x$ . Since most  $n$ 's are random,  $l(n)$  is mostly about  $\log n$ . In this case, about  $\log n$  bits of the shortest program  $p$  for  $x$  will be used to account for  $x$ 's length. For  $n$ 's that are easy to compute this is much less. This seems a minor problem at high complexities, but becomes an issue at low complexities, as follows. If the quantities of information related to the *pattern only* is low, say less than  $\log n$ , for two strings  $x$  and  $y$  of length  $n$ , then distinctions between these quantities for  $x$  and  $y$  may get blurred in the comparison between  $C(x)$  and  $C(y)$  if the quantity of information related to length  $n$  dominates in both. The  $C(x|l(x))$  complexity was meant to measure the information content of  $x$  apart from its length.

However, as the present exercise shows, in that case  $l(x)$  may contain already the complete description of  $x$  up to a constant number of bits.  
Source: D.W. Loveland, *Inform. Contr.*, 15(1969), 510–526.

**2.2.6.** [19] (a) Show that there is a constant  $d > 0$  such that for every  $n$  there are at least  $\lfloor 2^n/d \rfloor$  strings of length  $n$  with  $C(x|n), C(x) \geq n$ .

(b) Show that there are constants  $c, d' > 0$  such that for every large enough  $n$  there are at least  $\lfloor 2^n/d' \rfloor$  strings  $x$  of length  $n - c \leq l(x) \leq n$  with  $C(x|n), C(x) > n$ .

*Comments* Hint for Item (a): There is a constant  $c > 0$  such that for every  $n$  and every  $x$  of length  $l(x) \leq n - c$  we have  $C(x|n), C(x) \leq n$  (Theorem 2.1.2). Therefore, there are at most  $2^n - 2^{n-c+1}$  programs of length  $< n$  available as shortest programs for the strings of length  $n$ . Hint for Item (b): For every  $n$  there are equally many strings of length  $\leq n$  to be described and potential programs of length  $\leq n$  to describe them. Since some programs do not halt (Lemma 1.7.5 on page 34) for every large enough  $n$  there exists a string  $x$  of length at most  $n$  that has  $C(x|n), C(x) > n$  (and  $C(x|n), C(x) \leq l(x) + c$ ). Let there be  $m \geq 1$  such strings. Given  $m$  and  $n$  we can enumerate all  $2^{n+1} - m - 1$  strings  $x$  of length  $\leq n$  and complexity  $C(x|n) \leq n$  by dovetailing the running of all programs of length  $\leq n$ . The lexicographic first string of length  $\leq n$  not in the list satisfies  $\log m + O(1) \geq C(x|n) > n$ . The unconditional result follows by padding the description of  $x$  up to length  $n$ .

Source: H. Buhrman, P.M.B. Vitányi, June 1995. Also reported by M. Kummer and L. Fortnow.

**2.2.7.** [14] We can extend the notion of  $c$ -incompressibility as follows: (All strings are binary.) Let  $g : \mathcal{N} \rightarrow \mathcal{N}$  be unbounded. Call a string  $x$  of length  $n$   $g$ -incompressible if  $C(x) \geq n - g(n)$ . Let  $I(n)$  be the number of strings  $x$  of length at most  $n$  that are  $g$ -incompressible. Show that  $\lim_{n \rightarrow \infty} I(n)/2^{n+1} = 1$ .

*Comments.* Thus, the  $g$ -incompressible finite strings have uniform probability going to 1 in the set of strings of length  $n$  for  $n \rightarrow \infty$ .

**2.2.8.** [15] Prove that for each binary string  $x$  of length  $n$  there is a  $y$  equal to  $x$  but for one bit such that  $C(y|n) \leq n - \log n + O(1)$ .

*Comments.* Hint: the set of binary strings of length  $n$  constituting a Hamming code has  $2^n/n$  elements and is recursive. Source: personal communication Imre Csiszar, May 8, 1993.

**2.2.9.** [12] A Turing machine  $T$  computes an infinite sequence  $\omega$  if there is a program  $p$  such that  $T(p, n) = \omega_{1:n}$  for all  $n$ . Define  $C(\omega) = \min\{l(p) : U(p, n) = \omega_{1:n} \text{ for all } n\}$ , or  $\infty$  if such a  $p$  does not exist. Obviously, for all  $\omega$  either  $C(\omega) < \infty$  or  $C(\omega) = \infty$ .

- (a) Show that  $C(\omega) < \infty$  iff  $0.\omega$  is a *recursive real number*, Exercise 1.7.22, page 47. For the mathematical constants  $\pi$  and  $e$ ,  $C(\pi) < \infty$  and  $C(e) < \infty$ .
- (b) Show that the reals  $0.\omega$  with  $C(\omega) < \infty$  form a countably infinite set and that the reals  $0.\omega$  with  $C(\omega) = \infty$  have uniform measure one in the total set of reals in the interval  $[0, 1]$ .

**2.2.10.** [27] We consider how information about  $x$  can be dispersed. Let  $x \in \mathcal{N}$  with  $l(x) = n$  and  $C(x) = n + O(1)$ . Show that there are  $u, v, w \in \mathcal{N}$  such that

- (i)  $l(u) = l(v) = l(w) = n/2$ ,  $C(u) = C(v) = C(w) = n/2$  ( $+O(1)$ ), and they are pairwise independent:  $C(y|z) = n/2 + O(1)$  for  $y, z \in \{u, v, w\}$  and  $y \neq z$ ;
- (ii)  $x$  can be reconstructed from any two of them:  $C(x|y, z) = O(1)$ , where  $y, z \in \{u, v, w\}$  and  $y \neq z$ .

Can you give a general solution for finding  $m+k$  elements of  $\mathcal{N}$  such that each of them has length and complexity  $n/m$ , and  $x$  can be reconstructed from any  $m$  distinct elements?

*Comments.* It is surprising that  $x$  can be reconstructed from any two out of three elements, each of one-half the complexity of  $x$ . This shows that the identity of the individual bits is not preserved in the division. Hint: assume  $n = 2m$  and  $x = x_1 \dots x_{2m}$ ,  $u = u_1 \dots u_m$ ,  $v = v_1 \dots v_m$ , and  $w = w_1 \dots w_m$  with  $u_i = x_{2i-1}$ ,  $v_i = x_{2i}$ , and  $w_i = v_i \oplus u_i$ . (Recall that  $a \oplus b = 1$  iff  $a \neq b$ .) This solution apparently does not generalize. A general solution to distribute  $x$  over  $m+k$  elements such that any group of  $m$  elements determines  $x$  can be given as follows: Compute the least integer  $y \geq x^{1/m}$ . Let  $p_i$  be the  $i$ th prime, with  $p_1 = 2$ . Distribute  $x$  over  $u_1, \dots, u_{m+k}$ , where  $u_i \equiv x \pmod{p_i^{\alpha(i)}}$ , with  $\alpha(i) = \lceil y \log_{p_i} 2 \rceil$ . Using the Chinese Remainder Theorem we find that we can reconstruct  $x$  from any subset of  $m$  elements  $u_i$ . Source: A. Shamir, *Comm. ACM*, 22:11(1979), 612–613; M.O. Rabin, *J. ACM*, 36:2(1989), 335–348.

**2.2.11.** [18] Let  $A$  be the set of binary strings of length  $n$ . An element  $x$  in  $A$  is  $\delta$ -random if  $\delta(x|A) \leq \delta$ , where  $\delta(x|A) = n - C(x|A)$  is the randomness deficiency. Show that if  $x \in B \subseteq A$ , then

$$\log \frac{d(A)}{d(B)} - C(B|A) \leq \delta(x|A) + O(\log n).$$

*Comments.* That is, no random elements of  $A$  can belong to any subset  $B$  of  $A$  that is simultaneously pure (which means that  $C(B|A)$  is small), and not large (which means that  $d(A)/d(B)$  is large). Source: A.N. Kolmogorov and V.A. Uspensky, *Theory Probab. Appl.*, 32(1987), 389–412.

**2.2.12.** [27] Let  $x \in A$ , with  $d(A) < \infty$ . Then in Section 2.2 the *randomness deficiency* of  $x$  relative to  $A$  is defined as  $\delta(x|A) = l(d(A)) - C(x|A)$ . (Here  $C(x|A)$  is defined as  $C(x|\chi)$  with  $\chi$  the characteristic sequence of  $A$  and  $l(\chi) < \infty$ .) If  $\delta(x|A)$  is large, this means that there is a description of  $x$  with the help of  $A$  that is considerably shorter than just giving  $x$ 's serial number in  $A$ . Clearly, the randomness deficiency of  $x$  with respect to sets  $A$  and  $B$  can be vastly different. But then it is natural to ask whether there exist *absolutely nonrandom* objects, objects having large randomness deficiency with respect to any appropriate set.

Show the following: Let  $a$  and  $b$  be arbitrary constants; for every sufficiently large  $n$  there exists a binary string  $x$  of length  $n$  such that  $\delta(x|A) \geq b \log n$  for any set  $A$  containing  $x$  for which  $C(A) \leq a \log n$ .

*Comments.* Source: A.Kh. Shen', *Soviet Math. Dokl.*, 28(1983), 295–299. Compare with Kamae's Theorem, Exercise 2.7.5. Let us give some interpretation of such results bearing on statistical inference. Given an experimental result, the statistician wants to infer a statistical hypothesis under which the result is *typical*. Mathematically, given  $x$  we want to find a simple set  $A$  that contains  $x$  as a typical element. The above shows that there are outcomes  $x$  such that *no* simple statistical model of the kind described is possible. The question remains whether such objects occur in the real world.

**2.2.13.** [29] The Kolmogorov structure function  $C_k(x|n)$  was defined on page 114. Show that for every large enough  $n$  there are  $x$ 's of length  $n$  such that  $C_k(x|n) = n - k$  with  $k < C_k(x|n)$ .

*Comments.* These strings are strangely reluctant to reveal their structure. Source: V.V. V'yugin, *Theory Probab. Appl.*, 32(1987), 508-512.

## 2.3 *C* as an Integer Function

We consider  $C$  as an integer function  $C : \mathcal{N} \rightarrow \mathcal{N}$ , and study its behavior, Figure 2.1. First we observe that Theorem 2.1.2 gives an *upper bound* on  $C$ : there exists a constant  $c$  such that for all  $x$  in  $\mathcal{N}$  we have  $C(x) \leq l(x) + c$ , and by Theorem 2.2.1 this estimate is almost exact for the majority of  $x$ 's. This is a computable monotonic increasing upper bound that grows to infinity. It is also the least such upper bound. It turns out that the greatest monotonic nondecreasing lower bound also grows to infinity but does so noncomputably slowly.

- Theorem 2.3.1**
- (i) *The function  $C(x)$  is unbounded.*
  - (ii) *Define a function  $m$  by  $m(x) = \min\{C(y) : y \geq x\}$ . That is,  $m$  is the greatest monotonic increasing function bounding  $C$  from below. The function  $m(x)$  is unbounded.*

(iii) For any partial recursive function  $\phi(x)$  that goes monotonically to infinity from some  $x_0$  onward, we have  $m(x) < \phi(x)$  except for finitely many  $x$ . In other words, although  $m(x)$  goes to infinity, it does so more slowly than any unbounded partial recursive function.

**Proof.** (i) This follows immediately from (ii).

(ii) For each  $i$  there is at least  $x_i$  such that for all  $x > x_i$ , the smallest program  $p$  printing  $x$  has length greater than or equal to  $i$ . This follows immediately from the fact that there are only a finite number of programs of each length  $i$ . Clearly, for all  $i$  we have  $x_{i+1} \geq x_i$ . Now observe that the function  $m$  has the property that  $m(x) = i$  for  $x_i < x \leq x_{i+1}$ .

(iii) Assume the contrary: there is a monotonic nondecreasing unbounded partial recursive function  $\phi(x) \leq m(x)$  for infinitely many  $x$ . The domain  $A = \{x : \phi(x) < \infty\}$  of  $\phi$  is an infinite recursively enumerable set. By Lemma 1.7.4,  $A$  contains an infinite recursive subset  $B$ . Define

$$\psi(x) = \begin{cases} \phi(x) & \text{for } x \in B, \\ \phi(y) & \text{with } y = \max\{z : z \in B, z < x\}, \text{ otherwise.} \end{cases}$$

This  $\psi$  is total recursive, goes monotonically to infinity, and  $\psi(x) \leq m(x)$  for infinitely many  $x$ .

Now, define  $M(a) = \max\{x : C(x) \leq a\}$ . Then,  $M(a) + 1 = \min\{x : m(x) > a\}$ . It is easy to verify that

$$\max\{x : \psi(x) \leq a + 1\} \geq \min\{x : m(x) > a\} > M(a),$$

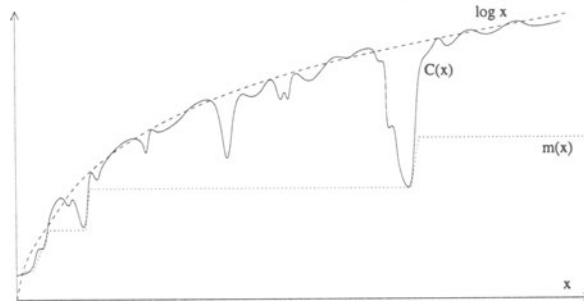
for infinitely many  $a$ 's, and the function  $F(a) = \max\{x : \psi(x) \leq a + 1\}$  is obviously total recursive. Therefore,  $F(a) > M(a)$  for infinitely many  $a$ 's. In other words,  $C(F(a)) > a$  for infinitely many  $a$ 's. But by Theorem 2.1.1 on page 97,

$$C(F(a)) \leq C_F(F(a)) + O(1) \leq l(a) + O(1).$$

This implies that there exists a constant  $c$  such that  $l(a) + c \geq a$  for infinitely many  $a$ , which is impossible.  $\square$

Notice that Items (ii) and (iii) of Theorem 2.3.1 do not hold for the length-conditional complexity  $C(x|l(x))$ . Namely, although  $C(x|l(x))$  is unbounded, it drops infinitely often to constant level. In other words, there is no unbounded monotonic function that is a lower bound on  $C(x|l(x))$  by Example 2.2.5. This phenomenon is further explored in the exercises.

The second cornerstone of the theory (millstone around its neck is probably more apt) is the *Noncomputability Theorem*.



**FIGURE 2.1.** The graph of the integer function  $C(x)$

**Theorem 2.3.2** *The function  $C(x)$  is not partial recursive. Moreover, no partial recursive function  $\phi(x)$  defined on an infinite set of points can coincide with  $C(x)$  over the whole of its domain of definition.*

**Proof.** This proof is related to that of Theorem 2.3.1, Item (iii). We prove that there is no partial recursive  $\phi$  as in the statement of the theorem. Every infinite recursively enumerable set contains an infinite recursive subset, Lemma 1.7.4. Select an infinite recursive subset  $A$  in the domain of definition of  $\phi$ . The function  $\psi(m) = \min\{x : C(x) \geq m, x \in A\}$  is (total) recursive (since  $C(x) = \phi(x)$  on  $A$ ), and takes arbitrarily large values, Theorem 2.3.1. Also, by definition of  $\psi$ , we have  $C(\psi(m)) \geq m$ . On the other hand,  $C(\psi(m)) \leq C_\psi(\psi(m)) + c_\psi$  by definition of  $C$ , and obviously  $C_\psi(\psi(m)) \leq l(m)$ . Hence,  $m \leq \log m$  up to a constant independent of  $m$ , which is false from some  $m$  onwards.  $\square$

That was the bad news; the good news is that we can approximate  $C(x)$ .

**Theorem 2.3.3** *There is a total recursive function  $\phi(t, x)$ , monotonic decreasing in  $t$ , such that  $\lim_{t \rightarrow \infty} \phi(t, x) = C(x)$ .*

**Proof.** We define  $\phi(t, x)$  as follows: For each  $x$ , we know that the shortest program for  $x$  has length at most  $l(x) + c$ , Theorem 2.1.2. Run the reference Turing machine  $U$  in the proof of Theorem 2.1.1 for  $t$  steps on each program  $p$  of length at most  $l(x) + c$ . If for any such input  $p$  the computation halts with output  $x$ , then define the value of  $\phi(t, x)$  as the length of the shortest such  $p$ , otherwise equal to  $l(x) + c$ . Clearly,  $\phi(t, x)$  is recursive, total, and monotonically nonincreasing with  $t$  (for all  $x$ ,  $\phi(t', x) \leq \phi(t, x)$  if  $t' > t$ ). The limit exists, since for each  $x$  there exists a  $t$  such that  $U$  halts with output  $x$  after computing  $t$  steps starting with input  $p$  with  $l(p) = C(x)$ .  $\square$

One cannot decide, given  $x$  and  $t$ , whether or not  $\phi(t, x) = C(x)$ . Since  $\phi(t, x)$  is nondecreasing and goes to the limit  $C(x)$  for  $t \rightarrow \infty$ , if there were a decision

procedure to test  $\phi(t, x) = C(x)$ , given  $x$  and  $t$ , then we could compute  $C(x)$ . But Theorem 2.3.2 tells us that  $C$  is not recursive.

Let  $g_1, g_2, \dots$  be a sequence of functions. We call  $f$  the *limit* of this sequence if  $f(x) = \lim_{t \rightarrow \infty} g_t(x)$  for all  $x$ . The limit is *recursively uniform* if for every rational  $\epsilon > 0$  there exists a  $t(\epsilon)$ , where  $t$  is a total recursive function, such that  $|f(x) - g_{t(\epsilon)}(x)| \leq \epsilon$ , for all  $x$ . Let the sequence of one-argument functions  $\psi_1, \psi_2, \dots$  be defined by  $\psi_t(x) = \phi(t, x)$ , for each  $t$  for all  $x$ . Clearly,  $C$  is the limit of the sequence of  $\psi$ 's. However, by Theorem 2.3.2, the limit is not recursively uniform. In fact, by the Halting Problem in Section 1.7, for each  $\epsilon > 0$  and  $t > 0$  there exist infinitely many  $x$  such that  $|C(x) - \psi_t(x)| > \epsilon$ . This means that for each  $\epsilon > 0$ , for each  $t$  there are many  $x$ 's such that our estimate  $\phi(t, x)$  overestimates  $C(x)$  by an *error* of at least  $\epsilon$ .

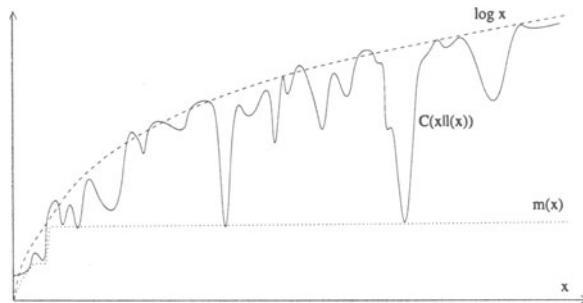
We describe some other characteristics of the function  $C$ .

**Continuity** The function  $C$  is continuous in the sense that there is a constant  $c$  such that  $|C(x) - C(x \pm h)| \leq 2l(h) + c$  for all  $x$  and  $h$ . (Hint: given a program that computes  $x$  we can change it into another program that adds (or subtracts)  $h$  from the output.)

**Logarithmic** The function  $C(x)$  mostly “hugs”  $\log x$ . It is bounded above by  $\log x + c$ , Theorem 2.1.2, page 100. On the other hand, by Theorem 2.2.1, page 109, for each constant  $k$ , the number of  $x$  of length  $n$  (about  $\log x$ ) such that  $C(x) < \log x - k$  is at most  $2^{n-k}$ .

**Fluctuation** The function  $C(x)$  fluctuates rapidly. Namely, for each  $x$  there exist two integers  $x_1, x_2$  within distance  $\sqrt{x}$  of  $x$  (that is,  $|x - x_i| \leq \sqrt{x}$  for  $i = 1, 2$ ) such that  $C(x_1) \geq l(x)/2 - c$  and  $C(x_2) \leq l(x)/2 + c$ . (Hint: change the low-order half of the bits of  $x$  to some incompressible string to obtain  $x_1$ , and change these bits to a very compressible string (such as all zeros) to obtain  $x_2$ .) Therefore, if  $x$  is incompressible with  $C(x) = l(x) - O(1)$ , then there is an  $x_2$  nearby where  $C(x_2)$  equals about  $C(x)/2$ , and if  $x$  is compressible with  $C(x) = o(l(x))$ , then there is an  $x_1$  nearby where  $C(x_1)$  equals about  $l(x)/2$ . These facts imply many fluctuations within small intervals because, for instance,  $C(x)$ ,  $C(x + \log x)$ ,  $C(x + \sqrt{x})$ ,  $C(cx)$ ,  $C(x^2)$ , and  $C(2^x)$ , all have about the same value.

**Long high-complexity runs** For each  $c$  there is a  $d$  such that there are no runs of  $d$  consecutive  $c$ -incompressible numbers. However, conversely, for each  $d$  there is a  $c$  such that there are runs of  $d$  consecutive  $c$ -incompressible numbers. (Hint: for the nonexistence part use numbers  $x$  of the form  $i2^j$  for which  $C(i2^j) \leq l(i) + l(j) + c < l(i2^j) - d$ ; for the existence part use the continuity property and the nearly logarithmic property above.)



**FIGURE 2.2.** The graph of the integer function  $C(x|l(x))$

**Example 2.3.1** It is not difficult to see that Theorems 2.3.1, Item (i), 2.3.2, and 2.3.3, Theorem 2.1.2 and above properties, hold for the length-conditional complexity measure  $C(x|l(x))$ . By the existence of  $n$ -strings, Example 2.2.5, the greatest monotonic lower bound on  $C(x|l(x))$  is a fixed constant, and therefore Items (ii) and (iii) of Theorem 2.3.1 do not hold for this complexity measure. Theorems 2.1.1, 2.2.1 are already proved for  $C(x|l(x))$  in their original versions. Namely, either they were proved for the conditional complexity in general, or the proof goes through as given for the length-conditional complexity. Thus, the general contour of the graph of  $C(x|l(x))$  looks *very roughly* similar to that of  $C(x)$ , except that there are dips below a fixed constant infinitely often, Figure 2.2.

Let us make an estimate of how often the dips occur. Consider the  $n$ -strings of Example 2.2.5. For each integer  $n$  there is an extension of the corresponding binary string with  $n - l(n)$  many 0's such that the resulting string  $x$  has complexity  $C(x|l(x)) \leq c$  for a fixed constant  $c$ . It is easy to see that  $\log x \approx n$ , and that for all  $n' < n$  the corresponding  $x' < x$ . Hence, the number of  $x' < x$  such that  $C(x'|l(x')) \leq c$  is at least  $\log x$ .  $\diamond$

## Exercises

**2.3.1.** [17] Let  $BB$  be the “busy beaver” function defined in Exercise 1.7.19, page 45. Show that  $C(m) = C(n) + O(1)$  for all  $m, n$  with  $m = BB(n)$ . Use this to provide an alternative proof for Theorem 2.3.1, Item (iii).

**2.3.2.** [15] Let  $\phi(t, x)$  be a recursive function and  $\lim_{t \rightarrow \infty} \phi(t, x) = C(x)$ , for all  $x$ . For each  $t$  define  $\psi_t(x) = \phi(t, x)$  for all  $x$ . Then  $C$  is the limit of the sequence of functions  $\psi_1, \psi_2, \dots$ . Show that for each error  $\epsilon$  and all  $t$  there are infinitely many  $x$  such that  $\psi_t(x) - C(x) > \epsilon$ .

*Comments.*  $C(x)$  is the uniform limit of the approximation above if for each  $\epsilon > 0$ , there exists a  $t$  such that for all  $x$ ,  $\psi_t(x) - C(x) \leq \epsilon$ . Item

(a) implies that  $C(x)$  is not the uniform limit of the above sequence of functions.

**2.3.3.** [23] Let  $\phi_1, \phi_2, \dots$  be the effective enumeration of partial recursive functions in Section 1.7. Define the *uniform complexity* of a finite string  $x$  of length  $n$  with respect to  $\phi$  (occurring in the above enumeration) as  $C_\phi(x; n) = \min\{l(p) : \phi(m, p) = x_{1:m} \text{ for all } m \leq n\}$  if such a  $p$  exists, and  $\infty$  otherwise. We can show an Invariance Theorem to the effect that there exists a universal partial recursive function  $\phi_0$  such that for all  $\phi$  there is a constant  $c$  such that for all  $x, n$  we have  $C_{\phi_0}(x; n) \leq C_\phi(x; n) + c$ . We choose a reference universal function  $\phi_0$  and define the *uniform Kolmogorov complexity* as  $C(x; n) = C_{\phi_0}(x; n)$ .

- (a) Show that for all finite binary strings  $x$  we have  $C(x) \geq C(x; l(x)) \geq C(x|l(x))$  up to additive constants independent of  $x$ .
- (b) Prove Theorems 2.1.1 to 2.3.3, with “ $C(x)$ ” replaced by “ $C(x; l(x))$ .”
- (c) Show that in contrast to measure  $C(x|l(x))$ , no constant  $c$  exists such that  $C(x; l(x)) \leq c$  for all  $n$ -strings (Example 2.2.5).
- (d) Show that in contrast to  $C(x|l(x))$ , the uniform complexity is *monotonic in the prefixes*: if  $m \leq n$ , then  $C(x_{1:m}; m) \leq C(x_{1:n}; n)$ , for all  $x$ .
- (e) Show that there exists an infinite binary sequence  $\omega$  and a constant  $c$  such that for infinitely many  $n$ ,  $C(\omega_{1:n}; n) - C(\omega_{1:n}|n) > \log n - c$ .

*Comments.* Item (b) shows that the uniform Kolmogorov complexity satisfies the major properties of the plain Kolmogorov complexity. Items (c) and (d) show that at least two of the objections to the length-conditional measure  $C(x|l(x))$  do not hold for the uniform complexity  $C(x; l(x))$ . Hint for Item (c): this is implied by the proof of Theorem 2.3.1 and Item (a). Item (e) shows as strong a divergence between the measures concerned as one could possibly expect. Source: D.W. Loveland, *Inform. Contr.*, 15(1969), 510–526.

**2.3.4.** • [35] Let  $\omega$  be an infinite binary string. We call  $\omega$  *recursive* if there exists a recursive function  $\phi$  such that  $\phi(i) = \omega_i$  for all  $i > 0$ . Show the following:

- (a) If  $\omega$  is recursive, then there is a constant  $c$  such that for all  $n$

$$\begin{aligned} C(\omega_{1:n}; n) &< c, \\ C(\omega_{1:n}|n) &< c, \\ C(\omega_{1:n}) - C(n) &< c. \end{aligned}$$

This is easy. The converses also hold but are less easy to show. They follow from Items (b), (e), (f),

- (b) For each constant  $c$ , there are only finitely many  $\omega$  such that for all  $n$ ,  $C(\omega_{1:n}; n) \leq c$ , and each such  $\omega$  is recursive.
- (c) For each constant  $c$ , there are only finitely many  $\omega$  such that for infinitely many  $n$ ,  $C(\omega_{1:n}; n) \leq c$ , and each such  $\omega$  is recursive.
- (d) There exists a constant  $c$  such that the set of infinite  $\omega$ , which satisfies  $C(\omega_{1:n}|n) \leq c$  for infinitely many  $n$ , has the cardinality of the continuum.
- (e) For each constant  $c$ , there are only finitely many  $\omega$  such that for all  $n$ ,  $C(\omega_{1:n}|n) \leq c$ , and each such  $\omega$  is recursive.
- (f) For each constant  $c$ , there are only finitely many  $\omega$  with  $C(\omega_{1:n}) \leq C(n) + c$  for all  $n$ , and each such  $\omega$  is recursive.
- (g) For each constant  $c$ , there are only finitely many  $\omega$  with  $C(\omega_{1:n}) \leq l(n) + c$  for all  $n$ , and each such  $\omega$  is recursive.
- (h) There exist nonrecursive  $\omega$  for which there exists a constant  $c$  such that  $C(\omega_{1:n}) \leq C(n) + c$  for infinitely many  $n$ .

*Comments.* Clearly Item (c) implies Item (b). In Item (d) conclude that not all such  $\omega$  are recursive. In particular, the analogue of Item (c) for  $C(\omega_{1:n}|n)$  does not hold. Namely, there exist nonrecursive  $\omega$  for which there exists a constant  $c$  such that for infinitely many  $n$  we have  $C(\omega_{1:n}|n) \leq c$ . Hint for Item (d): exhibit a one-to-one coding of subsets of  $\mathcal{N}$  into the set of infinite binary strings of which infinitely many prefixes are  $n$ -strings—in the sense of Example 2.2.5. Item (e) means that in contrast to the differences between the measures  $C(\cdot; l(\cdot))$  and  $C(\cdot|l(\cdot))$  exposed by the contrast between Items (c) and (d), Item (b) holds also for  $C(\cdot|l(\cdot))$ . Items (f) and (g) show a complexity gap, because  $C(n)$  can be much lower than  $l(n)$ . Hint for Item (h): use Item (d). Source: D.W. Loveland, *Inform. Contr.*, 15(1969), 510–526., for Items (b) through (e) and (h). Loveland attributes Item (e) to A.R. Meyer. The equivalence between bounded length-conditional complexity and bounded uniform complexity for prefixes of infinite strings is stated by A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124. Source of Items (f), (g) is G.J. Chaitin, *Theoret. Comput. Sci.*, 2(1976), 45–48. For the prefix complexity  $K$  introduced in Chapter 3, there are nonrecursive  $\omega$  such that  $K(\omega_{1:n}) \leq K(n) + O(1)$  for all  $n$  by a result of R.M. Solovay in Exercise 3.6.4.

**2.3.5.** [HM35] We want to show in some precise sense that the real line is computationally a fractal. (Actually, one is probably most interested in Item (a), which can be proved easily and elementarily from the following definition.) The required framework is as follows: Each infinite binary sequence  $\omega = \omega_1\omega_2\dots$  corresponds to a real number  $0 \leq 0.\omega < 1$ . Define the normalized complexity  $Cn(\omega) = \lim_{n \rightarrow \infty} C(\omega_{1:n})/n$ . If the limit does not exist, we set  $Cn(\omega)$  to half the sum of the upper and lower limit.

(a) Show that for all real  $\omega$  in  $[0, 1)$ , for every  $\epsilon > 0$  and all real  $r$ ,  $0 \leq r \leq 1$ , there exist real  $\zeta$  in  $[0, 1)$  such that  $|\omega - \zeta| < \epsilon$  and  $Cn(\zeta) = r$ . (For each real  $r$ ,  $0 \leq r \leq 1$ , the set of  $\omega$ 's with  $Cn(\omega) = r$  is dense on the real line  $[0, 1)$ .)

(b) Show that for all real  $\omega$ , all rational  $r$  and  $s$ , we have  $Cn(r\omega + s) = Cn(\omega)$  (both  $\omega$  and  $r\omega + s$  in  $[0, 1)$ ). Similarly, show that  $Cn(f(\omega)) = Cn(\omega)$  for all recursive functions  $f$ .

B. Mandelbrot defined a set to be a *fractal* if its Hausdorff dimension is greater than its topological dimension [B. Mandelbrot, *The Fractal Geometry of Nature*, W.H. Freeman, 1983; for definitions of the dimensions see W. Hurewicz and H. Wallman, *Dimension Theory*, Princeton Univ. Press, 1974].

(c) Show that for any real numbers  $0 \leq a < b \leq 1$ , the Hausdorff dimension of the set  $\{(\omega, Cn(\omega))\} \cap ([0, 1) \times [a, b])$  is  $1 + b$ .

(d) Show that the set  $G = \{(\omega, Cn(\omega)) : \omega \in [0, 1)\}$  has Hausdorff dimension 2 and topological dimension 1. (That is,  $G$  is a fractal.)

*Comments.* Source: J.-Y. Cai and J. Hartmanis, *J. Comput. System Sci.*, 49:3(1994), 605–619.

**2.3.6.** [M34] To investigate repeating patterns in the graph of  $C(x)$  we define the notion of a “shape.” Any function from the integers to the integers is a *shape*. A shape  $f$  matches the graph of  $C$  at  $j$  with *span*  $c$  iff for all  $x$  with  $j - c \leq x \leq j + c$  we have  $C(x) = C(j) + f(x - j)$ .

(a) Show that any matching shape has  $f(0) = 0$ . Thus, a shape is a template of which we align the center  $f(0)$  with  $j$  to see to what extent it matches  $C$ 's graph around the point of interest. We wish to investigate shapes that can be made to match arbitrarily far in each direction.

(b) A shape  $f$  is *recurrent* iff for all  $c$  there exists a  $j$  such that  $f$  matches the graph of  $C$  at  $j$  with span  $c$ . Show that there exists a recurrent shape.

(c) Show that there exists a constant  $c$  such that there are no runs  $C(n) = C(n + 1) = \dots = C(n + c)$  for any  $n$ .

(d) Prove that no recurrent shape is a recursive function.

*Comments.* Hints: for Item (b) use König's Infinity Lemma. Item (c) means the graph of  $C$  has no arbitrarily long flat spots. For Item (c), prove for sufficiently large  $c$  that for each integer  $i$ , for all  $n$  with  $C(n) = i$  the run  $C(n), C(n + 1), \dots, C(n + c)$  contains an element less than  $i$ . For Item (d) use a case analysis, and use in one case the proof of Item (c) and in the other cases the Recursion Theorem, Exercises 1.7.20, page 46. Source: H.P. Katseff and M. Sipser, *Theoret. Comput. Sci.*, 15(1981), 291–309.

## 2.4 Random Finite Sequences

---

One can consider those objects as nonrandom in which one can find sufficiently many regularities. In other words, we would like to identify “incompressibility” with “randomness”. This is proper if the sequences that are incompressible can be shown to possess the various properties of randomness (stochasticity) known from the theory of probability. That this is possible is the substance of the celebrated theory developed by the Swedish mathematician Per Martin-Löf.

There are many properties known that probability theory attributes to random objects. To give an example, consider sequences of  $n$  tosses with a fair coin. Each sequence of  $n$  zeros and ones is equiprobable as an outcome: its probability is  $2^{-n}$ . If such a sequence is to be random in the sense of a proposed new definition, then the number of ones in  $x$  should be near to  $n/2$ , the number of occurrences of blocks “00” should be close to  $n/4$ , and so on.

It is not difficult to show that each such single property separately holds for all incompressible binary strings. But we want to demonstrate that incompressibility implies all conceivable effectively testable properties of randomness (both the known ones and the as yet unknown ones). This way, the various theorems in probability theory about random sequences carry over automatically to incompressible sequences.

In the case of finite strings we cannot hope to distinguish sharply between random and nonrandom strings. For instance, considering the set of binary strings of a fixed length, it would not be natural to fix an  $m$  and call a string with  $m$  zeros random and a string with  $m + 1$  zeros nonrandom.

Let us borrow some ideas from statistics. We are given a certain sample space  $S$  with an associated distribution  $P$ . Given an element  $x$  of the sample space, we want to test the hypothesis “ $x$  is a typical outcome.” Practically speaking, the property of being typical is the property of belonging to any reasonable majority. In choosing an object at random, we have confidence that this object will fall precisely in the intersection of all such majorities. The latter condition we identify with  $x$  being random.

To ascertain whether a given element of the sample space belongs to a particular reasonable majority we introduce the notion of a test. Generally, a test is given by a prescription that, for every level of significance  $\epsilon$ , tells us for what elements  $x$  of  $S$  the hypothesis “ $x$  belongs to majority  $M$  in  $S$ ” should be rejected, where  $\epsilon = 1 - P(M)$ . Taking  $\epsilon = 2^{-m}$ ,  $m = 1, 2, \dots$ , we achieve this by saying that we have a description of the set  $V \subseteq \mathcal{N} \times S$  of nested *critical regions*

$$\begin{aligned} V_m &= \{x : (m, x) \in V\}, \\ V_m &\supseteq V_{m+1}, \quad m = 1, 2, \dots, \end{aligned}$$

while the condition that  $V_m$  be a critical region on the *significance level*  $\epsilon = 2^{-m}$  amounts to requiring, for all  $n$ ,

$$\sum_x \{P(x|l(x) = n) : x \in V_m\} \leq \epsilon.$$

The complement of a critical region  $V_m$  is called the  $(1 - \epsilon)$  *confidence interval*. If  $x \in V_m$ , then the hypothesis “ $x$  belongs to majority  $M$ ,” and therefore the stronger hypothesis “ $x$  is random,” is rejected with significance level  $\epsilon$ . We can say that  $x$  fails the test at the level of critical region  $V_m$ .

**Example 2.4.1** A string  $x_1x_2\dots x_n$  with many initial zeros is not very random. We can test this aspect as follows. The special test  $V$  has critical regions  $V_1, V_2, \dots$ . Consider  $x = 0.x_1x_2\dots x_n$  as a rational number, and each critical region as a half-open interval  $V_m = [0, 2^{-m})$  in  $[0, 1)$ ,  $m = 1, 2, \dots$ . Then the subsequent critical regions test the hypothesis “ $x$  is random” by considering the subsequent digits in the binary expansion of  $x$ . We reject the hypothesis on the significance level  $\epsilon = 2^{-m}$  provided  $x_1 = x_2 = \dots = x_m = 0$ , Figure 2.3.  $\diamond$

**Example 2.4.2** Another test for randomness of finite binary strings rejects when the relative frequency of ones differs too much from  $\frac{1}{2}$ . This particular test can be implemented by rejecting the hypothesis of randomness of  $x = x_1x_2\dots x_n$  at level  $\epsilon = 2^{-m}$  provided  $|2f_n - n| > g(n, m)$ , where  $f_n = \sum_{i=1}^n x_i$ , and  $g(n, m)$  is the least number determined by the requirement that the number of binary strings  $x$  of length  $n$  for which this inequality holds is at most  $2^{n-m}$ .  $\diamond$

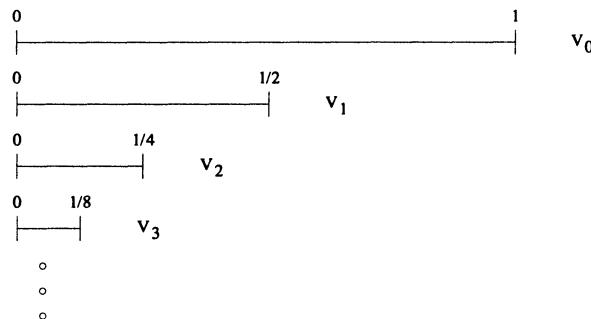


FIGURE 2.3. Test of Example 2.4.1

### 2.4.1 Randomness Tests

In practice, statistical tests are *effective* prescriptions such that we can compute, at each level of significance, for what strings the associated hypothesis should be rejected. It would be hard to imagine what use it would be in statistics to have tests that are not effective in the sense of computability theory (Section 1.7).

**Definition 2.4.1** Let  $P$  be a recursive probability distribution on sample space  $\mathcal{N}$ . A total function  $\delta : \mathcal{N} \rightarrow \mathcal{N}$  is a *P-test* (Martin-Löf test for randomness) if

1.  $\delta$  is enumerable (the set  $V = \{(m, x) : \delta(x) \geq m\}$  is recursively enumerable); and
2.  $\sum\{P(x|l(x) = n) : \delta(x) \geq m\} \leq 2^{-m}$ , for all  $n$ .

The critical regions associated with the common statistical tests are present in the form of the sequence  $V_1 \supseteq V_2 \supseteq \dots$ , where  $V_m = \{x : \delta(x) \geq m\}$ , for  $m \geq 1$ . Nesting is assured since  $\delta(x) \geq m + 1$  implies  $\delta(x) \geq m$ . Each set  $V_m$  is recursively enumerable because of Item 1.

Consider the important case of the uniform distribution, defined by  $L(x) = 2^{-2l(x)-1}$ . The restriction of  $L$  to strings of length  $n$  is defined by  $L_n(x) = 2^{-n}$  for  $l(x) = n$  and 0 otherwise. (By definition,  $L_n(x) = L(x|l(x) = n)$ .) Then Item 2 can be rewritten as  $\sum_{x \in V_m} L_n(x) \leq 2^{-m}$ , which is the same as

$$d(\{x : l(x) = n, x \in V_m\}) \leq 2^{n-m}.$$

In this case we often speak simply of a *test*, with the uniform distribution  $L$  understood.

In statistical tests membership of  $(m, x)$  in  $V$  can usually be determined in time polynomial in  $l(m) + l(x)$ .

**Example 2.4.3** The previous test examples can be rephrased in terms of Martin-Löf tests. Let us try a more subtle example. A real number such that all bits in odd positions in its binary representation are 1's is not random with respect to the uniform distribution. To show this we need a test that detects sequences of the form  $x = 1x_21x_41x_61x_8\dots$ . Define a test  $\delta$  by

$$\delta(x) = \max\{i : x_1 = x_3 = \dots = x_{2i-1} = 1\},$$

and  $\delta(x) = 0$  if  $x_1 = 0$ . For example:  $\delta(01111) = 0$ ;  $\delta(10011) = 1$ ;  $\delta(11011) = 1$ ;  $\delta(10100) = 2$ ;  $\delta(11111) = 3$ . To show that  $\delta$  is a test we have to show that  $\delta$  satisfies the definition of a test. Clearly,  $\delta$  is enumerable (even recursive). If  $\delta(x) \geq m$  where  $l(x) = n \geq 2m$ , then there are  $2^{m-1}$  possibilities for the  $(2m - 1)$ -length prefix of  $x$ , and  $2^{n-(2m-1)}$  possibilities for the remainder of  $x$ . Therefore,  $d\{x : \delta(x) \geq m, l(x) = n\} \leq 2^{n-m}$ .  $\diamond$

**Definition 2.4.2** A universal Martin-Löf test for randomness with respect to distribution  $P$ , a *universal P-test* for short, is a test  $\delta_0(\cdot|P)$  such that for each  $P$ -test  $\delta$ , there is a constant  $c$  such that for all  $x$  we have  $\delta_0(x|P) \geq \delta(x) - c$ .

We say that  $\delta_0(\cdot|P)$  (additively) majorizes  $\delta$ . Intuitively,  $\delta_0(\cdot|P)$  constitutes a test for randomness that incorporates all particular tests  $\delta$  in a single test. No test for randomness  $\delta$  other than  $\delta_0(\cdot|P)$  can discover more than a constant amount more deficiency of randomness in any string  $x$ . In terms of critical regions, a universal test is a test such that if a binary sequence is random with respect to that test, then it is random with respect to any conceivable test, neglecting a change in significance level. Namely, with  $\delta_0(\cdot|P)$  a universal  $P$ -test, let  $U = \{(m, x) : \delta_0(x|P) \geq m\}$ , and for any test  $\delta$ , let  $V = \{(m, x) : \delta(x) \geq m\}$ . Then, defining the associated critical zones as before, we find

$$V_{m+c} \subseteq U_m, \quad m = 1, 2, \dots,$$

where  $c$  is a constant (dependent only on  $U$  and  $V$ ).

It is a major result that there exists a universal  $P$ -test. The proof goes by first showing that the set of all tests is enumerable. This involves the first example of a type of construction we shall use over and over again in different contexts in Chapters 2, 3, and 4. The idea is as follows:

**Lemma 2.4.1** *We can effectively enumerate all  $P$ -tests.*

**Proof.** We start with the standard enumeration  $\phi_1, \phi_2, \dots$  of partial recursive functions from  $\mathcal{N}$  into  $\mathcal{N} \times \mathcal{N}$ , and turn this into an enumeration  $\delta_1, \delta_2, \dots$  of all and only  $P$ -tests. The list  $\phi_1, \phi_2, \dots$  enumerates all and only recursively enumerable sets of pairs of integers as  $\{\phi_i(x) : x \geq 1\}$  for  $i = 1, 2, \dots$ . In particular, for any  $P$ -test  $\delta$ , the set  $\{(m, x) : \delta(x) \geq m\}$  occurs in this list. The *only* thing we have to do is to eliminate those  $\phi_i$  whose range does not correspond to a  $P$ -test.

First, we effectively modify each  $\phi$  (we drop the subscript for convenience) to a function  $\psi$  such that range  $\phi$  equals range  $\psi$ , and  $\psi$  has the special property that if  $\psi(n)$  is defined, then  $\psi(1), \psi(2), \dots, \psi(n-1)$  are also defined. This can be done by dovetailing the computations of  $\phi$  on the different arguments: in the first phase do one step of the computation of  $\phi(1)$ , in the second phase do the second step of the computation of  $\phi(1)$  and the first step of the computation of  $\phi(2)$ . In general, in the  $n$ th phase we execute the  $n_1$ th step of the computation of  $\phi(n_2)$ , for all  $n_1, n_2$  satisfying  $n_1 + n_2 = n$ . We now define  $\psi$  as follows. If the first computation that halts is that of  $\phi(i)$ , then set  $\psi(1) := \phi(i)$ . If the second computation that halts is that of  $\phi(j)$ , then set  $\psi(2) := \phi(j)$ , and so on.

Secondly, use each  $\psi$  to construct a test  $\delta$  by approximation from below. In the algorithm, at each stage of the computation the local variable

array  $\delta(1 : \infty)$  contains the current approximation to the list of function values  $\delta(1), \delta(2), \dots$ . This is doable because the nonzero part of the approximation is always finite.

**Step 1** Initialize  $\delta$  by setting  $\delta(x) := 0$  for all  $x$ ; and set  $i := 0$ . {If the range of  $\psi$  is empty, then this assignment will not be changed in the remainder of the procedure. That is,  $\delta$  stays identically zero and it is trivially a test.}

**Step 2** Set  $i := i + 1$ ; compute  $\psi(i)$  and let its value be  $(x, m)$ .

**Step 3** If  $\delta(x) \geq m$  then go to Step 2 else set  $\delta(x) := m$ .

**Step 4** If  $\sum\{P(y|l(y) = l(x)) : \delta(y) \geq k\} > 2^{-k}$  for some  $k$ ,  $k = 1, \dots, m$  {since  $P$  is a recursive function we can effectively test whether the new value of  $\delta(x)$  violates Definition 2.4.1 on page 129} then set  $\delta(x) := 0$  and terminate {the computation of  $\delta$  is finished} else go to Step 2.

With  $P$  the uniform distribution, for  $i = 1$  the conditional in Step 4 simplifies to  $m > l(x)$ . In case the range of  $\psi$  is already a test, then the algorithm never finishes but forever approximates  $\delta$  from below. If  $\psi$  diverges for some argument then the computation goes on forever and does not change  $\delta$  any more. The resulting  $\delta$  is an enumerable test. If the range of  $\psi$  is not a test, then at some point the conditional in Step 4 is violated and the approximation of  $\delta$  terminates. The resulting  $\delta$  is a test, even a recursive one. Executing this procedure on all functions in the list  $\phi_1, \phi_2, \dots$ , we obtain an effective enumeration  $\delta_1, \delta_2, \dots$  of all  $P$ -tests (and only  $P$ -tests). We are now in the position to define a universal  $P$ -test.  $\square$

**Theorem 2.4.1** Let  $\delta_1, \delta_2, \dots$  be an enumeration of the above  $P$ -tests. Then  $\delta_0(x|P) = \max\{\delta_y(x) - y : y \geq 1\}$  is a universal  $P$ -test.

**Proof.** Note first that  $\delta_0(\cdot|P)$  is a total function on  $\mathcal{N}$  because of Item 2 in Definition 2.4.1, page 129.

(1) The enumeration  $\delta_1, \delta_2, \dots$  in Lemma 2.4.1 yields an enumeration of recursively enumerable sets:

$$\{(m, x) : \delta_1(x) \geq m\}, \{(m, x) : \delta_2(x) \geq m\}, \dots$$

Therefore,  $V = \{(m, x) : \delta_0(x|P) \geq m\}$  is recursively enumerable.

(2) Let us verify that the critical regions are small enough: for each  $n$ ,

$$\begin{aligned} & \sum_{l(x)=n} \{P(x|l(x)=n) : \delta_0(x|P) \geq m\} \\ & \leq \sum_{y=1}^{\infty} \sum_{l(x)=n} \{P(x|l(x)=n) : \delta_y(x) \geq m+y\} \\ & \leq \sum_{y=1}^{\infty} 2^{-m-y} = 2^{-m}. \end{aligned}$$

(3) By its definition,  $\delta_0(\cdot|P)$  majorizes each  $\delta$  additively. Hence, it is universal.  $\square$

By definition of  $\delta_0(\cdot|P)$  as a universal  $P$ -test, any particular  $P$ -test  $\delta$  can discover at most a constant amount more regularity in a sequence  $x$  than does  $\delta_0(\cdot|P)$ , in the sense that for each  $\delta_y$  we have  $\delta_y(x) \leq \delta_0(x|P) + y$  for all  $x$ .

For any two universal  $P$ -tests  $\delta_0(\cdot|P)$  and  $\delta'_0(\cdot|P)$ , there is a constant  $c \geq 0$  such that for all  $x$  we have  $|\delta_0(x|P) - \delta'_0(x|P)| \leq c$ .

#### 2.4.2 Explicit Universal Randomness Test

We started out with the objective to establish in what sense incompressible strings may be called random. In Section 2.2.1 we considered the randomness deficiency  $\delta(x|A)$  of a string  $x$  relative to a finite set  $A$ . With  $A$  the set of strings of length  $n$  and  $x \in A$  we find  $\delta(x|A) = \delta(x|n) = n - C(x|n)$ .

**Theorem 2.4.2** *The function  $\delta_0(x|L) = l(x) - C(x|l(x)) - 1$  is a universal  $L$ -test with  $L$  the uniform distribution.*

**Proof.** (1) We first show that  $f(x) = \delta_0(x|L)$  is a test with respect to the uniform distribution. The set  $\{(m, x) : f(x) \geq m\}$  is recursively enumerable by Theorem 2.3.3.

(2) We verify the condition on the critical regions. Since the number of  $x$ 's with  $C(x|l(x)) \leq l(x) - m - 1$  cannot exceed the number of programs of length at most  $l(x) - m - 1$ , we have  $d(\{x : f(x) \geq m\}) \leq 2^{l(x)-m} - 1$ .

(3) We show that for each test  $\delta$ , there is a constant  $c$  such that  $f(x) \geq \delta(x) - c$ . The main idea is to bound  $C(x|l(x))$  by exhibiting a description of  $x$ , given  $l(x)$ . Fix  $x$ . Let the set  $A$  be defined as

$$A = \{z : \delta(z) \geq \delta(x), l(z) = l(x)\}.$$

We have defined  $A$  such that  $x \in A$  and  $d(A) \leq 2^{l(x)-\delta(x)}$ . Let  $\delta = \delta_y$  in the standard enumeration  $\delta_1, \delta_2, \dots$  of tests. Given  $y$ ,  $l(x)$ , and  $\delta(x)$ ,

we have an algorithm to enumerate all elements of  $A$ . Together with the index  $j$  of  $x$  in enumeration order of  $A$ , this suffices to find  $x$ . We pad the standard binary representation of  $j$  with nonsignificant zeros to a string  $s = 00\dots0j$  of length  $l(x) - \delta(x)$ . This is possible since  $l(s) \geq l(d(A))$ . The purpose of changing  $j$  to  $s$  is that now the number  $\delta(x)$  can be deduced from  $l(s)$  and  $l(x)$ . In particular, there is a Turing machine that computes  $x$  from input  $\bar{y}s$ , when  $l(x)$  is given for free. Consequently, by Theorem 2.1.1,  $C(x|l(x)) \leq l(x) - \delta(x) + 2l(y) + 1$ . Since  $y$  is a constant depending only on  $\delta$ , we can set  $c = 2l(y) + 2$ .  $\square$

**Definition 2.4.3** Let us fix  $\delta_0(x|L) = l(x) - C(x|l(x)) - 1$  as the *reference universal test* with respect to the *uniform distribution L*. A string  $x$  is called *c-random*, if  $\delta_0(x|L) \leq c$ .

Randomness of a string is related to its incompressibility. It is easy to see that

$$C(x|l(x)) \leq C(x) \leq C(x|l(x)) + 2C(l(x) - C(x|l(x))),$$

up to fixed additive constants. (We can reconstruct  $l(x)$  from the length of a shortest program  $p$  for  $x$  and the quantity  $l(x) - l(p)$ .) This makes  $C(x)$  and  $C(x|l(x))$  about equal for the special case of  $x$  being incompressible. (However, for compressible  $x$ , like  $x = 0^n$ , the difference between  $C(x)$  and  $C(x|n)$  can rise to logarithmic in  $n$ .) Together with Theorem 2.4.2, this provides the relation between the outcome of the reference universal  $L$ -test and incompressibility. Fix a constant  $c$ . If string  $x$  is  $c$ -incompressible, then  $\delta_0(x|L) \leq c'$ , where  $c'$  is a constant depending on  $c$  but not on  $x$ . Similarly, if  $\delta_0(x|L) \leq c$ , then  $x$  is  $c'$ -incompressible, where  $c'$  is a constant depending on  $c$  but not on  $x$ . Roughly,  $x$  is *random*, or *incompressible*, if  $l(x) - C(x|l(x))$  is small with respect to  $l(x)$ .

**Example 2.4.4** It is possible to directly demonstrate a property of random binary strings  $x = x_1x_2\dots x_n$  related to Example 2.4.2: the number of ones,  $f_n = x_1 + \dots + x_n$ , must satisfy  $|2f_n - n| = O(\sqrt{n})$ . Assume that  $x$  is a binary string of length  $n$  that is random in the sense of Martin-Löf. Then, by Theorem 2.4.2,  $x$  is also  $c$ -incompressible for some fixed constant  $c$ . Let  $f_n = k$ . The number of strings satisfying this equality is  $\binom{n}{k}$ . By simply giving the index of  $x$  in the lexicographic order of such strings, together with  $n$ , and the excess of ones,  $d = |2k - n|$ , we can give a description of  $x$ . Therefore, using a short self-delimiting program for  $d$ :

$$C(x|n) \leq \log \binom{n}{k} + C(d) + 2l(C(d)).$$

For  $x$  given  $n$  to be  $c$ -incompressible for a constant  $c$ , we need  $C(x|n) \geq n - c$ . Then,

$$n - \log \binom{n}{k} - C(d) - 2l(C(d)) \leq c,$$

which can only be satisfied provided  $d = O(\sqrt{n})$  (estimate the binomial coefficient by Stirling's formula, page 17). Curiously, if  $d$  given  $n$  is easily describable (for example  $d = 0$  or  $d = \sqrt{n}$ ), then  $x$  given  $n$  is not random since it is not  $c$ -incompressible.  $\diamond$

Randomness in the sense of Martin-Löf means randomness insofar as it can be effectively certified. In other words, it is a negative definition. We look at objects through a special filter, which highlights some features but obscures others. We can perceive some qualities of nonrandomness through the limited sight of effective tests. Everything else we then call by definition random. This is a matter of a pragmatic, expedient, approach. It has no bearing on the deeper question about what properties real randomness, physically or mathematically, should have. It just tells us that an object is random as far as we will ever be able to tell, *in principle* and not just as a matter of *practicality*.

## Exercises

---

**2.4.1.** [20] For a binary string  $x$  of length  $n$ , let  $f(x)$  be the number of ones in  $x$ . Show that  $\delta(x) = \log(2n^{-1/2}|f(x) - n/2|)$  is a  $P$ -test with  $P$  the uniform measure.

*Comments.* Use Markov's inequality to derive that for each positive  $\lambda$ , the probability of  $2n^{-1/2}|f(x) - n/2| > \lambda$  is at most  $1/\lambda$ . Source: T.M. Cover, P. Gács, and R.M. Gray, *Ann. Probab.*, 17(1989), 840–865.

**2.4.2.** [23] Let  $x_1x_2\dots x_n$  be a random sequence with  $C(x|n) \geq n$ .

(a) Use a Martin-Löf test to show that  $x_10x_20\dots 0x_n$  is not random with respect to the uniform distribution.

(b) Use a Martin-Löf test to show that the ternary sequence  $y_1y_2\dots y_n$  with  $y_1 = x_n + x_1$  and  $y_i = x_{i-1} + x_i$  for  $1 < i \leq n$  is not random with respect to the uniform distribution.

*Comments.* Hint: in Item (b) in the  $y$ -string the blocks “02” and “20” do not occur. Extend the definition of random sequences from binary to ternary. Source: R. von Mises, *Probability, Statistics and Truth*, Dover, 1981.

**2.4.3.** [35] Let  $x$  be a finite binary sequence of length  $n$  with  $f_j = x_1 + x_2 + \dots + x_j$  for  $1 \leq j \leq n$ . Show that there exists a constant  $c > 0$  such that for all  $m \in \mathcal{N}$ , all  $\epsilon > 0$ , and all  $x$ ,

$$C(x|n, f_n) > \log \binom{n}{f_n} - \log(m\epsilon^4) + c$$

implies

$$\max_{m \leq j \leq n} \left| \frac{f_j}{j} - \frac{f_n}{n} \right| < \epsilon.$$

*Comments.* This result is called *Fine's Theorem*. This is an instance of the general principle that high probability of a computable property translates into the fact that high complexity implies that property. (For infinite sequences this principle is put in a precise and rigorous form in Theorem 2.5.5 on page 146.) Fine's Theorem shows that for finite binary sequences with high Kolmogorov complexity, given the length and the number of ones, the fluctuations of the relative frequencies in the initial segments is small. Since we deal with finite sequences this is called *apparent convergence*. Since virtually all finite binary strings have high complexity (Theorem 2.1.3), this "explains" why in a typical sequence produced by random coin throws the relative frequencies appear to converge or stabilize. Apparent convergence occurs because of, and not in spite of, the high irregularity (randomness or complexity) of a data sequence. Conversely, the failure of convergence forces the complexity to be less than maximal. Source: T.L. Fine, *IEEE Trans. Inform. Theory*, IT-16(1970), 251–257; also R. Heim, *IEEE Trans. Inform. Theory*, IT-25(1979), 557–566.

**2.4.4.** [36] (a) Consider a finite sequence of zeros and ones generated by independent tosses of a coin with probability  $p$  ( $0 < p < 1$ ) for "1." Let  $x = x_1 x_2 \dots x_n$  be a sequence of outcomes of length  $n$ , and let  $f_n = x_1 + x_2 + \dots + x_n$ . The probability of such an  $x$  is  $p^{f_n} (1-p)^{n-f_n}$ . If  $p$  is a recursive number, then the methods in this section can be used to obtain a proper definition of finite *Bernoulli sequences*, sequences that are random for this distribution. There is, however, no reason to suppose that in physical coins  $p$  is a recursive real. This prompts another approach where we disregard the actual probability distribution, but follow more closely the combinatorial spirit of Kolmogorov complexity. Define a finite Bernoulli sequence as a binary sequence  $x$  of length  $n$  whose only regularities are given by  $f_n$  and  $n$ . That is,  $x$  is a Bernoulli sequence iff  $C(x|n, f_n) = \log \binom{n}{f_n}$  up to a constant independent of  $x$ . Define a *Bernoulli test* as a test, with the condition that the number of sequences with  $f_n$  ones and  $n - f_n$  zeros in  $V_m$  is  $\leq 2^{-m} \binom{n}{f_n}$  for all  $m$ ,  $n$ , and  $f_n$ . Show that there exists a universal Bernoulli test  $\delta_0$ . Finite Bernoulli sequences are those sequences  $x$  such that  $\delta_0(x)$  is low. Show that up to a constant independent of  $x$

$$\delta_0(x) = \log \binom{n}{f_n} - C(x|n, f_n),$$

for all  $x$  (with  $n$  and  $f_n$  as above).

(b) We continue Item (a). In the current interpretation of probability the relative frequency of an event in a large number of experiments should not only be close to the probability, but there is an obscure secondary stipulation. If the probability of success is very small, we should be practically sure that the event should not occur in a single trial [A.N. Kolmogorov, *Grundbegriffe der Wahrscheinlichkeitsrechnung*, Berlin, 1933; English translation: Chelsea, 1956]. If  $x$  is a Bernoulli sequence (result of  $n$  independent coin tosses) with a very low relative success frequency  $f_n/n$  (the coin is heavily biased), then, necessarily,  $x_1 = 0$ . That is, the assumption that “1” occurs as the very first element implies substantial regularity of the overall sequence.

Show that there is a constant  $c$  such that  $\delta_0(x) \leq \log n/f_n - c$  implies  $x_1 = 0$ .

*Comments.* Hint for Item (b): construct the test that rejects at level  $2^{-m}$  when  $x_1 = 1$  and  $f_n \leq n2^{-m}$ . Show that this is a Bernoulli test. Compare this test with  $\delta_0$  in Item (a). For sequential Bernoulli tests for infinite sequences see Exercise 2.5.17 Source: P. Martin-Löf, *Inform. Contr.*, 9(1966), 602–619.

## 2.5 \*Random Infinite Sequences

### 2.5.1 Complexity Oscillations

Consider the question of how  $C$  behaves in terms of increasingly long initial segments of a fixed infinite binary sequence (or string)  $\omega$ . Is it monotone in the sense that  $C(\omega_{1:m}) \leq C(\omega_{1:n})$ , or  $C(\omega_{1:m}|m) \leq C(\omega_{1:n}|n)$ , for all infinite binary sequences  $\omega$  and all  $m \leq n$ ? We have already seen that the answer is negative in both cases. A similar effect arises when we try to use Kolmogorov complexity to solve the problem of finding a proper definition of random infinite sequences (*collectives*) according to the task already set by von Mises in 1919, Section 1.9.

It is seductive to call an infinite binary sequence  $\omega$  random if there is a constant  $c$  such that for all  $n$ , the  $n$ -length prefix  $\omega_{1:n}$  has  $C(\omega_{1:n}) \geq n - c$ . However, such sequences do not exist. We shall show that for high-complexity sequences, with  $C(\omega_{1:n}) \geq n - \log n - 2\log\log n$  for all  $n$ , this results in so-called *complexity oscillations*, where for every  $\epsilon > 0$

$$\frac{n - C(\omega_{1:n})}{\log n}$$

oscillates between 0 and  $1 + \epsilon$  for large enough  $n$ . First, we show that the  $C$  complexity of prefixes of each infinite binary sequence drops infinitely often unboundedly far below its own length.

**Theorem 2.5.1** *Let  $f : \mathcal{N}^+ \rightarrow \mathcal{N}$  be a total recursive function satisfying  $\sum_{n=1}^{\infty} 2^{-f(n)} = \infty$  (like  $f(n) = \log n$ ). Then for all infinite binary sequences  $\omega$ , we have  $C(\omega_{1:n}|n) \leq n - f(n)$  infinitely often.*

**Proof.** In order to get rid of some  $O(1)$  term in the final argument, we first change  $f$  into something that gets arbitrarily larger yet still diverges in the same way as  $f$ . Define

$$F(n) = \left\lfloor \log \left( \sum_{i=1}^n 2^{-f(i)} \right) \right\rfloor.$$

Note that  $\sum_{F(n)=m} 2^{-f(n)} \geq 2^m - 1$ . Now let  $g$ , also total recursive, be defined as  $g(n) = f(n) + F(n)$ . It follows that

$$\begin{aligned} \sum_{n=1}^{\infty} 2^{-g(n)} &= \sum_{m=1}^{\infty} \sum_{F(n)=m} 2^{-f(n)-m} \\ &\geq \sum_{m=1}^{\infty} 2^{-m}(2^m - 1) = \infty. \end{aligned}$$

Now we come to the real argument. Consider the unit interval  $[0, 1]$  laid out in a circle so that 0 and 1 are identified. The partial sums  $G_n = \sum_{i=1}^n 2^{-g(i)}$  mark off successive intervals  $I_n \equiv [G_{n-1}, G_n) \bmod 1$  on this circle. We exploit the fact that a point on the circle will be contained in many of these intervals.

For each  $x \in \{0, 1\}^*$ , the associated cylinder is the set

$$\Gamma_x = \{\omega \in \{0, 1\}^\infty : \omega_{1:l(x)} = x\}.$$

The geometric interpretation is  $\Gamma_x = [0.x, 0.x + 2^{-l(x)})$ . Let

$$A_n = \{x \in \{0, 1\}^n : \Gamma_x \bigcap I_n \neq \emptyset\}.$$

It follows from the divergence of  $G_n$  that for any  $\omega$  there is an infinite set  $N \subseteq \mathcal{N}$  consisting of the infinitely many  $n$  such that prefixes  $\omega_{1:n}$  belong to  $A_n$ . Describing a prefix  $\omega_{1:n} \in A_n$  by its index in the set, we have

$$\begin{aligned} C(\omega_{1:n}|n) &\leq \log |A_n| + O(1) \leq \log \frac{G_n - G_{n-1}}{2^{-n}} + O(1) \\ &= n - g(n) + O(1) \leq n - f(n). \end{aligned}$$

□

**Corollary 2.5.1**  $C(\omega_{1:n}) \leq n - f(n)$  infinitely often provided  $C(n|n - f(n)) = O(1)$  (like  $f(n) = \log n$ ).

**Proof.** This is a slightly stronger statement than Theorem 2.5.1. Let  $p$  be a description of  $\omega_{1:n}$ , given  $n$ , of length  $C(\omega_{1:n}|n)$ . Let  $f(), g()$  be

as in the proof of Theorem 2.5.1, and let  $q$  be an  $O(1)$ -length program that retrieves  $n$  from  $n - f(n)$ . Then  $l(\bar{q}p) \leq n - f(n)$  since  $l(p) \leq n - g(n)$  and  $g(n) - f(n)$  rises unboundedly. We pad  $\bar{q}p$  to length  $n - f(n)$ , obtaining  $\bar{q}1^{n-f(n)-l(\bar{q}p)-1}0p$ . This is a description for  $\omega_{1:n}$ . Namely, we first determine  $\bar{q}$  to find  $q$ . The length of the total description is  $n - f(n)$ . By assumption,  $q$  computes  $n$  from this. Given  $n$  we can retrieve  $\omega_{1:n}$  from  $p$ .  $\square$

In [P. Martin-Löf, *Z. Wahrsch. Verw. Geb.*, 19(1971), 225–230] Corollary 2.5.1 is stated to hold without the additional condition of  $n$  being retrievable from  $n - f(n)$  by an  $O(1)$  bit program. The proof of this fact is attributed to [P. Martin-Löf, On the oscillation of the complexity of infinite binary sequences (Russian), unpublished, 1965].

There is a simple proof that yields a result only slightly weaker than this. We first prove the result with the particular function  $\log n$  substituted for  $f(n)$  in the statement of the theorem. We then iterate the construction to prove a version of the theorem with a particular function  $g(n)$  substituted for  $f(n)$ . Our result is almost tight since for functions  $h(n)$  that are only slightly larger than  $g(n)$  the sum  $\sum 2^{-h(n)} < \infty$ . Moreover, the proof is explicit in that we exhibit  $g(n)$ .

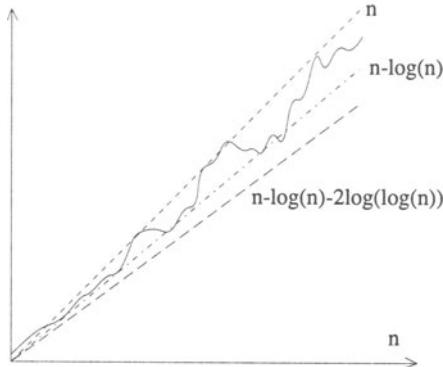
Let  $\omega$  be any infinite binary sequence, and  $\omega_{1:m}$  any  $m$ -length prefix of  $\omega$ . If  $\omega_{1:m}$  is the  $n$ th binary string in the lexicographical order  $0, 1, 00, \dots$ , that is,  $n = \omega_{1:m}$ ,  $m = l(n)$ , then  $C(\omega_{1:n}) = C(\omega_{m+1:n}) + c$ , with  $c$  a constant independent of  $n$  or  $m$ . Namely, with  $O(1)$  additional bits of information, we can trivially reconstruct the  $n$ th binary string  $\omega_{1:m}$  from the length  $n - l(n)$  of  $\omega_{m+1:n}$ . By Theorem 2.1.2 on page 100 we find that  $C(\omega_{m+1:n}) \leq n - l(n) + c$  for some constant  $c$  independent of  $n$ , whence the claimed result with  $f(n) = \log n$  follows.

It is easy to see that we get a stronger result by iteration of the above argument. There are infinitely many  $n$  such that the initial segment  $y = \omega_{1:n}$  can be divided as  $y = y_1y_2 \dots y_m$ , where  $l(y_1) = 2$ ,  $y_1 = l(y_2)$ ,  $y_2 = l(y_3)$ ,  $\dots$ ,  $y_{m-1} = l(y_m)$ . Use the usual pairing between natural numbers and binary strings. Clearly, given  $y_m$  we can easily compute all of  $\omega_{1:n}$ , by determining  $y_{m-1}$  as the binary representation of  $l(y_m)$ ,  $y_{m-2}$  as the binary representation of  $l(y_{m-1})$ , and so on until we find  $l(y_1) = y_1$ . (If  $\omega_{1:20} = 01001110110000011010$ , then  $y_1 = 01$ , which corresponds to natural number 4, so  $y_2 = 0011$ , which corresponds with the natural number 14, and finally  $y_3 = 10110000011010$ . Hence, given  $y_3$ , we can easily determine all of  $\omega_{1:20}$ .) Then, for infinitely many  $n$ ,

$$C(\omega_{1:n}) \leq m + c,$$

for  $m$  determined by  $n = m + l(m) + l(l(m)) + \dots + 2$ , all terms greater than or equal to 2. If we put  $g(n) = n - m$ , then it can be shown that  $\sum 2^{-g(n)} = \infty$ , but that for only slightly larger functions  $h(n) > g(n)$  the sum converges (for example  $h(n) = g(n) +$  the number of terms in  $g(n)$ ). There is an interesting connection with prefix codes, Section 1.11.1.

Our approach in this proof makes it easy to say something about the frequency of these complexity oscillations. Define a “wave” function  $w$  by  $w(1) = 2$  and  $w(i) = 2^{w(i-1)}$ , then the above argument guarantees that there are at least  $k$



**FIGURE 2.4.** Complexity oscillations of initial segments of infinite high-complexity sequences

values  $n_1, n_2, \dots, n_k$  less than  $n = w(k)$  such that  $C(\omega_{1:n_i}) \leq n_i - g(n_i) + c$  for all  $i = 1, 2, \dots, k$ . Obviously, this can be improved.

In Figure 2.4 we display the complexity oscillations of initial segments of high-complexity sequences as they must look according to Theorems 2.5.1, 2.5.4, 2.5.5. The upper bound on the oscillations,  $C(\omega_{1:n}) = n + O(1)$ , is reached infinitely often for almost all high-complexity sequences. Furthermore, the oscillations of all high-complexity sequences stay above  $n - \log n - 2 \log \log n$ , but dip infinitely often below  $n - \log n$ .

Having shown that the complexity of prefixes of each infinite sequence drops infinitely often unboundedly below the maximum value, we now want to show that Theorem 2.5.1 is optimal. But let's first discuss what this means. Clearly, Theorem 2.5.1 is only nontrivial for very irregular sequences  $x$ . It holds trivially for regular sequences like  $\omega = 0^\infty$ , where the complexity of the initial segments  $\omega_{1:n}$  is about  $\log n$ . We will prove that it is sharp for those  $\omega$  that are maximally random. To make this precise, we must define rigorously what we mean by a random infinite sequence. It is of major significance that in so doing we also succeed in completing in a satisfactory way the program outlined by von Mises.

Due to the complexity oscillations the idea of identifying random infinite sequences with those such that  $C(\omega_{1:n}) \geq n - c$ , for all  $n$ , is trivially infeasible. That is the bad news. In contrast, a similar approach in Section 2.4 for finite binary sequences turned out to work just fine. Its justification was found in Martin-Löf's important insight that to justify any proposed definition of randomness one has to show that the sequences that are random in the stated sense satisfy the several properties of stochasticity we know from the theory of probability. Instead of proving each such property separately, one may be able to show, once and for all, that

the random sequences introduced possess, in an appropriate sense, all possible properties of stochasticity.

The naive execution of the above ideas in classical mathematics is infeasible as shown by the following example: Consider as sample space  $S$  the set of all one-way infinite binary sequences. The cylinder  $\Gamma_x = \{\omega : \omega = x\ldots\}$  consists of all infinite binary sequences starting with the finite binary sequence  $x$ . For instance,  $\Gamma_\epsilon = S$ . The uniform distribution  $\lambda$  on the sample space is defined by  $\lambda(\Gamma_x) = 2^{-l(x)}$ . That is, the probability of an infinite binary sequence  $\omega$  starting with a finite initial segment  $x$  is  $2^{-l(x)}$ . In probability theory it is general practice that if a certain property, such as the Law of Large Numbers, or the Law of the Iterated Logarithm, has been shown to have probability one, then one calls this a *Law of Randomness*. For example, in our sample space the Law of Large Numbers says that  $\lim_{n \rightarrow \infty} (\omega_1 + \cdots + \omega_n)/n = \frac{1}{2}$ . If  $A$  is the set of elements of  $S$  that satisfy the Law of Large Numbers, then it can be shown that  $\lambda(A) = 1$ .

Generalizing this idea for  $S$  with measure  $\mu$ , one may identify *any* set  $B \subseteq S$ , such that  $\mu(B) = 1$ , with a Law of Randomness, namely, “to be an element of  $B$ .” Elements of  $S$  that do not satisfy the law “to be an element of  $B$ ” form a set of measure zero, a *null set*. It is natural to call an element of the sample space “random” if it satisfies all laws of randomness. Now we are in trouble. For each element  $\omega \in S$ , the set  $B_\omega = S - \{\omega\}$  forms a law of randomness. But the intersection of all these sets  $B_\omega$  of probability one is empty. Thus, no sequence would be random if we require that all laws of randomness that “exist” are satisfied by a random sequence.

It turns out that a constructive viewpoint enables us to carry out this program mathematically without such pitfalls. In practice, all laws that are proved in probability theory to hold with probability one are effective in the sense of Section 1.7. A straightforward formalization of this viewpoint is to require a law of probability to be partial recursive in the sense that we can effectively test whether it is violated. This suggests that the set of random infinite sequences should not be defined as the intersection of all sets of measure one, but as the intersection of all sets of measure one with a recursively enumerable complement. The latter intersection is again a set of measure one with a recursively enumerable complement. Hence, there is a single effective law of randomness that can be stated as the property “to satisfy all effective laws of randomness,” and the infinite sequences have this property with probability one.

### 2.5.2 Sequential Randomness Tests

As in Section 2.4, we define a test for randomness. However, this time the test will not be defined on the entire sequence (which is impossible for an effective test and an infinite sequence), but for each finite binary string. The value of the test for an infinite sequence is then defined as the maximum of the values of the test on all prefixes. Since this suggests an effective process of sequential approximations, we call it a sequential test. Below, we need to use notions of continuous sample spaces and measures as treated in Section 1.6.

**Definition 2.5.1** Let  $\mu$  be a recursive probability measure on the sample space  $\{0, 1\}^\infty$ . A total function  $\delta : \{0, 1\}^\infty \rightarrow \mathcal{N} \cup \{\infty\}$  is a *sequential  $\mu$ -test* (sequential Martin-Löf  $\mu$ -test for randomness) if

1.  $\delta(\omega) = \sup_{n \in \mathcal{N}} \{\gamma(\omega_{1:n})\}$ , where  $\gamma : \mathcal{N} \rightarrow \mathcal{N}$  is a total enumerable function ( $V = \{(m, y) : \gamma(y) \geq m\}$  is a recursively enumerable set); and
2.  $\mu\{\omega : \delta(\omega) \geq m\} \leq 2^{-m}$ , for each  $m \geq 0$ .

If  $\mu$  is the uniform measure  $\lambda$ , then we often use simply the term *sequential test*.

We can require  $\gamma$  to be a recursive function without changing the notion of a sequential  $\mu$ -test. By definition, for each enumerable function  $\gamma$  there exists a recursive function  $\phi$  with  $\phi(x, k)$  nondecreasing in  $k$  such that  $\lim_{k \rightarrow \infty} \phi(x, k) = \gamma(x)$ . Define a recursive function  $\gamma'$  by  $\gamma'(\omega_{1:n}) = \phi(\omega_{1:m}, k)$  with  $\langle m, k \rangle = n$ . Then,  $\sup_{n \in \mathcal{N}} \{\gamma'(\omega_{1:n})\} = \sup_{n \in \mathcal{N}} \{\gamma(\omega_{1:n})\}$ .

**Example 2.5.1** Consider  $\{0, 1\}^\infty$  with the uniform measure  $\lambda(x) = 2^{-l(x)}$ . An example of a sequential  $\lambda$ -test is to test whether there are 1's in even positions of  $\omega \in \{0, 1\}^\infty$ . Let

$$\gamma(\omega_{1:n}) = \begin{cases} n/2 & \text{if } \sum_{i=1}^{n/2} \omega_{2i} = 0, \\ 0 & \text{otherwise.} \end{cases}$$

The number of  $x$ 's of length  $n$  such that  $\gamma(x) \geq m$  is at most  $2^{n/2}$  for any  $m \geq 1$ . Therefore,  $\lambda\{\omega : \delta(\omega) \geq m\} \leq 2^{-m}$  for  $m > 0$ . For  $m = 0$ ,  $\lambda\{\omega : \delta(\omega) \geq m\} \leq 2^{-m}$  holds trivially. A sequence  $\omega$  is random with respect to this test if  $\delta(\omega) < \infty$ . Thus, a sequence  $\zeta$  with 0's in all even locations will have  $\delta(\zeta) = \infty$ , and it will fail the test and hence  $\zeta$  is not random with respect to this test. Notice that this is not a very strong test of randomness. For example, a sequence  $\eta = 010^\infty$  will pass  $\delta$  and be considered random with respect to this test. This test only filters out some nonrandom sequences with all 0's at the even locations and cannot detect other kinds of regularities.  $\diamond$

We continue the general theory of sequential testing. If  $\delta(\omega) = \infty$ , then we say that  $\omega$  *fails*  $\delta$ , or that  $\delta$  *rejects*  $\omega$ . Otherwise,  $\omega$  *passes*  $\delta$ . By definition, the set of  $\omega$ 's that is rejected by  $\delta$  has  $\mu$ -measure zero, and, conversely, the set of  $\omega$ 's that pass  $\delta$  has  $\mu$ -measure one.

Suppose that for a test  $\delta$  we have  $\delta(\omega) = m$ . Then there is a prefix  $y$  of  $\omega$ , with  $l(y)$  minimal, such that  $\gamma(y) = m$  for the  $\gamma$  used to define  $\delta$ . Then obviously, *each* infinite sequence  $\zeta$  that starts with  $y$  has  $\delta(\zeta) \geq m$ . The set of such  $\zeta$  is  $\Gamma_y = \{\zeta : \zeta = y\rho, \rho \in \{0, 1\}^\infty\}$ , the cylinder generated by

$y$ . Geometrically speaking,  $\Gamma_y$  can be viewed as the set of all real numbers  $0.y\dots$  corresponding to the half-open interval  $I_y = [0.y, 0.y + 2^{-l(y)})$ . For the uniform measure,  $\lambda(\Gamma_y)$  is equal to  $2^{-l(y)}$ , the common length of  $I_y$ .

In terms of common statistical tests, the critical regions are formed by the nested sequence  $V_1 \supseteq V_2 \supseteq \dots$ , where  $V_m$  is defined as  $V_m = \{\omega : \delta(\omega) \geq m\}$ , for  $m \geq 1$ . We can formulate the definition of  $V_m$  as

$$V_m = \bigcup \{\Gamma_y : (m, y) \in V\}.$$

In geometric terms,  $V_m$  is the union of a set of subintervals of  $[0, 1)$ . Since  $V$  is recursively enumerable, so is the set of intervals whose union is  $V_m$ . For each critical section we have  $\mu(V_m) \leq 2^{-m}$  (in the measure we count overlapping intervals only once).

Now we can reformulate the notion of passing a sequential test  $\delta$  with associated set  $V$ :

$$\delta(\omega) < \infty \text{ iff } \omega \notin \bigcap_{m=1}^{\infty} V_m.$$

**Definition 2.5.2** Let  $\mathbf{V}$  be the set of all sequential  $\mu$ -tests. An infinite binary sequence  $\omega$ , or the binary represented real number  $0.\omega$ , is called  $\mu$ -random if it passes *all* sequential  $\mu$ -tests:

$$\omega \notin \bigcup_{V \in \mathbf{V}} \bigcap_{m=1}^{\infty} V_m.$$

For each sequential  $\mu$ -test  $V$ , we have  $\mu(\bigcap_{m=1}^{\infty} V_m) = 0$ , by Definition 2.5.1. We call  $\bigcap_{m=1}^{\infty} V_m$  a *constructive  $\mu$ -null set*. Since there are only countably infinitely many sequential  $\mu$ -tests  $V$ , it follows from standard measure theory that

$$\mu \left( \bigcup_{V \in \mathbf{V}} \bigcap_{m=1}^{\infty} V_m \right) = 0,$$

and we call the set  $U = \bigcup_{V \in \mathbf{V}} \bigcap_{m=1}^{\infty} V_m$  the *maximal constructive  $\mu$ -null set*.

Similar to Section 2.4, we construct an enumerable function  $\delta_0(\omega|\mu)$ , the universal sequential  $\mu$ -test that incorporates (majorizes) all sequential  $\mu$ -tests  $\delta_1, \delta_2, \dots$  and that corresponds to  $U$ .

**Definition 2.5.3** A *universal sequential  $\mu$ -test*  $f$  is a sequential  $\mu$ -test such that for each sequential  $\mu$ -test  $\delta_i$  there is a constant  $c \geq 0$  such that for all  $\omega \in \{0, 1\}^\infty$ , we have  $f(\omega) \geq \delta_i(\omega) - c$ .

**Theorem 2.5.2** *There is a universal sequential  $\mu$ -test (denoted as  $\delta_0(\cdot|\mu)$ ).*

**Proof.** Start with the standard enumeration  $\phi_1, \phi_2, \dots$  of partial recursive functions from  $\mathcal{N}$  into  $\mathcal{N} \times \mathcal{N}$ . This way we enumerate all recursively enumerable sets of pairs of integers in the form of the ranges of the  $\phi_i$ 's. In particular, we have included *any* recursively enumerable set  $V$  associated with a sequential  $\mu$ -test. The *only* thing we have to do is to eliminate those  $\phi_i$ 's that do not correspond to a sequential  $\mu$ -test.

First, we effectively modify each  $\phi$  (we drop the subscript for convenience) to a function  $\psi$  such that range  $\phi$  equals range  $\psi$ , and  $\psi$  has the special property that if  $\psi(n) < \infty$ , then also  $\psi(1), \psi(2), \dots, \psi(n-1) < \infty$ .

Secondly, use each  $\psi$  to construct a function  $\delta : \{0, 1\}^\infty \rightarrow \mathcal{N}$  by approximation from below. In the algorithm, at each stage of the computation the arrays  $\gamma$  and  $\delta$  contain the current approximation to the function values of  $\gamma$  and  $\delta$ . This is doable because the nonzero part of the approximation is always finite.

**Step 1** Initialize  $\gamma$  and  $\delta$  by setting  $\delta(\omega) := \gamma(\omega_{1:n}) := 0$ , for all  $\omega \in \{0, 1\}^\infty, n \in \mathcal{N}$ , and set  $i := 0$ .

**Step 2** Set  $i := i + 1$ ; compute  $\psi(i)$  and let its value be  $(y, m)$ .

**Step 3** If  $\gamma(y) \geq m$  then go to Step 2 else set  $\gamma(y) := m$ .

**Step 4** If  $\mu(\bigcup_{\gamma(z) \geq k} \Gamma_z) > 2^{-k}$  for some  $k, k = 1, \dots, m$ , {since  $\mu$  is a recursive function we can effectively test whether the new assignment violates Definition 2.5.1} then terminate {the computation of  $\delta$  is finished} else set  $\delta(\omega) := \max\{m, \delta(\omega)\}$ , for all  $\omega \in \Gamma_y$ , and go to Step 2.

For the uniform measure  $\lambda(\Gamma_x) = 2^{-l(x)}$ , the conditional in Step 4 simplifies for  $i = 1$  to  $m > l(y)$ . In case the range of  $\psi$  is already a sequential  $\mu$ -test, then the algorithm never finishes but approximates  $\delta$  from below. If the range of  $\psi$  is not a sequential  $\mu$ -test, then at some point the conditional in Step 4 is violated and the computation of  $\delta$  terminates. The resulting  $\delta$  is a sequential  $\mu$ -test, even a recursive one. If the conditional in Step 4 is never violated, but the computation of  $\psi$  diverges for some argument, then  $\delta$  is trivially an enumerable sequential  $\mu$ -test.

Executing this procedure on all functions in the list  $\phi_1, \phi_2, \dots$ , we obtain an effective enumeration  $\delta_1, \delta_2, \dots$  of all sequential  $\mu$ -tests (and only sequential  $\mu$ -tests). The function  $\delta_0(\cdot|\mu)$  defined by

$$\delta_0(\omega|\mu) = \sup_{i \in \mathcal{N}} \{\delta_i(\omega) - i\}$$

is a universal sequential  $\mu$ -test.

First,  $\delta_0(\omega|\mu)$  is an enumerable function since the set  $\{(m, \omega) : \delta_0(\omega|\mu) \geq m\}$  is recursively enumerable. The proof that  $\delta_0(\cdot|\mu)$  is a sequential  $\mu$ -test, and majorizes all other sequential  $\mu$ -tests additively, is completely analogous to the proof for the similarly defined universal  $P$ -test in Section 2.4.  $\square$

Any other sequential  $\mu$ -test  $\delta_i$  can discover at most a constant amount more randomness in a sequence  $\omega$  than does  $\delta_0(\cdot|\mu)$ . That is,  $\delta_i(\omega) \leq \delta_0(\omega|\mu) + i$ , for all  $\omega$ .

The difference between any two universal sequential  $\mu$ -tests  $\delta_0(\cdot|\mu)$  and  $\delta'_0(\cdot|\mu)$  is bounded by a constant:  $|\delta_0(\omega|\mu) - \delta'_0(\omega|\mu)| \leq c$ , with  $c$  independent of  $\omega$ . We are now ready to separate the random infinite sequences from the nonrandom ones.

**Definition 2.5.4** Let the sample space  $\{0, 1\}^\infty$  be distributed according to  $\mu$ , and let  $\delta_0(\cdot|\mu)$  be a universal sequential  $\mu$ -test. An infinite binary sequence  $\omega$  is  $\mu$ -random in the sense of Martin-Löf if  $\delta_0(\omega|\mu) < \infty$ . We call such a sequence simply *random* where both  $\mu$  and Martin-Löf are understood. (This is particularly interesting for  $\mu$  is the uniform measure.)

Note that *this definition does not depend on the choice of the particular universal sequential  $\mu$ -test* with respect to which the level is defined. Hence, the line between random and nonrandom infinite sequences is drawn sharply without dependence on a reference  $\mu$ -test. It is easy to see that the set of infinite sequences that are not random in the sense of Martin-Löf forms precisely the maximal constructive  $\mu$ -null set of  $\mu$ -measure zero we have constructed above. Therefore,

**Theorem 2.5.3** *Let  $\mu$  be a recursive measure. The set of  $\mu$ -random infinite binary sequences has  $\mu$ -measure one.*

We say that the universal sequential  $\mu$ -test  $\delta_0(\cdot|\mu)$  rejects an infinite sequence with probability zero, and we conclude that a randomly selected infinite sequence passes all effectively testable laws of randomness with probability one.

The main question remaining is the following: Let  $\lambda$  be the uniform measure. Can we formulate a universal sequential  $\lambda$ -test in terms of complexity? In Theorem 2.4.2 on page 132 the universal (nonsequential) test is expressed that way. The most obvious candidate for the universal sequential test would be  $f(\omega) = \sup_{n \in \mathcal{N}} \{n - C(\omega_{1:n})\}$ , but it is improper. To see this, it is simplest to notice that  $f(\omega)$  would declare *all* infinite  $\omega$  to be nonrandom since  $f(\omega) = \infty$ , for all  $\omega$ , by Theorem 2.5.1. The same would be the case for  $f(\omega) = \sup_{n \in \mathcal{N}} \{n - C(\omega_{1:n}|n)\}$ , by about the same proof. It is difficult to express a

universal sequential test precisely in terms of  $C$ -complexity. But in Chapter 3 we show that it is easy to separate the random infinite sequences from the nonrandom ones in terms of “prefix” complexity.

### 2.5.3 Characterization of Random Sequences

How accurately can we characterize the set of infinite random sequences in complexity terms? It turns out that we can sandwich them between a proper superset in Theorem 2.5.4, page 145, and a proper subset in Theorem 2.5.5, page 146. First we bound the amplitude of the oscillations of random sequences.

**Definition 2.5.5** The infinite series  $\sum 2^{-f(n)}$  is *recursively convergent* if there is a recursive sequence  $n_1, n_2, \dots$  such that

$$\sum_{n=n_m}^{\infty} 2^{-f(n)} \leq 2^{-m}, \quad m = 1, 2, \dots$$

**Theorem 2.5.4** Let  $f(n)$  be a recursive function such that  $\sum_{n=1}^{\infty} 2^{-f(n)} < \infty$  is recursively convergent. If an infinite binary sequence  $\omega$  is random with respect to the uniform measure, then  $C(\omega_{1:n}|n) \geq n - f(n)$ , from some  $n$  onwards.

**Proof.** We define a sequential test that is passed only by the  $\omega$ 's satisfying the conditions in the theorem. For each  $m$ , let the critical section  $V_m$  consist of all infinite binary sequences  $\omega$  such that there exists an  $n \geq n_m$  for which  $C(\omega_{1:n}|n) < n - f(n)$ . In other words,  $V_m$  consists of the union of all intervals  $[0.\omega_{1:n}, 0.\omega_{1:n} + 2^{-n}]$  satisfying these conditions. We have to show that this is a sequential test. We can recursively enumerate the intervals comprising  $V_m$ , and therefore  $V_1, V_2, \dots$  is a sequence of recursively enumerable sets. Obviously, the sequence is nested. For every large enough  $n$ , at most  $2^{n-f(n)}$  strings  $y$  of length  $n$  satisfy  $C(y|n) < n - f(n)$  (Theorem 2.2.1 on page 109). Hence, with  $\lambda$  the uniform measure, we have, for all  $m$ ,

$$\lambda(V_m) \leq \sum_{n=n_m}^{\infty} 2^{n-f(n)} 2^{-n} \leq 2^{-m}.$$

Therefore, the sequence of critical regions forms a sequential test, and  $\lambda(\bigcap_{m=1}^{\infty} V_m) = 0$ . That is,  $\bigcap_{m=1}^{\infty} V_m$  is a constructive  $\lambda$ -null set associated with a sequential test. Consequently, it is contained in the maximal constructive  $\lambda$ -null set, which consists precisely of the sequences that are not random according to Martin-Löf.  $\square$

We can say fairly precisely which functions  $f$  satisfy the condition in Theorem 2.5.4. Examples are  $f(n) = 2 \log n$  and  $f(n) = \log n + 2 \log \log n$ .

In fact, the function  $g(n)$  used in the proof of Theorem 2.5.1 is about the borderline. It is *almost* sufficient that  $f(n) - g(n)$  be unbounded for a recursive function  $f(n)$  to satisfy the condition in Theorem 2.5.4. A precise form of the borderline function is given in Exercise 3.6.9, page 223.

With Theorem 2.5.4 we have shown that Theorem 2.5.1 is optimal in the sense that it gives the deepest complexity dips to be expected from sequences that are random with respect to the uniform measure (in the sense of Martin-Löf). But also, we have found a property of infinite sequences in terms of  $C$  that is implied by randomness with respect to the uniform measure. Is there also a property of infinite sequences in terms of complexity  $C$  that implies randomness with respect to the uniform measure?

**Theorem 2.5.5** *Let  $\omega$  be an infinite binary sequence.*

- (i) *If there exists a constant  $c$  such that  $C(\omega_{1:n}) \geq n - c$ , for infinitely many  $n$ , then  $\omega$  is random in the sense of Martin-Löf with respect to the uniform measure.*
- (ii) *The set of  $\omega$  for which there exists a constant  $c$  and infinitely many  $n$  such that  $C(\omega_{1:n}) \geq n - c$  has uniform measure one.*

**Proof.** We first prove the following claim:

**Claim 2.5.1** Let  $\omega \in \{0, 1\}^\infty$ . There exists a positive constant  $c$  such that  $C(\omega_{1:n}|n) \geq n - c$  for infinitely many  $n$  iff there exists a positive constant  $c$  such that  $C(\omega_{1:n}) \geq n - c$  for infinitely many  $n$ .

**Proof. (ONLY IF)** This is the easy direction, since conditional information does not increase complexity. Hence, for all  $\omega, n$ , we have  $C(\omega_{1:n}|n) \leq C(\omega_{1:n})$  up to a fixed additive constant.

**(IF)** For some fixed constant  $c_1$ , we have for all  $\omega, n$  that  $C(\omega_{1:n}) \leq C(\omega_{1:n}|n) + 2l(n - C(\omega_{1:n}|n)) + c_1$ . (The right-hand side of the inequality is the length of a description of  $\omega_{1:n}$ .) Since in the “If” direction we assume that  $C(\omega_{1:n}) \geq n - c$  for some  $c$  and infinitely many  $n$ , we obtain  $n - C(\omega_{1:n}|n) \leq c + c_1 + 2l(n - C(\omega_{1:n}|n))$  for this infinite sequence of  $n$ ’s. But that is only possible if there is a constant  $c_2$  such that  $n - C(\omega_{1:n}|n) \leq c_2$  for the same infinite sequence of  $n$ ’s, which finishes the proof.  $\square$

(i) Below, a “test” is a “ $\lambda$ -test” with  $\lambda$  the uniform measure. We denote sequential tests by  $\delta$ ’s, and (nonsequential) tests of Section 2.4 by  $\gamma$ ’s. Let  $\delta_0(\cdot|\lambda)$  denote the universal sequential test with respect to the uniform measure  $\lambda$ , and let  $\gamma_0(\cdot|L)$  denote the universal test with respect to the uniform distribution  $L$ .

Since a sequential test is a fortiori a test, there is a constant  $c$  such that  $\delta_0(\omega_{1:n}|\lambda) \leq \gamma_0(\omega_{1:n}|L) + c$ , for all  $\omega$  and  $n$ . By choosing the specific universal test of Theorem 2.4.2, we have  $\delta_0(\omega_{1:n}|\lambda) \leq n - C(\omega_{1:n}|n)$  up to a constant. Since  $\delta_0$  is monotonic nondecreasing,

$$\lim_{n \rightarrow \infty} \delta_0(\omega_{1:n}|\lambda) \leq \liminf_{n \rightarrow \infty} (n - C(\omega_{1:n}|n)) + O(1).$$

For those  $\omega$ 's satisfying the assumption in the statement of the theorem, by Claim 2.5.1 the right-hand side of the inequality is finite. By Theorem 2.4.2 therefore such  $\omega$ 's are random with respect to the uniform measure.

(ii) For each  $c$  and  $n$ , let  $V_{c,n}$  denote the union of the set of intervals associated with prefixes  $\omega_{1:n}$  of infinite binary sequences  $\omega$  such that  $C(\omega_{1:n}|n) \geq n - c$ . Let  $\lambda$  be the uniform measure. There are at most  $2^{n-c}$  strings of length less than  $n - c$  and therefore at least  $2^n - 2^{n-c}$  strings  $x$  of length  $n$  satisfy  $C(x|n) \geq n - c$ . Hence, for each  $m$  and  $c$  we have

$$\begin{aligned} \lambda \left( \bigcup_{n=m}^{\infty} V_{c,n} \right) &\geq \lambda(V_{c,m}) \\ &\geq (2^m - 2^{m-c})2^{-m} = 1 - 2^{-c}. \end{aligned}$$

Since the right-hand term is independent of  $m$ , we also have

$$\lambda \left( \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c,n} \right) \geq 1 - 2^{-c}.$$

Since

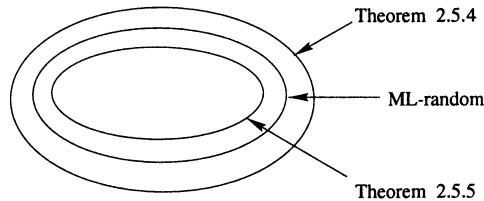
$$\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c,n} \subseteq \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c+1,n}$$

for all positive integers  $c$ , we obtain

$$\begin{aligned} \lambda \left( \bigcup_{c=1}^{\infty} \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c,n} \right) &= \lim_{c \rightarrow \infty} \lambda \left( \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c,n} \right) \\ &\geq \lim_{c \rightarrow \infty} (1 - 2^{-c}) = 1. \end{aligned}$$

The formula  $\bigcup_{c=1}^{\infty} \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} V_{c,n}$  denotes precisely the set of infinite sequences  $\omega$  for which there exists a positive integer constant  $c$  such that for infinitely many  $n$ ,  $C(\omega_{1:n}|n) \geq n - c$  holds. A fortiori, the same holds without the conditional up to an additive constant.  $\square$

We conclude that the set of sequences satisfying the condition in Theorem 2.5.4 contains the set of sequences that are random with respect



**FIGURE 2.5.** Three notions of “chaotic” infinite sequences

to the uniform measure (in Martin-Löf’s sense), and the latter contains the set of sequences satisfying the condition in Theorem 2.5.5; see Figure 2.5. It can be shown that these three sets are distinct, so the nesting is proper. Although the differences between each pair of the three sets are nonempty, they are not large, since all three sets have uniform measure one. For instance, the set of random sequences not satisfying the condition of Theorem 2.5.5 has uniform measure zero.

The combination of Theorems 2.5.4 and 2.5.5 enables us to give a relation between the upwards and downwards oscillations of the complexity of prefixes of the random sequences satisfying the property in Theorem 2.5.5 as follows:

**Corollary 2.5.2** If  $f$  is a recursive function such that  $\sum 2^{-f(n)}$  converges recursively and  $C(\omega_{1:n}) \geq n - c$  for some constant  $c$  and for infinitely many  $n$ , then  $C(\omega_{1:n}) \geq n - f(n)$  from some  $n$  onwards.

The universal sequential  $\mu$ -test characterizes the set of infinite random sequences. There are other ways to do so. We give an explicit characterization of infinite random sequences with respect to the uniform measure in Theorem 3.6.1 and its corollary, page 212. This characterization is an exact expression in terms of the prefix complexity developed in Chapter 3.

Apart from sequential tests as developed above there are other types of tests for randomness of individual infinite sequences. The extended theory of randomness tests can only be given after treating enumerable semimeasures in Chapter 4. This is done in Sections 4.5.7 and 4.5.6. There we give exact expressions for  $\mu$ -tests for randomness, for arbitrary recursive  $\mu$ .

We recall von Mises’s classic approach to obtaining infinite random sequences  $\omega$  as treated in Section 1.9, which formed a primary inspiration to the work reported in this section. It is of great interest whether one can, in his type of formulation, capture the intuitively and mathematically satisfying notion of infinite random sequence in the sense of Martin-Löf. According to von Mises an infinite binary sequence  $\omega$  was random (a collective) if

1.  $\omega$  has the property of frequency stability with limit  $p$ ; that is, if  $f_n = \omega_1 + \omega_2 + \dots + \omega_n$ , then the limit of  $f_n/n$  exists and equals  $p$ .
2. Any subsequence of  $\omega$  chosen according to an admissible place-selection rule has frequency stability with the same limit  $p$  as in condition 1.

One major problem was how to define “admissible,” and one choice was to identify it with Church’s notion of selecting a subsequence  $\zeta_1\zeta_2\dots$  of  $\omega_1\omega_2\dots$  by a partial recursive function  $\phi$  by  $\zeta_n = \omega_m$  if  $\phi(\omega_{1:r}) = 0$  for precisely  $n-1$  instances of  $r$  with  $r < m-1$  and  $\phi(\omega_{1:m-1}) = 0$ . We called these  $\phi$  “place-selection rules according to Mises-Church,” and the resulting sequences  $\zeta$  Mises-Wald-Church random.

In Section 1.9 we stated that there are Mises-Wald-Church random sequences with limiting frequency  $\frac{1}{2}$  that do not satisfy effectively testable properties of randomness like the Law of the Iterated Logarithm or the Recurrence Property. (Such properties are by definition satisfied by sequences that are Martin-Löf random.) In fact, the distinction between the two is quite large, since there are Mises-Wald-Church collectives  $\omega$  with limiting frequency  $\frac{1}{2}$  such that  $C(\omega_{1:n}) = O(f(n) \log n)$  for every unbounded, nondecreasing, total recursive function  $f$ , see also Exercise 2.5.13 on page 153. Such collectives are very nonrandom sequences from the viewpoint of Martin-Löf randomness where  $C(\omega_{1:n})$  is required to be asymptotic to  $n$ . See R.P. Daley, *Math. Systems Theory*, 9(1975), 83–94. It is interesting to point out that although a Mises-Wald-Church random sequence may have very low Kolmogorov complexity, it in fact has very high time-bounded Kolmogorov complexity. See Exercise 7.1.6 on page 473.

If we consider also sequences with limiting frequencies different from  $\frac{1}{2}$ , then it becomes even easier to find sequences that are random according to Mises-Wald-Church, but not according to Martin-Löf. Namely, *any* sequence  $\omega$  with limiting relative frequency  $p$  has complexity  $C(\omega_{1:n}) \leq Hn + o(n)$  for  $H = -(p \log p + (1-p) \log(1-p))$  ( $H$  is Shannon’s entropy). This means that for each  $\epsilon > 0$  there are Mises-Wald-Church random sequences  $\omega$  with  $C(\omega_{1:n}) < \epsilon n$  for all but finitely many  $n$ .

On the other hand, clearly all Martin-Löf random sequences are also Mises-Wald-Church random (each admissible selection rule is an effective sequential test).

This suggests that we have to liberate our notion of admissible selection rule somewhat in order to capture the proper notion of an infinite random sequence using von Mises’s approach. A proposal in this direction was given by A.N. Kolmogorov [*Sankhyā*, Ser. A, 25(1963), 369–376] and D.W. Loveland [*Z. Math. Logik Grundl. Math.* 12(1966), 279–294].

A Kolmogorov-Loveland admissible selection function to select an infinite subsequence  $\zeta_1\zeta_2\dots$  from  $\omega = \omega_1\omega_2\dots$  is any one-to-one recursive function  $\phi : \{0,1\}^* \rightarrow \mathcal{N} \times \{0,1\}$  from binary strings to (index, bit) pairs. The index gives the next position in  $\omega$  to be scanned, and the bit indicates whether the scanned bit of  $\omega$  must be included in  $\zeta_1\zeta_2\dots$  More precisely,

$$\phi(z) = \begin{cases} (i_1, a_1) & \text{if } z = \epsilon, \\ (i_m, a_m) & \text{if } z = z_1z_2\dots z_{m-1} \text{ with } z_j = \omega_i, 1 \leq j < m. \end{cases}$$

The sequence  $\zeta_1\zeta_2\dots$  is called a *Kolmogorov-Loveland random sequence*.

Because  $\phi$  is one-to-one, it must hold that  $i_m \neq i_1, i_2, \dots, i_{m-1}$ . The described process yields a sequence of  $\phi$ -values  $(z_1, a_1), (z_2, a_2), \dots$ . The selected subsequence  $\zeta_1\zeta_2\dots$  consists of the ordered sequence of those  $z_i$ 's whose associated  $a_i$ 's equal one.

As compared to the Mises-Wald-Church approach, the liberation is contained in the fact that the order of succession of the terms in the subsequence chosen is not necessarily the same as that of the original sequence. In comparison, it is not obvious (and it seems to be unknown) whether a subsequence  $\zeta_1\zeta_2\dots$  selected from a Kolmogorov-Loveland random sequence  $\omega_1\omega_2\dots$  by a Kolmogorov-Loveland place-selection rule is itself a Kolmogorov-Loveland random sequence. Note that the analogous property necessarily holds for Mises-Wald-Church random sequences.

The set of Kolmogorov-Loveland random sequences is contained in the set of Mises-Wald-Church random sequences and contains the set of Martin-Löf random sequences. If  $\omega_1\omega_2\dots$  is Kolmogorov-Loveland random then clearly  $\zeta_1\zeta_2\dots$  defined by  $\zeta_i = \omega_{\sigma(i)}$  with  $\sigma$  being a recursive permutation, is also Kolmogorov-Loveland random. The Mises-Wald-Church notion of randomness does not have this important property of randomness of staying invariant under recursive permutation. Loveland gave the required counterexample in the cited reference. Hence, the containment of the set of Kolmogorov-Loveland random sequences in the set of Mises-Wald-Church random sequences is proper. This leaves the question of whether the containment of the set of Martin-Löf random sequences in the set of Kolmogorov-Loveland random sequences is proper. Kolmogorov has stated in [*Problems Inform. Transmission*, 5(1969), 3–4] without proof that there exists a Kolmogorov-Loveland random sequence  $\omega$  such that  $C(\omega_{1:n}) = O(\log n)$ . But A.N. Muchnik (not to be confused with A.A. Muchnik) showed that this is false since no  $\omega$  with  $C(\omega_{1:n}) \leq cn + O(1)$  for a constant  $c < 1$  can be Kolmogorov-Loveland random. Nonetheless, containment is proper since A.Kh. Shen' [*Soviet Math. Dokl.*, 38:2(1989), 316–319] has shown there exists a Kolmogorov-Loveland random sequence which is not random in Martin-Löf's sense. Therefore, the problem of giving a satisfactory definition of infinite Martin-Löf random sequences in the form proposed by von Mises has not yet been solved. See also [A.N. Kolmogorov and V.A. Uspensky, *Theory Probab. Appl.*, 32(1987), 389–412; V.A. Uspensky, A.L. Semenov, and A.Kh. Shen', *Russ. Math. Surveys*, 45:1(1990), 121–189].

## Exercises

**2.5.1.** [13] Consider  $\{0, 1\}^\infty$  under the uniform measure. Let  $\omega = \omega_1\omega_2\dots \in \{0, 1\}^\infty$  be random in the sense of Martin-Löf.

- (a) Show that  $\zeta = \omega_n\omega_{n+1}\dots$  is Martin-Löf random for each  $n$ .
- (b) Show that  $\zeta = x\omega$  is Martin-Löf random for each finite string  $x$ .

*Comments.* Source: C. Calude, I. Chitescu, *Bulletino U.M.I.*, (7) 3-B(1989), 229–240.

**2.5.2.** [21] Consider  $\{0, 1\}^\infty$  under the uniform measure. Let  $\omega = \omega_1\omega_2\dots \in \{0, 1\}^\infty$ .

(a) Show that if there is an infinite recursive set  $\{(i, \omega_i)\}$ , then  $\omega$  is not random in the sense of Martin-Löf.

(b) Show that if the set  $\{i : \omega_i = 0\}$  contains an infinite recursively enumerable subset, then  $\omega$  is not random in the sense of Martin-Löf.

*Comments.* Source: C. Calude and I. Chitescu, *Ibid.*.

**2.5.3.** [21] Let  $\omega = \omega_1\omega_2\dots$  be any infinite binary sequence. Define  $\zeta = \zeta_1\zeta_2\dots$  by  $\zeta_i = \omega_i + \omega_{i+1}$ ,  $i \geq 1$ . This gives a sequence over the alphabet  $\{0, 1, 2\}$ . Show that  $\zeta$  is not random in the sense of Martin-Löf under the uniform measure (extend the definition from binary to ternary sequences).

*Comments.* Hint: the blocks “02” and “20” do not occur in  $\zeta$ . Source: R. von Mises, *Probability, Statistics and Truth*, Dover, 1981.

**2.5.4.** [23] Let  $\omega$  be any infinite binary sequence. Show that for all constants  $c$  there are infinitely many  $m$  such that for all  $n$  with  $m \leq n \leq 2m$ ,  $C(\omega_{1:n}) \leq n - c$ .

*Comments.* We are guaranteed to find long complexity oscillations (of length  $m$ ) in an infinite binary sequence  $\omega$  relatively near the beginning (namely  $\omega_{m:2m}$ ), even if  $\omega$  is Martin-Löf random. Source: H.P. Katseff and M. Sipser, *Theoret. Comput. Sci.*, 15(1981), 291–309.

**2.5.5.** [M19] Let  $f$  be such that  $\sum 2^{-f(n)} < \infty$ . Show that the set of infinite binary sequences  $\omega$  satisfying  $C(\omega_{1:n}|n) \geq n - f(n)$  for all but finitely many  $n$  has uniform measure 1.

*Comments.* Hint: the number of  $y$  with  $l(y) = n$  such that  $C(y) < n - f(n)$  is less than  $2^n - f(n)$ . This implies that the probability that this inequality is satisfied is less than  $2^{-f(n)}$ , and the result follows by the Borel-Cantelli Lemmas, see Exercise 1.10.3, page 64. Source: P. Martin-Löf, *Z. Wahrsch. Verw. Geb.*, 19(1971), 225–230.

**2.5.6.** [09] Consider infinite binary sequences  $\omega$  with respect to the uniform measure. Show that with probability one there exists a constant  $c$  such that  $C(\omega_{1:n}|n) \geq n - c$  for infinitely many  $n$ .

*Comments.* Hint: use Theorem 2.5.5 Item (ii) and Claim 2.5.1. Source: P. Martin-Löf, *Z. Wahrsch. Verw. Geb.*, 19(1971), 225–230. .

**2.5.7.** [19] Consider infinite binary sequences  $\omega$  with respect to the uniform measure. Show that if  $f$  is a recursive function and  $\sum 2^{-f(n)}$  converges recursively and  $C(\omega_{1:n}) \geq n - c$  for some constant  $c$  and infinitely many  $n$ , then  $C(\omega_{1:n}) \geq n - f(n)$  for all but finitely many  $n$ .

*Comments.* This formulation establishes a connection between upward and downward oscillations of the complexity of prefixes of almost all (random) infinite binary sequences. For  $f$  we can take  $f(n) = 2 \log n$ , or

$f(n) = \log n + 2 \log \log n$ , and so on. Hint: combine Theorems 2.5.4, 2.5.5.  
Source: P. Martin-Löf, *Z. Wahrscheinlichkeitstheorie verw. Geb.*, 19(1971), 225–230.

**2.5.8.** [19] Show that there exists an infinite binary sequence  $\omega$  and a constant  $c > 0$  such that  $\liminf_{n \rightarrow \infty} C(\omega_{1:n}|n) \leq c$ , but for any unbounded function  $f$  we have  $\limsup_{n \rightarrow \infty} C(\omega_{1:n}|n) \geq n - f(n)$ .

*Comments.* The oscillations can have amplitude  $\Omega(n)$ . Hint: use the  $n$ -strings defined above. Construct  $\omega = y_1 y_2 \dots$  from finite strings  $y_i$ . Let  $c$  be a fixed independent constant. For odd  $i$  choose  $y_i$  such that  $y_1 \dots y_i$  is an  $n$ -string, which implies that  $C(y_{1:i}|l(y_{1:i})) < c$ . For even  $i$  choose  $y_i$  as a long enough random string so that  $C(y_{1:i}|l(y_{1:n})) \geq l(y_{1:n}) - f(l(y_{1:n}))$ .  
Source: H.P. Katseff and M. Sipser, *Theoret. Comput. Sci.*, 15(1981), 291–309.

**2.5.9.** [39] Consider a measure on the set of intervals contained in  $[0, 1]$ . For example, the Lebesgue measure  $\lambda$  defined by  $\lambda(\Gamma_y) = 2^{-l(y)}$ . (Recall that for each finite binary string  $y$  the cylinder  $\Gamma_y$  is the set of all infinite strings  $\omega$  starting with  $y$ .) Let  $\omega$  be an infinite binary sequence such that for every recursively enumerable sequence  $A_1, A_2, \dots$  of sets of intervals with the property that the series  $\sum_i \lambda(A_i) < \infty$  converges,  $\omega$  is contained in only finitely many  $A_i$ . Show that for the  $\omega$ 's defined this way, the *Solovay random sequences* are precisely the infinite binary sequences that are random in the sense of Martin-Löf with respect to the uniform measure.

*Comments.* In Martin-Löf's definition of randomness (with respect to the uniform measure) of infinite binary sequences, he required that  $\lambda(A_i) \leq 2^{-i}$ . That definition is equivalent to stipulating the existence of some “regulator of convergence”  $f(i) \rightarrow \infty$  that is recursive and nondecreasing such that  $\lambda(A_i) \leq 2^{-f(i)}$ . Solovay's definition has the advantage that it does not require such a regulator. Source: R.M. Solovay, *Lecture Notes*, 1975, unpublished; and G.J. Chaitin, *Algorithmic Information Theory*, Cambridge University Press, 1987; A. Kh. Shen', *Soviet Math. Dokl.*, 38:2(1989), 316–319.

**2.5.10.** [35] (a) Show that for every positive constant  $c$  there is a positive constant  $c'$  such that  $\{\omega_{1:n} : C(\omega_{1:n}; n) \geq n - c\} \subseteq \{\omega_{1:n} : C(\omega_{1:n}|n) \geq n - c'\}$ .

(b) Use the observation in Item (a) to show that Theorem 2.5.4 holds for the uniform complexity measure  $C(\cdot; l(\cdot))$ .

(c) Show that if  $f$  is a recursive function and  $\sum 2^{-f(n)} = \infty$ , then for all infinite  $\omega$  we have  $C(\omega_{1:n}; n) \leq n - f(n)$  for infinitely many  $n$ . Hence, Theorem 2.5.1 holds for uniform complexity.

*Comments.* Hint for Item (a): define the notion of a (universal) uniform test as a special case of Martin-Löf's (universal) test. Compare this

result with the other exercises to conclude that whereas the uniform complexity tends to be higher in the low-complexity region, the length-conditional complexity tends to be higher in the high-complexity region. Source: D.W. Loveland, *Inform. Contr.*, 15(1969), 510–526.

**2.5.11.** • [31] Show that the following statements are equivalent for an infinite binary sequence  $\omega$ : For some constant  $c$  and infinitely many  $n$ ,

$$\begin{aligned} C(\omega_{1:n}|n) &\geq n - c, \\ C(\omega_{1:n}; n) &\geq n - c, \\ C(\omega_{1:n}) &\geq n - c. \end{aligned}$$

*Comments.* Hint: use Claim 2.5.1 and Exercise 2.5.10. In view of Theorem 2.5.5 these conditions equivalently imply that  $\omega$  is random in the sense of Martin-Löf. Source: R.P. Daley, pp. 113–122 in: *Computational Complexity*, ed. R. Rustin, *Courant Comput. Sci. Symp.* 7(1971).

**2.5.12.** [30] (a) Show that the following statements are equivalent for an infinite binary sequence  $\omega$ : There exists a  $c$  such that for infinitely many  $n$ ,

$$\begin{aligned} C(\omega_{1:n}|n) &\leq c, \\ C(\omega_{1:n}; n) &\leq l(n) + c, \\ C(\omega_{1:n}) &\leq l(n) + c. \end{aligned}$$

The sequences thus defined are called *pararecursive sequences*.

(b) Show that no pararecursive sequence is random in the sense of Martin-Löf.

*Comments.* Comparison with the other exercises shows that the recursive sequences are contained in the pararecursive sequences. It also shows that the pararecursive sequences have the cardinality of the continuum, so this containment is proper. R.P. Daley [*J. Symb. Logic*, 41(1976), 626–638] has shown that the enumerable sequences (characteristic sequences of recursively enumerable sets) are pararecursive, and this containment is proper by the same argument as before. Hint for Item (b): use Theorem 2.5.4. Despite Item (b), there are pararecursive sequences that are close to being random. For any unbounded function  $f$ , there is a pararecursive sequence  $\omega$  such that for infinitely many  $n$  we have  $C(\omega_{1:n}|n) \geq n - f(n)$ . Source: H.P. Katseff and M. Sipser, *Theoret. Comput. Sci.*, 15(1981), 291–309.

**2.5.13.** [43] Let  $A$  be the set of Mises-Wald-Church random sequences with  $p = \frac{1}{2}$ . The admissible place-selection rules are the partial recursive functions.

(a) Show that there is an  $\omega \in A$  such that for each unbounded, nondecreasing, total recursive function  $f$ ,  $C(\omega_{1:n}; n) \leq f(n) \log n$  from some  $n$  onwards.

(b) Show that for all  $\omega \in A$ , there is a constant  $c$  such that  $C(\omega_{1:n}; n) \geq \log n - c$  from some  $n$  onwards.

Consider the larger class  $B \supset A$  that is defined just as  $A$  but with the admissible place-selection rules are restricted to the total recursive functions.

(c) Show that there is an  $\omega \in B$  such that  $C(\omega_{1:n}; n) \leq f(n)$  from some  $n$  onwards, for each  $f$  as in Item (a).

(d) Show that for each  $\omega \in B$ , for each constant  $c$ , we have  $C(\omega_{1:n}; n) \geq c$  from some  $n$  onwards.

*Comments.* Compare page 149. This shows that there are Mises-Wald-Church random sequences of quite low complexity, and that it makes a difference whether the admissible place-selection rules are partial recursive or total recursive. Source: R.P. Daley, *Math. Systems Theory*, 9 (1975), 83–94. Item (a) is proved using Item (c), which is attributed to D.W. Loveland, and uses a construction (LMS algorithm) in D.W. Loveland, *Z. Math. Logik*, 12(1966), 279–294. Compare with Exercise 7.1.6, page 473, to see what happens when we impose a total recursive time bound on the decoding process.

**2.5.14.** [O37] Is there a Mises-Wald-Church random sequence  $\omega$  with limiting frequency  $\frac{1}{2}$  with  $C(\omega_{1:n}) = O(\log n)$ ?

*Comments.* Compare Exercise 2.5.13.

**2.5.15.** [33] (a) Show that there exists an infinite binary sequence  $\omega$  that is random with respect to the uniform measure, but for each constant  $c$  there are only finitely many  $n$  such that  $C(\omega_{1:n}|n) > n - c$  (the condition of Theorem 2.5.5 does not hold.)

(b) Show that there exists an infinite binary sequence  $\omega$  satisfying (i)  $C(\omega_{1:n}) > n - f(n)$  from some  $n$  onwards and  $\sum 2^{-f(n)}$  converges recursively (the condition in Theorem 2.5.4 holds), and (ii)  $\omega$  is not random with respect to the uniform measure.

*Comments.* Thus, each containment in the nested sequence of sets of infinite binary sequences that satisfy the condition in Theorem 2.5.4, randomness according to Martin-Löf, and the condition in Theorem 2.5.5, as in Figure 2.5, is proper. Source, C.P. Schnorr, *Math. Systems Theory*, 5(1971), 246–258.

**2.5.16.** [21] (a) Show that none of the variants of algorithmic complexity, such as  $C(x)$ ,  $C(x|l(x))$ , and  $C(x;l(x))$ , is invariant with respect to cyclic shifts of the strings.

- (b) Show that all these variants coincide to within the logarithm of the minimum of all these measures.

*Comments.* Hint: use the idea in the proof of Theorem 2.5.4. This invariance cannot be expected from any complexity measure in this book at all. Source: C.P. Schnorr, pp. 193–211 in: R.E. Butts and J. Hintikka (Eds.), *Basic Problems in Methodology and Linguistics*, D. Reidel, 1977.

**2.5.17.** [36] (a) Consider an infinite sequence of zeros and ones generated by independent tosses of a coin with probability  $p$  ( $0 < p < 1$ ) for “1.” If  $p$  is a recursive number, the methods in this section can be applied to obtain a proper definition for infinite Bernoulli sequences. There is, however, no reason to suppose that in physical coins  $p$  is a recursive real. This prompts another approach where we disregard the actual probability distribution, but follow more closely the combinatorial spirit of Kolmogorov complexity. Define sequential Bernoulli tests (in analogy with Section 2.5 and Exercise 2.4.4 on page 135). Show that there exists a universal sequential Bernoulli test  $\delta_0$ . An infinite binary sequence  $\omega$  is a Bernoulli sequence if  $\delta_0(\omega) < \infty$ . Show that the set of Bernoulli sequences has measure one with respect to the measure induced in the set of infinite binary sequences interpreted as independent  $(p, 1 - p)$  Bernoulli trials.

(b) In our definition of infinite Bernoulli sequences no restrictions were laid on the limiting behavior of the relative frequency, such as, for instance, required in Condition 1 of Definition 1.9.1 of an infinite random sequence (collective)  $\omega$ , where we require that  $\lim_{n \rightarrow \infty} f_n/n = p$  for some  $p$  ( $0 < p < 1$ ) (Section 1.9). The relative frequency  $f_n/n$  of ones in increasingly longer prefixes  $\omega_{1:n}$  of a collective  $\omega$  does not oscillate indefinitely, but converges to a definite limit. Show that remarkably, this is also the case for an infinite Bernoulli sequence  $\omega$ .

(c) By the law of large numbers, *all* real numbers  $p$  in  $[0, 1]$  occur as limit frequencies  $\lim_{n \rightarrow \infty} f_n/n$  for infinite random binary sequences  $\omega$ , and not only the recursive ones. Show that in contrast, for infinite Bernoulli sequences  $\omega$  the limit relative frequency cannot vanish,  $\lim_{n \rightarrow \infty} f_n/n = 0$ , unless  $\omega_n = 0$  for all  $n$ .

*Comments.* Hint for Item (b): for an arbitrary rational  $\epsilon > 0$  construct a sequential Bernoulli test that rejects at level  $2^{-m}$  if  $|f_i/i - f_j/j| > \epsilon$ , for some  $i, j \geq h(m)$ , for some suitable nondecreasing total recursive function. Compare with the universal Bernoulli test of Exercise 2.4.4 on page 135.. Hint for Item (c): this is the infinite analogue of the phenomenon in Item (b) of Exercise 2.4.4 on page 135. From Item (c) we conclude that an event with vanishing limit frequency is actually impossible. This is in stark contrast with von Mises’s explicit statement of the opposite for his conception of random sequences (collectives). [R.

von Mises, *Probability, Statistics and Truth*, Dover, 1981 (Reprinted). Source: P. Martin-Löf, *Inform. Contr.*, 9(1966), 602–619. Additionally we mention the following result [L.A. Levin, *Sov. Math. Dokl.*, 14(1973), 1413–1416]. Suppose we are given a constructively closed family  $\mathcal{M}$  of measures (this notion is defined naturally on the space of measures). Let a test  $f$  be called *uniform* for  $\mathcal{M}$  if for all measures in  $\mathcal{M}$ , for all positive integers  $k$ , the measure of outcomes  $\omega$  where  $f(\omega) > k$  is at most  $2^{-k}$ . There exists a universal uniform test.

**2.5.18.** [37] Let  $\mu$  be a recursive measure on the sample space  $\{0, 1\}^\infty$ . Recall from Section 2.5 Martin-Löf's construction of a constructive  $\mu$ -null set using a notion of sequential test  $V$  with associated critical regions  $V_1 \supseteq V_2 \supseteq \dots$  of measures  $\mu(V_i) \leq 2^{-i}$ , for  $i \geq 1$ . A constructive  $\mu$ -null set is a *total recursive  $\mu$ -null set* if additionally,  $f(i) = \mu(V_i)$  is a total recursive function. Call an infinite sequence  $\mu$ -random in the sense of Schnorr if it is not contained in any total recursive  $\mu$ -null set.

- (a) Show that there is no universal total recursive  $\mu$ -null set that contains all others.
- (b) Show that the set of sequences that are  $\mu$ -random in the sense of Martin-Löf is a subset of the set of sequences that are  $\mu$ -random in the sense of Schnorr.
- (c) An  $\omega \in \{0, 1\}^\infty$  is an *atom* with respect to  $\mu$  if  $\mu(\omega) > 0$ . A measure  $\mu$  is called *discrete* if the set of atoms of  $\mu$  has  $\mu$ -measure one. (Obviously all atoms of recursive  $\mu$  are recursive sequences.) Show that the Schnorr- $\mu$ -random sequences coincide with the Martin-Löf- $\mu$ -random sequences iff  $\mu$  is discrete.

*Comments.* Item (c) implies that for nondiscrete  $\mu$  (like the uniform measure) Schnorr randomness is weaker than Martin-Löf randomness. The notion of total recursive null sets is the recursive analogue of the intuitionistic notion of sets of measure zero by L.E.J. Brouwer [A. Heyting, *Intuitionism, an Introduction*, North-Holland, 1956]. Sometimes the following statement is called *Schnorr's Thesis*: “A sequence behaves within all effective procedures like a  $\mu$ -random sequence iff it is  $\mu$ -random in the sense of Schnorr.” Source: C.P. Schnorr, pp. 193–211 in: R.E. Butts and J. Hintikka (Eds.), *Basic Problems in Methodology and Linguistics*, D. Reidel, 1977.

**2.5.19.** [M42] We abstract away from levels of significance and concentrate on the arithmetical structure of statistical tests. Statistical tests are just  $\Pi_n^0$  null sets, for some  $n$  (Exercise 1.7.21, page 46). The corresponding definition of randomness is defined as, “an infinite binary sequence is  $\Pi_n^0$ -random with respect to a recursive measure  $\mu$  if it is not contained in any  $\Pi_n^0$  set  $V$  with  $\mu$ -measure zero.” Has the set of  $\Pi_n^0$ -random sequences  $\mu$ -measure one?

*Comments.* Source: H. Gaifman and M. Snir, *J. Symb. Logic*, 47(1982), 495–548.

**2.5.20.** [M43] We assume familiarity with the unexplained notions below. We leave the arithmetic hierarchy of Exercise 2.5.19 and consider hyperarithmetic sets. Define an infinite binary sequence to be *hyperarithmetically random* if it belongs to the intersection of all hyperarithmetic sets of measure one. (A hyperarithmetic set can be regarded as a constructive version of the restriction to Borel sets that is usually accepted in probability theory—Section 1.6. The specific Borel sets considered there are always obtained by applying the Borelian operations to recursive sequences of previously defined sets, which means precisely that they are hyperarithmetical.)

- (a) Show that the set of sequences that are hyperarithmetically random is a  $\Sigma_1^1$  set of measure one ( $\Sigma_1^1$  in the *analytic hierarchy*).
- (b) Show that a hyperarithmetic sequence is not hyperarithmetically random.
- (c) Show that the set of hyperarithmetically random sequences is not hyperarithmetical.

*Comments.* Already A. Wald proposed to sharpen von Mises's notion of randomness by defining a sequence to be random if it possesses all properties of probability one that are expressible within a certain formalized logic such as *Principia Mathematica*. This exercise is a variation on this idea. Just as here, Wald's proposal can be expected to define a set of random strings that is no longer expressible in the language with which we started. (This does not happen for Martin-Löf random sequences as defined in Section 2.5 because of the existence of a universal sequential test.) However, with the present proposal the resulting class of random strings, while escaping the hyperarithmetic hierarchy, does not escape us completely but belongs to a class of sets that can still be handled constructively. Source: P. Martin-Löf, pp. 73–78 in: *Intuitionism and Proof Theory*, A. Kino et al. (Eds.), North-Holland, 1970. For related work more in the direction of Wald's ideas, see [P.A. Benioff, *J. Math. Phys.*, 17:5(1976), 618–628, 629–640; L. Longpré and V. Kreinovich, 'Randomness as incompressibility: a non-algorithmic analogue,' Manuscript, 1996].

## 2.6 Statistical Properties of Finite Sequences

Each individual infinite sequence generated by a  $(\frac{1}{2}, \frac{1}{2})$  Bernoulli process (flipping a fair coin) has (with probability 1) the property that the relative frequency of zeros in an initial  $n$ -length segment goes to  $\frac{1}{2}$  as  $n$  goes to infinity. Such randomness related statistical properties of individual (high) complexity finite binary sequences are often required in applications of incompressibility arguments. The situation for infinite random sequences is better studied, and therefore we look there first.

### Definition 2.6.1

E. Borel has called an infinite sequence of zeros and ones *normal* in the scale of two if for each  $k$ , the frequency of occurrences of each block  $y$  of length  $k$  in the initial segment of length  $n$  goes to limit  $2^{-k}$  as  $n$  grows unboundedly.

It is known that normality is not sufficient for randomness, since Champernowne's sequence 123456789101112... is normal in the scale of ten. On the other hand, it is universally agreed that a random infinite sequence must be normal. (If not, then some blocks occur more frequently than others, which can be used to obtain better than fair odds for prediction.)

We know from Section 2.5 that each infinite sequence that is random with respect to the uniform measure satisfies all effectively testable properties of randomness: it is normal, it satisfies the so-called Law of the Iterated Logarithm, the number of 1's minus the number of 0's in an initial  $n$ -length segment is positive for infinitely many  $n$  and negative for another infinitely many  $n$ , and so on. While the statistical properties of infinite sequences are simple corollaries of the theory of Martin-Löf randomness, for finite sequences the situation is less simple. Here, we determine the frequencies of occurrence of substrings in high Kolmogorov complexity strings. In Section 6.4.1 the similar question is treated for subgraphs of high Kolmogorov complexity graphs.

### 2.6.1 Statistics of 0's and 1's

In the finite case, randomness is a matter of degree, because it would be clearly unreasonable to say that a sequence  $x$  of length  $n$  is random and to say that a sequence  $y$  obtained by flipping the first "1" bit of  $x$  is nonrandom. What we can do is to express the degree of incompressibility of a finite sequence in the form of its Kolmogorov complexity, and then analyze the statistical properties of the sequence—for example, the number of 0's and 1's in it.

Since almost all finite sequences have about maximal Kolmogorov complexity, each individual maximal complexity sequence must possess approximately the expected (average) statistical properties of the overall set. For example, we can a priori state that each high-complexity finite binary sequence is "normal" in the sense that each binary block of length  $k$  occurs about equally frequently for  $k$  relatively small. In particular,

this holds for  $k = 1$ . However, in many applications we need to know exactly what “about” and the “relatively small” in this statement mean. In other words, we are interested in the extent to which “Borel normality” holds in relation to the complexity of a finite sequence.

Let  $x$  have length  $n$ . By Example 2.4.4, if  $C(x|n) = n + O(1)$ , then the number of zeros it contains is

$$\frac{n}{2} + O(\sqrt{n}).$$

**Notation 2.6.1** The quantity  $K(x|y)$  in this section satisfies

$$C(x|y) \leq K(x|y) \leq C(x|y) + 2 \log C(x|y) + 1.$$

We can think of it as roughly the length of a self-delimiting version of a program  $p$  of length  $l(p) = C(x|y)$ . In Chapter 3 it is defined as “prefix complexity.”

**Definition 2.6.2** The class of *deficiency* functions is the set of functions  $\delta : \mathcal{N} \rightarrow \mathcal{N}$  satisfying  $K(n, \delta(n)|n - \delta(n)) \leq c_1$  for all  $n$ . (Hence,  $C(n, \delta(n)|n - \delta(n)) \leq c_1$  for all  $n$ .)

This way we can retrieve  $n$  and  $\delta(n)$  from  $n - \delta(n)$  by a self-delimiting program of at most  $c_1$  bits. We choose  $c_1$  so large that each monotone sublinear recursive function that we are interested in, such as  $\log n$ ,  $\sqrt{n}$ ,  $\log \log n$ , is such a deficiency function. The constant  $c_1$  is a benchmark that stays fixed throughout this section.

We denote the number of 1's in a binary string  $x \in \{0, 1\}^*$  by  $\#\text{ones}(x)$ .

**Lemma 2.6.1** *There is a constant  $c$  such that for all deficiency functions  $\delta$ , for each  $n$  and  $x \in \{0, 1\}^n$ , if  $C(x) \geq n - \delta(n)$ , then*

$$\left| \#\text{ones}(x) - \frac{n}{2} \right| \leq \sqrt{\frac{3}{2}(\delta(n) + c)n / \log e}. \quad (2.3)$$

**Proof.** A general estimate of the tail probability of the binomial distribution, with  $s_n$  the number of successful outcomes in  $n$  experiments with probability of success  $0 < p < 1$ , is given by Chernoff's bounds, Lemma 1.10.1 on page 61:

$$\Pr(|s_n - pn| > m) \leq 2e^{-m^2/3pn}. \quad (2.4)$$

Let  $s_n$  be the number of 1's in the outcome of  $n$  fair coin flips, which means that  $p = \frac{1}{2}$ . Define  $A = \{x \in \{0, 1\}^n : |\#\text{ones}(x) - n/2| > m\}$  and apply Equation 2.4 to obtain

$$d(A) \leq 2^{n+1} e^{-2m^2/3n}.$$

We choose  $m$  such that for some constant  $c$  to be determined later,

$$\frac{2m^2 \log e}{3n} = \delta(n) + c.$$

We can compress any  $x \in A$  in the following way:

1. Let  $s$  be a self-delimiting program to retrieve  $n$  and  $\delta(n)$  from  $n - \delta(n)$ , of length at most  $c_1$ .
2. Given  $n$  and  $\delta(n)$ , we can effectively enumerate  $A$ . Let  $i$  be the index of  $x$  in such an effective enumeration of  $A$ . The length of the (not necessarily self-delimiting) description of  $i$  satisfies

$$\begin{aligned} l(i) &\leq \log d(A) \leq n + 1 - (2m^2 \log e)/3n \\ &= n + 1 - \delta(n) - c. \end{aligned}$$

The string  $si$  is padded to length  $n + 1 - \delta(n) - c + c_1$ . From  $si$  we can reconstruct  $x$  by first using  $l(si)$  to compute  $n - \delta(n)$ , then compute  $n$  and  $\delta(n)$  from  $s$  and  $n - \delta(n)$ , and subsequently enumerate  $A$  to obtain the  $i$ th element. Let  $T$  be the Turing machine embodying the procedure for reconstructing  $x$ . Then by Theorem 2.1.1,

$$C(x) \leq C_T(x) + c_T \leq n + 1 - \delta(n) - c + c_1 + c_T.$$

Choosing  $c = c_1 + c_T + 2$  we find  $C(x) < n - \delta(n)$ , which contradicts the condition of the theorem. Hence,  $|\#ones(x) - n/2| \leq m$ .  $\square$

It may be surprising at first glance, but there are no maximally complex sequences with about an equal number of zeros and ones. Equal numbers of zeros and ones is a form of regularity, and therefore lack of complexity. That is, for  $x \in \{0, 1\}^n$ , if  $|\#ones(x) - n/2| = O(1)$ , then the randomness deficiency  $\delta(n) = n - C(x)$  is nonconstant (order  $\log n$ ).

**Lemma 2.6.2** *There is a constant  $c$  such that for all  $n$  and all  $x \in \{0, 1\}^n$ , if*

$$\left| \#ones(x) - \frac{n}{2} \right| \leq 2^{-\delta(n)-c} \sqrt{n},$$

*then  $C(x) \leq n - \delta(n)$ .*

**Proof.** Let  $m = 2^{-\delta(n)-c} \sqrt{n}$ , with  $c$  a constant to be determined later. Let  $A = \{x \in \{0, 1\}^n : |\#ones(x) - n/2| \leq m\}$ . There is a constant  $c_2$  such that there are only

$$d(A) \leq (2m + 1) \binom{n}{n/2} \leq c_2 \frac{2^n m}{\sqrt{n}} \tag{2.5}$$

elements in  $A$  (use Stirling's approximation, Exercise 1.5.4 on page 17). Thus, for each  $x \in A$ , we can encode  $x$  by its index in an enumeration of  $A$ . We can find  $A$  from  $n$  and  $\delta(n)$ . We can find  $n$  and  $\delta(n)$  from  $n - \delta(n)$  by a self-delimiting program of size at most  $c_1$ . Altogether, this description takes  $\log d(A) + c_1 = n - \delta(n) - c + c_1 + \log c_2$  bits. Let this process of reconstructing  $x$  be executed by Turing machine  $T$ . Choosing  $c = c_1 + \log c_2 + c_T$  we obtain by Theorem 2.1.1

$$C(x) \leq C_T(x) + c_T \leq n - \delta(n).$$

□

**Example 2.6.1** We consider some particular values of  $\delta(n)$ . Set  $\delta_1(n) = \frac{1}{2} \log n - \log \log n$ . If  $|\#\text{ones}(x) - n/2| = O(\log n)$ , then  $C(x) \leq n - \delta_1(n) + O(1)$ . Set  $\delta_2(n) = \frac{1}{2} \log n$ . If

$$\left| \#\text{ones}(x) - \frac{n}{2} \right| = O(1),$$

then  $C(x) \leq n - \delta_2(n) + O(1)$ . That is, if the number of 1's is too close to the number of 0's, then the complexity of the string drops significantly below its maximum. ◇

An incompressible string of length  $n$  cannot have precisely or almost  $n/2$  ones by Lemma 2.6.2. Then how many ones should an incompressible string contain? The next lemma shows that for an incompressible  $x$  having  $j + n/2$  ones,  $K(j|n)$  must be at least about order  $\log n$ .

**Lemma 2.6.3** *There is a constant  $c$  such that for all  $n$  and all  $x \in \{0, 1\}^n$ , if*

$$\left| \#\text{ones}(x) - \frac{n}{2} \right| = j,$$

*then  $C(x|n) \leq n - \frac{1}{2} \log n + K(j|n) + c$ .*

**Proof.** Let  $A = \{x \in \{0, 1\}^n : |\#\text{ones}(x) - n/2| = j\}$ . There is a constant  $c_3$  such that there are

$$d(A) \leq \binom{n}{n/2} \leq c_3 \frac{2^n}{\sqrt{n}} \tag{2.6}$$

elements in  $A$  (use Stirling's approximation, Exercise 1.5.4 on page 17). In order to enumerate elements in  $A$ , we only need to describe  $j$  and  $n$ . Thus, for any  $x \in A$ , we can encode  $x$  by its index  $i$  (in  $\log d(A)$  bits) in an enumeration of  $A$ . With  $n$  given, we can recover  $x$  from an encoding of  $j$  in  $K(j|n)$  bits, followed by  $i$ . This description of  $x$ , given  $n$ , takes  $\log d(A) + K(j|n) \leq n - \frac{1}{2} \log n + \log c_3 + K(j|n)$  bits. Let  $T$

be the Turing machine embodying this procedure to recover  $x$  given  $n$ . Choosing  $c = \log c_3 + c_T$ , we have

$$C(x|n) \leq C_T(x|n) + c_T \leq n - \frac{1}{2} \log n + K(j|n) + c.$$

□

**Example 2.6.2** For  $j = O(1)$  we have  $C(x|n) \leq n - \frac{1}{2} \log n + O(1)$  which is slightly stronger than the statement about unconditional  $C(x)$  in Example 2.6.1. For  $j = \sqrt{n}$  and  $j$  incompressible ( $K(j|n) \geq \frac{1}{2} \log n$ ), we have  $C(x|n) \leq n - O(1)$ . Only for such  $j$ 's is it possible that a number  $x$  is incompressible. ◇

## 2.6.2 Statistics of Blocks

The analysis up till now has been about the statistics of 0's and 1's. But in a normal infinite binary sequence according to Definition 2.6.2 each block of length  $k$  occurs with limiting frequency of  $2^{-k}$ . That is, blocks 00, 01, 10, and 11 should occur about equally often, and so on. Finite sequences will generally not be exactly normal, but normality will be a matter of degree. We investigate the block statistics for finite binary sequences.

**Definition 2.6.3** Let  $x = x_1 \dots x_n$  be a binary string of length  $n$ , and  $y$  a much smaller string of length  $l$ . Let  $p = 2^{-l}$  and  $\#y(x)$  be the number of (possibly overlapping) distinct occurrences of  $y$  in  $x$ . For convenience, we assume that  $x$  "wraps around" so that an occurrence of  $y$  starting at the end of  $x$  and continuing at the start also counts.

**Theorem 2.6.1** *Assume the notation of Definition 2.6.3 with  $l \leq \log n$ . There is a constant  $c$  such that for all  $n$  and  $x \in \{0, 1\}^n$ , if  $C(x) \geq n - \delta(n)$ , then*

$$|\#y(x) - pn| \leq \sqrt{\alpha pn},$$

with  $\alpha = [K(y|n) + \log l + \delta(n) + c]3l/\log e$ .

**Proof.** We prove by contradiction. Assume that  $n$  is divisible by  $l$ . (If it is not, then we can put  $x$  on a Procrustean bed to make its length divisible by  $l$  at the cost of having the above frequency estimate  $\#y(x)$  plus or minus an error term of at most  $l/2$ .) There are  $l$  ways of dividing (the ring)  $x$  into  $N = n/l$  contiguous equal-sized blocks, each of length  $l$ . For each such division  $i \in \{0, 1, \dots, l-1\}$ , let  $\#y(x, i)$  be the number of (now nonoverlapping) occurrences of block  $y$ . We apply the Chernoff bound, Equation 2.4, again. With  $A = \{x \in \{0, 1\}^n : |\#y(x, i) - pN| > m\}$  this

gives  $d(A) \leq 2^{n+1}e^{-m^2/3pN}$ . We choose  $m$  such that for some constant  $c$  to be determined later,

$$\frac{m^2 \log e}{3pN} = K(\langle y, i \rangle | n) + \delta(n) + c.$$

To describe an element  $x$  in  $A$ , we now need only to enumerate  $A$  and indicate the index of  $x$  in such an enumeration. The description contains the following items:

1. *A description used to enumerate A.* Given  $n - \delta(n)$ , we can retrieve  $n$  and  $\delta(n)$  using a self-delimiting description of at most  $c_1$  bits. To enumerate  $A$ , we also need to know  $i$  and  $y$ . Therefore, given  $n - \delta(n)$ , the required number of bits to enumerate  $A$  is at most

$$K(\langle y, i, \delta(n), n \rangle | n - \delta(n)) \leq K(\langle y, i \rangle | n) + c_1.$$

2. *A description of the index of x.* The number of bits to code the index  $j$  of  $x$  in  $A$  is

$$\begin{aligned} \log d(A) &\leq \log \left( 2^{n+1}e^{-m^2/3pN} \right) \\ &= n + 1 - \frac{m^2 \log e}{3pN} \\ &= n + 1 - K(\langle y, i \rangle | n) - \delta(n) - c. \end{aligned}$$

This total description takes at most  $n + 1 - \delta(n) - c + c_1$  bits. Let  $T$  be a Turing machine reconstructing  $x$  from these items. According to Theorem 2.1.1, therefore

$$C(x) \leq C_T(x) + c_T \leq n + 1 - \delta(n) - c + c_1 + c_T.$$

With  $c = c_1 + c_T + 2$  we have  $C(x) < n - \delta(n)$ , which contradicts the assumption of the theorem. Therefore,  $|\#y(x, i) - pN| \leq m$ , which in turn implies

$$|\#y(x, i) - pN| \leq \sqrt{\frac{K(\langle y, i \rangle | n) + \delta(n) + c}{\log e}} 3pN.$$

For each division  $i$  ( $0 \leq i \leq l - 1$ ) this inequality follows from the  $\delta(n)$  incompressibility of  $x$ . Notwithstanding the fact that occurrences of substrings in different divisions are dependent, the inequality holds for each division separately and independently. The theorem now follows by noting that  $|\#y(x) - pn| = \sum_{i=0}^{l-1} |\#y(x, i) - pN|$ ,  $K(\langle y, i \rangle | n) \leq K(y | n) + K(i | n) + O(1)$  and  $K(i | n) \leq \log l + O(1)$ .  $\square$

Similar to the analysis of blocks of length 1, the complexity of a string drops below its maximum in case some block  $y$  of length  $l$  occurs in one of the  $l$  block divisions, say  $i$ , with frequency exactly  $pN$  ( $p = 1/2^l$ ). Then we can point out  $x$  by giving  $n, y, i$ , and its index in a set of cardinality

$$\binom{N}{pN} (2^l - 1)^{N-pN} = O\left(\frac{2^{Nl}}{\sqrt{p(1-p)N}}\right).$$

Therefore,

$$C(x|\langle n, y \rangle) \leq n - \frac{1}{2} \log n + \frac{1}{2}(l + 3 \log l) + O(1).$$

### 2.6.3 Length of Runs

It is known from probability theory that in a randomly generated finite sequence the *expectation* of the length of the longest run of zeros or ones is pretty high. For each individual finite sequence with high Kolmogorov complexity we are *certain* that it contains each block (say, a run of zeros) up to a certain length.

**Theorem 2.6.2** *Let  $x$  of length  $n$  satisfy  $C(x) \geq n - \delta(n)$ . Then for sufficiently large  $n$ , each block  $y$  of length*

$$l = \log n - \log \log n - \log(\delta(n) + \log n) - O(1)$$

*occurs at least once in  $x$ .*

**Proof.** We are sure that  $y$  occurs at least once in  $x$  if  $\sqrt{\alpha p n}$  in Theorem 2.6.1 is less than  $p n$ . This is the case if  $\alpha < p n$ , that is,

$$\frac{K(y|n) + \log l + \delta(n) + O(1)}{\log e} 3l < p n.$$

$K(y|n)$  is majorized by  $l + 2 \log l + O(1)$  (since  $K(y|n) \leq K(y) + O(1)$ ) and  $p = 2^{-l}$  with  $l$  set at

$$l = \log n - \log(3\delta(n) \log n + 3 \log^2 n)$$

(which equals  $l$  in the statement of the theorem up to an additive constant). Substitution yields

$$\frac{l + 3 \log l + \delta(n) + O(1)}{\log e} 3l < 3(\delta(n) \log n + \log^2 n),$$

and it is easy to see that this holds for sufficiently large  $n$ .  $\square$

**Corollary 2.6.1** *If  $\delta(n) = O(\log n)$ , then each block of length  $\log n - 2 \log \log n - O(1)$  occurs at least once in  $x$ .*

In Lemma 6.9.1 we show that if  $C(x|n, p) \geq n$  then no substring of length greater than  $2\log n$  occurs (possibly overlapping) twice in  $x$ . Here,  $n = l(x)$ , and  $p$  is some fixed program used to reconstruct  $x$  from a description of length  $C(x|n, p)$  and  $n$ .

Analyzing the proof of Theorem 2.6.2 we can improve the corollary for low values of  $K(y|n)$ .

**Corollary 2.6.2** If  $\delta(n) = O(\log \log n)$ , then for each  $\epsilon > 0$  and every large enough  $n$ , every string  $x$  of length  $n$  contains an all-zero run  $y$  (for which  $K(y|n) = O(\log l)$ ) of length  $l = \log n - (1 + \epsilon) \log \log n + O(1)$ .

Since there are  $2^n(1 - O(1/\log n))$  strings  $x$  of length  $n$  with  $C(x) \geq n - \log \log n + O(1)$ , the *expected length of the longest run of consecutive zeros* if we flip a fair coin  $n$  times is at least  $l$  as in Corollary 2.6.2.

We show in what sense Theorem 2.6.2 is sharp. Let  $x = uvw$ ,  $l(x) = n$ , and  $C(x) \geq n - \delta(n)$ . We can describe  $x$  by giving

1. A description of  $v$  in  $K(v)$  bits;
2. The literal representation of  $uw$ ;
3. A description of  $l(u)$  in  $\log n + \log \log n + 2 \log \log \log n + O(1)$  bits.

Then, since we can find  $n$  by  $n = l(v) + l(uw)$ ,

$$\begin{aligned} C(x) &\leq n - l(v) + K(v) + \log n \\ &\quad + (1 + o(1)) \log \log n + O(1). \end{aligned} \tag{2.7}$$

Substitute  $C(x) = n - \delta(n)$  and  $K(v) = o(\log \log n)$  (choose  $v$  to be very regular) in Equation 2.7 to obtain

$$l(v) \leq \delta(n) + \log n + (1 + o(1)) \log \log n.$$

This means that for instance, for each  $\epsilon > 0$ , no maximally complex string  $x$  with  $C(x) = n + O(1)$  contains a run of zeros (or the initial binary digits of  $\pi$ ) of length  $\log n + (1 + \epsilon) \log \log n$  for  $n$  large enough and regular enough. By Corollary 2.6.2, on the other hand, such a string  $x$  *must* contain a run of zeros of length  $\log n - (1 + \epsilon) \log \log n + O(1)$ .

## Exercises

---

**2.6.1.** [24] The great majority of binary strings of length  $n$  have a number of 0's in between  $n/2 - \sqrt{n}$  and  $n/2 + \sqrt{n}$ . Show that there are  $x$ 's of length  $n$  with an order  $\sqrt{n}$  excess of 0's with  $C(x) = n + O(1)$ .

**2.6.2.** [29] Let  $\omega = \omega_1\omega_2\dots$  be an infinite binary sequence with

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \omega_i = p,$$

for some  $p$  between 0 and 1 (compare Section 1.9).

- (a) Show that if  $C(\omega_{1:n}) \sim n$ , then  $p = \frac{1}{2}$ .
- (b) Show that if  $p = \frac{1}{4}$ , then  $C(\omega_{1:n}) \leq 0.80n$  asymptotically.
- (c) Show that in general, if  $c = -p \log p - (1-p) \log(1-p)$ , then  $C(\omega_{1:n}) \leq cn + o(n)$ . If  $p$  is about  $\frac{1}{4}$ , then  $C(\omega_{1:n}) \leq 0.80n + o(n)$ .

*Comments.* Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987; attributed to A.N. Kolmogorov. Hint: Item (a), use Lemma 2.6.2.

**2.6.3.** [M35] A finite binary string  $x$  of length  $n$  is called  $\delta$ -random iff  $C(x|n) \geq n - \delta$ . A Turing machine place-selection rule  $R$  is a Turing machine that “selects” and outputs a (not necessarily consecutive) substring  $R(x)$  from its input  $x$ . If  $R$  is the  $k$ th Turing machine in the standard enumeration, then  $C(R) = C(k)$ .

Show that for any  $\epsilon > 0$ , there exist numbers  $n_0$  and  $\mu > 0$  such that if  $l(x) = n$ ,  $l(R(x)) = r \geq n_0$ ,  $\sum_{1 \leq i \leq r} R(x)_i = m$ , and  $(\delta + C(R|n))/r < \mu$ , then

$$\left| \frac{m}{r} - \frac{1}{2} \right| < \left( \frac{\delta + C(R|n) + 2.5 \log r}{(2 \log e - \epsilon)r} \right)^{1/2}.$$

*Comments.*  $\delta$ -random sequences were introduced by A.N. Kolmogorov, *Lecture Notes in Math.*, vol. 1021, Springer-Verlag, 1983, 1–5. He noted that “sequences satisfying this condition have for sufficiently small  $\delta$  the particular property of frequency stability in passing to subsequences.” In this exercise we supply some quantitative estimates of frequency stability. Source: E.A. Asarin, *SIAM Theory Probab. Appl.*, 32(1987), 507–508. This exercise is used by Asarin to show that  $\delta$ -random elements of certain finite sets obey familiar probability-theoretic distribution laws. In particular, for some particular finite sets each  $\delta$ -random element is  $\epsilon$ -normal; see also E.A. Asarin, *Soviet Math. Dokl.*, 36(1988), 109–112.

## 2.7 Algorithmic Properties of $C$

---

By algorithmic properties of  $C$  we mean properties with a recursion-theoretic flavor as in Section 1.7. We have already met a few of these. By Theorem 2.3.1, the greatest monotonic lower bound on  $C(n)$  is unbounded, but goes to infinity more slowly than any monotonic unbounded partial recursive function. By Theorem 2.3.2 the integer function  $C(n)$  is not recursive. Nonetheless, by Theorem 2.3.3 the integer function  $C(n)$  can be approximated arbitrarily closely from above (is co-enumerable). Unfortunately, at each stage of such an approximation process, for each size of error, there are infinitely many  $x$  such that the approximation of  $C(x)$  and its real value are at least this error apart.

In fact, a much stronger statement holds: For each total recursive function  $f$  with  $\lim_{x \rightarrow \infty} f(x) = \infty$  the set of  $x$  for which we can prove  $C(x) > f(x)$  is finite (Theorem 2.7.1 Item (iii) below). Thus, if we choose  $f(x) \ll l(x)$ , then we know that  $C(x) \gg f(x)$  for almost all  $x$  of each length, Theorem 2.2.1, yet we can prove this only for finitely many  $x$ .

- Theorem 2.7.1**
- (i) *The set  $A = \{(x, a) : C(x) \leq a\}$  is recursively enumerable, but not recursive.*
  - (ii) *Every partial recursive function  $\phi(x)$  that is a lower bound on  $C(x)$  is bounded.*
  - (iii) *Let  $f(x)$  be a total recursive function with  $g(x) \leq f(x) \leq l(x)$  for all  $x$  and some unbounded monotonic function  $g$ . Then the set  $B = \{x : C(x) \leq f(x)\}$  is simple. That is,  $B$  is recursively enumerable and the complement of  $B$  is infinite but does not contain an infinite recursively enumerable subset.*

**Proof.** (i) That  $A$  is recursively enumerable follows immediately from Theorem 2.3.3. However,  $A$  is not recursive. Namely, If  $A$  is recursive, then we can compute  $C(x)$  by asking consecutive questions “is  $C(x) \leq a$ ” for  $a := 0, 1, \dots$ , contradicting Theorem 2.3.2.

(ii) Let  $\phi$  be a partial recursive function and define  $D = \{x : \phi(x) \leq C(x)\}$ . If  $D$  is finite, there is nothing to prove. Assume  $D$  is infinite and  $\phi$  is unbounded, by way of contradiction. Recursively enumerate the domain of definition of  $\phi$  without repetition, and define a total recursive function  $g$  by  $g(n)$  that equals the least  $x$  in this enumeration such that  $\phi(x) \geq n$ . For each  $n$  there is such an  $x$ , by the contradictory assumption. If  $\phi = \phi_k$  in the standard effective enumeration  $\phi_1, \phi_2, \dots$  of the partial recursive functions, as in Section 1.7, then  $n \leq C(x) \leq l(n) + l(\bar{k})$ , up to a constant. For  $n$  large enough we have a contradiction.

(iii) That  $B$  is recursively enumerable follows from Item (i). The complement of  $B$  is infinite by Theorem 2.2.1. We prove that  $B$  is simple.

Let  $D$  be an infinite recursively enumerable set contained in the complement of  $B$ . The restriction  $f_D(x)$  of  $f(x)$  to  $D$  is a partial recursive lower bound for  $C(x)$ . By Item (ii) therefore  $f_D(x)$  is bounded. Since  $f(x)$  rises unboundedly with  $x$  this is possible only if  $D$  is finite.  $\square$

**Corollary 2.7.1** The set RAND defined by  $\{x : C(x) \geq l(x)\}$  is immune—it is infinite and has no infinite recursively enumerable subset. In fact, the proofs support stronger results in that the set  $B$  above is *effectively* simple and RAND is *effectively* immune (Exercise 2.7.6).

### 2.7.1 Undecidability by Incompressibility

This approach allows us to give a result similar to Lemma 1.7.6 on page 35, but with examples of undecidable statements that differ from the ones given by Gödel. Namely, for each formal system  $T$ , there is a constant  $c_T$  such that no formula of form “ $C(x) \geq c_T$ ” is provable in  $T$ .

#### Example 2.7.1

If  $T$  is an axiomatizable sound theory whose axioms and rules of inference require about  $k$  bits to describe, then  $T$  cannot be used to prove the randomness of any number much longer than  $k$  bits. If the system could prove randomness for a number much longer than  $k$  bits, then the *first* such proof (first in an unending enumeration of all proofs obtainable by repeated application of axioms and rules of inference) could be used to derive a contradiction: an approximately  $k$ -bit program to find and print out the specific random number mentioned in this proof, a number whose smallest program is by assumption considerably larger than  $k$  bits. Therefore, even though most strings are random, we will never be able to explicitly exhibit a string of reasonable size that demonstrably possesses this property. Formally,

- Let  $T$  be an axiomatizable theory ( $T$  is a recursively enumerable set consisting of axioms and provable formulas). Hence, there is a  $k$  such that  $T$  is describable in  $k$  bits:  $C(T) \leq k$ .
- Let  $T$  be sound: all formulas in  $T$  are true (in the standard model of the natural numbers).
- Let  $S_c(x)$  be a formula in  $T$  with the meaning “ $x$  is the lexicographically least binary string of length  $c$  with  $C(x) \geq c$ .” Here  $x$  is a formal parameter and  $c$  an explicit constant, so  $C(S_c) \leq \log c$  up to a fixed constant independent of  $T$  and  $c$ .

For each  $c$ , there exists an  $x$  such that  $S_c(x) = \text{true}$  is a true statement by a simple counting argument (Theorem 2.3.1). Moreover,  $S_c$  expresses that this  $x$  is unique. It is easy to see that combining the descriptions of  $T$ ,  $S_c$ , we obtain a description of this  $x$ . Namely, for each candidate string  $y$  of length  $c$ , we can decide  $S_c(y) = \text{true}$  (which is the case

for  $y = x$ ) or  $\neg S_c(y) = \text{true}$  (which is the case for  $y \neq x$ ) by simple enumeration of all proofs in  $T$ . (Here we use the soundness of  $T$ .) We need to distinguish the descriptions of  $T$  and  $S_c$ . We can do this by coding  $T$ 's description in self-delimiting format (see page 13). This takes not more than  $2k$  bits. Hence, for some fixed constant  $c'$  independent of  $T$  and  $c$ , we find  $C(x) \leq 2k + \log c + c'$ , which contradicts  $C(x) > c$  for all  $c > c_T$ , where  $c_T = 3k + c'$  for another constant  $c'$ . (A minor improvement of the argument shows that  $c_T = k + 2 \log k + c'$  suffices.)  $\diamond$

- Corollary 2.7.2** There is a recursively enumerable set  $B$  with an infinite complement, such that for every axiomatizable sound theory  $T$  there are only finitely many  $n$  for which the formula " $n \notin B$ " is both true and provable in  $T$ . (But with finitely many exceptions all infinitely many such formulas are true.)

**Proof.** Let  $B$  be the simple set in Theorem 2.7.1, Item (iii), and let  $\bar{B}$  be its complement. Clearly, the set  $D \subseteq \bar{B}$  of elements  $n$  that can be proved in  $T$  to belong to  $\bar{B}$  is recursively enumerable. Since  $B$  is simple its complement  $\bar{B}$  does not contain an infinite recursively enumerable subset. Therefore,  $D$  is finite, which proves the theorem.  $\square$

We have formulated Corollary 2.7.2 so as to bring out some similarities and differences with Lemma 1.7.6 as clearly as possible. As pointed out in Section 1.7, the set  $K_0$  used in Lemma 1.7.6 is "complete," whereas the set  $B$  used in Corollary 2.7.2 is "simple." According to generally accepted viewpoints in recursion theory, the set  $K_0$  is different from the set  $B$  in an essential way. Therefore, we can regard the proofs of the existence of undecidable statements in sufficiently rich axiomatizable theories by Lemma 1.7.6 and Corollary 2.7.2 as essentially different.

The set  $K_0$  is not only complete in the sense of Turing reducibility, it is also complete in the sense of "many-to-one reducibility." A set  $A$  is *many-to-one* reducible to a set  $B$  if there exists a recursive function  $f$  such that for all  $x$ ,  $x \in A$  iff  $f(x) \in B$ . A set  $A$  is complete in the sense of many-to-one reducibility, *m-complete* for short, if  $A$  is recursively enumerable and all recursively enumerable sets  $B$  are many-to-one reducible to  $A$ . As it turns out, *m*-complete sets are nonrecursive. This raises the question: are all nonrecursive recursively enumerable sets *m*-complete? The answer was given by E. Post in 1944 by introducing simple sets as the first examples of nonrecursive recursively enumerable sets that are not *m*-complete. (For the notions of reducibility, completeness, and simple sets see the exercises in Section 1.7, in particular, Exercises 1.7.16, 1.7.15.)

The set  $B$  used in Corollary 2.7.2 is a simple set and hence not *m*-complete. Therefore, although  $B$  is many-to-one reducible to  $K_0$ , the set  $K_0$  is not many-to-one reducible to  $B$ . This shows that  $K_0$  is of a so-called higher degree of unsolvability with respect to many-to-one reducibility than  $B$ , which is the substance of the viewpoint that they differ in an essential way.

In less formal terms the approaches are different because the first one can be viewed as a form of *Russell's Paradox* and the other one as a form of the *Richard-Berry Paradox*. Both paradoxes are described in [B. Russell and A.N. Whitehead, *Principia Mathematica*, Oxford, 1917]. While the first paradox formed the original incentive for the authors to supply the sophisticated logical foundation for set theory in the *Principia*, in a footnote they state that the second paradox “was suggested to us by Mr. G.G. Berry of the Bodleian Library.”

The paradox due to Bertrand Russell (1872–1970) arises when the collection of all sets that are not members of themselves is considered as a set. If this collection is a member of itself, then it contradicts the set definition, but if it is not a member of itself, then by the set definition it is a member of itself (which is a contradiction as well). There is a close connection between Russell's Paradox and the result of Gödel cited as Lemma 1.7.6. We have seen that this result was proved by reducing the halting problem in the form of  $K_0$  to the decision problem in a sufficiently strong, sound, axiomatized theory. Since  $K_0$  is  $m$ -complete, this shows that any problem shown to be unsolvable *this way* must have a degree of unsolvability at least as high as the maximal degree of unsolvability with respect to many-to-one reducibility as any recursively enumerable set.

The Richard-Berry Paradox is the definition of a number as “the least number that cannot be defined in fewer than twenty words.” Formalizing the notion of “definition” as the shortest program from which a number can be computed by the reference machine  $U$ , it turns out that the quoted statement (reformulated appropriately) is not an *effective description*. This was essentially what we did in the proof of Corollary 2.7.2, by reducing the set  $B$  to the decision problem in a sufficiently strong, sound, axiomatizable theory. But  $B$  is of a lesser degree of unsolvability with respect to many-to-one reducibility than is  $K_0$ . Therefore, showing undecidability of sufficiently rich axiomatizable theories using Kolmogorov complexity *this way* is essentially different from Gödel's original approach.

Gödel's first incompleteness theorem entails an *explicit* construction of a statement  $s$ , associated with each sufficiently strong, sound, axiomatized theory  $T$  that is undecidable in  $T$ . Formula  $s$  simply says of itself “I am unprovable in  $T$ .” In contrast, the construction in Corollary 2.7.2 says that for any sound axiomatized system  $T$  there is a constant  $c_T < \infty$  such that *all* true statements with the meaning “ $C(x) \geq c_T$ ” are unprovable in  $T$ . By Theorem 2.3.1 there are infinitely many such statements. Now suppose we have an effective procedure to find such constants for given theories, that is, a total recursive function  $\phi$  such that  $\phi(T) \geq c_T$  for all  $T$ . Then, unfortunately, Theorem 2.7.1, Item (ii), tells us that *no* effective procedure can determine for more than finitely many pairs  $(x, T)$  whether or not  $C(x) \geq \phi(T)$ . This shows that in general, although the undecidable statements based on Kolmogorov complexity are plentiful for each theory, we cannot explicitly construct them. Thus, the new approach entails loss of constructivity.

### 2.7.2 Barzdins's Lemma

Using Kolmogorov complexity one can quantify the distinction between recursively enumerable sets and recursive sets. Let  $A$  be a set of natural numbers.

**Definition 2.7.1** The *characteristic sequence* of  $A \subseteq \mathcal{N}$  is an infinite binary sequence  $\chi = \chi_1\chi_2\dots$  defined by

$$\chi_i = \begin{cases} 1 & \text{if } i \in A, \\ 0 & \text{otherwise.} \end{cases}$$

If  $A$  is recursively enumerable, and also its complement consisting of the  $i$ 's such that  $\chi_i = 0$  is recursively enumerable, then  $f(i) = \chi_i$  is recursive, and the conditional complexity  $C(\chi_{1:n}|n)$  is bounded by a fixed constant for all  $n$ . (The converse also holds, Exercise 2.3.4 on page 124.) But in the general case of recursively enumerable sets  $A$  the complexity  $C(\chi_{1:n}|n)$  can grow unboundedly with  $n$ . However, this growth is at best logarithmically slow, which shows that such characteristic sequences are very nonrandom. For instance, they are not random in the sense of Martin-Löf according to Theorem 2.5.4. The result below is known as *Barzdins's Lemma*. (Actually, J.M. Barzdins proved the sharper version of Exercise 2.7.2 on page 173.)

**Theorem 2.7.2**

- (i) Any characteristic sequence  $\chi$  of a recursively enumerable set  $A$  satisfies  $C(\chi_{1:n}|n) \leq \log n + c$  for all  $n$ , where  $c$  is a constant dependent on  $A$  (but not on  $n$ ).
- (ii) Moreover, there is a recursively enumerable set such that its characteristic sequence  $\chi$  satisfies  $C(\chi_{1:n}) \geq \log n$  for all  $n$ .

**Proof.** (i) Since  $A$  is recursively enumerable, there is a partial recursive function  $\phi$  such that  $A = \{x : \phi(x) < \infty\}$ . Dovetail the computations of  $\phi(1), \phi(2), \dots$ . This way we enumerate  $A$  without repetitions in the order in which the computations of the  $\phi(i)$ 's terminate. The prefix  $\chi_{1:n}$  can be reconstructed from the number  $m$  of 1's it contains. For if we know  $m$ , then it suffices to use  $\phi$  to enumerate the elements of  $A$  until we have found  $m$  distinct such elements less than or equal to  $n$ . If the set of these elements is  $B = \{a_1, a_2, \dots, a_m\}$ , then by assumption these are all elements in  $A$  that do not exceed  $n$ . Hence, from  $B$  we can reconstruct all 1's in  $\chi_{1:n}$ , and the remaining positions must be the 0's. This way we can reconstruct  $\chi_{1:n}$ , given  $n$ , from a description of  $\phi$  and  $m$ . Since  $m \leq n$  and  $C(\phi) < \infty$  we have proved (i).

(ii) Let  $\phi_0$  be the additively optimal function  $\phi_0$  of Theorem 2.1.1. Define  $\chi = \chi_1\chi_2\dots$  by

$$\chi_i = \begin{cases} 1 & \text{if } \phi_0(i, i) = 0, \\ 0 & \text{if } \phi_0(i, i) \neq 0 \text{ or } \phi_0(i, i) = \infty. \end{cases}$$

Obviously,  $\chi$  is the characteristic sequence of a recursively enumerable subset of the natural numbers. We prove that  $\chi$  satisfies the property stated in the theorem. For suppose to the contrary that  $C(\chi_{1:n}) < \log(n)$  for some  $n$ . This means there is a short program  $p$ , of length less than  $\log(n)$  that computes  $\chi_{1:n}$ . Since  $p < n$  this implies that  $\phi_0(p, p) = \chi_p$ , which contradicts the definition of  $\chi_p$ .  $\square$

The converse of Theorem 2.7.2 Item (i) does not hold in general. This follows by the construction of a meager set that is not recursively enumerable. For instance, let  $\chi$  be the characteristic sequence of a set  $A$  and  $C(\chi_{1:n}|n) \geq n - c$  for infinitely many  $n$  and a fixed constant  $c$ . By Theorem 2.5.5 such strings are abundant, and by Theorem 2.7.2 Item (i) we find that  $A$  is not recursively enumerable. Construct a sequence  $\zeta$  by  $\zeta = \chi_1\alpha_1\chi_2\alpha_2\dots$  with  $\alpha_i = 0^{f(i)}$ , where  $f(i)$  is some fast-growing total recursive function with an inverse. Obviously, if  $\zeta$  is the characteristic sequence of set  $B$ , then  $B$  is not recursively enumerable. But also there is now another constant  $c$  such that  $C(\zeta_{1:n}|n) \leq f^{-1}(n) + c$  for all  $n$ . Choosing  $f$  such that  $\log \log f(n) = n$  gives

$$C(\zeta_{1:n}|n) \leq \log \log n + O(1).$$

Theorem 2.7.2 Item (i) cannot be improved to the unconditional “ $C(\chi_{1:n}) \leq \log n + c$  for all  $n$  and some  $c$ ,” since all  $\chi$ ’s satisfying this are recursive (and hence the corresponding sets  $A$  are recursive) by Exercise 2.3.4 on page 124. Theorem 2.7.2 Item (ii) cannot be improved to the conditional “ $C(\chi_{1:n}|n) \geq \log n$  for all  $n$ ” by Exercise 2.7.3, page 173.

**Example 2.7.2** Diophantine equations are algebraic equations of the form  $X = 0$ , where  $X$  is built up from nonnegative integer variables and nonnegative integer constants by a finite number of additions ( $A + B$ ) and multiplications ( $A \times B$ ). The best known examples are  $x^n + y^n = z^n$ , where  $n = 1, 2, \dots$ .

Pierre de Fermat (1601–1665) has stated that this equation has no solution in positive integers  $x$ ,  $y$ , and  $z$  for  $n$  an integer greater than 2. (For  $n = 2$  there exist solutions, for instance  $3^2 + 4^2 = 5^2$ .) However, he did not supply a proof of this assertion, often called *Fermat’s Last Theorem*. After 350 years of withstanding concerted attempts to come up with a proof or disproof, the problem had become a celebrity among unsolved mathematical problems. (At this time of writing, November 1996, it appears that Andrew Wiles has settled the problem affirmatively.) Suppose we substitute all possible values for  $x, y, z$  with  $x + y + z \leq n$ , for  $n = 3, 4, \dots$ . This way, we recursively enumerate all solutions of Fermat’s equation. Hence, such a process will eventually give a counterexample to Fermat’s conjecture *if one exists*, but the process will never yield conclusive evidence if the conjecture happens to be true.

In his famous address to the International Mathematical Congress in 1900, D. Hilbert proposed twenty-three mathematical problems as a program to direct the mathematical efforts in the twentieth century.

The tenth problem asks for an algorithm that given an arbitrary Diophantine equation, produces either an integer solution for this equation or indicates that no such solution exists. After a great deal of preliminary work by other mathematicians, the Russian mathematician Yu.V. Matijasevich finally showed that no such algorithm exists. Suppose we weaken the problem as follows. First, effectively enumerate all Diophantine equations, and consider the characteristic sequence  $\Delta = \Delta_1 \Delta_2 \dots$ , defined by  $\Delta_i = 1$  if the  $i$ th Diophantine equation is solvable, and 0 otherwise. Then  $C(\Delta_{1:n}) \leq n + O(1)$ . But the theorem above shows that  $C(\Delta_{1:n}|n) \leq \log n + c$ , for some fixed constant  $c$ . The nonrandomness of the characteristic sequence means that the solvability of Diophantine equations is highly interdependent—it is impossible for a random sequence of them to be solvable and the remainder unsolvable. ◇

- Example 2.7.3** In the proof of Theorem 2.7.2, Item (ii), we used a set recursively isomorphic to the halting set  $K_0 = \{\langle x, y \rangle : \phi_x(y) < \infty\}$ . In fact, almost every recursively enumerable set that is complete under the usual reducibilities would have done. We can use this to obtain natural examples for incompressible finite strings. Let  $d$  be the number of elements in  $K_0$  that are less than  $2^n$ . Then by running all the computations of all  $\phi_x(y)$  with  $z < 2^n$ ,  $z = \langle x, y \rangle$ , in parallel until  $d$  of them have halted, we effectively find all computations among them that halt. That is, if  $\chi$  is the characteristic sequence of  $K_0$ , then we can effectively compute  $\chi_{1:m}$  (where  $m = 2^n$ ) from  $d$ . The program  $p$  for this computation has length  $l(p) \leq l(d) + c \leq n + c$ , for some fixed constant  $c$  independent of  $n$ . By Theorem 2.7.2, Item (ii), the shortest program from which we can compute  $\chi_{1:m}$  has length at least  $n$ . Hence, there is another constant  $c$  such that  $p$  is  $c$ -incompressible. ◇

## Exercises

- 2.7.1.** [10] Show that there exists a constant  $c$  such that  $C(0^n|n) \leq c$  for all  $n$ , and  $C(0^n) \geq \log n - c$  for infinitely many  $n$ .

*Comments.* Hint: use Barzdins's Lemma (Theorem 2.7.2).

- 2.7.2.** • [26] Let  $A \subseteq \mathcal{N}$  be a recursively enumerable set, and let  $\chi = \chi_1 \chi_2 \dots$  be its characteristic sequence. We use the “uniform complexity”  $C(\chi_{1:n}; n)$  of Exercise 2.3.3 on page 124.

(a) Show that  $C(\chi_{1:n}; n) \leq \log n + O(1)$  for all  $A$  and  $n$ .

(b) Show that there exists an  $A$  such that  $C(\chi_{1:n}; n) \geq \log n$  for all  $n$ .

*Comments.* This implies Theorem 2.7.2. It is the original Barzdins's Lemma. Source: J.M. Barzdins, *Soviet Math. Dokl.*, 9(1968), 1251–1254.

- 2.7.3.** [27] Is there a symmetric form of Theorem 2.7.2 (Barzdins's Lemma) using only conditional complexities? The answer is negative.

Show that there is no recursively enumerable set such that its characteristic sequence  $\chi$  satisfies  $C(\chi_{1:n}|n) \geq \log n + O(1)$  for all  $n$ .

*Comments.* Hint: let  $\chi$  be the characteristic sequence of a recursively enumerable set  $A$ . Consider  $C(\chi_{1:f(n)}|f(n))$ , with  $\chi_{1:f(n)}$  containing exactly  $2^{2^n}$  ones. Then,  $\log f(n) \geq 2^n$ . But using the partial recursive function enumerating  $A$ , we can compute  $\chi_{1:f(n)}$ , given  $f(n)$ , from just the value of  $n$ . Hence, we have a program of  $\log n + O(1)$  bits for  $\chi_{1:f(n)}$ . Compare this to Barzdins's Lemma (Theorem 2.7.2) and Exercise 2.7.2. Source: R.M. Solovay, sci.logic electronic newsgroup, 24 November, 1989.

**2.7.4.** [34] Is there a symmetric form of Theorem 2.7.2 (Barzdins's Lemma) using only unconditional complexities? The answer is negative. Show that there is a recursively enumerable set  $A \subseteq \mathcal{N}$  and a constant  $c$  such that its characteristic sequence  $\chi$  satisfies  $C(\chi_{1:n}) \geq 2 \log n - c$  for infinitely many  $n$ .

*Comments.* First note the easy fact that the Kolmogorov complexity of  $\chi_{1:n}$  is at most  $2 \log n + O(1)$  for all  $n$  ( $\leq \log n$  bits to specify  $n$  and  $\leq \log n$  bits to specify  $k = \sum_{i=1}^n \chi_i$ ). Hint: partition  $\mathcal{N}$  into exponentially increasing half-open intervals  $I_k = (t_k, 2^{t_k}]$  with  $t_0 = 0$ . Note that  $\log(2^{t_k} - t_k) = 2 \log t_{k+1} - 2 - o(1)$  for  $k \rightarrow \infty$ . Use increasingly precise approximations of  $C(\chi_{1:n})$  for  $n \in I_k$  for increasing  $k$  to enumerate  $A$ . Source: R.M. Solovay, sci.logic electronic newsgroup, 24 November, 1989; posed as open problem [O39] in the first printing of this book; solved by M. Kummer [*SIAM J. Comput.*, 25:6(1996), 1123–1143].

**2.7.5.** [25] Prove the following strange fact (Kamae's Theorem). For each natural number  $m$  there is a string  $x$  such that for all but finitely many strings  $y$ ,  $C(x) - C(x|y) \geq m$ . That is,

*Comments.* That is, there exist finite objects such that almost all finite objects contain a large amount of information about them. Hint:  $x$  must be such that almost all large numbers contain much information about  $x$ . Let  $c$  be a large enough fixed constant. Let  $A$  be a recursively enumerable set of integers, and let  $\alpha_1\alpha_2\dots$  be the characteristic sequence of  $A$ . Set  $x = x(k) = \alpha_1\alpha_2\dots\alpha_h$ , where  $h = 2^k$ . By Barzdins's Lemma, Theorem 2.7.2 on page 171, we can assume  $C(x(k)) \geq k$ . Enumerate  $A$  without repetition as  $b_1, b_2, \dots$ . Let  $m(k) = \max\{i : b_i \leq 2^k\}$ . Then for any integer  $y \geq m(k)$  we have  $C(x(k)|y) \leq \log k + c$ . Namely, using  $y$  we can enumerate  $b_1, b_2, \dots, b_t$  and with  $\log k$  extra information describing  $k$  we can find  $x(k)$ . Therefore,  $C(x) - C(x|y) \geq k - \log k - c$ . Source: T. Kamae, *Osaka J. Math.*, 10(1973), 305–307. See also Exercise 2.2.12.

**2.7.6.** [25] Consider an enumeration  $W_1, W_2, \dots$  of all recursively enumerable sets. A simple set  $A$  is *effectively simple* if there is a recursive function  $f$  such that  $W_i \subseteq \bar{A}$  implies that  $d(W_i) \leq f(i)$  (where  $\bar{A}$  is the complement of  $A$ ). The set  $\bar{B}$  is called *effectively immune*.

- (a) Show that the set  $B$  defined in Theorem 2.7.1 is effectively simple.
- (b) Show that the set RAND defined by  $\{x : C(x) \geq l(x)\}$  is effectively immune.
- (c) Show that Item (a) implies that  $B$  is Turing-complete for the recursively enumerable sets.

*Comments.* Hint: the proof of Theorem 2.7.1 Item (iii) showing that  $B$  is simple actually shows that  $B$  is also effectively simple, which demonstrates Item (a). For Item (c) see for example P. Odifreddi, *Classical Recursion Theory*, North-Holland, 1989.

**2.7.7.** [32] Let  $\phi_1, \phi_2, \dots$  be the standard enumeration of partial recursive functions. The *diagonal halting set* is  $\{x : \phi_x(x) < \infty\}$  (also denoted as  $K$ ). The *Kolmogorov set* is  $\{(x, y) : C(x) \leq y\}$ . We assume familiarity with notions in Exercise 1.7.16. To say a set  $A$  is *recursive in* a set  $B$  is the same as saying  $A$  is Turing reducible to  $B$ .

- (a) Show that the diagonal halting set is recursive in the Kolmogorov set.
- (b) Show that the Kolmogorov set is recursive in the diagonal halting set.
- (c) Show that the Kolmogorov set is Turing-complete for the recursively enumerable sets.

*Comments.* This means that if we can solve the halting problem, then we can compute  $C$ , and conversely. Hint for Item (a): given  $x$  we want to know if  $x \in K$ , that is, whether  $T_x(x)$  halts. Let  $l(\langle x, T_x \rangle) = n$ . Now use the Kolmogorov set to recursively find the least number  $t$  such that for all  $y$  with  $l(y) = 2n$  and  $C(y) < 2n$  the reference universal machine  $U$  computes  $y$  from some program of length less than  $2n$  in at most  $t$  steps. Note that  $t$  is found with some organized dovetailing. Claim:  $T_x(x)$  halts iff  $T_x(x)$  halts within  $t$  steps (hence we can see whether  $T_x(x)$  halts). If not, then we can use  $T_x(x)$  as a “clock” and run the same dovetailing process as above, but now we produce a string of complexity  $2n$  via a description of length  $n$ . Source: Attributed to P. Gács by W. Gasarch, personal communication February 13, 1992.

**2.7.8.** [42] We can express the nonrecursivity of  $C(x)$  in terms of  $C(C(x)|x)$ , which measures what we may call the *complexity of the complexity function*. Denote  $l(x)$  by  $n$ .

- (a) Prove the *upper bound*  $C(C(x)|x) \leq \log n + O(1)$ .
- (b) Show the following *lower bound*: For each length  $n$  there are strings  $x$  such that

$$C(C(x)|x) \geq \log n - \log \log n - O(1).$$

*Comments.* This means that  $x$  only marginally helps to compute  $C(x)$ ; most information in  $C(x)$  is extra information related to the halting problem. Hint for Item (b): same proof as in Section 3.8. Source: P. Gács, *Soviet Math. Dokl.*, 15(1974), 1477–1480.

**2.7.9.** [44] Show that every infinite sequence is Turing-reducible (Exercise 1.7.16, page 43, with sets replaced by characteristic sequences of sets) to an infinite sequence that is random with respect to the uniform measure.

*Comments.* Charles Bennett raised the question of whether every infinite binary sequence can be obtained from an incompressible one by a Turing machine. He proved this for a special case. Philosophically, the result implied in the exercise allows us to view even very pathological sequences as the result of two relatively well-understood notions, to wit, the completely chaotic outcome of coin-tossing and a Turing machine transducer algorithm. Source: P. Gács, *Inform. Contr.*, 70(1986), 186–192.

**2.7.10.** [30] This exercise assumes knowledge of the notion of *Turing degree*, Exercise 1.7.16. Every Turing degree contains a set  $A$  such that if  $\chi$  is the characteristic sequence of  $A$ , then  $C(\chi_{1:n}|n) \leq \log n$  for all  $n$ .

*Comments.* Hence, high degree of unsolvability of a set does not imply high Kolmogorov complexity of the associated characteristic sequence. Hint: call a set  $B$  *semirecursive* if there exists a recursive linear ordering  $<_B$  of  $\mathcal{N}$  such that there exists a lower cut element  $y$  such that  $B = \{x : x \leq_B y\}$ . For any set  $A$  there is a semirecursive set  $B$  such that  $B \equiv_T A$  [C.G. Jockusch, *Trans. AMS* 131(1968), 420–436]. Every semirecursive set  $B$  has a characteristic sequence  $\chi$  of  $(\mathcal{N}, <_B)$  such that  $C(\chi_{1:n}|n) \leq \log n + c$ , by the same proof as Theorem 2.7.2, Item (i). Since  $<_B$  is recursive, the same property holds for the usual characteristic sequence of  $B$ . Source: W. Gasarch, Letter, August 1988. See also R.P. Daley, *J. Comput. System Sci.*, 9(1974), 151–163; *Math. Systems Theory*, 9(1975), 83–94; *Inform. Contr.*, 44(1980), 236–244.

**2.7.11.** [42] Use Kolmogorov complexity to show the existence of Turing degrees of unsolvability (Exercise 1.7.16) in between the recursive sets and Turing complete sets (like  $K_0$ ).

*Comments.* Source: R.P. Daley, *J. Symb. Logic*, 46(1981), 460–474; *Inform. Contr.*, 52(1982), 52–67.

**2.7.12.** [39] We assume familiarity with the notion of truth-table reducibility. Let  $\chi$  be the characteristic sequence of a recursively enumerable set  $A$ . Here  $C(\chi_{1:n}; n)$  is the uniform complexity of Exercise 2.3.3.

(a) Show that  $A$  is complete under weak truth-table reducibility iff for some unbounded total recursive function  $f(n)$ , we have  $C(\chi_{1:n}; n) \geq f(n)$ .

- (b) Show that  $A$  is complete under Turing reducibility iff  $C(\chi_{1:n}; n) \geq f(n)$  for some unbounded total function  $f$  recursive in  $A$ .

*Comments.* For resource-bounded versions of Kolmogorov complexity the situation is quite different. Source: M.I. Kanovich, *Soviet Math. Dokl.*, 10(1969), 700–701; 11(1970), 1224–1228.

**2.7.13.** [20] Define the *state complexity*  $S(x)$  of a finite binary string  $x$  as the least  $n$  such that there is a Turing machine with  $n$  states that started in the standard initial conditions of empty tape and distinguished start state will eventually halt with  $x$  on its output tape. All machines considered are of the original model as in Section 1.7. Define  $B = \{\langle x, y \rangle : S(x) \leq y\}$ .

- (a) Prove that  $B$  is recursively enumerable but not recursive.  
 (b) Prove that  $B$  is Turing complete (in the sense of Exercise 1.7.16).

*Comments.* Suppose our Turing machines use an  $m$ -letter alphabet. Let  $T_m(x)$  denote the complexity of  $x$  in terms of the minimal number of internal states of a Turing machine. Then

$$T_m(x) \sim C(x)/(m - 1) \log C(x).$$

Source: problem by J. Andrews, electronic news, June 24, 1988; solutions by V. Pratt, R.M. Solovay, electronic news, June 1988.

**2.7.14.** [22] Show that the set  $K_0$  used in Lemma 1.7.6 is not many-to-one reducible to the set  $B$  featuring in Corollary 2.7.2 on page 169, while  $B$  is many-to-one reducible to  $K_0$ .

*Comments.* Hint: use Exercise 1.7.16.  $K_0$  is  $m$ -complete, while  $B$  is simple and hence not  $m$ -complete. The set  $K_0$  is of a *higher degree of unsolvability with respect to many-to-one reducibility* than  $B$ .

**2.7.15.** [32] Show that there exists an immune set  $I$  (a set without infinite recursively enumerable subset, for instance the complement of a set  $B$  as in Theorem 2.7.1, Item (iii)), such that there is a probabilistic machine that computes the characteristic function of some infinite subset of  $I$ .

*Comments.* Hint: use Theorems 2.5.4, 2.7.1, and the following framework. A *probabilistic* machine is just like a deterministic machine except that at some steps there are several actions (instead of a single action) that the machine can perform with given probabilities. For simplicity assume that there are exactly two possible actions, each with probability  $\frac{1}{2}$ . (That is, at each such choice the machine flips a coin.) A probabilistic machine computes a function  $\phi$  with probability  $p$  means that the machine with input  $x$  halts with output  $\phi(x)$  with probability  $p$ . We usually assume  $p > \frac{1}{2}$ . It can be shown that a machine with any

value  $p$  in between zero and one can be simulated by a machine with value  $p$  close to one. It turns out that  $\phi$  is computable by a probabilistic machine iff  $\phi$  is partial recursive [K. de Leeuw, et al., pp. 183–212 in: *Automata Studies*, C.E. Shannon and J. McCarthy (Eds), Princeton Univ. Press, 1956]. This result is often interpreted as showing that probabilistic machines cannot perform tasks that are impossible for deterministic machines. But a task may not consist of only finding an unambiguous value, but may consist in finding any one of a set of possible values. In this form there are obviously tasks that deterministic machines cannot do but that probabilistic machines can do, like the construction of a non-recursive sequence. Source: A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124, attributed to J.M. Barzdins.

**2.7.16.** [37] A set  $H$  of natural numbers is called *hyperimmune* if there is no total recursive function  $f$  such that  $f(i) > h_i$  for all  $i$ , where  $h_i$  is the  $i$ th element of  $H$  in increasing order. That is,  $H$  is immune (Exercise 2.7.15, page 177) but the variety of immunity of  $H$  is due to the fact that the function that enumerates  $H$ 's elements in increasing order of size grows faster than any recursive function. Prove the following:

- (a) Every hyperimmune set  $H$  contains an infinite subset whose characteristic sequence is not computable by a probabilistic machine.
- (b) However, there is a probabilistic machine that computes the characteristic sequence of some hyperimmune set.

*Comments.* Source: A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124, attribute Item (a) to V.N. Agafonov and L.A. Levin, and Item (b) to N.V. Petri. See also Theorems 2.2 and 2.3 in P. Gács, *Theoret. Comput. Sci.*, 22(1983), 71–93.

**2.7.17.** [19] We define a variant of the Busy Beaver function  $BB(n)$  in Exercise 1.7.19. Let  $BC(n)$  be the largest natural number  $m$  such that  $C(m) \leq n$ . Let  $\phi_1, \phi_2, \dots$  be the standard effective enumeration of partial recursive functions.

- (a) Show that  $BC(n) > \phi(n)$ , where  $\phi = \phi_k$ , for all  $n \geq 2C(k) + O(1)$ .
- (b) Show that  $BC(n)$  is not a recursive function.
- (c) Show that the nonrecursiveness of  $BC(n)$  is equivalent to the unsolvability of the halting problem (Lemma 1.7.5, page 34).
- (d) Let  $F$  be an axiomatizable sound formal theory that can be described completely (axioms, inference rules, ...) in  $m$  bits. Show that  $F$  cannot be used to determine the value of  $BC(n)$  for any  $n > m + O(1)$ .

*Comments.* Hint for Item (a):  $C(\phi(n) + 1) \leq C(k) + C(n - C(k)) + 2 \log \min(C(k), n - C(k)) + O(1)$ . Then,  $C(\phi(n) + 1) \leq n/2 + 2 \log n + O(1)$ . Hence,  $BC(n) \geq \phi(n) + 1$ , for all  $n$  from  $2C(k) + O(1)$  onwards. Hint for

Item (b): it grows faster than any recursive function. Hint for Item (c): if the halting problem were solvable, we could compute  $BC(n)$  from the outputs of all halting programs of length at most  $n$ . Conversely, every halting program  $p$  halts within  $BC(n)$  steps, for  $n \geq l(p) + O(1)$ . So recursiveness of  $BC$  implies the solvability of the halting problem. This exercise is an application of Theorem 2.3.1, Item (iii). In fact,  $BC(n)$  is some sort of inverse function of  $m(n)$ , the greatest monotonic increasing function bounding  $C(n)$  from below. Source: G.J. Chaitin, pp. 108–111 in: *Open Problems in Communication and Computation*, T.M. Cover, B. Gobinath (Eds.), Springer-Verlag, 1988.

## 2.8 Algorithmic Information Theory

One interpretation of the complexity  $C(x)$  is as the quantity of information needed for the recovery of an object  $x$  from scratch. Similarly, the conditional complexity  $C(x|y)$  quantifies the information needed to recover  $x$  given only  $y$ . Hence the complexity is “absolute information” in an object. Can we obtain similar laws for complexity-based “absolute information theory” as we did for the probability-based information theory of Section 1.11?

If  $C(x|y)$  is much less than  $C(x)$ , then we may interpret this as an indication that  $y$  contains a lot of information about  $x$ .

**Definition 2.8.1** The *algorithmic information* about  $y$  contained in  $x$  is defined as

$$I_C(x : y) = C(y) - C(y|x).$$

Choosing reference function  $\phi_0$  in Theorem 2.1.1 such that  $\phi_0(x, \epsilon) = x$ ,

$$C(x|x) = 0 \text{ and } I_C(x : x) = C(x).$$

By the additive optimality of  $\phi_0$ , these equations hold up to an additive constant independent of  $x$ , for any reference function  $\phi_0$ . In this way we can view the complexity  $C(x)$  as the algorithmic information contained in an object about itself. For applications, this definition of the quantity of information has the advantage that it refers to individual objects, and not to objects treated as elements of a set of objects with a probability distribution given on it, as in Section 1.11.

Does the new definition have the desirable properties that hold for the analogous quantities in classic information theory? We know that equality and inequality can hold only up to additive constants, according to the indeterminacy in the Invariance Theorem 2.1.1. Intuitively, it is reasonable to require that

$$I_C(x : y) \geq 0,$$

up to an additive fixed constant independent of  $x$  and  $y$ . Formally, this follows easily from the definition of  $I_C(x : y)$ , by noting that  $C(y) \geq C(y|x)$  up to an independent additive constant.

### 2.8.1 Entropy, Information, and Complexity

The major point we have to address is the relation between the Kolmogorov complexity and Shannon's entropy as defined in Section 1.11. Briefly, classic information theory says a random variable  $X$  distributed according to  $P(X = x)$  has entropy (complexity)  $H(X) = -\sum P(X = x) \log P(X = x)$ , where the interpretation is that  $H(X)$  bits are on the average sufficient to describe an outcome  $x$ . Algorithmic complexity says that an object  $x$  has complexity, or algorithmic information,  $C(x)$  equal to the minimum length of a binary program for  $x$ . It is a beautiful fact that these two notions turn out to be much the same. The statement below may be called the Theorem of Equality between Stochastic Entropy and Expected Algorithmic Complexity. (The theorem actually gives an inequality, but together with the simple argument in Example 2.8.1 on page 181 this turns into an asymptotic equality.)

**Theorem 2.8.1** *Let  $x = y_1y_2\dots y_m$  be a finite binary string with  $l(y_1) = \dots = l(y_m) = n$ . Let the frequency of occurrence of the binary representation of  $k = 1, 2, 3, \dots, 2^n$  as a  $y$ -block be denoted by  $p_k = d(\{i : y_i = k\})/m$ . Then up to an independent additive constant,*

$$C(x) \leq m(H + \epsilon(m)),$$

*with  $H = -\sum p_k \log p_k$ , the sum taken for  $k$  from 1 to  $2^n$ , and  $\epsilon(m) = 2^{n+1}l(m)/m$ . Note that  $\epsilon(m) \rightarrow 0$  as  $m \rightarrow \infty$  with  $n$  fixed.*

**Proof.** Denote  $2^n$  by  $N$ . To reconstruct  $x$  it suffices to know the number  $s_k = p_k m$  of occurrences of  $k$  as a  $y_i$  in  $x$ ,  $k = 1, 2, \dots, N$ , together with  $x$ 's serial number  $j$  in the ordered set of all strings satisfying these constraints. That is, we can recover  $x$  from  $s_1, \dots, s_N, j$ . Therefore, up to an independent fixed constant,  $C(x) \leq 2l(s_1) + \dots + 2l(s_N) + l(j)$ . By construction,

$$j \leq \binom{m}{s_1, \dots, s_N},$$

a multinomial coefficient (page 10). Since also each  $s_k \leq m$ , we find

$$C(x) \leq 2^{n+1}l(m) + l\left(\binom{m}{s_1, \dots, s_N}\right).$$

Writing the multinomial coefficient in factorials, and using Stirling's approximation, Exercise 1.5.4 on page 17, to approximate  $j$ , the theorem is proved.  $\square$

In Theorem 2.8.1 we have separated the frequency regularities from the remaining regularities. The entropy component  $mH$  measures only the frequency regularities, while the remaining component  $m\epsilon(m)$  accounts for all remaining factors.

**Example 2.8.1** For  $x$  representing the sequence of outcomes of independent trials, the inequality in Theorem 2.8.1 can be replaced by asymptotic equality with high probability.

We give a simple example to show the relation between the entropy  $H$  of a stochastic source  $X$ , emitting  $n$ -length binary words with probability  $P(X = x)$  of outcome  $x$ , and the complexity  $C$ . Let  $P(X = x) = 2^{-n}$  be the uniform probability distribution on the outcomes of length  $n$ . The entropy  $H$  in Theorem 2.8.1 is, according to Section 1.11, especially designed to measure frequency regularities. We show that it is asymptotically equal to the expected complexity of a string. By Theorem 2.2.1 almost all  $x$  are  $c$ -incompressible, that is, there are  $2^n(1 - 2^{-c+1})$  many  $x$ 's that have  $C(x) \geq n - c$ . A simple computation shows that the entropy  $H(X) = -\sum_{l(x)=n} P(X = x) \log P(X = x)$  is asymptotically equal to the *expected complexity*  $\mathbf{E} = \sum_{l(x)=n} P(X = x)C(x)$  of an  $n$ -length word. Namely, we find

$$\frac{n}{n + O(1)} \leq \frac{H(X)}{\mathbf{E}} < \frac{n}{(1 - 2^{-c+1})(n - c)}.$$

Substitute  $c = \log n$  to obtain

$$\lim_{n \rightarrow \infty} \frac{H(X)}{\mathbf{E}} = 1.$$

◇

In Section 4.3.5 we show the following generalization: Let  $p$  be a *recursive* probability distribution on  $\mathcal{N}$ , that is, there is a Turing machine computing the function  $p$ . Let moreover  $K(x)$  be the prefix complexity as defined in Section 1.11.1. Then  $-\log p(x)$  and  $K(x)$  are close to each other with high probability. Since  $|K(x) - C(x)| \leq 2\log C(x)$  by Example 3.1.3 on page 194, also  $-\log p(x)$  and  $C(x)$  are close to each other with high probability.

In particular, the entropy  $H = -\sum_{l(x)=n} P(x) \log P(x)$  of the distribution  $P$  is asymptotically equal to the expected complexity  $\sum_{l(x)=n} P(x)K(x)$  of words of length  $n$ , see Section 8.1.

Because we saw a few lines above that  $K(x)$  and  $C(x)$  are equal up to a logarithmic additive term, the expected plain complexity  $C(\cdot)$  is also asymptotically equal to the entropy,

$$\sum_x P(x)C(x) \sim -\sum_x P(x) \log P(x).$$

Thus, the intended interpretation of complexity  $C(x)$  as a measure of the information content of an individual object  $x$  is supported by a tight quantitative relationship to Shannon's probabilistic notion.

### 2.8.2 Symmetry of Information

Is algorithmic information symmetric? In Section 1.11 we have noted that in Shannon's information theory the mutual information in one random variable about another one is symmetric. While the equation  $I_C(x : y) = I_C(y : x)$  cannot be expected to hold exactly, a priori it can be expected to hold up to a constant related to the choice of reference function  $\phi_0$  in Theorem 2.1.1. However, with the current definitions, information turns out to be symmetric only up to a logarithmic additive term.

**Example 2.8.2** By Theorem 2.2.1, there is a binary string  $x$  of each length  $n$  such that  $C(x|n) \geq n$ . Similarly, there are infinitely many  $n$  such that  $C(n) \geq l(n)$ . Choosing  $x$  such that its length  $n$  is random in this sense yields, up to independent constants,

$$\begin{aligned} I_C(x : n) &= C(n) - C(n|x) \geq l(n), \\ I_C(n : x) &= C(x) - C(x|n) \leq n - n = 0. \end{aligned}$$

◇

This example shows that the difference (the asymmetry of algorithmic information)  $|I_C(x : y) - I_C(y : x)|$  can be of the order of the logarithm of the complexities of  $x$  and  $y$ . However, it cannot be greater, as we proceed to show now. This may be called the Theorem of Symmetry of Algorithmic Information for  $C$ -complexity. As usual,  $C(x, y) = C(\langle x, y \rangle)$  is the length of the least program of  $U$  that prints out  $x$  and  $y$  and a way to tell them apart.

**Theorem 2.8.2** *For all  $x, y \in \mathcal{N}$ ,  $C(x, y) = C(x) + C(y|x) + O(\log C(x, y))$ .*

Since  $C(x, y) = C(y, x)$  up to an additive constant term, the following Symmetry of Information property follows immediately.

**Corollary 2.8.1** Up to an additive term  $O(\log C(x, y))$ ,

$$C(x) - C(x|y) = C(y) - C(y|x).$$

Therefore,

$$|I_C(x : y) - I_C(y : x)| = O(\log C(x, y)).$$

Theorem 2.8.2 cannot be improved in general, since in Example 2.8.2 we have seen that the difference  $|I_C(x : y) - I_C(y : x)| \geq \log C(x)$  for some nontrivial  $x$  and  $y$ . The proof of Theorem 2.8.2 follows.

**Proof.** ( $\leq$ ) We can describe  $\langle x, y \rangle$  by giving a description of  $x$ , a description of  $y$  given  $x$ , and an indication of where to separate the two

descriptions. If  $p$  is a shortest program for  $x$  and  $q$  is a shortest program for  $y$ , with  $l(p) \leq l(q)$ , then there is a Turing machine for which  $\overline{l(p)pq}$  is a program to compute  $\langle x, y \rangle$ . Invoking the Invariance Theorem 2.1.1,  $C(x, y) \leq C(x) + C(y|x) + 2l(C(x)) + O(1)$ .

( $\geq$ ) Assume to the contrary: for every constant  $c$ , there are  $x$  and  $y$  such that

$$C(y|x) > C(x, y) - C(x) + cl(C(x, y)). \quad (2.8)$$

Let  $A = \{\langle u, z \rangle : C(u, z) \leq C(x, y)\}$ . Given  $C(x, y)$ , the set  $A$  can be recursively enumerated. Let  $A_x = \{z : C(x, z) \leq C(x, y)\}$ . Given  $C(x, y)$  and  $x$ , we have a simple algorithm to recursively enumerate the set  $A_x$ . One can describe  $y$ , given  $x$ , using its serial number in enumeration order of  $A_x$  and  $C(x, y)$ . Therefore,

$$C(y|x) \leq l(d(A_x)) + 2l(C(x, y)) + O(1). \quad (2.9)$$

By Equations 2.8, 2.9, for every constant  $c$ , there are  $x$  and  $y$  such that

$$d(A_x) > 2^e, \quad e = C(x, y) - C(x) + cl(C(x, y)). \quad (2.10)$$

But now we can obtain a too short description for  $x$  as follows. Given  $C(x, y)$  and  $e$ , we can recursively enumerate the strings  $u$  which are candidates for  $x$  by satisfying

$$\begin{aligned} A_u &= \{z : C(u, z) \leq C(x, y)\}, \\ 2^e &< d(A_u). \end{aligned} \quad (2.11)$$

Denote the set of such  $u$  by  $U$ . Clearly,  $x \in U$ . Also,

$$\{\langle u, z \rangle : u \in U, z \in A_u\} \subseteq A. \quad (2.12)$$

The number of elements in  $A$  cannot exceed the available number of programs that are short enough to satisfy its definition:

$$d(A) \leq 2^{C(x, y)+O(1)}. \quad (2.13)$$

Using Equations 2.11, 2.12, and 2.13,

$$d(U) < \frac{d(A)}{2^e} \leq \frac{2^{C(x, y)+O(1)}}{2^e}.$$

Hence, we can reconstruct  $x$  from  $C(x, y)$ ,  $e$ , and the serial number of  $x$  in enumeration order of  $U$ . Therefore,

$$C(x) < 2l(C(x, y)) + 2l(e) + C(x, y) - e + O(1).$$

Substituting  $e$  as given in Equation 2.10, this yields a contradiction,  $C(x) < C(x)$ , for large enough  $c$ .  $\square$

## Exercises

**2.8.1.** [17] The following equality and inequality seem to suggest that the shortest descriptions of  $x$  contain some extra information besides the description of  $x$ .

- (a) Show that  $C(x, C(x)) = C(x) + O(1)$ .
- (b) Show that  $C(x|y, i - C(x|y, i)) \leq C(x|y, i)$ .

*Comments.* These (in)equalities are in some sense pathological and may not hold for all reasonable descriptional complexities. However, these phenomena also hold for the prefix complexity  $K$  introduced in Chapter 3. Source: P. Gács, *Lecture Notes on Descriptional Complexity and Randomness*, Manuscript, Boston University, 1987.

**2.8.2.** [15] Let  $x$  be a string of length  $n$ .

- (a) Show that  $C(x, C(x)) = C(x)$ .
- (b) Show that the equality  $C(x, C(x)) = C(C(x)|x) + C(x)$  can only be satisfied to within an additive term of about  $\log n$ .
- (c) Prove that  $C(x, y) = C(x|y) + C(y)$  can only hold to within an additive logarithmic term without using Item (a) and Exercise 2.7.8.

*Comments* Hint for Item (b): use Item (a) and Exercise 2.7.8. Hint for Item (c): additivity is already violated on random strings of random length. Source: P. Gács, *Lecture Notes on Descriptional Complexity and Randomness*, Manuscript, Boston University, 1987; A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124.

**2.8.3.** [28] Let  $\omega = \omega_1\omega_2\dots$  be an infinite binary sequence. The entropy function  $H(p)$  is defined by  $H(p) = p \log p + (1-p) \log(1-p)$ . Let  $\lim_{n \rightarrow \infty} \sum_{i=1}^n \omega_i = p$ .

- (a) Show that

$$C(\omega_{1:n}|n) \leq nH\left(\frac{1}{n} \sum_{i=1}^n \omega_i\right) + \log n + c.$$

(b) Show the following: If the  $\omega_i$ 's are generated by coin flips with probability  $p$  for outcome 1 (a Bernoulli process with probability  $p$ ), then for all  $\epsilon > 0$ ,

$$\Pr \left\{ \omega : \left| \frac{C(\omega_{1:n}|n)}{n} - H(p) \right| > \epsilon \right\} \rightarrow 0,$$

as  $n$  goes to infinity.

## 2.9 History and References

---

The confluence of ideas leading to Kolmogorov complexity is analyzed in Section 1.8 through Section 1.12. Who did what, where, and when, is exhaustively discussed in Section 1.13. The relevant documents are dated R.J. Solomonoff, 1960/1964, A.N. Kolmogorov, 1965, and G.J. Chaitin, 1969. According to L.A. Levin, Kolmogorov in his talks used to give credit also to A.M. Turing (for the Universal Turing machine). The notion of nonoptimal complexity (as a complexity based on shortest descriptions but lacking the Invariance Theorem) can be attributed, in part, also to A.A. Markov [*Soviet Math. Dokl.*, 5(1964), 922–924] and G.J. Chaitin [*J. ACM*, 13(1966), 547–569], but that is not a very crucial step from Shannon’s coding concepts.

The notion “sophistication” was introduced by M. Koppel and H. Atlan [*Unpublished Manuscript*] and M. Koppel [*Complex Systems*, 1 (1987), 1087–1091; *The Universal Turing Machine: A Half-Century Survey*, R. Herken, ed., Oxford Univ. Press, 1988, 435–452].

The connection between incompressibility and randomness was made explicit by Kolmogorov and later by Chaitin. Theorem 2.2.1 is due to Kolmogorov. The idea to develop an algorithmic theory of information is due to Kolmogorov, as is the notion of deficiency of randomness. Universal a priori probability (also based on the Invariance Theorem) is due to Solomonoff. This is treated in more detail in Chapter 4. (Solomonoff did not consider descriptive complexity itself in detail.)

The material in Section 2.2.2 is based on A.N. Kolmogorov, *Lecture Notes in Mathematics*, vol. 1021, Springer-Verlag, 1983, 1–5; A.N. Kolmogorov, V.A. Uspensky, *Theory Probab. Appl.*, 32(1987), 389–412; a. Kh. Shen’, *Soviet Math. Dokl.*, 28(1983), 295–299; and V.V. V'yugin, *Theory Probab. Appl.*, 32(1987), 508–512. For the Kolmogorov minimal sufficient statistic we mainly followed T.M. Cover and J.A. Thomas, *Elements of Information Theory*, Wiley, 1991, 169–182.

In his 1965 paper, Kolmogorov mentioned the uncomputability of  $C(x)$  in a somewhat vague form: “[...] the function  $C_\phi(x|y)$  cannot be effectively calculated (generally recursive) even if it is known to be finite for all  $x$  and  $y$ .” Also Solomonoff suggests this in his 1964 paper: “it is clear that many of the individual terms of Eq. (1) are not ‘effectively computable’ in the sense of Turing [...] but can be used] as the heuristic basis of various approximations.” Related questions were considered by L. Löfgren [*Automata Theory*, E. Caianiello, Ed., Academic Press, 1966, 251–268; *Computer and Information Sciences II*, J. Tou (Ed.), Academic Press, 1967, 165–175]. Theorem 1 in the latter reference demonstrates in general that for *every* universal function  $\phi_0$ ,  $C_{\phi_0}(x)$  is not recursive in  $x$ . (In the Invariance Theorem we only considered universal functions using a special type of coding.)

Despite the depth of the main idea of Kolmogorov complexity, the technical expression of the basic quantities turned out to be inaccurate in the sense that many important relationships hold only to within an error term such as the logarithm of complexity. For instance, D.W. Loveland introduced  $n$ -strings in [*Inform. Contr.*, 15(1969), 510–526; *Proc. ACM 1st Symp. Theory Comput.*, 1969, 61–65] and proved that the length-conditional  $C(x_{1:n}|n)$  measure is not monotonic in  $n$ , Example 2.2.5, page 112. He proposed a uniform complexity to solve this problem, and relationships between these complexities are the subject of Exercises 2.3.3, 2.3.4, 2.5.10, 2.5.11, 2.5.12, and 2.5.13.

In the subsequent development of this chapter we have used time and again the excellent 1970 survey by L.A. Levin and A.K. Zvonkin [*Russ. Math. Surv.*, 25:6(1970), 83–124] which describes mainly the research performed in the former USSR. We have drawn considerably on and profited from the point of view expressed in P. Gács's [*Komplexität und Zufälligkeit*, Ph.D. Thesis, J.W. Goethe Univ., Frankfurt am Main, 1978, unpublished; *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987]. Another source for the Russian school is the survey by V.V. V'yugin, *Selecta Mathematica* formerly *Sovietica*, 13:4(1994), 357–389, (Translated from the Russian *Semiotika and Informatika*, 16(1981), 14–43).

In [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124], Theorems 2.2.1 through 2.3.2 are attributed to Kolmogorov. The result on meager sets in Section 2.2 is from [M. Sipser, *Lecture Notes on Complexity Theory*, MIT Lab Computer Science, 1985, unpublished]. We avoided calling such sets sparse sets because we need to reserve the term “sparse set” for sets that contain a polynomial number of elements for each length. The approximation Theorem 2.3.3 is stated in some form in [R.J. Solomonoff, *Inform. Contr.*, 7(1964), 1–22, 224–254], and is attributed also to Kolmogorov by Levin and Zvonkin. Some other properties of the integer function  $C$  we mentioned were observed by H.P. Katseff and M. Sipser [*Theoret. Comput. Sci.*, 15(1981), 291–309].

The material on random strings (sequences) in Sections 2.4, 2.5 is primarily due to P. Martin-Löf [*Inform. Contr.*, 9(1966), 602–619; *Z. Wahrsch. Verw. Geb.*, 19(1971), 225–230]. The notions of random finite strings and random infinite sequences, complexity oscillations, enumerable (sequential) Martin-Löf tests and the existence of universal (sequential) tests, the use of constructive measure theory, and Theorems 2.4.2, 2.5.3, 2.5.1 through 2.5.5 are taken from P. Martin-Löf's papers. Weaker oscillations are mentioned by G.J. Chaitin [*J. ACM*, 16(1969), 145–159]. We also used [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124; P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987].

As noted in the main text, the complexity oscillations of infinite sequences prevent a clear expression of randomness in terms of complexity. This problem was investigated by L.A. Levin in [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124] and independently by C.P. Schnorr [*Lecture Notes in Mathematics*, vol. 218, Springer-Verlag, 1971]. As a part of the wider issue of (pseudo) random number generators and (pseudo) randomness tests, the entire issue of randomness of individual finite and infinite sequences is thoroughly reviewed by D.E. Knuth, *Seminumerical Algorithms*, Addison-Wesley, 1981, pp. 142–169; summary, history, and references: pp. 164–166. The whole matter of randomness of individual finite and infinite sequences of zeros and ones is placed in a wider context of probability theory and stochastics, and is exhaustively analyzed in [A.N. Kolmogorov and V.A. Uspensky, *Theory Probab. Appl.*, 32(1987), 389–412; V.A. Uspensky, A.L. Semenov and A.Kh. Shen', *Russ. Math. Surv.*, 45:1(1990), 121–189].

Section 2.6, which analyzes precisely the relative frequencies of 0's and 1's and  $k$ -length blocks in individual infinite and finite sequences in terms of their Kolmogorov complexity, is from [M. Li and P.M.B. Vitányi, *Math. Systems Theory*, 27(1994), 365–376].

The recursion-theoretic properties we treat in Section 2.7 are related to Gödel's famous Incompleteness Theorem. Theorem 2.7.1 is attributed to J.M. Barzdins in [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124]. The proof of Corollary 2.7.2 was given by G.J. Chaitin [*J. ACM*, 21(1974), 403–423; *Scientific American*, 232:5(1975), 47–52]. This application and some philosophical consequences have been advocated with considerable eloquence by G.J. Chaitin and C.H. Bennett [C.H. Bennett and M. Gardner, *Scientific American*, 241:5(1979), 20–34].

We also used the insightful discussion in [P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987]. These results are analyzed and critically discussed from a mathematical logic point of view by M. van Lambalgen [*J. Symb. Logic*, 54(1989), 1389–1400]. Theorem 2.7.2, Barzdins's Lemma, occurs both in [J.M. Barzdins, *Soviet Math. Dokl.*, 9(1968), 1251–1254, and D.W. Loveland, *Proc. ACM 1st Symp. Theory Comput.*, 1969, 61–65]. Examples in Section 2.7 are due to Kolmogorov, 1970, published much later as [*Russ. Math. Surv.*, 38:4(1983), 27–36] and a footnote in [L.A. Levin, *Problems Inform. Transmission*, 10:3(1974), 206–210].

The treatment of the relation between plain Kolmogorov complexity and Shannon's entropy in Section 2.8 is based on the work of A.N. Kolmogorov [*Problems Inform. Transmission*, 1:1(1965), 1–7; *IEEE Trans. Inform. Theory*, IT-14(1968), 662–665; *Russ. Math. Surv.*, 38:4(1983), 27–36; *Lecture Notes in Mathematics*, vol. 1021, Springer-Verlag, 1983,

1–5] and on [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124]. The latter reference attributes Theorem 2.8.1 to Kolmogorov. Theorem 2.8.2 and its Corollary 2.8.1, establishing the precise error term in the additivity of complexity and symmetry of information as logarithmic in the complexity, is due to Levin and Kolmogorov [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124; A.N. Kolmogorov, *Russ. Math. Surv.*, 38:4(1983), 27–36].

# Algorithmic Prefix Complexity

While the development of an algorithmic theory of complexity according to the original definitions (plain Kolmogorov complexity) in Chapter 2 was very fruitful, for certain goals the mathematical framework is not yet satisfactory. This has resulted in a plethora of proposals of modified measures to get rid of one or the other problem. Let us list a few conspicuous inconveniences.

- The plain complexity is *not subadditive*: the inequality  $C(x, y) \leq C(x) + C(y)$  holds only to within a term logarithmic in  $C(x)$  or  $C(y)$ —Example 2.1.4 on page 101 and Example 2.2.3 on page 111. An attempt to solve this problem is to use conditional complexity, in which case we indeed find  $C(x, y|C(x)) \leq C(x) + C(y)$  up to an additive constant.
- Another problem is *nonmonotonicity over prefixes*: it would be pleasing if the complexity of  $xy$  were never less than the complexity of  $x$ . The complexity measure  $C(x)$  does not have this property. In contrast to the subadditivity question, here use of conditional complexity  $C(x|l(x))$  instead of  $C(x)$  does not help. Not only is this measure nonmonotonic, but also it has another counterintuitive property: it drops infinitely often below a fixed constant as  $x$  runs through the natural numbers. A first proposal to remedy this defect was the “uniform complexity” variant  $C(x; l(x))$ , see Exercise 2.3.3, page 124. Informally, this measure gives the length of the shortest program  $p$  that computes  $x_i$  for all inputs  $i$ ,  $i = 1, 2, \dots, l(x)$ .
- In the development of the theory of *random infinite sequences* it would be natural to identify infinite random sequences with infinite sequences of which all finite initial segments are random. As it

turned out, no infinite sequence satisfies this criterion, due to complexity oscillations. This holds for any of the  $C(x)$ ,  $C(x|l(x))$ , and  $C(x;l(x))$  measures. It is quite complicated to express Martin-Löf's notion of randomness in terms of the  $C$ -complexity of prefixes.

- The original motivation of Solomonoff to introduce algorithmic complexity was as a device through which to assign a universal *prior probability* to each finite binary string (Sections 1.6, 1.10). Choosing the reference Turing machine  $U$  as in the Invariance Theorem 2.1.1, page 97, this induces a distribution  $P$  over  $\mathcal{N}$  (equivalently  $\{0,1\}^*$ ) defined by  $P(x) = \sum 2^{-l(p)}$ , the sum taken over all inputs  $p$  for which  $U$  computes output  $x$  and halts.

This approach is different from the one in Example 1.1.3 on page 5, where we used  $P'(x) = 2^{-l(x^*)}$ , where  $x^*$  is a shortest program for  $x$ . Anticipating the sequel, if we allow only self-delimiting programs as developed in this chapter, the two approaches turn out to be the same, see Theorem 4.3.3 on page 253, namely,  $P(x) = \Theta(P'(x))$  if we allow only Turing machines with self-delimiting programs (no program for which the machine halts is the prefix of another program for which the machine halts).

Unfortunately,  $P$  is not a probability distribution, since the series  $\sum_x P(x)$  diverges. Worse, for each individual  $x$  we have  $P(x) = \infty$ . Namely, for each  $x \in \mathcal{N}$  there is a Turing machine  $T$  that computes the constant  $x$  for all inputs. If  $T$  denotes the description of  $T$  used by  $U$ , then

$$P(x) \geq 2^{-l(T)} \sum_{p \in \{0,1\}^*} 2^{-l(p)} = \infty.$$

Our next try is to redefine  $P(x)$  by not considering all programs that output  $x$ , but only the shortest program that outputs  $x$ . This yields  $P(x) = 2^{-C(x)}$ . However, we still have that  $\sum_x P(x)$  diverges, so again  $P$  is not a probability distribution. This holds also with respect to  $C(x|l(x))$  and  $C(x;l(x))$ , the simple reason being that since all of these measures are close to  $\log x$  for almost all  $x$ , the corresponding  $P(x)$  will be close to  $1/x$  for almost all  $x$ . Divergence of  $\sum_x P(x)$  follows then from divergence of the harmonic series  $\sum_x 1/x$ .

- There is a close relation between the complexity of a finite string and its *Shannon entropy*, Section 2.8 and the later Section 8.1. Indeed, it would be troublesome if it were not so, since both notions have been introduced as a measure of information content: in the one case of individual objects, and in the other case of the event space. Therefore, we could hope that classical information-theoretic identities as derived in Section 1.11 would have complexity analogues that are satisfied up to an additive constant (reflecting the

choice of reference machine). However, in Section 2.8 we have established that the complexity analogues of some information-theoretic identities are satisfied only to within an additive logarithmic term and that this cannot be improved in general.

Looking at this list of deficiencies with the wisdom of hindsight, it is not too difficult to transcend a decade of strenuous investigation and alternative proposals by proceeding immediately to what now seems a convenient version of algorithmic complexity. This is the complexity induced by Turing machines with a set of programs in which no program is a proper prefix of another program. This proposal gets rid of all problems above (except for nonmonotonicity). Let us consider this matter in some detail.

The study of the algorithmic complexity of descriptions asks in effect for a code in the sense of Section 1.11.1. In his original paper Shannon restricted attention to those codes for which no value is the prefix of another value, the so-called prefix-codes. This restriction is highly motivated by the implicit assumption that descriptions will be concatenated and must be uniquely decodable.

Recall from Section 1.11.1 that uniquely decodable codes and prefix-codes share the same sets of code-word lengths. Moreover, the minimal average code-word length  $L$  of a prefix-code encoding source word sequences emitted by a random variable with probability distribution  $P$  satisfies

$$H(P) \leq L \leq H(P) + 1,$$

where  $H(P)$  is the entropy of the random variable (Theorem 1.11.2 on page 75). This is obtained (up to an additive constant term) by assigning code-word lengths  $l_i = \lceil \log 1/p_i \rceil$  to the  $i$ th outcome of the random variable, where  $p_i$  is the probability of that outcome. Viewed this way,  $H(P)$  may be interpreted as the minimal expected length of a description of a single outcome in prefix-free codes. Thus, the average additive fudge term in complexity equations based on prefix-codes is  $O(1)$  rather than logarithmic. However, for the individual complexity equations the situation is more complicated.

The divergence of the series  $\sum 2^{-C(x)}$  was a major blow for Solomonoff's program. But if the set of programs of the reference machine is prefix-free, then convergence of this series as required by the probability interpretation is a property ensured by Kraft's Inequality on page 74.

At present, prefix-code based complexity is often considered as some sort of a standard algorithmic complexity. Lest the reader be deluded into the fallacy that this is the most perfect of all possible worlds, we state that different applications turn out to require different versions of complexity, and all of these are natural for their own purposes.

### 3.1 The Invariance Theorem

#### Definition 3.1.1

In the theory of algorithmic complexity we have a more refined goal than in classic information theory, but not essentially different objectives. Here too, it is natural to restrict the effective descriptions to uniquely decodable codes, and since our interest is only in the length of code-words, to prefix-codes.

A *partial recursive prefix function*  $\phi : \{0, 1\}^* \rightarrow \mathcal{N}$  is a partial recursive function such that if  $\phi(p) < \infty$  and  $\phi(q) < \infty$ , then  $p$  is not a proper prefix of  $q$ .

Let  $T_1, T_2, \dots$  be the standard enumeration of Turing machines, and let  $\phi_1, \phi_2, \dots$  be the corresponding enumeration of partial recursive functions (Section 1.7). Obviously, all partial recursive prefix functions occur in this list. We change every Turing machine  $T$  computing  $\phi$  to a machine  $T'$  computing a partial recursive prefix function  $\psi$ , where  $\psi = \phi$  if  $\phi$  was already a partial recursive prefix function. The machine  $T'$  executes the algorithm below using machine  $T$ . In contrast to  $T$ , which is presented with each input out of  $\{0, 1\}^*$  delimited by endmarkers, machine  $T'$  is presented with each potentially infinite binary input sequence  $b_1 b_2 \dots$ , which it reads from left-to-right without backing up.

A *halting input*, or *program*, of  $T'$  is an initial segment  $b_1 b_2 \dots b_m$  such that  $T'$  halts after reading  $b_m$  and before reading  $b_{m+1}$ . There are no other programs of  $T'$ . This way no program is a proper prefix of any other program, that is, the set of programs is prefix free. Moreover,  $T'$  determines the end of each program without reading the next symbol. Such programs are called *self-delimiting*; see Definition 1.11.4 on page 76. Each such program  $p$  is the code word for the source word  $T'(p)$  consisting of the word written on the output tape by the time  $T'$  halts its computation with program  $p$ . Machine  $T'$  executes the following algorithm:

**Step 1** Set  $p := \epsilon$ .

**Step 2** Dovetail all computations of  $T$  computing  $\phi(p, q)$ , for all  $q \in \{0, 1\}^*$ . {Let  $q_j$  be the  $j$ th string of  $\{0, 1\}^*$ . Dovetailing means executing consecutive stages 1, 2, ..., such that in the  $i$ th stage we do one next step of the computation of  $\phi(pq_j)$  for all  $j$  with  $j \leq i$ .} If  $\phi(pq) < \infty$  is the first halting computation, then go to Step 3.

**Step 3** If  $q = \epsilon$  then output  $\phi(p)$  and halt else read  $b :=$  next input bit; set  $p := pb$ ; go to Step 2.

This construction produces an effective enumeration

$$T'_1, T'_2, \dots$$

of Turing machines, called *prefix machines*, computing all, and only, partial recursive prefix functions  $\psi_1, \psi_2, \dots$ . This brings us to the *Invariance Theorem for prefix complexity*.

**Example 3.1.1** An alternative way of defining a prefix machine is as follows. It is a Turing machine with a separate one-way input tape, a separate one-way output tape, and a two-way work tape, all of them one-way infinite. In the start state the input is written on the input tape, and the read-only head on that tape scans the first symbol. Initially, both the work tape and the output tape are empty and the read/write heads are at leftmost squares. At the start of each step the state of the finite control and the symbols under scan on the input tape and the work tape determine what the Turing machine does in this step. It does all of the following: It either moves the reading head on the input tape to the next input symbol or does not move the reading head at all, followed by a computation step consisting of either changing the scanned symbol on the work tape or moving the work tape head one square left, right, or not at all, followed by a change of state of the finite control and either writing a symbol to the next empty square on the output tape or not writing on the output tape.

Let the infinite input tape contain only 0's or 1's (no blanks). Clearly, with this definition every prefix machine  $T$  scans one maximal prefix, say  $p$ , of any infinite binary input. The output is the contents of the output tape when the machine halts. We leave it to the reader to show that such prefix machines can be effectively enumerated. Moreover, every partial recursive prefix function is computed by a prefix machine and every prefix machine computes a partial recursive prefix function.  $\diamond$

Recall Definition 2.0.1, page 95, of a function that is universal (or additively optimal) for a class of functions.

**Theorem 3.1.1** *There exists a universal partial recursive prefix function  $\psi_0$  such that for any partial recursive prefix function  $\psi$  there is a constant  $c_\psi$  such that  $C_{\psi_0}(x|y) \leq C_\psi(x|y) + c_\psi$ , for all  $x, y \in \mathcal{N}$ .*

**Proof.** Using standard techniques, like those in Section 1.7, we can show that there exists a universal prefix machine  $U$  such that  $U(\langle y, \langle n, p \rangle \rangle) = T_n(y, p)$  for all  $y, p \in \mathcal{N}$ . The remainder of the proof is analogous to the proof of Theorem 2.1.1, page 97.  $\square$

For each pair of additively optimal partial recursive prefix  $\psi$  and  $\psi'$ ,

$$|C_\psi(x|y) - C_{\psi'}(x|y)| \leq c_{\psi, \psi'},$$

for all  $x$  and some constant  $c_{\psi, \psi'}$ . We fix one additively optimal partial recursive prefix function  $\psi_0$  as the standard reference, and  $U$  in the proof

of Theorem 3.1.1 as the *reference prefix machine*, , and define the *prefix complexity* of  $x$ , conditional to  $y$ , by  $K(x|y) = C_{\psi_0}(x|y)$  for all  $x$ . The (unconditional) prefix complexity of  $x$  is defined as  $K(x) = K(x|\epsilon)$ .

**Example 3.1.2** Define  $K(x, y) = K(\langle x, y \rangle)$ . Requiring the decoding algorithm to be prefix-free has advantages, since now we can concatenate descriptions without marking where one description ends and the other one begins. Previously, endmarkers disappeared in concatenation of descriptions, which was responsible for logarithmic fudge terms in our formulas. More formally, in contrast to  $C$  (Example 2.2.3 on page 111) the measure  $K$  is *subadditive*, that is,

$$K(x, y) \leq K(x) + K(y) + O(1).$$

Namely, let  $U$  be the reference machine of Theorem 3.1.1. Let  $U(x^*) = x$  with  $l(x^*) = K(x)$ , and  $U(y^*) = y$  with  $l(y^*) = K(y)$ . Since the set of programs of  $U$  is a prefix-code, we can modify  $U$  to a machine  $V$  that first reads  $x^*$  and computes  $x$ , then reads  $y^*$  and computes  $y$ , and finally computes and outputs  $\langle x, y \rangle$ . If  $V = T_n$  in the enumeration of prefix Turing machines, then  $U(\langle n, \langle x^*, y^* \rangle \rangle) = \langle x, y \rangle$ . A similar argument shows that  $K(xy) \leq K(x) + K(y) + O(1)$ .  $\diamond$

**Example 3.1.3** The functions  $C$  and  $K$  are asymptotically equal. For all  $x$  and  $y$ , we have, up to additive constant terms,

$$C(x|y) \leq K(x|y) \leq C(x|y) + 2 \log C(x|y).$$

Since the prefix partial recursive functions are a restriction on the notion of partial recursive functions, it is obvious that  $C(x|y) \leq K(x|y)$ . To prove the second inequality, we recall that for each  $x$  and  $y$  the reference machine of the  $C$ -complexity of Definition 2.1.2 on page 98 computes output  $x$  from some input  $\langle y, p \rangle$  with  $l(p) = C(x|y)$ . We know that  $\overline{l(p)}p$  is a self-delimiting encoding for  $p$ . Therefore,  $K(x|y) \leq C(x|y) + 2l(C(x|y)) + O(1)$ .  $\diamond$

It is straightforward to extend this idea by having a prefix machine  $V$  compute  $x$  from input  $p_r 0 p_{r-1} 0 \dots 0 p_0 1$ . Here  $p_{i+1}$  is the shortest program for the length of  $p_i$ , and  $p_0$  is the shortest program to compute  $x$  given  $y$ . If  $V = T_n$ , then  $U(n, y, p_r 0 p_{r-1} 0 \dots 0 p_0 1) = x$ , and

$$\begin{aligned} K(x|y) &\leq C(x|y) + C(C(x|y)) + \dots + O(1) + r \\ &\leq C(x|y) + C(C(x|y)) + O(\log C(C(x|y))) + O(1), \end{aligned}$$

where the  $\dots$  indicates all  $r$  positive terms, and the  $O(1)$  term comprises the cost of encoding  $n$ . Using more efficient prefix-codes, we can improve the above estimate. Let  $l^*$  be as defined by Equation 1.17; we have

$$C(x|y) \leq K(x|y) \leq C(x|y) + l^*(C(x|y)) + O(1). \quad (3.1)$$

A more precise relation between  $C(x)$  and  $K(x)$  was shown by R.M. Solovay:

$$\begin{aligned} K(x) &= C(x) + C(C(x)) + O(C(C(C(x)))), \\ C(x) &= K(x) - K(K(x)) - O(K(K(K(x))). \end{aligned}$$

The content of this observation is that in order to make the minimal  $C$ -style program self-delimiting, we must prefix it with a self-delimiting encoding of its own length. To within the cited error term, this simple procedure is always optimal.

In Theorem 2.1.2, page 100, we showed the upper bound  $C(x) \leq n + O(1)$  for all  $x$  of length  $n$ . We have  $K(x) \leq n + 2 \log n + O(1)$  by encoding  $x$  as  $\overline{l(x)}x$ . We give a concrete upper bound on the implied constant of the  $O(1)$  term in Section 3.2. We can improve the estimate by more efficient prefix-codes (Section 1.11.1) to

$$K(x) \leq \log^* n + n + l(n) + l(l(n)) + \cdots + O(1), \quad (3.2)$$

where the sum is taken over all positive terms.

**Example 3.1.4** One reason to introduce the original conditional complexity in Chapter 2 was to obtain length-conditional complexity. The analogous length-conditional prefix complexity satisfies, with  $l(x) = n$ ,

$$\begin{aligned} K(x|n) &\leq K(x) + O(1) \\ &\leq K(x|n) + K(n) + O(1) \\ &\leq K(x|n) + \log^* n + l(n) + l(l(n)) + \cdots + O(1). \end{aligned}$$

The second inequality holds since the concatenation of a self-delimiting program for  $l(x)$  and a self-delimiting program to compute  $x$  given  $l(x)$  constitutes a self-delimiting program for  $x$ . The third inequality follows from Equation 3.2 since  $K(n) \leq \log^* n + l(n) + l(l(n)) + \cdots + O(1)$ .  $\diamond$

## Exercises

**3.1.1.** [12] Use the Kraft Inequality, page 74, to show that  $K(x) \geq \log x + \log \log x$  for infinitely many  $x$ .

**3.1.2.** [15] We investigate transformations of complexity under the arithmetic operation  $*$ . Show that

- (a)  $K(x * y) \leq K(x) + K(y) + O(1)$ ;
- (b)  $K(x * y) + \log \log(x * y) \geq K(x) + K(y) + O(1)$ ;
- (c) If  $x$  and  $y$  are relatively prime, then  $K(x * y) = K(x, y) + O(1)$ ;

*Comments.* Item (b): order all pairs  $\langle u, v \rangle$  satisfying  $u * v = x * y$  lexicographically. Let  $\langle x, y \rangle$  be the  $j$ th element. Use  $j$  and  $x * y$ .

**3.1.3.** [16] Let us backtrack to the original definition of a Turing machine in Section 1.7. Instead of a tape alphabet consisting of 0, 1, and the blank symbol  $B$ , we want to consider Turing machines with a tape alphabet consisting only of 0 and 1. We can modify the original effective enumeration  $T_1, T_2, \dots$  in Section 1.7 to an effective enumeration  $T'_1, T'_2, \dots$  of Turing machines using a purely binary tape alphabet, without blanks  $B$ . Furthermore, assume that the set of inputs for halting computations of each Turing machine form a uniquely decodable code.

- (a) Define  $K'(x) = \min\{l(p) + i + 1 : T'_i(p) = x, i \geq 1\}$ . Show that  $K'(x) = K(x) + O(1)$ .
- (b) Show that restricting the Turing machines to a one-letter tape alphabet (with or without delimiting blanks) does not yield a proper complexity measure.

*Comments.* By reformulating the notion of Turing machine with a purely binary tape alphabet, without distinguished blank symbol  $B$ , we naturally end up with the prefix complexity  $K(x)$ . Hint for Item (a): use the McMillan-Kraft Theorem, Exercise 1.11.9 on page 82. Further restrictions of the tape alphabet do not yield proper complexity measures. Hint for Item (b): show that there is no additively optimal partial recursive function in this enumeration of partial recursive functions.

- 3.1.4.** • [19] Do Exercise 2.1.9 on page 107 for  $K$ -complexity.

**3.1.5.** • [31] We can derive  $K(x)$  in another way. Define a *complexity* function  $F(x)$  to be a function on the natural numbers with the property that the series  $\sum 2^{-F(x)} \leq 1$ , and such that  $F(x)$  is co-enumerable. That is, the set  $\{(m, x) : F(x) \leq m\}$  is recursively enumerable. (If  $F(x) = \infty$ , then  $2^{-F(x)} = 0$ .)

- (a) Show that the number of  $x$ 's for which  $F(x) \leq m$  is at most  $2^m$ .
- (b) Show that there is an additively optimal (minimal) complexity  $F$  such that for each complexity  $F'$  there is a constant  $c$  depending only on  $F$  and  $F'$  but not on  $x$  such that  $F(x) \leq F'(x) + c$ . This is the Invariance Theorem corresponding to Theorem 3.1.1.
- (c) Show that  $F(x) = K(x) + O(1)$ , where  $F(x)$  is the optimal complexity in Item (b).

*Comments.* Hint for Item (b): define  $F(x) = \min_k \{1, F_k(x) + k\}$ , where  $F_k$  is the complexity resulting from the  $k$ th partial recursive function after modifying it to satisfy the requirements above. This approach has been proposed by L.A. Levin in a sequence of papers: [*Russ. Math. Surv.*, 25:6(1970), 83–124, *Sov. Math. Dokl.*, 14(1973), 1413–1416, and *Problems Inform. Transmission*, 10(1974), 206–210 (where the equivalence with the prefix machine approach is established)]. See also [P.

Gács, *Sov. Math. Dokl.*, 15(1974), 1477–1480]. G.J. Chaitin [*J. ACM*, 22(1975), 329–340] used the prefix machine approach to define  $K(x)$  but gave a different interpretation to the conditional complexity, resulting in the  $Kc$  version in Example 3.9.1 and Exercise 3.9.3. There we will observe that unconditionally  $Kc(x) = K(x)$ . Let  $\langle x, y \rangle$  denote as usual a fixed recursive one-to-one correspondence between natural numbers and pairs of natural numbers. Formally, define  $Kc(x|y) = Kc(\langle x, y \rangle) - Kc(x)$  and  $I(x; y) = Kc(y) - Kc(y|x)$ . Now we can derive all (in)equalities of Exercise 3.9.3 without giving an algorithmic interpretation like  $K(x|\langle y, K(y) \rangle)$  to  $Kc(x|y)$ .

**3.1.6.** [39] Giving a definition of complexity of description, we use recursive decoding functions from the set of finite binary strings to itself. However, we can consider this set with different topologies. If we consider distinct binary strings as incomparable, then this set corresponds to the natural numbers  $\mathcal{N}$ . If we consider the distinct binary words as compared by the prefix relation, then we view them as nodes in a binary tree  $\mathcal{B}^*$ ,  $\mathcal{B} = \{0, 1\}$ , in the obvious way. Then there are four possible definitions of recursive decoding mappings: from  $\mathcal{N}$  to  $\mathcal{N}$ , from  $\mathcal{N}$  to  $\mathcal{B}^*$ , from  $\mathcal{B}^*$  to  $\mathcal{N}$ , and from  $\mathcal{B}^*$  to  $\mathcal{B}^*$ . Exploit this idea to formalize the following:

- (a) Using recursive decoding from  $\mathcal{N}$  to  $\mathcal{N}$  we obtain the plain complexity  $C(x)$ .
- (b) Using recursive decoding from  $\mathcal{N}$  to  $\mathcal{B}^*$  we obtain the uniform complexity  $C(x; n)$  as defined in Exercise 2.3.3, page 124.
- (c) Using recursive decoding from  $\mathcal{B}^*$  to  $\mathcal{N}$  we obtain the prefix complexity  $K$ .
- (d) Using recursive decoding from  $\mathcal{B}^*$  to  $\mathcal{B}^*$  we obtain the monotone complexity as defined by Levin and Schnorr (Section 4.5.4).

*Comments.* Source: A. Kh. Shen', *Sov. Math. Dokl.*, 29:3(1984), 569–574; V.A. Uspensky, pp. 85–101 in: *Kolmogorov Complexity and Computational Complexity*, O. Watanabe, Springer-Verlag, 1992.

## 3.2 \*Sizes of the Constants

The concrete Kolmogorov complexity and prefix complexity of an object depend on the particular choice of universal machine fixed as the reference machine. A minor modification of the universal Turing machine  $U$  of Example 1.7.4 on page 30 yields small implied constants—1 or 2 instead of “ $O(1)$ ”—in relations such as  $C(x) \leq l(x) + O(1)$  (Theorem 2.1.2 on page 100) and  $K(x) \leq l(x) + 2 \log l(x) + O(1)$  (page 195). For example, modify  $U$  to the reference universal machine  $U'$  such that  $U'(0p) = p$  and  $U'(1p) = U(p)$ . But it is quite tedious to exhibit the explicit encoding of such a machine. On the other hand, if we aim for a “small” universal

Turing machine, then it becomes difficult to determine the concrete implied constants—which also may become quite large. With this in mind, we want to find a machine model in which a concrete universal machine is easily, and shortly, implemented from weak elementary operations and the implied constants above are small and can be easily determined. R. Penrose in [*The Emperor's New Mind*, Oxford University press, 1989] encoded a universal Turing machine in 5495 bits. G.J. Chaitin [*Complexity*, 1:4(1995/1996), 55–59] used LISP and powerful elementary operations to define a universal machine to concretely determine prefix complexity and some associated constants.

Lambda calculus is simpler and more elegant than LISP. Pure lambda calculus with combinators  $S$  and  $K$ , is elegant, but slow to evaluate. Since we are concerned with simplicity of definition rather than with efficiency of execution, we base a concrete definition of Kolmogorov complexity on combinatory logic. What follows is a very brief exposition on combinatory logic (see [H.P. Barendregt, *The Lambda Calculus, its Syntax and Semantics*, North-Holland, Amsterdam, 1984]), which develops into a discussion of how to represent binary strings in it.

The *primitive combinators*  $S$  and  $K$  are “defined” by the rewrite rules

$$\begin{aligned} Sxyz &= xz(yz), \\ Kxy &= x. \end{aligned}$$

Other *combinators* are constructed from these two by application. We assume that application of  $x$  on  $y$ , written  $(xy)$ , associates to the left, and omit parentheses accordingly. Thus,  $Sxyz$  should be interpreted as  $((Sx)y)z$ . Perhaps surprisingly,  $S$  and  $K$  suffice to define any function definable in standard lambda calculus. As an example, we can check that  $I = SKz$ , for every  $z$ , is the identity function:

$$SKzx = Kx(zx) = x.$$

The combinator  $\infty \equiv SSK(S(SSK))$  is another interesting example—when applied to any combinator, it leads to an infinite rewriting chain (abbreviating  $S(SSK)$  to  $L$ ):

$$\begin{aligned} \infty x &\equiv SSKLx = SL(KL)x = Lx(KLx) = S(SSK)xL \\ &= SSKL(xL) \equiv \infty(xL) = \infty(xLL) = \infty(xLLL) \dots. \end{aligned}$$

In our machine model we want to consider each combinator to be a machine, capable of processing input bits and producing output bits. This entails representing an input binary string by a combinator, to which the machine can be applied, and a way of interpreting the result, if possible, as an output binary string.

A rewrite of  $Sxyz$  or  $Kxy$  will form the basic “machine cycle.” Our combinator machines employ a *lazy* reduction strategy, in which only the “frontal” primitive combinator can be reduced. This is to avoid endless rewrites that might be inconsequential for interpreting the result as a string. Lazy reduction ends when a combinator of the form  $K$ ,  $Kx$ ,  $S$ ,  $Sx$ , or  $Sxy$  is reached, where we lack the required number of arguments for reducing the front  $S/K$  (as in  $S(KKK)K$ ).

### 3.2.1 Encoding Combinators as Binary Strings

The nice thing about combinators is that they have a wonderfully simple encoding as binary strings that is also self-delimiting: encode  $S$  as 00,  $K$  as 01, and application as 1. Formally, we define the encoding  $\langle C \rangle$  of a combinator  $C$  as

$$\begin{aligned}\langle S \rangle &\equiv 00 \\ \langle K \rangle &\equiv 01 \\ \langle C_0 C_1 \rangle &\equiv 1 \langle C_0 \rangle \langle C_1 \rangle\end{aligned}$$

For instance, the combinator  $S(KSS)$ ,  $(S((KS)S))$  in full, is encoded as 10011010000. The length of this encoding is  $3n - 1$ , where  $n$  is the number of primitive combinators ( $S$  and  $K$ ) it contains.

### 3.2.2 Encoding Booleans, Pairs, and Binary Strings as Combinators

Define

$$\begin{aligned}0 &\equiv K \quad (\text{false}) \\ 1 &\equiv SK \quad (\text{true})\end{aligned}$$

Then, “if  $B$  then  $P$  else  $Q$ ” can be represented by “ $BQP$ .” The simplest way to pair two combinators  $x$  and  $y$  is to have something that when applied to 0 gives  $x$ , and when applied to 1 gives  $y$ . What we want then is a combinator, call it  $P$ , that satisfies  $Pxyz = zxy$ , for example,

$$P \equiv S(S(KS)(S(KK)(S(KS)(S(S(KS)(SK))K))))(KK).$$

Lists are then constructed by repeated pairing as usual:  $Pl_0(Pl_1(Pl_2\dots))$  represents the list  $[l_0, l_1, l_2, \dots]$ . A special element to signal the end of a list, like the null character signaling the end of a string in the C programming language, will be helpful. Since we only consider lists of bits, we require an element that is easily distinguished from 0 and 1. The simple  $KS$  is such an element, since  $KS00 = 1$  while  $000 = 100 = 0$ . An empty list is then a list whose first element is  $KS$ , and we define this to be  $\$ \equiv K(KS)$ . This way we can make arbitrary lists of 0s and 1s. For instance, the string  $s = 0111001$  maps to the list  $P0(P1(P1(P1(P0(P0(P1\$))))))$ , conveniently—and a little confusingly, since  $s$  itself doesn’t represent a combinator—denoted by  $(s\$)$ . One can check that  $(s\$)0$  gives the first bit 0 of  $s$ , while  $(s\$)1$  gives the tail  $P1(P1(P1(P0(P0(P1\$)))) \equiv (111001\$)$ .

### 3.2.3 Output Conventions

If the machine halts, then how do we interpret the result as a binary string? It is unreasonable to expect the result to be exactly identical to  $(s\$)$  for some string  $s$ , due to the lazy reduction. Instead, we should focus on the individual bits  $s_0 s_1 \dots s_{n-1}$  making up a string  $s$ . For a combinator  $C$ , let  $C11\dots10$  with  $i$  1's be denoted by  $C^i$ , that is, the  $i$ th element of the list  $C$ . We say that combinator  $C$  produces output  $s = s_0 \dots s_{n-1}$  iff  $C^i x_0 x_1$  reduces to  $x_s$ , ( $= s_i x_0 x_1$ ) for variables  $x_0$  and  $x_1$  and  $i < n$ , while  $C^n x_0 x_1$  reduces to  $Sx_1$  ( $= \$0 x_0 x_1$ ). In this way a single reduction distinguishes among the three possibilities.

Of course, not all combinators produce proper output in this sense. For instance,  $S^0 x_0 x_1 \equiv S0x_0x_1 = 0x_1(x_0x_1) = x_1$  and  $S^1 x_0 x_1 \equiv S10x_0x_1 = 1x_0(0x_0)x_1 = 0x_0x_1 = x_0$ , but  $S^2 x_0 x_1 \equiv S110x_0x_1 = 10(10)x_0x_1 = 10x_0x_1 = x_0x_1$ . So the primitive combinator  $S$  can be interpreted as a list starting with 2 bits, but the remainder cannot be interpreted as either starting with a bit or behaving like the empty list. Such computations must then be considered as invalid. Having defined now how a combinator can be seen as a machine taking binary strings as input and producing binary strings as output, we return to our main objective.

### 3.2.4 The Universal Combinator

A machine in our model is a combinator that is applied to the input combinator—the list that the input string maps to—in order to produce a proper output. For machines that are supposed to be able to detect the end of their input, an input string  $s$  is mapped to the list  $(s\$)$ .

For *prefix machines*, which need to decide when to stop reading more input based on what they have read so far, we can present them lists of the form  $(s\infty)$  for increasingly longer  $s$  until they halt. Recall that  $\infty$  leads to an infinite computation. Therefore, halting of the machine on  $(s\infty)$  implies it has not “read past” string  $s$ .

### 3.2.5 Decoding a Combinator Encoding

We require combinators that given a list (i) produce the combinator encoded on the initial part of that list and (ii) return the remaining part of the list.

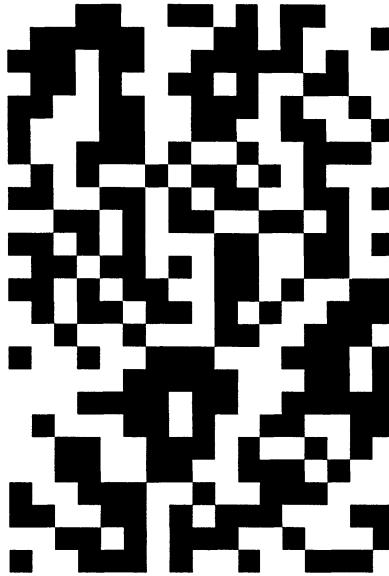
The following combinators do just that:

$$\begin{aligned} Qx &= x0(x10SK)(Q(x1)(Q(R(x1)))) \quad (\text{initial combinator}) \\ Rx &= x0(x11)(R(R(x1))) \quad (\text{skip combinator}) \end{aligned}$$

It may look like cheating to define these combinators recursively, but the fixpoint operator  $Y \equiv SSK(S(K(SS(S(SSK))))K)$  helps out.  $Y$  has the nice property that  $Yx = x(Yx)$ , allowing us to take  $Q \equiv YQ'$  and  $R \equiv YR'$ , where  $R'xy = ya(ybb)(x(x(yb)))$  and  $R'xy = ya(ybaSK)(x(yb)(x(R(yb))))$ .

Putting the pieces together, we form the universal combinator

$$U \equiv SQR,$$



**FIGURE 3.1.** The 425-bit universal combinator  $U'$  in pixels

which has the property that

$$U(\langle C \rangle x) \equiv SQR(\langle C \rangle x) = Q(\langle C \rangle x)(R(\langle C \rangle x)) = Cx.$$

### 3.2.6 Concrete Complexities

At long last, we can define Kolmogorov complexity and prefix complexity concretely.

The Kolmogorov complexity  $C(x)$  of a string  $x$  is the length of a minimal string  $p$  such that  $U(p\$)$  outputs  $x$ . The prefix complexity  $K(x)$  of a string  $x$  is the length of a minimal string  $p$  such that  $U(p\infty)$  outputs  $x$ .

To define Kolmogorov complexity conditional to a string  $y$ , we simply, after applying  $U$  to the input, apply the result to  $(y\$)$  in turn. (It is clear how this approach can be extended to an arbitrary number of conditional arguments.)

The Kolmogorov complexity  $C(x|y)$  of a string  $x$  conditional to a string  $y$  is the length of a minimal string  $p$  such that  $U(p\$(y\$))$  outputs  $x$ . The prefix complexity  $K(x|y)$  of a string  $x$  conditional to a string  $y$  is the length of a minimal string  $p$  such that  $U(p\infty(y\$))$  outputs  $x$ . This way combinatory logic makes an excellent vehicle for defining Kolmogorov complexity concretely, and we obtain the following:

**Theorem 3.2.1**

$$C(x) \leq l(x) + 8, \quad C(\epsilon) = 2,$$

$$\begin{aligned} C(x|y) &\leq l(x) + 2, \\ K(x) &\leq 2l(x) + 231, \quad K(\epsilon) = 11, \\ K(x|l(x)) &\leq l(x) + 689, \\ K(x) &\leq l(x) + 2l(l(x)) + 1066. \end{aligned}$$

For instance, to prove the result about  $K(x)$ , note that a string  $s$  can be encoded in a self-delimiting way as  $1s_01s_1\dots1s_{n-1}0$  of length  $2n+1$ . The combinator  $D$  defined by  $Dx = x0\$(P(x10)(D(x11)))$ , expressible as 77 primitive combinators, decodes this. The universal combinator  $U$  takes 518 bits. A mere 425 bits suffices to encode the currently most parsimonious universal combinator  $U'$  nearly 13 times smaller than the—admittedly less optimized—universal Turing machine of Penrose referred to above. Figure 3.1 represents  $U'$  graphically by a  $17 \times 25$  matrix of pixels, where a white pixel denotes “1” and a black pixel denotes “0.” The matrix should be read in row-major order.

### 3.3 Incompressibility

The quantitative relations between  $K$  and  $C$  show that there must be many incompressible strings with respect to  $K$ . For the  $C$ -complexity, for each  $n$ , the maximal complexity of strings of length  $n$  equals  $n$ , up to a fixed constant. With  $K$ -complexity life is not so simple.

#### Example 3.3.1

We first observe that since the range of  $K(\cdot)$  is the code-word length set of a prefix-code, and it follows from the Kraft Inequality, page 74, that  $\sum 2^{-K(x)} \leq 1$ . This implies the following:

If  $f(x)$  is a function such that the series  $\sum_x 2^{-f(x)}$  diverges, then  $K(x) > f(x)$  infinitely often. Otherwise  $K(x)$  could not satisfy the Kraft Inequality. For instance, choose  $f = f_k$  for each fixed  $k$ , with  $f_k(x) = \log x + \log \log x + \dots + \log \log \dots \log x$ , the last term consisting of the  $k$ -fold iterated logarithm.

Let  $f$  be a function such that the series  $\sum_x 2^{-f(x)} \leq 1$ . By the Kraft Inequality, the set  $\{f(x) : x \in \mathcal{N}\}$  is the length set of the code word set of a prefix-code. If, moreover, this prefix-code is decoded by a prefix partial recursive function, then by Theorem 3.1.1, we have  $K(x) \leq f(x) + O(1)$ . An easy example is  $f(x) = \log(x) + 2 \log \log(x)$ .  $\diamond$

We cannot even give a recursive upper bound  $f(n)$  on the maximum of  $K(x)$  for strings of length  $n$  that is always sharp to within a fixed constant. Nonetheless, we can give a precise expression for the maximal complexity of a string of length  $n$ . As we can expect, almost all  $x$  have nearly maximal complexity. This is expressed in Theorem 3.3.1, which we may regard as a prefix complexity analogue of Theorem 2.2.1 on page 109. Item (ii) gives the distribution of description lengths.

- Theorem 3.3.1**
- (i) For each  $n$ ,  $\max\{K(x) : l(x) = n\} = n + K(n) + O(1)$ .
  - (ii) For each fixed constant  $r$ , the number of  $x$  of length  $n$  with  $K(x) \leq n + K(n) - r$  does not exceed  $2^{n-r+O(1)}$ .

**Proof.** (i) ( $\leq$ ) Let  $U$  be the reference prefix machine of Theorem 3.1.1. Consider a prefix machine  $T$  that on input  $qx$ , where  $U(q) = l(x)$  computes  $T(qx) = x$ . Let  $T_1, T_2, \dots$  be the standard enumeration of prefix machines. Since  $T$  is a prefix machine we have  $T = T_m$  for some  $m$ . Then,  $U(\bar{m}qx) = x$ . Hence,  $K(x) \leq l(x) + K(l(x)) + 2l(m)$  with  $m$  independent of  $x$ .

( $\geq$ ) Since there are  $2^n$  strings  $x$  of length  $n$ , this follows from Item (ii).

(ii) Assume a string  $x$  of length  $n$  satisfies

$$K(x) \leq n + K(n) - r.$$

We now use a remarkable equality of which the proof is omitted at this point:

$$K(x) + K(n|x, K(x)) = K(n) + K(x|n, K(n)) + O(1).$$

(The proof uses Theorem 3.9.1—rather its corollary, Theorem 3.9.2 on page 233—which in turn depends on Theorem 4.3.4 on page 255. We did not find a satisfactory way to avoid dependence on later material.) Substitute  $K(n|x, K(x)) = O(1)$ , since  $n = l(x)$ , to obtain

$$K(x|n, K(n)) \leq n - r + O(1). \tag{3.3}$$

By simple counting, there are fewer than  $2^{n-r+O(1)}$  strings  $x$  of length  $n$  satisfying Equation 3.3, which is the required result.  $\square$

By almost the same proof, the analogue of Theorem 3.3.1 holds for the conditional complexity  $K(x|y)$ . The notion of  $c$ -incompressible strings can be formulated for  $K$ -complexity. But as the upper bound in Theorem 3.3.1 suggests, this is more subtle than for  $C$ -complexity.

**Example 3.3.2** By Theorem 3.3.1, the great majority of all  $x$ 's of length  $n$  cannot be compressed. In fact, for the majority of  $x$ 's the complexity  $K(x)$  significantly exceeds  $l(x)$ . The maximal randomness achievable by a string  $x$  of length  $n$  is  $K(x) = n + K(n) + O(1)$ . It is already natural to call  $x$  incompressible if the length of a shortest program for  $x$  exceeds  $x$ 's length.  $\diamond$

**Definition 3.3.1** Let  $x$  be a string of length  $n$ . If  $K(x) \geq n$ , then  $x$  is *incompressible*.

**Example 3.3.3** Consider the strings that are shortest programs themselves. Are they compressible with respect to  $K$ -complexity? Denote a *shortest program* for  $x$  by  $x^*$ . Clearly,  $K(x^*) \leq l(x^*) + O(1)$ , since a modification of the reference machine  $U$  changes it into an identity machine  $V$  using the same prefix-free set of programs as does  $U$ . The inequality follows if we simulate  $V$  by  $U$ . Moreover,  $K(x^*) \geq l(x^*) + O(1)$ , since otherwise  $x^*$  is not a shortest program for  $x$ , by the usual argument. Consequently, the prefix complexity of a shortest program does not rise above its length:

$$K(x^*) = l(x^*) + O(1).$$

By Theorem 3.3.1 therefore, the number of shortest programs occurring in the set of all strings of length  $n$  does not exceed  $2^{n-K(n)+O(1)}$ . That is, the fraction of shortest programs among the strings of length  $n$  is at most  $2^{-K(n)+O(1)}$ , which goes to zero as  $n$  rises unboundedly. In fact, it can be shown that the number of short programs is small, and that the number of shortest programs for every string  $x$  is  $O(1)$  (Exercise 4.3.5, page 266).  $\diamond$

## Exercises

---

**3.3.1.** [22] Show that neither  $K(x)$  nor  $K(x|l(x))$  is invariant with respect to cyclic shifts. For example,  $K(x_{1:n}) = K(x_{m+1:n}x_{1:m}) + O(1)$  is not satisfied for all  $m$ ,  $1 \leq m \leq n$ .

*Comments.* Hint: choose  $x = 10\dots 0$ ,  $l(x) = 2^k$ , and  $K(m|l(x)) \geq l(l(x)) + O(1)$ .

**3.3.2.** [32] Show that Kamae's result, Exercise 2.7.5 on page 174, does not hold for  $K(x|y)$ .

*Comments.* Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

**3.3.3.** [17] Show that  $K(x) \leq K(x|n, K(n)) + K(n) + O(1)$ , for  $n = l(x)$ .

*Comments.* Hint:  $K(x) \leq K(x, n) + O(1)$  and apply the addition identity to  $K(x, n)$ . Source: G.J. Chaitin, *J. ACM*, 22(1975), 329–340.

**3.3.4.** [15] Show that with  $n = l(x)$ , we have  $C(x|n) \leq K(x) \leq C(x|n) + l^*(C(x|n)) + l^*(n) + O(1)$ .

*Comments.* Hint: this is an easy consequence of Equation 3.1. Source: S.K. Leung-Yan-Cheong and T.M. Cover, *IEEE Trans. Inform. Theory*, IT-24(1978), 331–339.

**3.3.5.** [15] Let  $\phi(x, y)$  be a recursive function.

(a) Show that  $K(\phi(x, y)) \leq K(x) + K(y) + c_\phi$ , where  $c_\phi$  is a constant depending only on  $\phi$ .

- (b) Show that (a) does not hold for  $C$ -complexity.

*Comments.* Hint: in Item (b) use the fact that the logarithmic error term in Theorem 2.8.2, page 182, cannot be improved.

- 3.3.6.** [11] Show that  $K(x) \leq C(x) + C(C(x)) + O(\log C(C(x)))$ .

- 3.3.7.** [23] Let  $f(x, n)$  be the number of binary strings  $p$  of length  $n$  such that  $U(p) = x$ . ( $p$  is a program for  $x$  on the reference machine  $U$  of Theorem 3.1.1.)

- (a) Show that for every  $n \geq K(x)$  we have  $\log f(x, n) = n - K(x, n) + O(1)$ .

- (b) Show that the number of shortest programs of any object is bounded by a universal constant.

- (c) Show that  $K(x, K(x)) = K(x) + O(1)$ .

*Comments.* Hint: in Item (b), use  $K(x) \leq K(x, n) + O(1)$ ; substitute  $n = K(x)$  in the expression in Item (a) to obtain  $\log f(x, K(x)) = K(x) - K(x, K(x)) + O(1) \leq O(1)$ . Hint for Item (c): use the hint of Item (b) and the number  $\log f(x, K(x)) > 0$ . Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

- 3.3.8.** [17] The following equality and inequality seem to suggest that the shortest descriptions of  $x$  contain some extra information besides the description of  $x$ .

- (a) Show that  $K(x, K(x)) = K(x) + O(1)$ .

- (b) Show that  $K(x|y, i - K(x|y, i)) \leq K(x|y, i)$ .

*Comments.* These (in)equalities are in some sense pathological. But they hold also for  $C(\cdot)$ . Source: P. Gács, *Ibid*.

- 3.3.9.** [25] Let  $f(i, x)$  be the number of strings  $y$  such that  $K(y|x) = i$ . Show that  $\log f(i, x) = i - K(i|x) \pm O(1)$ .

- 3.3.10.** [46] The following formulas relate  $C$  and  $K$ :

- (a) Show that  $K(x) = C(x) + C(C(x)) + O(C(C(C(x))))$ .

- (b) Show that  $C(x) = K(x) - K(K(x)) - O(K(K(K(x))))$ .

- (c) Show that  $C(C(x)) - K(K(x)) = O(K(K(K(x))))$ .

- (d) Show that  $K(K(K(x))) \sim C(C(C(x)))$ .

*Comments.* Granted Items (c) and (d), it follows that Items (a) and (b) are equivalent. These formulas express the number of extra bits needed to convert a minimal length program for the reference machine for  $C$ -complexity into a minimal length program for the reference machine of

$K$ -complexity. Source is R.M. Solovay, *Lecture Notes*, 1975, unpublished, and quoted in R.M. Solovay, *Non-Classical Logics, Model Theory and Computability*, N.C.A. da Costa and R. Chuaqui (Eds.), North-Holland, 1977, 283–307.

**3.3.11.** [24] Let  $\phi : \{0, 1\}^* \rightarrow \mathcal{N}$  be a prefix algorithm, that is, a partial recursive function with a prefix-free domain. Then the *extension complexity* of  $x$  with respect to  $\phi$  is defined by  $E_\phi(x) = \min\{l(p) : x \text{ is a prefix of } \phi(p)\}$ , or  $E_\phi(x) = \infty$  if there is no such  $p$ .

(a) Show that there is an additively optimal prefix algorithm  $\phi_0$  such that for any other prefix algorithm  $\phi$  there is a constant  $c$  such that  $E_{\phi_0}(x) \leq E_\phi(x) + c$  for all  $x$ . Select one such  $\phi_0$  as reference and set the *extension complexity*  $E(x) = E_{\phi_0}(x)$ . Similarly, we can define the conditional extension complexity  $E(x|y)$ .

(b) Show that for the relation between the extension complexity  $E$  and the prefix complexity  $K(x)$ , with  $n = l(x)$ ,

$$E(x) \leq K(x) \leq E(x) + K(n) + O(1) \leq E(x) + l^*(n) + O(1).$$

*Comments.* Source: S.K. Leung-Yan-Cheong, T.M. Cover, *IEEE Trans. Inform. Theory*, IT-24(1978), 331–339.

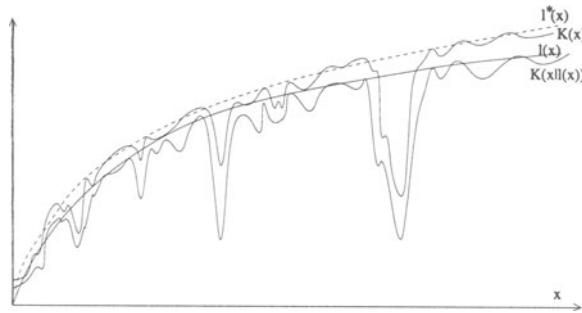
### 3.4 $K$ as an Integer Function

Consider  $K$  as an integer function  $K : \mathcal{N} \rightarrow \mathcal{N}$  and determine its behavior. Most of the properties stated for  $C$  and  $C(\cdot|l(\cdot))$  hold for  $K$  and  $K(\cdot|l(\cdot))$ . We look at the items in the same order, and indicate the similarities and differences.

All of Theorem 2.3.1, page 119, holds with  $K$  substituted for  $C$ . In particular, therefore,  $K(x)$  is unbounded and it (its limit inferior) goes to infinity more slowly than any unbounded partial recursive function. Although the conditional complexity  $K(\cdot|l(\cdot))$  is unbounded, there is no unbounded monotonic lower bound on it. That is,  $K(\cdot|l(\cdot))$  drops below a fixed constant infinitely often, Figure 3.2. This is just like the case of  $C(\cdot|l(\cdot))$  and has a similar proof.

In fact, even the least monotonic upper bounds on the length conditional complexities are similar.  $C(x|l(x))$  reaches upper bound  $\log(x) + O(1)$  for infinitely many  $x$  just as  $C(x)$  does. But, while  $K(x)$  exceeds  $\log x + \log \log x$  infinitely often by Theorem 3.3.1, it is easy to show that  $K(x|l(x)) \leq l(x) + O(1)$  for all  $x$  by Exercise 3.4.2 on page 3.4.2.

Both Theorem 2.3.2, page 121, and Theorem 2.3.3, page 121, hold with  $C$  replaced by  $K$ . That is,  $K$  is not partial recursive and can be computed only for finitely many  $x$ . However,  $K$  is co-enumerable, that is, the set



**FIGURE 3.2.** The graphs of  $K(x)$  and  $K(x|l(x))$

$\{(m, x) : K(x) \leq m\}$  is recursively enumerable. The same properties hold for the conditional complexity  $K(x|y)$ .

There are small differences with respect to the four properties discussed in Section 2.3. The function  $K$  is *continuous* in the sense that  $|K(x) - K(x \pm h)| \leq K(h) + O(1)$ . Function  $K(x)$  mostly “hugs”  $\log^* x + l(x) + l(l(x)) + \dots + O(1)$  in the sense that this is a good upper bound. The function  $K(x)$  has many *fluctuations* for the same reasons as  $C(x)$ . For each  $c$  there is a bound on the length of a run of consecutive  $c$ -incompressible numbers, however, the length of runs of  $c$ -incompressible numbers rises unboundedly with  $c$ .

## Exercises

---

**3.4.1.** [10] Let  $C^+(x) := \max\{C(y) : y \leq x\}$ ,  $K^+(x) := \max\{K(y) : y \leq x\}$ . (a) Show that  $C^+(x) = \log x + O(1)$ .

(b) Show that  $K^+(x) = \log x + K(\lceil \log x \rceil) + O(1)$ .

**3.4.2.** [20] Analyze the integer function  $K(x|n)$  with  $n = l(x)$ .

(a) Show that there is a constant  $c$  such that there are infinitely many  $x$  such that  $K(x|n) \leq c$ .

(b) Let  $h = n - C(x|n)$ . Show that  $K(x|n) \leq C(x|n) + K(h|n) + O(1)$ .

(c) Use Item (b) to show that  $K(x|n) \leq n + O(1)$  for all  $x$ .

(d) Show that  $K(x|n) \leq C(x|n) + \log n + O(1)$ .

**3.4.3.** [12] Show that  $\sum_x 2^{-K(x|l(x))}$  does not converge.

**3.4.4.** [39] (a) Show that  $l^*(x) = \log x + \log \log x + \dots$  (all positive terms) satisfies  $\sum_x 2^{-l^*(x)} < \infty$ .

(b) Show that for all  $x$  we have  $K^+(x) \leq l^*(x) + O(1)$  ( $K^+$  as in Exercise 3.4.1).

(c) Show that for most  $x$  we have  $K^+(x) = l^*(x) + O(1)$ .

*Comments.* Hint for Item b: use Item (a). Source: attributed to T.M. Cover [P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987].

### 3.5 Random Finite Sequences

Recall that  $c$ -incompressible strings with respect to  $C$  coincide with  $c'$ -random strings under the uniform distribution. Since a  $C$ -incompressible string can be equated with its shortest program, Example 3.3.3 shows that for appropriate constants  $c$  and  $c'$ , the  $c'$ -incompressible strings with respect to  $K$ -complexity are a superset of the  $c$ -incompressible strings with respect to  $C$ -complexity. But the following lemma shows that the converse does not hold: there are incompressible strings with respect to  $K$ -complexity that are not  $O(1)$ -random with respect to the uniform distribution.

**Lemma 3.5.1** *For infinitely many  $n$  there are strings  $x$  of length  $n$  that have  $K(x) \geq n$  and  $C(x) \leq n - \log n$ .*

**Proof.** Consider infinite sequences  $\omega$ . By Corollary 2.5.1, for each  $\omega$  there are infinitely many  $n$  such that  $C(\omega_{1:n}) < n - \log n$ . On the other hand, in Theorem 3.6.1 we will show that  $\omega$  is (Martin-Löf) random with respect to the uniform measure iff  $K(\omega_{1:n}) \geq n + O(1)$  for all  $n$ . Hence, the  $n$ -length prefixes of any random  $\omega$ , for the sequence of values  $n$  denoting points where complexity oscillations reach  $\log n$  below  $n$ , is an example demonstrating the lemma.  $\square$

**Corollary 3.5.1** For each  $c$ , there are finite strings  $x$  that are incompressible according to prefix complexity ( $K(x) \geq l(x)$ ) and that are not  $c$ -random finite strings with respect to the uniform distribution.

**Example 3.5.1** How much can  $K(x)$  exceed  $C(x)$  in effective terms? For most  $x$  we have

$$K(x) = C(x) + K(C(x)) + O(1).$$

Namely, on the one hand, for all  $x$  we have  $K(x) \leq C(x) + K(C(x)) + O(1)$ , since if  $p$  is a shortest program for  $x$  (on an ordinary machine) with  $l(p) = C(x)$ , and  $q$  is a shortest program for  $l(p)$  (on a prefix machine) with  $l(q) = K(l(p))$ , then  $qp$  is a program for  $x$  on some prefix machine. On the other hand, by Theorem 3.3.1 for most  $x$  we have  $K(x) \geq C(x) + K(C(x)) + O(1)$ . Hence, we can conclude that the difference between  $K$  and  $C$  satisfies

$$K(x) - C(x) \leq \log^*(n) + l(n) + l(l(n)) + \dots + O(1),$$

for all  $x$  of length  $n$ , and that the inequality is nearly sharp for infinitely many  $x$ , namely, for those  $x$  with  $C(x) = n + O(1)$  and  $K(n) \geq l(n) + l(l(n))$ .  $\diamond$

We have seen that  $C(x)$  is less than  $K(x)$ , since in the definition of  $C(x)$  we did not take the information into account that is needed to make the program prefix-free. We can express  $C$  precisely in terms of  $K$ , namely,  $C$  is the unique (to within an additive constant) function satisfying

$$C(x) = \min\{i : K(x|i) \leq i\} = K(x|C(x)) + O(1). \quad (3.4)$$

(This equation can be generalized in the natural way to  $C(x|y)$ .)

( $\geq$ ) We prove  $C(x) \geq K(x|C(x)) + O(1)$ . The shortest Turing machine program  $x^*$  for  $x$  has length  $C(x)$ . There is a prefix machine that given  $C(x)$ , computes  $x$  from  $x^*$ . Clearly, this also implies  $C(x) \geq \min\{i : K(x|i) \leq i\} + O(1)$ .

( $\leq$ ) We prove  $C(x) \leq \min\{i : K(x|i) \leq i\} + O(1)$ . This is the same as, if  $K(x|i) \leq i$ , then  $C(x) \leq i + O(1)$ . Assume  $x$  is computed by reference prefix machine  $U$ , given  $i$ , from input  $p$ ,  $l(p) \leq i$ . A Turing machine (not necessarily prefix machine)  $T$ , when presented input  $0^{i-l(p)-1}1p$ , in case  $i-l(p)-1 > 0$ , or input  $p$  otherwise, can extract  $i$  from the input length and simulate  $U$  on  $p$ , using the extracted  $i$  as the conditional value. By Theorem 2.1.1, page 97,  $C(x) \leq C_T(x) \leq i$ , up to additive constants. This proves the required inequality. The same proof, with  $C(x)$  substituted for  $i$ , proves  $C(x) \leq K(x|C(x)) + O(1)$ . This finishes the proof of Equation 3.4. Source: L.A. Levin, *Soviet Math. Dokl.*, 17:2(1976), 522–526.

Denoting  $l(x)$  by  $n$ , Theorem 2.4.2 on page 132 identifies

$$\delta_0(x|L) = n - C(x|n) - 1$$

as a universal Martin-Löf test for the uniform distribution  $L$ . Corollary 3.5.1 shows that we cannot simply express randomness of a finite string under the uniform distribution in terms of its incompressibility in the sense of prefix complexity.

Substituting the length-conditional version of Equation 3.4,

$$C(x|n) = K(x|C(x), n) + O(1),$$

we obtain an expression that involves both  $K$  and  $C$ . This can be avoided, in a somewhat contrived manner, by introducing an auxiliary complexity based on  $K$ . Define

$$\bar{K}(x; k) = \min\{i : K(x|k-i) \leq i\}.$$

Similar to Equation 3.4 we obtain

$$\bar{K}(x; k) = K(x|k - \bar{K}(x; k)).$$

Then we can define the conditional auxiliary complexity in a similar way and denote it by  $\overline{K}(x; k|y)$ . Using the same arguments as above, we find

$$\overline{K}(x; k|k) = C(x|k) + O(1).$$

Then we can express the randomness deficiency found by a universal  $P$ -test in terms of the auxiliary complexity. For  $P$  a recursive function it can be proven that

$$\delta_0(x|P) = -\log P(x) - \overline{K}(x; -\log P(x)|n)$$

is a universal  $P$ -test. If  $P$  is the uniform distribution on strings of length  $n$ , then  $-\log P(x) = n$  for all  $x$  of length  $n$  and

$$\delta_0(x|L) = n - \overline{K}(x; n|n) = n - C(x|n) + O(1)$$

is the familiar universal test for the uniform distribution  $L$ . Thus, while Martin-Löf's universal test for randomness of *finite* strings under the uniform distribution is easily expressed in terms of  $C$ -complexity, Section 2.4, we have to go through quite some contortions to express it in terms of (a variant of)  $K$ -complexity. These contortions are necessary, since it is only the form using  $K$ -complexity that is generalizable to universal  $P$ -tests for arbitrary computable distributions  $P$ .

The extended theory of  $P$ -tests, and exact expressions for universal  $P$ -tests, are treated in Section 4.3.5. These constructions show that the presence of  $C$  in Martin-Löf's expression  $n - C(x | n)$  is a lucky freak of the uniform distribution. See P. Gács, *Komplexität und Zufälligkeit*, Ph.D. Thesis, Mathematics Department, J.W. Goethe Universität, Frankfurt am Main, 1978.

In contrast, while in Section 2.5 we did not succeed at all in expressing Martin-Löf's universal sequential test for randomness of *infinite* sequences under the uniform distribution in terms of  $C(x)$ , in the next Section 3.6 this will turn out to be straightforward in terms of  $K(x)$ .

## Exercises

**3.5.1.** [43] Recall the universal Martin-Löf test for randomness of finite strings  $x$  of length  $n$  under the uniform distribution:  $\delta_0(x|L) = n - C(x) + O(1)$ . This was based on Theorems 2.1.2, 2.2.1 on pages 100, 109. The analogous facts of Theorem 3.3.1 for  $K(x)$  suggest a test  $\theta_0(x|L) = n + K(n) - K(x) + O(1)$ , with the same intuitive interpretation.

(a) Show that  $\theta_0(x|L) \geq \delta_0(x|L) + O(\log \delta_0(x|L) + 1)$ . Hence, if  $\theta_0$  is small, then  $\delta_0$  is small (and this is the sense in which strings that are “ $K$ -random” are “ $C$ -random”).

(b) Construct an infinite sequence of finite strings  $x_m$  with the following properties: (i)  $l(x_m) \rightarrow \infty$  for  $m \rightarrow \infty$ ; (ii)  $C(x_m) = l(x_m) + O(1)$ ; and

$\lim_{m \rightarrow \infty} \theta_0(x_m | L) / \log^2 l(x_m) = 1$ . In this sense “ $C$ -randomness” does not entail “ $K$ -randomness.”

(c) Use Exercise 3.3.10 to show that

$$\limsup_{n \rightarrow \infty} \frac{K(n) - C(n) - K(K(n))}{K(K(K(n)))} \leq 1.$$

Use Item (b) to improve this to an equality. The same method provides a counterexample to various improvements in the error terms of the formulas relating  $C(x)$  and  $K(x)$  in Exercise 3.3.10.

(d) Use Item (c) to show that the seductive equality  $K(x) = C(x) + K(C(x)) + O(1)$  is *false* in general (although it obviously holds for all  $x$  that are random enough).

*Comments.* Source is R.M. Solovay, *Lecture Notes*, UCLA, 1975, unpublished.

## 3.6 \*Random Infinite Sequences

Let  $\omega$  be an infinite binary sequence. We are interested in the behavior of  $K(\omega_{1:n})$  as a function of  $n$ . The complexity  $K(\omega_{1:n})$  will turn out to be nonmonotonic in  $n$ , just like  $C$ -complexity. We will find that an infinite sequence  $\omega$  is random under the uniform measure iff a simple condition on  $K(\omega_{1:n})$  is satisfied. For the  $C$ -complexity we were unable to find such a condition in Section 2.5.

### Example 3.6.1

We show that  $K(0^n)$  as a function of  $n$  is not monotonic. Choose  $n$  with  $K(n) \geq l(n)$ . Let  $m = 2^k \geq n$  with  $k$  minimal. Then  $0^n$  is a prefix of  $0^m$ . Firstly,  $K(0^n) \geq \log n + O(1)$ . Secondly,  $K(m) = K(k) + O(1) \leq 2 \log \log n + O(1)$ . Therefore,  $K(0^n)$  is exponential in  $K(0^m)$ .  $\diamond$

Theorems 2.5.1, 2.5.4 on pages 136, 145 show, for almost all infinite binary sequences, that the oscillations of the  $C$ -complexity of initial segments of almost all sequences, as a function of their length, are confined to a thin wedge below the identity function  $n + O(1)$ . Both random sequences and some nonrandom sequences oscillate in this wedge. For  $K$ -complexity the wedge

$$n + O(1) < K(\omega_{1:n}) \leq n + K(n) + O(1)$$

contains all and only infinite random sequences. (The upper bound holds for all infinite sequences by Theorem 3.3.1, the lower bound for random sequences is proved in Theorem 3.6.1 on page 212.) The complexity oscillations are in some form still there, but  $K$  exceeds  $C$  by so much that the complexity of the prefix does not drop below the length of the prefix itself (for random infinite  $\omega$ ).

### 3.6.1 Explicit Universal Randomness Tests

The idea that random infinite sequences are those sequences such that the complexity of the initial  $n$ -length segments is at most a fixed additive constant below  $n$ , for all  $n$ , is one of the first-rate ideas in the area of Kolmogorov complexity. In fact, this was one of the motivations for Kolmogorov to invent Kolmogorov complexity in the first place, see page 55. We have seen in Section 2.5 that this does not work for the plain Kolmogorov complexity  $C(\cdot)$ , due to the complexity oscillations. The next result is important, and is a culmination of the theory. For prefix-complexity  $K(\cdot)$  it is indeed the case that random sequences are those sequences for which the complexity of each initial segment is at least its length.

A.N. Kolmogorov suggested the above relation between the complexity of initial segments and randomness of infinite sequences [A.N. Kolmogorov, *IEEE Trans. Inform. Theory*, IT-14:5(1968), 662–664]. This approach being incorrect using  $C(\cdot)$  complexity, P. Martin-Löf [*Inform. Contr.*, 9(1966), 602–619] developed the theory as put forth in Section 2.5. Nevertheless, Kolmogorov did not abandon the general outlines of his original idea of connecting randomness of infinite sequences with complexity, see pp. 405–406 of [V.A. Uspensky, *J. Symb. Logic*, 57:2(1992), 385–412].

C.P. Schnorr [*J. Comput. System Sci.*, 7(1973), 376–388] for the uniform distribution, and L.A. Levin [*Sov. Math. Dokl.*, 14(1973), 1413–1416] for arbitrary computable distributions, introduced simultaneously and independently similar but unequal versions of complexity to characterize randomness. These versions are “process complexity” and the “monotone” variant of complexity  $Km(x)$  (Definition 4.5.9, page 282), respectively. They gave the first realizations of Kolmogorov’s idea by showing Corollary 4.5.3, page 295: an infinite sequence  $\omega$  is random iff  $|Km(\omega_{1:n}) - n| = O(1)$  (Levin) and a similar statement for process complexity (Schnorr).

G.J. Chaitin [*J. ACM*, 22(1975), 329–340], proposed calling an infinite sequence  $\omega$  random iff  $K(\omega_{1:n}) \geq n - O(1)$  for all  $n$ . Happily, this proposal characterizes once again precisely those infinite binary sequences that are random in Martin-Löf’s sense (without proof attributed to C.P. Schnorr, 1974, in Chaitin’s paper). This important result was widely circulated, but the first published proof appears, perhaps, only as Corollary 3.2 in [V.V. V'yugin, *Semiotika i Informatika*, 16(1981), 14–43, in Russian]. See for historical notes [A.N. Kolmogorov and V.A. Uspensky, *SIAM J. Theory Probab. Appl.*, 32(1987), 387–412].

Another equivalent proposal in terms of constructive measure theory was given by R.M. Solovay, Exercise 2.5.9, page 152. The fact that such different effective formalizations of infinite random sequences turn out to define the same mathematical object constitutes evidence that our intuitive notion of infinite sequences that are effectively random coincides with the precise notion of Martin-Löf random infinite sequences.

**Theorem 3.6.1** *An infinite binary sequence  $\omega$  is random with respect to the uniform measure iff there is a constant  $c$  such that for all  $n$ ,  $K(\omega_{1:n}) \geq n - c$ .*

**Proof.** (ONLY IF) Assume  $\omega$  is a random infinite sequence. Then, for each sequential test  $\delta$ , we have  $\delta(\omega) < \infty$ , Section 2.5. We construct a particular sequential test, say  $\delta$ , such that if  $\delta(\omega) < \infty$ , then there is a constant  $c$  such that  $K(\omega_{1:n}) \geq n - c$ , for all  $n$ .

Recall the usual preliminaries. If  $y$  is a finite string, then  $\Gamma_y$  denotes the set of infinite sequences  $\omega$  that start with  $y$ . Let  $\lambda$  denote the uniform measure, so that with  $l(y) = n$  we have  $\lambda(\Gamma_y) = 2^{-n}$ . Geometrically speaking,  $\Gamma_y$  corresponds to the half-open interval  $[0.y, 0.y + 2^{-n})$ .

Define  $\delta$  as follows. Consider a sequence  $A_0 \supset A_1 \supset \dots$  of sets of finite strings such that

$$A_k = \bigcup_{n \geq 1} \{y : K(y) \leq n - k - c', n = l(y)\},$$

with  $c'$  a fixed constant that is large enough to make the remainder of the proof go through. Define a total function  $\gamma$  by  $\gamma(y) = \sup_{k \in \mathcal{N}} \{k : y \in A_k\}$ . Since  $K$  is a co-enumerable function (it can be approximated from above),  $\gamma$  is an enumerable function (it can be approximated from below). Finally, for each infinite sequence  $\omega$  define  $\delta(\omega) = \sup_{n \in \mathcal{N}} \{\gamma(\omega_{1:n})\}$ .

To prove that  $\delta$  is a sequential test, it suffices to show that  $\lambda\{\omega : \delta(\omega) \geq k\} \leq 2^{-k}$ , for each  $k \geq 0$ .

There are fewer than  $a(k, n) = 2^{n-K(n)-k}$  strings  $y$  of length  $n$  of complexity  $K(y) \leq n - k - c'$  (Theorem 3.3.1). This bounds the number of strings of length  $n$  in  $A_k$ . Overestimating  $\lambda\{x : \delta(\omega) \geq k\}$ , using  $a(k, n)$ , rearranging terms, and employing Kraft's Inequality (page 74) to argue that  $\sum 2^{-K(n)} \leq 1$ , we have

$$\begin{aligned} \lambda\{\omega : \delta(\omega) \geq k\} &\leq \sum_{y \in A_k} \lambda\{\Gamma_y\} \\ &\leq \sum_{n \in \mathcal{N}} a(k, n) 2^{-n} \\ &= 2^{-k} \sum_{n \in \mathcal{N}} 2^{-K(n)} \leq 2^{-k}. \end{aligned}$$

It is now proved that  $\delta$  is a sequential test. If  $\omega$  is random in the sense of Martin-Löf, then  $\delta(\omega) < \infty$ . That is, for each random  $\omega$  there is a constant  $c < \infty$ , such that  $K(\omega_{1:n}) \geq n - c$ , for all  $n$ .

(IF) Suppose that  $\omega$  is not random, that is, there is a sequential test  $\delta$  such that  $\delta(\omega) = \infty$ . We show that this implies that  $n - K(\omega_{1:n})$  is positively unbounded. Let  $\gamma$  be the defining enumerable function of  $\delta$  as in Definition 2.5.1. For all  $k \geq 1$ , define  $A_k$  as the maximal prefix-free subset of  $\{y : \gamma(y) \geq k\}$ . Therefore,  $\sum_{y \in A_k} 2^{-l(y)} \leq 1$  by Kraft's

Inequality on page 74. Then

$$L = \{l(y) - k : y \in A_{2k}, k > 1\}$$

satisfies Kraft's Inequality

$$\sum_{k>1} \sum_{y \in A_{2k}} 2^{k-l(y)} \leq \sum_{k>1} 2^{-k} \leq 1.$$

This means that  $L$  is the length set of a prefix-code. Use  $\gamma$  to enumerate  $\bigcup\{A_{2k} : k > 1\}$ . Use this enumeration to construct a prefix machine  $T$  such that  $K_T(y) = l(y) - k$  for all  $y \in A_{2k}$ ,  $k > 1$ .

We have assumed  $\delta(\omega) = \infty$ . Hence, for each  $k$ , there is an  $n$  such that  $\omega_{1:n} \in A_{2k}$ , which means  $K_T(\omega_{1:n}) \leq n - k$ . By Theorem 3.1.1, there is a constant  $c_T$  such that  $K(\omega_{1:n}) \leq K_T(\omega_{1:n}) + c_T$ . Therefore,

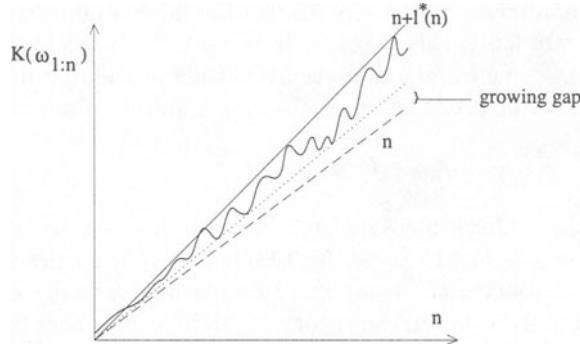
$$\limsup_{n \rightarrow \infty} (n - K(\omega_{1:n})) = \infty.$$

□

**Corollary 3.6.1** The function  $\rho_0(\omega|\lambda) = \sup_{n \in \mathcal{N}} \{n - K(\omega_{1:n})\}$  characterizes the random infinite sequences by  $\rho_0(\omega|\lambda) < \infty$  iff  $\omega$  is random with respect to the uniform measure  $\lambda$ . (In terms of our previous notation  $\rho_0(\omega|\lambda) < \infty$  iff  $\delta_0(\omega|\lambda) < \infty$  where  $\delta_0(\cdot|\lambda)$  is the universal sequential  $\lambda$ -test with  $\lambda$  the uniform measure.)

There are different types of tests that characterize precisely the same class of random infinite sequences. The sequential tests are but one type. The test  $\rho_0$  is an example of an integral test (a universal one) with respect to the uniform measure as defined in Section 4.5.6. The introduction of different types of tests awaits the machinery developed in Chapter 4. The theory of integral tests and martingale tests of infinite sequences is developed in Sections 4.5.6 and 4.5.7, respectively. There we also give exact expressions for tests that separate the random infinite sequences from the nonrandom ones with respect to any recursive measure.

**Example 3.6.2** The proof of Theorem 3.6.1 shows that the separation of a random infinite sequence from the nonrandom ones involves a complexity gap. That is, for random  $\omega$  the amount by which  $K(\omega_{1:n})$  exceeds  $n$  goes to  $\infty$  in the limit; for nonrandom  $\omega$  the amount by which  $n$  exceeds  $K(\omega_{1:n})$  may fluctuate but is also unbounded. So the value of  $K(\omega_{1:n})$  as a function of  $n$  is bounded away from the diagonal in either direction for all  $\omega$ : nonconstant far below for the nonrandom sequences, and nonconstant far above for the random sequences, Figure 3.3. See also [G.J. Chaitin, *Algorithmic Information Theory*, Cambridge Univ. Press, 1987].



**FIGURE 3.3.** Complexity oscillations of a typical random sequence  $\omega$

$\omega$ IS RANDOM iff	$\omega$ IS NOT RANDOM iff
$\rho_0(\omega \lambda) < \infty$ iff $\exists c \forall n [K(\omega_{1:n}) \geq n - c]$ iff $\lim_{n \rightarrow \infty} K(\omega_{1:n}) - n = \infty$	$\rho_0(\omega \lambda) = \infty$ iff $\forall c \exists n [K(\omega_{1:n}) < n - c]$ iff $\limsup_{n \rightarrow \infty} n - K(\omega_{1:n}) = \infty$

**FIGURE 3.4.**  $K$ -complexity criteria for randomness of infinite sequences

In the  $C$ -complexity version, random infinite sequences have oscillations of the complexity of initial segments below the diagonal. With  $K$  complexity similar oscillations must take place above the diagonal. Some relevant properties are surveyed in Figure 3.4.

For example, the maximal complexity of some  $\omega_{1:n}$  is  $K(\omega_{1:n}) = n + K(n) + O(1)$ , Theorem 3.3.1. But similar to the case for  $C(\omega)$  (Theorem 2.5.1 on page 136), no single  $\omega$  satisfies  $K(\omega_{1:n}) \geq n + K(n) + O(1)$  for all  $n$ . In fact, in analogy to the proof of Theorem 2.5.1, we find that for all  $\omega$ , for infinitely many  $n$ ,

$$\begin{aligned} K(\omega_{1:n}) &\leq K(C(\omega_{1:n})) + C(\omega_{1:n}) + O(1) \\ &\leq n - g(n) + K(n - g(n)) + O(1) \\ &\leq n - g(n) + K(n) + O(1), \end{aligned}$$

where  $g(n)$  is as defined in the proof of Theorem 2.5.1. Actually,

$$\limsup_{n \rightarrow \infty} K(n) - g(n) = \infty,$$

since the series  $\sum_n 2^{-g(n)}$  diverges and the series  $\sum_n 2^{-K(n)}$  converges. The infinity on the right-hand side means something very slowly increasing, like  $\log^* n$  and  $l^*(n)$ . Generally, the complexity oscillations of a random sequence  $\omega$  will look as in Figure 3.3.

In analogy to Theorem 2.5.4, page 145, we observe the following: By Theorem 2.5.5, page 146, the set of infinite binary sequences that are

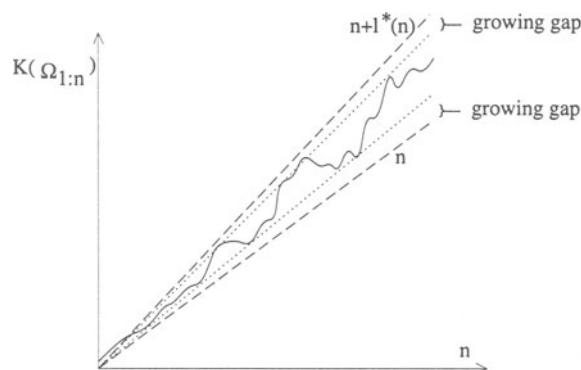
random in the sense of Martin-Löf have (uniform) measure one. By Theorem 3.6.1, this is precisely the set of  $\omega$ 's with  $K(\omega_{1:n}) \geq n + O(1)$ . To get some insight into the amplitude of the upward oscillations, we note that the set of  $\omega$ 's such that for infinitely many  $n$ ,

$$K(\omega_{1:n}) \geq n + K(n) + O(1)$$

has uniform measure one, Exercise 3.6.3 on page 219. M. van Lambalgen [*J. Symb. Logic*, 52(1987), 725–755] suggests that far from being a nuisance, the complexity oscillations actually enable us to discern a fine structure in the theory of random sequences (see also Exercise 3.6.5, page 221). For all sequences relatively low in the recursive hierarchy, like  $\Delta_2^0$  definable sequences (for a definition see Exercise 3.6.5 on page 221), the upward oscillations are not maximal, since

$$\lim_{n \rightarrow \infty} n + K(n) - K(\omega_{1:n}) = \infty.$$

There are Martin-Löf random sequences that are  $\Delta_2^0$  definable, such as Chaitin's halting probability number  $\Omega$  of Section 3.6.2. The complexity oscillations of this restricted type of random sequences, of which  $\Omega$  is a typical example, are confined to a narrower wedge, as in Figure 3.5, than the general random sequences of Figure 3.3. The monotone complexity, Chapter 4, was developed to smooth out the oscillations and to characterize randomness. But monotone complexity also obliterates all quantitative differences between Martin-Löf random sequences, and hence does not allow us to distinguish stronger randomness properties.  $\diamond$



**FIGURE 3.5.** Complexity oscillations of  $\Omega$

### 3.6.2 Halting Probability

It is impossible to construct an infinite random sequence by algorithmic means. But, using the reference prefix machine  $U$  of Theorem 3.1.1, we can define a particular random infinite binary sequence in a more or less natural way. The *halting probability* is the real number  $\Omega$  defined by

$$\Omega = \sum_{U(p) < \infty} 2^{-l(p)},$$

the sum taken over all inputs  $p$  for which the reference machine  $U$  halts. Since  $U$  halts for some  $p$ , we have  $\Omega > 0$ . Because  $U$  is a prefix machine, the set of its programs forms a prefix-code, and by Kraft's Inequality, page 74, we find  $\Omega \leq 1$ . Actually,  $\Omega < 1$  since  $U$  does not always halt.

We call  $\Omega$  the halting probability because it is the probability that  $U$  halts if its program is provided by a sequence of fair coin flips. The number  $\Omega$  has interesting properties. In the first place, the binary representation of the real number  $\Omega$  encodes the halting problem very compactly. Denote the initial  $n$ -length segment of  $\Omega$  after the decimal dot by  $\Omega_{1:n}$ . If  $\Omega$  were a terminating binary rational number, then we use the representation with infinitely many zeros, so that  $\Omega < \Omega_{1:n} + 2^{-n}$ . We shall show that the binary expansion of  $\Omega$  is an *incompressible* sequence.

**Claim 3.6.1** Let  $p$  be a binary string of length at most  $n$ . Given  $\Omega_{1:n}$ , it is decidable whether the reference prefix machine  $U$  halts on input  $p$ .

**Proof.** Clearly,

$$\Omega_{1:n} \leq \Omega < \Omega_{1:n} + 2^{-n}. \quad (3.5)$$

Dovetail the computations of  $U$  on all inputs as follows: The first phase consists of  $U$  executing one step of the computation on the first input. In the second phase,  $U$  executes the second step of the computation on the first input and the first step of the computation of the second input. Phase  $i$  consists of  $U$  executing the  $j$ th step of the computation on the  $k$ th input, for all  $j$  and  $k$  such that  $j + k = i$ . We start with an approximation  $\Omega' := 0$ . Execute phases 1, 2,  $\dots$ . Whenever any computation of  $U$  on some input  $p$  terminates, we improve our approximation of  $\Omega$  by executing

$$\Omega' := \Omega' + 2^{-l(p)}.$$

This process eventually yields an approximation  $\Omega'$  of  $\Omega$ , such that  $\Omega' \geq \Omega_{1:n}$ . If  $p$  is not among the halted programs that contributed to  $\Omega'$ , then  $p$  will never halt. With a new  $p$  halting we add a contribution of  $2^{-l(p)} \geq 2^{-n}$  to the approximation of  $\Omega$ , contradicting Equation 3.5 by

$$\Omega \geq \Omega' + 2^{-l(p)} \geq \Omega_{1:n} + 2^{-n}.$$

□

**Claim 3.6.2** There is a constant  $c$  such that  $K(\Omega_{1:n}) \geq n - c$  for all  $n$ .

That is,  $\Omega$  is a *particular random real*, and one that is naturally defined to boot. That  $\Omega$  is random implies that it is not computable, and therefore *transcendental*. Namely, if it were computable, then  $K(\Omega_{1:n}|n) = O(1)$ , which contradicts Claim 3.6.2. By the way, irrationality of  $\Omega$  implies that both inequalities in Equation 3.5 are strict.

**Proof.** From Claim 3.6.1 it follows that given  $\Omega_{1:n}$ , one can calculate all programs  $p$  of length not greater than  $n$  for which the reference prefix machine  $U$  halts. For every  $x$  that is not computed by any of these halting programs, the shortest program  $x^*$  has size greater than  $n$ , that is,  $K(x) > n$ . Hence, we can construct a recursive function  $\phi$  computing such high complexity  $x$ 's from initial segments of  $\Omega$ , such that for all  $n$ ,

$$K(\phi(\Omega_{1:n})) > n.$$

Given a description of  $\phi$  in  $c$  bits, for each  $n$  we can compute  $\phi(\Omega_{1:n})$  from  $\Omega_{1:n}$ , which means

$$K(\Omega_{1:n}) + c \geq n,$$

which was what we had to prove. □

**Corollary 3.6.2** By Theorem 3.6.1, we find that  $\Omega$  is random with respect to the uniform measure.

It is possible to determine the first couple of bits of the halting probability. For example, with the concrete Kolmogorov complexity fixed in Section 3.2, J. Tromp has computed  $0.00106502 < \Omega < 0.217643$ . This is because many short strings are either not syntactically correct programs for the concrete universal machine, or they halt quickly, or looping is easily detected. But knowing, say, the first 10,000 bits of  $\Omega$  enables us to solve the halting of all programs of fewer than 10,000 bits. This includes programs looking for counterexamples to Fermat's Last Theorem, Goldbach's Conjecture, Riemann's Hypothesis, and most other conjectures in mathematics that can be refuted by single finite counterexamples. Moreover, for all axiomatic mathematical theories that can be expressed compactly enough to be conceivably interesting to human beings, say in fewer than 10,000 bits,  $\Omega_{10,000}$  can be used to decide for every statement in the theory whether it is true, false, or independent. Finally, knowledge of  $\Omega_{1:n}$  suffices to determine whether  $K(x) \leq n$  for each finite binary string  $x$ . Thus,  $\Omega$  is truly the number of Wisdom, and "can be known of, but not known, through human reason" [C.H. Bennett and M. Gardner, *Scientific American*, 241:11(1979), 20–34]. But even if you possess  $\Omega_{1:10,000}$ , you cannot use it except by spending time of a thoroughly unrealistic nature. (The time  $t(n)$  it takes to find all halting programs of length less than  $n$  from  $\Omega_{1:n}$  grows faster than any recursive function.)

## Exercises

---

**3.6.1.** [21] Let  $A$  be an infinite recursively enumerable set of natural numbers.

(a) Show that if we define  $\theta = \sum_{n \in A} 2^{-K(n)}$ , then  $K(\theta_{1:n}|n) \geq n - O(1)$  for all  $n$ . (Therefore,  $\theta$  is a random infinite sequence in the sense of Martin-Löf by Theorem 3.6.1.)

(b) Show that  $\theta$  is not a recursive number.

*Comments.* Note that for  $A$  the set of natural numbers we have  $\theta = \Omega$ , the halting probability in Section 3.6.2. Because  $\theta$  is not a recursive real number it is irrational and even transcendental. Source: G.J. Chaitin, *Algorithmic Information Theory*, Cambridge University Press, 1987.

**3.6.2.** [27] Show that a real number in the unit interval  $[0, 1]$  (equivalently, infinite sequence over  $r$  letters) is Martin-Löf random with respect to the uniform distribution in base  $r$  iff it is random with respect to the uniform distribution in base  $s$ .

*Comments* Hint: by Exercise 3.1.4 we have for each pair of integers  $r, s \geq 2$  that  $|K_r(x) - K_s(x)| \leq c_{r,s}$  for all  $x \in \mathcal{N}$  and some constant  $c_{r,s}$  independent of  $x$ . Use the fact that an infinite binary sequence  $\omega$  is random with respect to the uniform distribution iff  $K(\omega_{1:n}) \geq n - c$  for some constant  $c$  and all  $n$ , Theorem 3.6.1. (Note  $K \equiv K_2$ .) The argument cannot be too trivial since a total recursive 1-1 function  $f$  can map  $\omega = \omega_1\omega_2\dots$  to  $\zeta = \omega_10\omega_20\dots$ . Then  $f^{-1}(\zeta) = \omega$ . But  $\zeta$  is not random, even if  $\omega$  is. Source: C. Calude and H. Jürgenson, pp. 44–66 in: H. Maurer, J. Karhumaki, and G. Rozenberg (Eds.), *Results and Trends in Theoretical Computer Science*, Springer-Verlag, Berlin, 1994. This result is a special case of the simple and much more general approach of Exercise 4.5.14 on page 304.

**3.6.3.** [29] We investigate the complexity oscillations of  $K(\omega_{1:n})$  for infinite binary sequences  $\omega$ . If  $\omega$  is an infinite sequence that is random in the sense of Martin-Löf, then these oscillations take place above the identity line,  $K(\omega_{1:n}) \geq n$ . The maximal complexity of  $\omega_{1:n}$  is  $K(\omega_{1:n}) = n + K(n) + O(1)$ , Theorem 3.3.1. But similar to the case for  $C(\omega)$ , Theorem 2.5.1, page 136, no  $\omega$  satisfies the following: there is a constant  $c$  such that for all  $n$ , we have  $K(\omega_{1:n}) \geq n + K(n) - c$ .

(a) Show that for all  $\omega$  infinitely many  $n$ ,  $K(\omega_{1:n}) \leq n + K(n) - g(n) + O(1)$ , where  $g(n)$  as defined in the proof of Theorem 2.5.1. Can you generalize this to obtain the analogue of Theorem 2.5.1?

(b) Let the series  $\sum 2^{-f(n)} < \infty$  converge recursively. Show that for almost all infinite binary sequences  $\omega$  we have  $K(\omega_{1:n}) \geq n + K(n) - f(n)$  ( $\geq n + O(1)$ ) for all  $n$ . This gives a lower bound of on the dips of the downward oscillations for almost all  $\omega$ .

(c) In analogy with Theorem 2.5.4, page 145, we can note the following: By Theorem 2.5.5, page 146, the set of infinite binary sequences that are random in the sense of Martin-Löf have uniform measure one. By Theorem 3.6.1, therefore, the superset of  $\omega$ 's with  $K(\omega_{1:n}) \geq n + O(1)$  has also uniform measure one. Show that the set of  $\omega$ 's with

$$K(\omega_{1:n}) \geq n + K(n) + O(1)$$

for infinitely many  $n$  has uniform measure one. This gives some insight into the amplitude of the upward oscillations. Not all Martin-Löf random sequences achieve this.

(d) Let  $f(x)$  be a recursive function such that the series  $\sum 2^{-f(x)}$  diverges. Show that the set of  $\omega$ 's with

$$K(\omega_{1:n}) \geq n + f(n) + O(1)$$

for infinitely many  $n$  has uniform measure one.

*Comments.* Hint for Item (a):  $K(\omega_{1:n}) \leq C(\omega_{1:n}) + K(C(\omega_{1:n})) + O(1)$ , and use Theorem 2.5.1, page 136. Note that  $\limsup_{n \rightarrow \infty} K(n) - g(n) = \infty$ . This gives an upper bound on the dips of the downwards oscillations: the least monotonic upper bound on  $K(n) - g(n)$  rises very, very, slowly—more slowly than the  $k$ -fold iterated logarithm for any fixed  $k$ . Hint for Item (b):  $\sum 2^{-f(n)-d} \leq 1$  for some nonnegative constant  $d$ . By the Kraft Inequality, page 74, there is an effective prefix-code  $E$  such that  $l(E(n)) = f(n) + d$  for all  $n$ . By Theorem 3.1.1,  $K(n) \leq l(E(n))$  for every effective prefix-code  $E$ , up to a constant depending only on  $E$ . Thus,  $K(n) \leq f(n) + d$  for a different constant  $d$  depending only on  $f$ . The result now follows from Section 3.6. Source: M. van Lambalgen [*Random Sequences*, Ph.D. Thesis, University of Amsterdam, 1987; *J. Symb. Logic*, 52(1987), 725–755]. Items (c) and (d) are attributed to R.M. Solovay.

**3.6.4.** [35] Recall that an infinite sequence  $\omega$  is recursive iff  $C(\omega_{1:n}) \leq C(n) + O(1)$ , Example 2.3.4. If  $\omega$  is recursive, then for all  $n$  we have  $K(\omega_{1:n}|n) = O(1)$  and  $K(\omega_{1:n}) \leq K(n) + O(1)$ .

(a) Show that if  $K(\omega_{1:n}|n) = O(1)$ , then  $\omega$  is recursive.

(b) Show that if  $K(\omega_{1:n}) \leq K(n) + O(1)$ , then  $\omega$  is recursive in  $0'$ , the Turing degree of the halting problem.

(c) Show that there are nonrecursive  $\omega$ 's satisfying  $K(\omega_{1:n}) \leq K(n) + O(1)$ , for all  $n$ .

(d) Show that for each constant  $c$ , there are at most  $2^c$  many  $\omega$  such that for all  $n$ ,  $K(\omega_{1:n}) \leq K(n) + c$ .

*Comments.* Item (b) is attributed to G.J. Chaitin, Item (c) to R.M. Solovay. Source is R.M. Solovay, *Lecture Notes*, UCLA, 1975, unpublished,

and personal communication. For Item (d) and other discussion on this topic, see [D. Zambella, *On sequences with simple initial segments*, ITLI Tech. Rept. ML-90-05, Fac. Wiskunde en Informatica, University of Amsterdam, 1990].

**3.6.5.** [27] Far from being a nuisance, the complexity oscillations actually enable us to discern a fine structure in the theory of random sequences. A sequence  $\omega$  is  $\Delta_2^0$  definable iff the set  $\{n : \omega_n = 1\}$  is  $\Delta_2^0$  definable, Exercise 1.7.21 on page 46. We consider infinite binary sequences that are  $\Delta_2^0$  definable (like the halting probability  $\Omega$ , Section 3.6.2).

- (a) Show that if  $\omega$  is  $\Delta_2^0$  definable, then  $\lim_{n \rightarrow \infty} n - K(\omega_{1:n}|n) = \infty$ . (This is, of course, only interesting for random  $\omega$ 's.)
- (b) Show that if  $\omega$  is  $\Delta_2^0$  definable, then  $\lim_{n \rightarrow \infty} n + K(n) - K(\omega_{1:n}) = \infty$ .
- (c) Show that if there is a constant  $c$  such that for infinitely many  $n$  we have  $n + K(n) - K(\omega_{1:n}) \leq c$ , then  $\omega$  is not  $\Delta_2^0$  definable. That is, if such an  $\omega$  is random, then it is not a “simply definable” random sequence.

*Comments.* Hint for Item (b): use Item (a). Items (a) and (b) delimit the upswing of the complexity oscillations for  $\Delta_2^0$  definable  $\omega$ . Hint for Item (c): use Exercise 3.6.3. The  $\Delta_2^0$  definable random sequences are rather atypical random sequences. (An example is the halting probability  $\Omega$ .) The  $K$ -complexity allows us to distinguish between easily definable random sequences and those that are not so easily definable. Source: M. van Lambalgen [*Random Sequences*, Ph.D. Thesis, University of Amsterdam, 1987; *J. Symb. Logic*, 54(1989), 1389–1400].

**3.6.6.** [31] Let  $\omega$  be an infinite binary sequence. Show that if there exists a constant  $c$  such that  $K(\omega_{1:n}) \geq n - c$  for all  $n$ , then for all  $k$  we have  $K(\omega_{1:n}) - n \geq k$  from some  $n$  onwards.

*Comments.* In Theorem 3.6.1 we showed that an infinite binary sequence is random with respect to the uniform measure if  $K(\omega_{1:n}) \geq n - c$  for some  $c$  and almost all  $n$ . There is not only a sharp dividing line but a wide complexity gap that separates the random infinite sequences from the nonrandom infinite sequences. With respect to the uniform measure, we have in Exercise 2.5.9 defined the notion of a *Solovay test* for randomness, and shown that  $\omega$  is Martin-Löf random iff it is Solovay random. Call  $\omega$  *weakly Chaitin* random if there is a constant  $c$  such that for all  $n$ , we have  $K(\omega_{1:n}) \geq n - c$ . Call  $\omega$  *strongly Chaitin* random if for all constants  $c$ , for almost all  $n$ , we have  $K(\omega_{1:n}) \geq n + c$ . The results referred to and this exercise show that the sets of infinite sequences defined by all of these definitions coincide precisely. That different points of departure attempting to formalize an intuitive notion (like randomness) turn out to define the same objects is commonly viewed as evidence that the

formalization correctly represents our intuition. Source: G.J. Chaitin, *Algorithmic Information Theory*, Cambridge Univ. Press, 1987; attributed to C.P. Schnorr.

**3.6.7.** [38] Let  $\omega$  be an infinite binary sequence that is random with respect to the uniform measure. Let  $g$  be a recursive function such that the series  $\sum_n 2^{-g(n)} = \infty$ , for example  $g(n) = \log n$ . Let  $h$  be a recursive function that is monotone and unbounded, like  $h(n) = \log \log n$ . Show that for infinitely many  $n$  we have  $K(n) \geq g(n)$  and  $K(\omega_{1:n}) \leq n + h(n)$ .

*Comments.* Note that this does not imply that the monotonic lower bound on  $K(\Omega_{1:n}) - n$  is unbounded. This is nonetheless true. Source: R. Solovay, *Lecture Notes*, UCLA, 1975, unpublished.

**3.6.8.** [41] Consider an infinite binary sequence  $\omega$  as a real number  $0.\omega$ . We call  $\omega$  an  *$\Omega$ -like number* if (i) there is an enumerable function  $f$  and an argument  $n$  such that  $f(n) = 0.\omega$ ; and (ii)  $K(\omega_{1:n}) = K(\Omega_{1:n}) + O(1)$ . A real number  $\omega$  is *arithmetically random* (with the uniform measure understood) if any arithmetic property of  $\omega$  (viewed as an infinite sequence of zeros and ones) holds for some set of reals of uniform measure 1. In other words,  $\omega$  is arithmetically random iff  $\omega$  does not belong to any arithmetic set of reals of uniform measure zero. Since the arithmetic sets are Borel sets, and the union of countably many Borel sets of uniform measure zero is again a Borel set of uniform measure zero, the set of arithmetically random reals has uniform measure one. Let  $m(n)$  be the greatest monotonic lower bound on  $K(n)$ , that is,  $m(n) = \min\{K(x) : x \geq n\}$ . By the same argument as in the proof of Theorem 2.3.2, page 121, for any recursive function  $\phi(n)$  that goes to infinity monotonically with  $n$ , we have  $m(n) < \phi(n)$  from some  $n$  onwards.

- (a) Show that  $\Omega$  is  $\Omega$ -like.
- (b) Show that  $\Omega$ -like reals are Martin-Löf random.
- (c) Show that arithmetically random reals are Martin-Löf random.
- (d) Show that no  $\Omega$ -like real is arithmetically random, so the set of  $\Omega$ -like reals has uniform measure zero.
- (e) Show that if  $\omega$  is arithmetically random, then there is a constant  $c$  (depending on  $\omega$ ) such that for infinitely many  $n$  we have  $K(\omega_{1:n}) \geq n + K(n) - c$ .
- (f) Show that if  $\omega$  is  $\Omega$ -like then  $K(\omega_{1:n}) \leq n + K(n) - m(n) + O(\log m(n))$ .
- (g) Show that if  $\omega$  is arithmetically random, then  $K(\omega_{1:n}) \geq n + \alpha(n) + O(\log \alpha(n))$ .

*Comments.* Hint for Item (f): compute  $\Omega_{m(n)}$  effectively from  $n$  by doing  $n$  steps in the computation of  $\Omega$ . R.M. Solovay credits this idea to C.P. Schnorr. Source is R. Solovay, *Lecture Notes*, UCLA, 1975, unpublished.

**3.6.9.** [39] (a) Show that an infinite binary sequence  $\omega$  is random with respect to the uniform measure iff  $C(\omega_{1:n} \mid n) \geq n - K(n)$  for all  $n$ .

(b) Show that  $\gamma_0(\omega_{1:n} \mid L) = n - C(\omega_{1:n} \mid n) - K(n)$  is finite iff  $\omega$  is random with respect to the universal measure.

*Comments.* Therefore, the formula for  $\gamma_0$  expresses concisely and precisely how the complexity oscillations of  $C(\omega_{1:n} \mid n)$  of random infinite sequences behave. This is *Levin's test*—the characterization for random sequences from which Theorem 2.5.4 on page 145 follows. Source: P. Gács, *Ph.D. Thesis*, Frankfurt am Main, 1978, Theorem 5.4, Corollary 2; *Z. Math. Logik Grundl. Math.*, 26(1980), 385–394.

**3.6.10.** [29] Let  $\omega$  be an infinite binary sequence that is Martin-Löf random with respect to the uniform distribution (like the halting probability  $\Omega$  of Section 3.6.2).

(a) Show that  $\omega$  is von Mises-Wald-Church random (Section 1.9).

(b) Show that  $\omega$  is *effectively unpredictable*. That is, let  $f$  be a recursive function that given an arbitrary finite binary string predicts “the next bit is one,” “the next bit is zero,” or “no prediction.” Then if  $f$  predicts infinitely many bits of  $\omega$ , it does no better than chance because in the limit the relative frequency of both correct predictions and incorrect predictions is  $\frac{1}{2}$ .

(c) Use Item (b) to show that the zeros and ones have limiting relative frequency  $\frac{1}{2}$ .

(d) Show that  $\omega$  is *normal* in base two in the sense of Borel. That is, each of the  $2^k$  possible blocks of  $k$  bits in  $\omega$  has limiting relative frequency  $2^{-k}$ .

(e) Show that the Law of the Iterated Logarithm (Section 1.10) applies to the relative frequency of correct and incorrect predictions of bits of  $\omega$ .

*Comments.* Hint for Item (c): predict the next bit of  $\omega$  by a total recursive function that always has value one. Source: G.J. Chaitin, *Algorithmic Information Theory*, Cambridge Univ. Press, 1987.

## 3.7 Algorithmic Properties of $K$

According to Section 3.4 the function  $K$  is not recursive but it is co-enumerable (it can be approximated from above). This is also the case for the two variable function  $K(x|y)$ . Also, Theorem 2.7.1 on page 167 and Corollary 2.7.2 on page 169 hold in precisely the same way by the same proofs with  $K$  replacing  $C$ . Barzdins's Lemma, Theorem 2.7.2, page 171, holds by the same proof with the following obvious modification:

**Definition 3.7.1** Define  $K^+(n) = \max\{K(m) : m \leq n\}$ .

Any characteristic sequence  $\chi$  of a recursively enumerable set  $A$  satisfies  $K(\chi_{1:n}|n) \leq K^+(n) + c$  for all  $n$ , where  $c$  is a constant depending on  $A$  but not on  $n$ . There is a recursively enumerable set such that its characteristic sequence satisfies  $K(\chi_{1:n}) \geq \log n - c$  for all  $n$ , where  $c$  is a constant that does not depend on  $n$ .

### 3.7.1 Randomness in Diophantine Equations

There is an algorithm to decide the solvability of the first  $n$  Diophantine equations, given about  $\log n$  bits extra information (Example 2.7.2). Namely, given the number  $m \leq n$  of soluble equations in the first  $n$  equations, we can recursively enumerate solutions to the first  $n$  equations in the obvious way until we have found  $m$  solvable equations. The remaining equations are unsolvable. This shows that the solubility of the enumerated Diophantine equations is interdependent in some way.

Suppose we replace the question of mere solubility by the question of whether there are finitely many or infinitely many different solutions. That is, no matter how many solutions we find for a given equation, by itself this can give no information on the question to be decided. It turns out that the set of indices of the Diophantine equations with infinitely many different solutions is not recursively enumerable. In particular, in the characteristic sequence each initial segment of length  $n$  has Kolmogorov complexity of about  $n$ .

**Claim 3.7.1** There is an (exponential) Diophantine equation

$$A(n, x_1, x_2, \dots, x_m) = 0$$

that has only finitely many solutions  $x_1, x_2, \dots, x_m$  if the  $n$ th bit of  $\Omega$  is zero and that has infinitely many solutions  $x_1, x_2, \dots, x_m$  if the  $n$ th bit of  $\Omega$  is one.

The role of *exponential* Diophantine equations should be clarified. Yu.V. Matijasevich [*Soviet Math. Dokl.*, 11(1970), 354–357] proved that every recursively enumerable set has a polynomial Diophantine representation. J.P. Jones and Yu.V. Matijasevich [*J. Symbol. Logic*, 49(1984), 818–829] proved that every recursively enumerable set has a singlefold exponential Diophantine representation. It is not known whether singlefold representation (which is important in our application) is always possible without exponentiation. See also G.J. Chaitin, *Algorithmic Information Theory*, Cambridge University Press, 1987.

**Proof.** By dovetailing the running of all programs of the reference prefix machine  $U$  in the obvious way, we find a recursive sequence of rational numbers  $\omega_1 \leq \omega_2 \leq \dots$  such that  $\Omega = \lim_{n \rightarrow \infty} \omega_n$ . The set

$$R = \{(n, k) : \text{the } n\text{th bit of } \omega_k \text{ is a one}\}$$

is a recursively enumerable (even recursive) set. The main step is to use a theorem due to J.P. Jones and Yu.V. Matijasevich [*J. Symbol. Logic* 49(1984), 818–829] to the effect that “every recursively enumerable set  $R$  has a singlefold exponential Diophantine representation  $A(\cdot, \cdot)$ .” That is,  $A(p, y) = 0$  is an exponential Diophantine equation, and the singlefoldedness consists in the property that  $p \in R$  iff there is a  $y$  such that  $A(p, y) = 0$  is satisfied, and moreover, there is only a single such  $y$ . (Here both  $p$  and  $y$  can be multiples of integers, in our case  $p$  represents  $\langle n, x_1 \rangle$ , and  $y$  represents  $\langle x_2, \dots, x_m \rangle$ . For technical reasons we consider as proper solutions only solutions  $x$  involving no negative integers.) It follows that there is an exponential Diophantine equation  $A(n, k, x_2, \dots, x_m) = 0$  that has exactly one solution  $x_2, \dots, x_m$  if the  $n$ th bit of the binary expansion of  $\omega_k$  is a one, and it has no solution  $x_2, \dots, x_m$  otherwise. Consequently, the number of different  $m$ -tuples  $x_1, x_2, \dots, x_m$  that are solutions to  $A(n, x_1, x_2, \dots, x_m) = 0$  is infinite if the  $n$ th bit of the binary expansion of  $\Omega$  is a one, and this number is finite otherwise.  $\square$

This can be interpreted as follows: Consider the sequence of Diophantine equations  $D_1, D_2, \dots$  such that  $D_n(x_1, x_2, \dots, x_m) = A(n, x_1, \dots, x_m)$ . Let us say that a formal theory “settles  $k$  cases” if it enables one to prove that the number of solutions of  $D_n$  is finite or is infinite for  $k$  specific values (not necessarily consecutive) of  $n$ . It is not difficult to show that no formal theory in which one can prove only true theorems and which is completely describable in  $n$  bits can settle more than  $n + K(n) + O(1)$  cases. (Hint: use arguments about how many (scattered) bits of  $\Omega$  can be determined (together with their positions) in a formal theory of given complexity.)

“This is a region in which mathematical truth has no discernible structure or pattern and appears to be completely random. These questions are completely beyond the power of human reasoning. Mathematics cannot deal with them. Quantum physics has shown that there is randomness in nature. I believe that we have demonstrated [...] that randomness is already present in pure Mathematics. This does not mean that the universe and Mathematics are completely lawless, it means that laws of a different kind apply: statistical laws. [...] Perhaps number theory should be pursued more openly in the spirit of an experimental science! To prove more one must assume more.” [Chaitin]

## Exercises

- 3.7.1.** [42] (a) Show that  $K^+$  (Exercise 3.4.1) is nonrecursive in the sense that there is no total recursive function  $f$  such that  $K^+(x) = f(x) + O(1)$  for all  $x$ .

- (b) Show that there is a recursive upper bound  $f(x)$  of the function  $K(x)$  with the property that  $f(x) = K(x)$  holds for infinitely many  $x$ .
- (c) Show that we cannot effectively find infinitely many  $x$ 's where some recursive upper bound on  $K(x)$  (as in Item (b)) is sharp. (The same statement for  $C(x)$  follows from Theorem 2.3.2, page 121.)
- (d) Let  $f(x)$  be a recursive upper bound on  $K(x)$  as in Item (b). Show that, in addition to Item (c), there is no recursive function  $g(x)$  mapping each  $x$  to a finite set  $X \subseteq \{y : y > x\}$  such that each  $X$  contains some  $y$  for which  $f(y) \leq C(y)$ . (Notice that the function  $\log x$ , or one almost equal to it, has this property for  $C(x)$ .)

*Comments.* The sharp monotonic estimate on  $\log x + K(\lceil \log n \rceil)$  on  $K(x)$  is less satisfying than the sharp monotonic estimate  $\log x$  on  $C(x)$ , because it is not a recursive function. We can derive several recursive upper bounds on  $K(x)$  such as  $\log x + 2 \log \log x$  or  $l^*(x)$ , but none of these is sharp for large  $x$ . The present Item (b) shows that, nonetheless, we can find a recursive upper bound on  $K(x)$  that is sharp infinitely often. However, we cannot find infinitely many  $x$ 's for which this is the case, since Item (c) shows that every infinite set of  $x$ 's where the recursive upper bound coincides with  $K(x)$  is nonrecursive. This also holds for  $C(x)$ . However, Item (d) shows a difference between  $K(x)$  and  $C(x)$  in effectiveness of a recursive upper bound that is sharp infinitely often. Source: R. Solovay, *Lecture Notes*, 1975, unpublished; P. Gács, *Ph.D. Thesis*, Mathematics Department, Frankfurt am Main, 1978.

**3.7.2.** [35] Show that there is a recursive upper bound  $f$  on  $K$  and a constant  $c$  with the property that there are infinitely many  $x$  such that for all  $k > 0$ , the quantity of numbers  $y \leq x$  with  $K(y) < f(y) - k$  is less than  $cx2^{-k}$ .

*Comments.* Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

## 3.8 \*Complexity of Complexity

The complexity function  $K$  (or  $C$  for that matter) is nonrecursive (Theorem 2.3.2, page 121). If the number of bits one needs to know to compute  $f(x)$  from  $x$  is not constant for all  $x$ , then  $f$  is nonrecursive. For instance, let  $f$  be defined by  $f(x) = \Omega_{1:x}$ , with  $\Omega$  the halting probability in Section 3.6.2. Then  $K(f(x)|x) \geq l(f(x)) + O(1)$ . That is, knowledge of  $x$  does not help to describe  $f(x)$  at all.

If  $f$  is recursive, then  $K(f(x)|x) = O(1)$ . But if  $f$  is 0, 1-valued, then  $K(f(x)|x) = O(1)$  too (since the programs for the constant functions of value 0 and 1 provide a bound). Thus, the complexity  $K(f(x)|x)$  is a property of  $f$  that tells something about the degree of nonrecursiveness of  $f$ , but the converse is not necessarily the case.

**Definition 3.8.1** The *complexity*  $K(f)$  of function  $f$  is defined as the function  $K(f(x)|x)$ .

We analyze  $K(K(x)|x)$ , the *complexity of the complexity function*  $K$ . There is a simple *upper bound* for all  $x$ . If  $l(x) = n$ , then  $K(x) \leq n + 2 \log n$  and therefore

$$K(K(x)|x) \leq \log n + 2 \log \log n + O(1).$$

It turns out that  $K(K(x)|x)$  can be very close to this upper bound.

**Theorem 3.8.1** For  $n$  there are strings  $x$  of length  $n$  such that  $K(K(x)|x) \geq \log n - \log \log n + O(1)$ . The same lower bound holds for  $C(C(x)|x)$ .

**Proof.** The proof is a little tricky. Let  $U$  be the reference machine of Theorem 3.1.1. Fix a large enough  $n$ . In the sequel all  $x$ 's considered have length  $n$ . Define the maximum value of  $K(K(x)|x)$  by

$$s = \max_{l(x)=n} \min\{l(p) : U(p, x) = K(x)\}.$$

To prove the theorem we only need to show

$$s \geq \log n - \log \log n + O(1). \quad (3.6)$$

Since  $K(x) \leq n + 2 \log n + O(1)$  for all  $x$  of length  $n$ ,

$$K(K(x)|x) \leq s \leq \log n + 2 \log \log n + O(1).$$

Let  $U$  be the reference prefix machine. A binary string  $p$  is called a *suitable* program for  $x$  if for some  $q$  we have

- $l(p) \leq s$ ;
- $U$  computes  $l(q)$  from  $p$ , given  $x$ ; and
- $U$  computes  $x$  from  $q$ .

In particular, there is a suitable program  $p$  for  $x$ , of length  $l(p) = K(K(x)|x) \leq s$  (for instance, with corresponding  $q$  of length  $l(q) = K(x)$ ). We denote by  $M_i$  the set of strings  $x$  for which there are at least  $i$  different suitable  $p$ . Now consider the sequence

$$\emptyset = M_{j+1} \subseteq M_j \subseteq \cdots \subseteq M_0 = \{0, 1\}^n, \quad M_j \neq \emptyset.$$

There are  $2^{s+1} - 1$  strings of length at most  $s$ . Therefore,

$$j \leq 2^{s+1}.$$

To prove the theorem, it suffices to show, for all  $i \leq j$ ,

$$l(d(M_i)) \leq l(d(M_{i+1})) + 5 \log n. \quad (3.7)$$

Namely, this implies that  $j \geq n/(5 \log n)$ , which together with  $j \leq 2^{s+1}$ , shows Equation 3.6.

We prove Equation 3.7. There exists an  $x \in M_i - M_{i+1}$ , which is found by the following procedure with input  $i, s, n, d(M_{i+1}), l(d(M_i))$ , for all  $i$  with  $0 \leq i \leq j$ :

**Step 1** Recursively enumerate all of  $M_{i+1}$ . {We know when we are done by the time we have found  $d(M_{i+1})$  elements.}

**Step 2** Recursively enumerate enough of  $M_i - M_{i+1}$  to make the sequel meaningful. For each  $z$  in  $M_i - M_{i+1}$  we obtain, find by enumeration *all*  $i$  suitable programs for  $z$ . A suitable program  $p$  with  $U(p, z)$  minimal satisfies  $U(p, z) = K(z)$ . We may assume that

$$\log d(M_i - M_{i+1}) \geq \log d(M_i) - 1, \quad (3.8)$$

since otherwise Equation 3.7 holds trivially. From Equation 3.8 it follows, by Theorem 3.3.1, that there exists a  $z_{\max}$  among the  $z$ 's for which

$$K(z_{\max}) \geq l(d(M_i)) - 1. \quad (3.9)$$

So we keep on enumerating  $z$ 's until we find the first such  $z_{\max}$ .

**Step 3** Set  $x := z_{\max}$ .

This algorithm provides a description of  $x$  containing the following items, each item in self-delimiting code:

- A description of this discussion on  $O(1)$  bits;
- a description of  $d(M_{i+1})$  in at most  $l^*(n) + l(d(M_{i+1}))$  bits (since  $K(l(d(M_{i+1}))) \leq l^*(n)$ );
- a description of  $l(d(M_i))$  in  $l^*(n)$  bits;
- descriptions of  $i, n$ , and  $s$ , in  $l^*(n), l^*(n)$ , and  $l^*(\log n + 2 \log \log n)$  bits, respectively.

The size of the description of  $x$  gives an upper bound on  $K(x)$ ,

$$K(x) \leq 4l^*(n) + O(l^*(\log n)) + l(d(M_{i+1})) + O(1). \quad (3.10)$$

Equations 3.9 and 3.10 imply Equation 3.7, and hence the theorem. Precisely the same proof shows  $C(C(x)|x) \geq \log n - \log \log n + O(1)$ .  $\square$

**Example 3.8.1** Since  $C(C(x)) \leq \log n + O(1)$  for all  $x$  of length  $n$ , Theorem 3.8.1 indicates that for some  $x$ , knowledge of  $x$  only marginally helps to compute  $C(x)$ ; most information in  $C(x)$  is extra information. It turns out that these  $x$ 's have lower complexity than random  $x$ 's. Consider an  $x$  of length  $n$  with  $C(x) \geq n - k$ . Then  $C(C(x)|x) \leq C(k) + c'$  for some fixed constant  $c'$  independent of  $x$ . If Theorem 3.8.1 holds for  $x$ , then it follows that  $C(k) \geq \log n - \log \log n - c' + O(1)$ . Also,  $C(k) \leq \log k + c$  for another constant  $c$  independent of  $x$ . Therefore,  $k = \Omega(n/\log n)$ . Therefore, if Theorem 3.8.1 holds for  $x$ , then

$$C(x) \leq n - \Omega\left(\frac{n}{\log n}\right).$$

In fact, the same argument shows the following: Fix a constant  $c'$ . If Theorem 3.8.1 holds for  $x$ , then for each  $k$ ,  $0 \leq k \leq n$ , such that  $C(k|x) \leq c'$  (for example  $k = n/2$  or  $k = \sqrt{n}$ ), we have

$$C(x) \notin [k - \delta, k + \delta],$$

with  $\delta = O(n/\log n)$ . ◇

## Exercises

**3.8.1.** [20] Show that for Theorem 3.8.1 to hold, we must have  $K(x) \leq n - \Omega(n/\log n)$ .

## 3.9 \*Symmetry of Algorithmic Information

Recall the *symmetry of algorithmic information* question raised in Sections 1.11 and 2.8. Up to what precision does the equality

$$K(x, y) = K(y|x) + K(x)$$

hold? It turns out that the precision is low, just as for  $C$  complexity, albeit for more fundamental reasons. It is at once obvious that

$$K(x, K(x)) = K(x) + O(1). \quad (3.11)$$

Namely, we can reconstruct both  $x$  and  $K(x)$  from the shortest program for  $x$ . (Similarly  $C(x, C(x)) = C(x) + O(1)$ .) Now consider the difference of  $K(x, y)$  from  $K(y|x) + K(x)$  for  $y = K(x)$ . Combining Theorem 3.8.1 and Equation 3.11, for each  $n$ , there is a string  $x$  of length  $n$ ,

$$K(x, K(x)) \leq K(x) + K(K(x)|x) - \log n + \log \log n + O(1). \quad (3.12)$$

For Equation 3.12 it does not really matter whether we use  $K(x)$  or  $C(x)$ . However,  $C(x)$  is nonadditive for reasons much simpler than  $K(x)$ : additivity is violated already on random strings as shown in Section 2.8. But for  $K(x)$ , the possibly high-complexity  $K(K(x)|x)$  is the only reason. Equation 3.11 is an obvious property that can be used for most variants of complexity to prove the corresponding version of Equation 3.12, which amounts to *asymmetry of information*.

The above shows that analogues of the information-theoretic identities can hold only up to an additive term logarithmic in the complexity. In the case of  $C(x)$  this is primarily caused by the randomness (incompressibility) of the strings concerned, Section 2.4. The overwhelming majority of strings is random enough for this effect to occur. One reason for focusing on the  $K(x)$  measure was to eliminate this randomness-based effect. But it turns out we still cannot obtain sharp analogues of the information-theoretic identities, due to the extra information besides  $x$  contained in  $K(x)$ . This is the sole reason that the information-theoretic identities do not hold for  $K(x)$  precisely. However, there are only relatively few  $x$ 's with large  $K(K(x)|x)$ . This follows immediately from Example 3.8.1. Another argument is as follows:

Let  $\chi$  be the characteristic function (equivalently, sequence) of a recursively enumerable set such that  $K(\chi_{1:m}|n) \geq n$  where  $m = 2^n$ . Such sets exist by Barzdins's Lemma, Theorem 2.7.2 on page 171. Let  $I(\chi : x) = K(x) - K(x|\chi)$ , that is, the information in  $\chi$  about  $x$ . For instance, we can choose  $\chi$  as the characteristic function of the halting set  $K_0 = \{\langle x, y \rangle : \phi_x(y) < \infty\}$ . We can take also any other complete set. Then

$$K(K(x)|x) \leq I(\chi : x) + 3 \log I(\chi : x) + O(1). \quad (3.13)$$

This means that the extra information contained in  $K(x)$ , apart from  $x$ , is less than the information that any complete set contains about  $x$ . The a priori probability of  $x$  with high  $K(K(x)|x)$  is low. Namely, it can be shown (but we omit the proof) that for any recursive or enumerable probability distribution (measure) on the  $x$ 's,  $k$  bits of information concerning  $\chi$  occur in  $x$  with probability  $2^{-k}$ . Thus, we can derive, from the displayed inequality,

$$P\{x : K(K(x)|x) \geq i\} \leq i^3 2^{-i}, \quad (3.14)$$

for all recursive and enumerable probability distributions  $P$ . Equation 3.13 tells us that the extra information in  $K(x)$ , apart from  $x$ , is less than the information that any complete set contains about  $x$ . Equation 3.14 tells us that the a priori probability of all  $x$  with high  $K(K(x)|x)$  is low. Source: attributed to L.A. Levin [P. Gács, *Soviet Math. Dokl.*, 15(1974), 1477–1480, and Ph.D. Thesis, J.W. Goethe Univ., Frankfurt, 1978].

### 3.9.1 Algorithmic Information and Entropy

Let us refresh some relations between Shannon entropy and algorithmic information, or complexity. Let  $P$  be a recursive probability distribution. In Example 2.8.1 on page 181 we have shown that for recursive probability distributions  $P(\cdot)$  the expected value  $C(\cdot)$  is asymptotic to  $H(P)$  (in case both values grow unboundedly we have shown that the quotient of the compared quantities goes to 1). This implies the similar relation between expected  $K(\cdot)$  value and  $H(P)$ , see Section 8.1.

A more direct proof uses the universal distribution in Chapter 4. Since the set of  $K(x)$ 's is the length set of a prefix code, the first inequality of the Noiseless Coding Theorem, Theorem 1.11.2, page 75, shows that  $H(P) \leq \sum_x P(x)K(x)$ . Moreover, an effective version of the Shannon-Fano code in Example 1.11.2

guarantees that  $K(x) \leq -\log P(x) + O(1)$  (a formal proof is given later in Lemma 4.3.3 on page 253). Together this shows that the *entropy*

$$H(P) = - \sum_x P(x) \log P(x)$$

of the distribution  $P$  is asymptotically equal to the *expected complexity*

$$\sum_x P(x)K(x)$$

with respect to probability  $P(\cdot)$ . However, we can make more precise calculations. The equality between expected prefix complexity and entropy is treated in detail in Section 8.1. It holds extremely precisely—up to an additive constant.

Instead of requiring  $P(\cdot)$  to be recursive, it suffices to require that  $P$  be an enumerable function. Together with  $\sum P(x) = 1$ , the latter requirement implies that  $P(\cdot)$  is recursive (Example 4.3.2 on page 246).

The fact that the  $P$ -expectations of  $-\log P(x)$  and  $K(x)$  are close does not necessarily mean that those quantities are close together for all arguments. However, in Example 4.3.10 on page 261 we also show that for recursive  $P(\cdot)$  the values  $-\log P(x)$  and  $K(x)$  are close to each other with high probability.

This establishes a tight quantitative relation between Shannon's statistical conception of entropy of a probability distribution and our intended interpretation of  $K(x)$  as the information content of an individual object.

Let  $X$  and  $Y$  be two discrete random variables with a joint distribution. In Section 1.11 we defined the following notions: the *conditional entropy*  $H(Y|X)$  of  $Y$  with respect to  $X$ , the *joint entropy*  $H(X, Y)$  of  $X$  and  $Y$ , and the *information*  $I(x : y)$  in  $x$  about  $y$ . The relations between these quantities are governed by Equations 1.11, 1.12, 1.13. The crucial one is Equation 1.11 on page 69: the additivity rule  $H(X, Y) = H(X) + H(Y|X)$ . In general we would like to derive the same relations for the complexity analogues. In Section 2.8 it turned out that the complexity analogue of Equation 1.11 holds in terms of  $C$ , within a logarithmic additive term (Theorem 2.8.2, page 182). The proof that Theorem 2.8.2 is sharp used strings whose lengths were random. For  $K$  we also find that the exact analogue of Equation 1.11 on page 69 does not hold (Example 3.9), not because it is violated on random  $x$ 's with random lengths as is  $C$ , but for more subtle reasons.

### 3.9.2 Exact Symmetry of Information

By Equation 3.12, for each length  $n$  there are  $x$  of length  $n$  and  $y$  such that

$$|K(x, y) - K(x) - K(y|x)| = \Omega(\log K(x)),$$

showing that additivity can only be satisfied to within a logarithmic term. Since the complexity of the complexity function as expressed by

Theorem 3.8.1 holds for all proper variants of complexity, additivity corresponding to Equation 1.11 cannot hold exactly for any proper variant of complexity.

While the complexity of the complexity function prevents an exact analogue to Equation 1.11, page 69, it turns out that nonetheless we can find an exact additivity property for  $K$ , by replacing the conditional  $x$  by  $\langle x, K(x) \rangle$  (equivalently, by the shortest program for  $x$ ).

**Theorem 3.9.1** *Let  $x$  and  $y$  be finite binary strings. Then up to a fixed additive constant,  $K(x, y) = K(x) + K(y|x, K(x))$ .*

**Proof.** ( $\leq$ ) Let  $p$  be a shortest program from which the reference prefix machine computes  $x$ , and let  $q$  be a shortest program for which it computes  $y$  given  $x$  and  $K(x)$ . But then we can find another prefix machine that with input  $pq$  uses  $p$  to compute  $x$  and  $K(x)$  ( $= l(p)$ ) and subsequently uses  $x$  and  $K(x)$  to compute  $y$  from  $q$ .

( $\geq$ ) Consider  $x$  and  $K(x)$  as fixed constants. We need the following results, which are proven only later in Chapter 4 (the conditional version Theorem 4.3.4 of Theorem 4.3.3 on page 255):

1. Recall Definition 1.7.3, page 35, of enumerable functions  $f : \mathcal{N} \rightarrow \mathcal{R}$ . Let  $X$  consist of enumerable functions  $f_x(y)$  with  $\sum_y f_x(y) < \infty$ . There exists a  $g \in X$ , such that for all  $f_x \in X$ ,  $f_x(y) = O(g(y))$ .
2. We can set this  $g(y) = 2^{-K(y|x, K(x))}$ .

Clearly, with  $x$  fixed the function  $h_x(y) = 2^{K(x)-K(x,y)}$  is enumerable. Moreover, with  $x$  fixed the set  $\{K(x, y) : y \in \mathcal{N}\}$  is the length set of a prefix-free set of programs from which the reference prefix machine  $U$  computes  $\langle x, y \rangle$ . By the Kraft Inequality, page 74, we have

$$\sum_y 2^{-K(x,y)} \leq 1.$$

Therefore,

$$\sum_y h_x(y) \leq 2^{K(x)} < \infty.$$

Hence,  $h_x \in X$ , which implies by Items 1 and 2 that  $h_x(y) = O(g(y))$ . Substituting the expressions for  $g$  and  $h_x$ , taking the logarithm of both sides, and rearranging, we obtain

$$K(x, y) \geq K(x) + K(y|x, K(x)) + O(1).$$

□

This implies immediately the *subadditive property*

$$K(x, y) \leq K(x) + K(y|x) + O(1) \leq K(x) + K(y) + O(1),$$

which does not hold for  $C$  by Example 2.2.3 on page 111. This is a straightforward consequence of the fact that a prefix machine does not require extra information to see where one program ends and the other one begins—information that needs to be provided to the plain decoding algorithms in Chapter 2.

We cannot replace  $\langle x, K(x) \rangle$  by  $K(x)$  in Theorem 3.9.1. But we can replace  $x$  by  $\langle x, K(x) \rangle$ , as follows. Note that

$$K(x, y) = K(x, K(x), y) + O(1). \quad (3.15)$$

Substitution in Theorem 3.9.1 shows the following:

**Corollary 3.9.1**  $K$ -complexity is additive in the form

$$K(\langle x, K(x) \rangle, y) = K(\langle x, K(x) \rangle) + K(y|\langle x, K(x) \rangle) + O(1).$$

**Definition 3.9.1** The  $K$ -complexity of information in  $x$  about  $y$  is

$$I(x : y) = K(y) - K(y|x). \quad (3.16)$$

Information is symmetric if the information in  $x$  about  $y$  equals (up to additive constants) the information in  $y$  about  $x$ . Rewriting  $K(x, y)$  in two different ways, it follows from Theorem 3.9.1 that

$$K(y) - K(y|x, K(x)) = K(x) - K(x|y, K(y)) + O(1).$$

**Theorem 3.9.2** *Symmetry of information for  $K$ -complexity holds in the following form:*  
 $I(\langle x, K(x) \rangle : y) = I(\langle y, K(y) \rangle : x) + O(1).$

We cannot replace  $\langle x, K(x) \rangle$  by  $K(x)$  and  $\langle y, K(y) \rangle$  by  $K(y)$  in Theorem 3.9.2. The complexity version of individual information is asymmetric, in contrast to the expectation in the statistical entropy version of Equation 1.13, which is symmetric.

**Lemma 3.9.1** *The error in symmetry of information using  $K$ -complexity is given by*

$$\begin{aligned} & \log l(x) - \log \log l(x) + O(1) \\ & \leq |I(x : y) - I(y : x)| \\ & \leq \log K(x) + \log K(y) + 2 \log \log K(x) + 2 \log \log K(y) + O(1). \end{aligned}$$

**Proof.** ( $\leq$ ) Define the conditional

$$I(x : y|z) = K(y|z) - K(y|x, z). \quad (3.17)$$

Since  $K(y|z) \leq K(y|x, z) + K(x|z) + O(1)$ ,

$$I(x : y|z) \leq K(x|z) + O(1). \quad (3.18)$$

We find, rewriting using Equations 3.15, 3.17, then Equation 3.18, and the usual estimates of the relevant quantities,

$$\begin{aligned} K(x) + K(y|x) - K(x, y) &= I(K(x) : y|x) + O(1) \\ &\leq K(K(x)|x) + O(1) \\ &\leq \log K(x) + 2 \log \log K(x) + O(1). \end{aligned} \quad (3.19)$$

Using Equation 3.19 we find

$$\begin{aligned} I(x : y) - I(y : x) &= I(K(y) : x|y) - I(K(x) : y|x) + O(1) \\ &\leq \log K(x) + \log K(y) + 2 \log \log K(x) \\ &\quad + 2 \log \log K(y) + O(1). \end{aligned}$$

( $\geq$ ) Compute the error in additivity for the special case  $y = \langle x, K(x) \rangle$ . Rewrite using Definition 3.16, and apply Equation 3.11 (denoting  $x^* = \langle x, K(x) \rangle$ ):

$$\begin{aligned} I(x : x^*) - I(x^* : x) &= K(x^*) - K(x^*|x) - K(x) + K(x|x^*) \\ &= K(x, x^*) - K(x) - K(x^*|x) + O(1) \\ &= -K(K(x)|x) + O(1). \end{aligned}$$

By Theorem 3.8.1, for each  $n$  there is an  $x$  of length  $n$  such that

$$-K(K(x)|x) \leq -\log n + \log \log n + O(1).$$

□

We can use this to show that in fact  $I(x : y)$  is *not even asymptotically symmetric*. Let  $l(x) = n$ . From  $K(K(x)) \leq \log n + 2 \log \log n + O(1)$  and Theorem 3.8.1 it follows that, up to an additive constant,

$$I(x : K(x)) = K(K(x)) - K(K(x)|x) \leq 3 \log \log n. \quad (3.20)$$

By Theorem 3.8.1, using  $K(K(x)|\langle x, K(x) \rangle) = O(1)$ , there exist  $x$  of each length  $n$ , such that

$$I(\langle x, K(x) \rangle : K(x)) = K(K(x)) + O(1) \geq \log n - \log \log n + O(1).$$

By writing out the definitions, we have

$$I(K(x) : x) = I(K(x) : \langle x, K(x) \rangle) + O(1).$$

Consequently,

$$\begin{aligned} \log n - 4 \log \log n + O(1) \\ \leq |I(\langle x, K(x) \rangle : K(x)) - I(K(x) : \langle x, K(x) \rangle)| \\ + |I(x : K(x)) - I(K(x) : x)|. \end{aligned}$$

This shows that the symmetry of information is violated strongly (in exponentially greater measure since  $I(x : K(x)) = O(\log \log n)$ ) on at least one of the pairs  $x, K(x)$  or  $x^*, K(x)$ . Source: attributed to L.A. Levin [P. Gács, *Soviet Math. Dokl.* 15(1974), 1477–1480].

**Example 3.9.1** If a problem does not have a satisfactory solution as it is posed originally, it is a common mathematical ploy to change the definitions to accommodate the problem. Given  $\langle x, K(x) \rangle$ , we can enumerate all shortest programs for  $x$ . The first one we find is denoted by  $x^*$ . From  $x^*$  we can compute  $\langle x, K(x) \rangle$ . Hence,  $x^*$  and  $\langle x, K(x) \rangle$  contain the same information although they are not identical strings. Below, we can replace  $x^*$  by  $\langle x, K(x) \rangle$ . Define  $Kc(x|y) = K(x|y^*)$ . Of course,  $Kc(x) = K(x)$  and  $Kc(x, y) = K(x, y)$ . We can formulate the additivity property as

$$Kc(x, y) = Kc(x) + Kc(y|x) + O(1).$$

Define the *mutual information* of two objects  $x$  and  $y$  by  $I(x; y) = Kc(y) - Kc(y|x)$  (which we may view as the  $Kc$  analogue of  $I(x : y)$ ). Using Theorem 3.9.1, we find that the mutual information is a *symmetric quantity* as desired:

$$I(x; y) = K(x) + K(y) - K(x, y) + O(1) = I(y; x) + O(1).$$

In general, the world looks prettier using  $Kc$ , and lots of basic identities can be derived (Exercise 3.9.3, page 236). One drawback about the conditional complexity  $Kc(x|y)$  is that it is *not co-enumerable*. Namely, suppose the contrary and a function  $Kc(x|y) = K(x, y) - K(y)$  is co-enumerable. Then for fixed  $y$ , the function  $2^{-Kc(x|y)}$  of variable  $x$  is enumerable and satisfies  $\sum_x 2^{-Kc(x|y)} \leq 1$ . Hence, by the Kraft Inequality, for each fixed  $y$  the set  $\{Kc(x|y) : y \geq 1\}$  is the length set of a prefix-code. But for fixed  $y$ , the set  $\{K(x|y) : y \geq 1\}$  is also the length set of a prefix-code, and by the conditional variant of the Invariance Theorem, the shortest one among the co-enumerable length sets. That is, for each  $y$ , there exists a constant  $c_{Kc(\cdot|y)}$  such that for all  $x$ ,

$$K(x|y) \leq Kc(x|y) + c_{Kc(\cdot|y)}.$$

Substituting  $x = K(y)$ , we find  $Kc(K(y)|y) = O(1)$ , while  $K(K(y)|y)$  can be quite large by Theorem 3.9.1.  $\diamond$

## Exercises

---

**3.9.1.** [22] Show that  $\{K(x) : x = 1, 2, \dots\}$  is the length set of an additively optimal universal code in the sense of Section 1.11.1.

**3.9.2.** [12] Let  $x^*$  be the first enumerated shortest program for  $x$ . Show that  $x^*$  and  $\langle x, K(x) \rangle$  contain the same information:  $K(x^*) = K(\langle x, K(x) \rangle) + O(1)$ .

**3.9.3.** [20] Define conditional complexity  $Kc(x|y) = K(x|y^*)$  with  $y^* = \langle y, K(y) \rangle$ , or  $y^*$  is the first enumerated shortest program for  $y$ . Define  $Kc(x) = K(x)$  and  $Kc(x, y) = K(\langle x, y \rangle) = K(x, y)$ . Prove the following identities (up to an additive constant):

- (a)  $Kc(x, y) = Kc(y, x)$ .
- (b)  $Kc(x|x) = 0$ .
- (c)  $Kc(Kc(x)|x) = 0$ . (Contrast this with Theorem 3.8.1.)
- (d)  $Kc(x) \leq Kc(x, y)$ .
- (e)  $Kc(x|y) \leq Kc(x)$ .
- (f)  $Kc(x, y) = Kc(x) + Kc(y|x)$ . (Contrast this with the case for  $K(x, y)$ , Theorem 3.9.1.)
- (g) We have defined  $I(x; y) = Kc(x) - Kc(x|y)$ . Show that  $I(x; y) \geq 0$ ,  $I(x; x) = Kc(x)$ , and  $I(\epsilon; x) = I(x; \epsilon) = 0$ .
- (h)  $I(x; y) = Kc(x) + Kc(y) - Kc(x, y) + O(1) = I(y; x) + O(1)$ . (In fact,  $I(x; y)$  is the symmetric quantity of mutual information of objects  $x$  and  $y$  on page 235.)
- (i)  $Kc(x, Kc(x)) = Kc(x)$ . (Hint: use (f) and (c). Contrast this derivation with the identical but differently formulated Equation 3.11.)
- (j)  $Kc(x, y) = Kc(x, y, Kc(x), Kc(y|x))$ .
- (k)  $Kc(Kc(x), Kc(y|x)|x, y) = 0$ .
- (l)  $Kc(Kc(x), Kc(y), Kc(y|x), Kc(x, y), Kc(x, y)|x, y) = 0$ .
- (m)  $Kc(I(x; y)|x, y) = 0$ .
- (n)  $Kc(x) \leq Kc(x|y) + Kc(y|z) + Kc(z)$ .
- (o)  $Kc(x, y, z) = Kc(x|y, z) + Kc(y|z) + Kc(z)$ .
- (p) Show that  $Kc(x|y)$  is not co-enumerable. (Hint: use Theorem 3.8.1.) This contrasts with  $K(x|y)$  which is co-enumerable.

*Comments.* See also Example 3.9.1 in Section 3.9.1. Essentially, we have now at our disposal the entire calculus of information theory. In fact, the  $Kc$  calculus is somewhat richer since it contains rules like Items (c) and (i) that have no counterpart in classical information theory. Note that the difference between  $K$  and  $Kc$  is the conditional form  $K(x|y)$

versus  $Kc(x|y)$ .  $Kc$  complexity was introduced by G.J. Chaitin in *J. ACM* 22(1975), 329–340, and a recent exposé appears in *Algorithmic Information Theory*, Cambridge Univ. Press, 1987.

**3.9.4.** [12] Let  $n = l(x)$ , and let  $x$  and  $n$  be incompressible. Show that  $K(x, n) = K(x) + K(n|x) = K(n) + K(x|n)$  up to additive constants.

*Comments.* This relates to the symmetry of information issue for  $K$ . The proof we gave that Theorem 2.8.2 on page 182 is sharp for  $C$  does not hold for  $K$ . Hence, to obtain that the analogue of Theorem 2.8.2 for  $K$  is sharp one has to use another argument. This argument uses the complexity of the complexity function, Theorem 3.8.1, which in fact holds for all variants of Kolmogorov complexity. Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

## 3.10 History and References

Prefix complexity was first introduced in [L.A. Levin, *Problems Inform. Transmission*, 10:3(1974), 206–210; P. Gács, *Soviet Math. Dokl.*, 15(1974), 1477–1480; G.J. Chaitin, *J. ACM*, 22(1975), 329–340]. It resolves the technical problems of Solomonoff’s original proposal for a universal a priori probability [R.J. Solomonoff, *A preliminary report on a general theory of inductive inference*, Tech. Rept. ZTB-138, Zator Company, Cambridge, Mass., November 1960; *Inform. Contr.*, 7(1964), 1–22, 224–254]. L.A. Levin [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surv.*, 25:6(1970), 83–124] identified universal a priori probability with the maximal enumerable semimeasure  $M$  over sample space  $\{0, 1\}^\infty$ , and it turns out that the negative logarithm of the  $m$  version of  $M$  coincides with complexity  $K$  (Theorem 4.3.3 on page 253 in Chapter 4). In some cases  $K$  is a more convenient complexity measure than  $C$ , or conversely, but for many applications one can use both equally well because they coincide to within a logarithmic additive term.

The material in Section 3.2 on the concrete implementation of a universal partial recursive function in combinatory logic to obtain explicit constants bounding how far  $K(x)$  can exceed  $l(x) + 2l(l(x))$ , how far  $K(x|l(x))$  can exceed  $l(x)$ , how far  $C(x)$  can exceed  $l(x)$ , and so on, is due to J. Tromp [‘A CL-based definition of Kolmogorov complexity,’ Manuscript, CWI, Amsterdam, November 1996]. Earlier, G.J. Chaitin [*Complexity*, 1:4(1995/1996), 55–59] used a LISP implementation of a reference universal prefix-machine to define concrete prefix complexity. He explicitly calculated constants involved in upper bounds on prefix complexity and with the “halting probability” (Section 3.6.2).

Estimates for the quantitative relation between  $C$  and  $K$  in Section 3.1 are from [L.A. Levin, *Problems Inform. Transmission*, 10:3(1974), 206–210; S.K. Leung-Yan-Cheong and T.M. Cover, *IEEE Trans. Inform.*

*Theory*, IT-24(1978), 331–339; R. Solovay, *Lecture Notes*, UCLA, 1975, unpublished]. The numerical estimates on  $K$ -complexity and compressibility, like the upper bound on  $K$  and the distribution of description lengths in Theorem 3.3.1, are from [G.J. Chaitin, *J. ACM*, 22(1975), 329–340]. In the original submission of that paper Chaitin proposed to call an infinite binary sequence  $x$  random iff  $K(x_{1:n}) \geq n - O(1)$ . This was shown to be equivalent to Martin-Löf’s notion of randomness, Theorem 3.6.1, by C.P. Schnorr, acting as a referee of that paper.

The halting probability  $\Omega$ , Section 3.6.2, was popularized by Chaitin [*J. ACM*, 22(1975), 329–340] and C.H. Bennett [C.H. Bennett and M. Gardner, *Scientific American*, 241:11(1979), 20–34]. The relation between the halting probability and the solvability of whether Diophantine equations have finitely many solutions or infinitely many solutions, Section 3.7.1, is due to Chaitin [*Adv. Appl. Math.*, 8(1987), 119–146; *Algorithmic Information Theory*, Cambridge Univ. Press, 1987]. The latter book also surveys prefix complexity and randomness of infinite sequences. It incorrectly attributes the presented results only to G.J. Chaitin, P. Martin-Löf, C.P. Schnorr, and R.M. Solovay; see also the book review by P. Gács [*J. Symb. Logic*, 54(1989), 624–627]. Related surveys are [A.N. Kolmogorov and V.A. Uspensky, *Theory Probab. Appl.*, 32(1987), 389–412; V.A. Uspensky, A.L. Semenov and A.Kh. Shen’, *Russ. Math. Surv.*, 45:1(1990), 121–189; V.A. Uspensky, *J. Symb. Logic*, 57:2(1992), 385–412].

Theorem 3.8.1 on the complexity of the complexity function is due to P. Gács [*Soviet Math. Dokl.*, 15(1974), 1477–1480] and was later found independently by R.M. Solovay [*Lecture Notes*, UCLA, 1975, unpublished]. This crucial result establishes the lower bound on the error term up to which information can be symmetrical, for all possible variants of Kolmogorov complexity: Theorem 3.9.1, attributed to L.A. Levin in [P. Gács, *Soviet Math. Dokl.*, 15(1974), 1477–1480]. The equality of the expected algorithmic information (both the Kolmogorov complexity and prefix complexity variants) with Shannon’s entropy is treated in detail in Section 8.1.

## Algorithmic Probability

P.S. Laplace (1749–1827) has pointed out the following reason why intuitively a regular outcome of a random event is unlikely:

"We arrange in our thought all possible events in various classes; and we regard as *extraordinary* those classes which include a very small number. In the game of heads and tails, if heads comes up a hundred times in a row then this appears to us extraordinary, because the almost infinite number of combinations that can arise in a hundred throws are divided in regular sequences, or those in which we observe a rule that is easy to grasp, and in irregular sequences, that are incomparably more numerous."

If we define a "regular object" as an "object with significantly less than maximal complexity," then the number of all regular events is small. This implies that the event that any one of them occurs has small probability (in the uniform distribution). Yet, the classical calculus of probabilities tells us that 100 heads are just as probable as any other sequence of heads and tails, even though our intuition tells us that it is less "random" than some others. Listen to the redoubtable Dr. Samuel Johnson:

"Dr. Beattie observed, as something remarkable which had happened to him, that he chanced to see both the No. 1 and the No. 1000, of the hackney-coaches, the first and the last; 'Why, Sir,' said Johnson, 'there is an equal chance for one's seeing those two numbers as any other two.' He was clearly right; yet the seeing of two extremes, each of which is in some degree more conspicuous than the rest, could not but strike one in a stronger manner than the sight of any other two numbers." [Boswell's *Life of Johnson*]

Laplace distinguishes between the object itself and a cause of the object:

"The regular combinations occur more rarely only because they are less numerous. If we seek a cause wherever we perceive symmetry, it is not that we

regard the symmetrical event as less possible than the others, but, since this event ought to be the effect of a regular cause or that of chance, the first of these suppositions is more probable than the second. On a table we see letters arranged in this order C o n s t a n t i n o p l e, and we judge that this arrangement is not the result of chance, not because it is less possible than others, for if this word were not employed in any language we would not suspect it came from any particular cause, but this word being in use among us, it is incomparably more probable that some person has thus arranged the aforesaid letters than that this arrangement is due to chance."

Let us try to turn Laplace's argument into a formal one. Suppose we observe a binary string  $x$  of length  $n$  and want to know whether we must attribute the occurrence of  $x$  to pure chance or to a cause. "Chance" means that the literal  $x$  is produced by fair coin tosses. "Cause" means that the reference prefix machine  $U$  computes  $x$  when its program is provided by fair coin tosses. The chance of generating  $x$  literally is about  $2^{-n}$ . But the chance of generating  $x$  in the form of a short program from which  $U$  computes  $x$  is at least  $2^{-K(x)}$ . In other words, if  $x$  is regular, then  $K(x) \ll n$ , and it is about  $2^{n-K(x)}$  times more likely that  $x$  arose as the result of computation from some simple cause (like a short program) than literally by a random process.

This gives an objective and absolute definition of "simplicity" as "low Kolmogorov complexity." Consequently, one obtains an objective and absolute version of the classic maxim of William of Ockham (1290?–1349?), known as Occam's razor: "If there are alternative explanations for a phenomenon, then, all other things being equal, we should select the simplest one." One identifies "simplicity of an object" with "an object having a short effective description." In other words, a priori we consider objects with short descriptions more likely than objects with only long descriptions. That is, objects with low complexity have high probability while objects with high complexity have low probability. Pursuing this idea leads to the remarkable probability distribution  $2^{-K(x)}$  below.

## 4.1 Enumerable Functions Revisited

We continue the treatment of enumerable functions where we left it in Section 1.7.3. Co-enumerable and recursive functions were defined in Definition 1.7.3 on page 35. Nontrivial examples of functions that are co-enumerable but not recursive are  $C(x)$ ,  $C(x|y)$ ,  $K(x)$ , and  $K(x|y)$  (Theorems 2.3.2, 2.3.3 on pages 121, 121, respectively). Examples of functions that are enumerable but not recursive are  $-C(x)$ ,  $-K(x)$ ,  $2^{-K(x)}$ , and the universal Martin-Löf test  $\delta_0(x|L) = l(x) - C(x|l(x)) - 1$  with respect to the uniform distribution  $L$ .

**Definition 4.1.1** An enumerable function  $f$  is *universal* if there is an effective enumeration  $f_1, f_2, \dots$  of enumerable functions such that  $f(i, x) = f_i(x)$ , for all  $i, x \in \mathcal{N}$ . (Define  $f(i, x) = f(\langle i, x \rangle)$ .)

**Lemma 4.1.1** *There is a universal enumerable function.*

**Proof.** Let  $\phi_1, \phi_2, \dots$  be the effective enumeration of partial recursive functions in Section 1.7. Consider each partial recursive  $\phi$  as a function  $\phi : \mathcal{N} \times \mathcal{N} \rightarrow \mathcal{Q}$  by interpreting  $\phi(\langle x, k \rangle) = \langle p, q \rangle$  as  $\phi(x, k) = p/q$ . Define for each  $i$  the function  $f_i$  by

$$f_i(x) = \max_{k \in \mathcal{N}} \{\phi_i(x, k)\}, \quad (4.1)$$

or  $\infty$  if such a maximum does not exist. Each such function  $f_i$  is enumerable since we can dovetail the computations of  $\phi_i(x, k)$  for all  $k \geq 1$ . That is, the dovetailed computation proceeds by stages 1, 2,  $\dots$ . At each stage  $j$  the overall computation executes step  $j - k$  of the particular subcomputation of  $\phi_i(x, k)$ , for each  $k$  such that  $j - k > 0$ . On the other hand, if  $f$  is an enumerable function, then there exists a partial recursive function  $\phi$  as in Equation 4.1 by Definition 1.7.3. This way we obtain an enumeration  $f_1, f_2, \dots$  of all and only partial functions that are enumerable.

Let  $\phi_0(i, y)$  be a universal partially recursive function such that for each  $i$ ,  $\phi_0(i, y) = \phi_i(y)$ , for all  $y$ . Interpret  $y$  as a pair  $\langle x, k \rangle$  above. For some index  $j$  we have  $\phi_0(i, y) = \phi_j(\langle i, y \rangle)$ . For this index  $j$  we define  $f_0(i, x) = f_j(\langle i, x \rangle)$ . Then,  $f_0(i, x) = f_i(x)$ , for all  $i$  and  $x$ .  $\square$

An analogous argument shows how to construct a function that is *universal co-enumerable*. In contrast, there is no universal recursive function.

**Example 4.1.1** If  $f(x, y) \geq C(x|y)$ , for all  $x$  and  $y$ , then we call  $f(x, y)$  a *majorant* of  $C(x|y)$ . The minimum of any finite number of majorants is again a majorant. This is the way we combine different heuristics for recognizing patterns in strings. All co-enumerable majorants of  $C(x|y)$  share an interesting and useful property:

**Lemma 4.1.2** *Let  $f(x, y)$  be co-enumerable. For all  $x, y$  we have  $C(x|y) \leq f(x, y) + O(1)$  iff  $d(\{x : f(x, y) \leq m\}) = O(2^m)$ , for all  $y$  and  $m$ .*

**Proof. (ONLY IF)** Suppose to the contrary that for each constant  $c$ , there exist  $y$  and  $m$  such that the number of  $x$ 's with  $f(x, y) \leq m$  exceeds  $c2^m$ . Then by counting, there is an  $x$  such that  $C(x|y, m) \geq m + \log c$ . Together with  $C(x|y) \geq C(x|y, m) + O(1)$ , this yields  $C(x|y) \geq f(x, y) + \log c + O(1)$ . Letting  $c \rightarrow \infty$  contradicts  $C(x|y) \leq f(x, y) + O(1)$ .

(IF) Without loss of generality we can choose the minimal  $m$  satisfying the “if” assumption, for each  $x$  and  $y$ . That is, given  $x$  and  $y$ , set  $m := f(x, y)$ . Let  $g$  be a partial recursive function such that  $g(k, x, y) \geq g(k+1, x, y)$  and  $\lim_{k \rightarrow \infty} g(k, x, y) = f(x, y)$ . Then we can describe  $x$ , given  $y$  and  $m$ , by the recursive function  $g$  approximating  $f$  from above, together with the index of  $x$  in enumeration order, namely, by enumerating all  $x$ 's that satisfy  $f(x, y) \leq m$ .

Hence,  $C(x|y, m) \leq m + O(1)$ . Choose  $h$  such that  $C(x|y, m) = m - h$ . Using first the fact that we can reconstruct  $m$  from  $C(x|y, m)$  and  $h$  and then substituting according to the definition of  $h$  gives

$$\begin{aligned} C(x|y) &\leq C(x|y, m) + 2 \log h + O(1) \\ &\leq m - h + 2 \log h + O(1). \end{aligned}$$

Since we have chosen  $f(x, y) = m$  this proves the “if” part.  $\square \diamond$

A real number  $x = 0.x_1x_2\dots$  is recursively enumerable, or *enumerable* for short, if the set of rationals below  $x$  is recursively enumerable. A number  $x$  is *co-enumerable* if  $-x$  is enumerable. A number  $x$  is *recursive* if it is both enumerable and co-enumerable.

It is easy to show that  $x$  is enumerable (recursive) iff there is an enumerable (recursive) function  $f$  such that  $f(i) = x_i$  for all  $i$ . The halting probability  $\Omega = \sum_{U(p) < \infty} 2^{-l(p)}$  (Section 3.6.2 on page 217) is an enumerable real. Let us explicitly construct the approximation. Define  $\phi(n) = \sum 2^{-l(p)}$ , the sum taken over all programs  $p$  of the reference prefix machine  $U$  of Theorem 3.1.1 on page 193 with  $l(p) \leq n$  that halt within  $n$  steps. Obviously,  $\phi$  is a recursive function, and  $\phi(1), \phi(2), \dots$  is a monotonic nondecreasing sequence of rational numbers with limit  $\Omega$ . Similarly,  $\sum_x 2^{-K(x)} < \Omega$  and moreover, the entire class of  $\Omega$ -like reals introduced in Exercise 3.6.8 on page 222, are enumerable reals.

## 4.2 Nonclassical Notation of Measures

In this section, and in fact in the sequel of this entire chapter, we assume some knowledge of measure theory in the classical sense, say, the basics in Section 1.6. In algorithmic probability theory it is customary to use a nonstandard approach to measures. For better or worse we will follow this usage. Let us first look at standard measure theory.

Classically, the framework is as follows: Let  $\mathcal{B}$  be a finite or countably infinite set of *basic elements*. For example,  $\mathcal{B} = \{0, 1\}$ ,  $\mathcal{B} = \{0, 1\}^*$ , or  $\mathcal{B} = \mathbb{N}$  (the natural numbers). Consider the continuous sample space  $S = \mathcal{B}^\infty$ , that is, all one-way infinite sequences over  $\mathcal{B}$ .

We want to extend the idea of probability from finite sample spaces like the outcomes {head, tail} for fair coin tosses to continuous sample spaces

like  $S$ . Probability cannot be properly defined for the individual elements of  $S$ . (The probability of selecting a particular real number  $r$  from the interval  $[0, 1]$  is necessarily 0 for all but countably many elements.) Therefore, one defines the probability for subsets of  $S$ . Since there are too many subsets to describe, one first defines probability for countably many sets that are easily described. These sets are called *cylinder* sets. Subsequently, by the operation of union, intersection, complement, and countable union, the probability definition is extended to many more subsets of  $S$  according to the Kolmogorov Axioms in Section 1.6. (These “Borel sets” are by no means all subsets of  $S$ .)

A *cylinder* is a set  $\Gamma_x \subseteq S$  defined by

$$\Gamma_x = \{\omega : \omega_{1:l(x)} = x, x \in \mathcal{B}^*\}.$$

Let  $\mathcal{G} = \{\Gamma_x : x \in \mathcal{B}^*\}$  be the set of all cylinders in  $S$ .

A function  $\mu : \mathcal{G} \rightarrow \mathcal{R}$  defines a *probability measure* if

$$\begin{aligned}\mu(\Gamma_\epsilon) &= 1, \\ \mu(\Gamma_x) &= \sum_{b \in \mathcal{B}} \mu(\Gamma_{xb}).\end{aligned}$$

(For general measures we can take  $\mu(\Gamma_\epsilon) \in \mathcal{R}$  or even  $\infty$ .) In this definition we have only defined the measures  $\mu(\Gamma)$  for all cylinders  $\Gamma \subseteq S$ . It induces measures for all subsets of  $S$  obtainable from the cylinders by intersection, union, complement, and countable union. It is a theorem of measure theory that  $\mu$  uniquely induces a measure on the Borel sets. However, in the sequel we are primarily interested in the measures of the cylinders. For convenience of notation we replace the set function on cylinder sets by the isomorphic function on the defining initial segments.

**Notation 4.2.1** A measure  $\mu$  is defined by the function  $\mu : \mathcal{G} \rightarrow \mathcal{R}$ . Now, consider the function  $\mu' : \mathcal{B}^* \rightarrow \mathcal{R}$  defined by  $\mu'(x) = \mu(\Gamma_x)$ . Trivially, from  $\mu'$  we can reconstruct  $\mu$ . From now on we call the  $\mu'$  functions “measures” and drop the primes. Formally, we use the definition of measure below. One should keep in mind that our notation is shorthand for the original measure.

**Definition 4.2.1** A function  $\mu : \mathcal{B}^* \rightarrow \mathcal{R}$  is a *probability measure* if

$$\begin{aligned}\mu(\epsilon) &= 1, \\ \mu(x) &= \sum_{b \in \mathcal{B}} \mu(xb),\end{aligned}$$

for all  $x \in \mathcal{B}^*$ . A semimeasure is a defective measure. A function  $\mu : \mathcal{B}^* \rightarrow \mathcal{R}$  is a *semimeasure* if for all  $x \in \mathcal{B}^*$ ,

$$\begin{aligned}\mu(\epsilon) &\leq 1, \\ \mu(x) &\geq \sum_{b \in \mathcal{B}} \mu(xb).\end{aligned}$$

We can transform any semimeasure  $\mu$  into a measure  $\rho$  by adding a distinguished element  $u$  not in  $\mathcal{B}$ , called the *undefined element*. We simply concentrate the surplus probability in Definition 4.2.1 on  $u$  by setting

$$\begin{aligned}\rho(\epsilon) &= 1, \\ \rho(xu) &= \rho(x) - \sum_{b \in \mathcal{B}} \mu(xb).\end{aligned}$$

It is straightforward that for all  $x \in \mathcal{B}^*$ , we have  $\rho(x) = \mu(x)$ .

**Example 4.2.1** For each  $x \in \{0, 1\}^*$  define the measure  $\lambda(x) = 2^{-l(x)}$ . This is the Lebesgue measure, or *uniform measure*, on the closed unit interval  $[0, 1]$ . It has a geometric interpretation. Consider the real numbers in  $[0, 1]$  as being represented by their binary representation after the decimal dot. A real like  $\frac{1}{2}$  has two representations, namely,  $0.100\dots$  and  $0.011\dots$ . We denote it by  $0.1$  and choose the representation with infinitely many zeros. The uniform measure  $\lambda(x)$  of the cylinder  $\Gamma_x$  is the length  $2^{-l(x)}$  of the interval  $[0.x, 0.x + 2^{-l(x)}]$ .  $\diamond$

This discussion leads to the central notion of this chapter: enumerable and recursive semimeasures.

**Definition 4.2.2** A semimeasure  $\mu$  is *enumerable (recursive)* if the function  $\mu$  is enumerable (recursive).

## Exercises

---

**4.2.1.** [18] Let  $\mu$  be a semimeasure over  $\mathcal{B}^*$ . Show that if  $\mu$  is recursive, then we can find an algorithm to compute  $\mu(x)$  and  $\sum_{b \in \mathcal{B}} \mu(xb)$ , for all  $x \in \mathcal{B}^*$ , to any degree of accuracy.

*Comments.* These properties are implicitly used throughout Section 4.5 on continuous semimeasures. Source: V.G. Vovk, *Soviet Math. Dokl.*, 35(1987), 656–660.

**4.2.2.** [14] (a) Let  $U$  be the reference prefix machine of Theorem 3.1.1 on page 193. Define  $P(x) = \sum_{U(p)=x} 2^{-l(p)}$ . Show that  $\sum_x P(x) \leq 1$ , so  $P(x)$  qualifies as a probability distribution over the integers. (We use the term “probability distribution” loosely here for nonnegative real-valued functions summing up to *at most* 1.)

- (b) Define  $P(x) = 2^{-K(x)}$ . Show that  $\sum_x P(x) \leq 1$ , the sum taken over all  $x$ , so  $P(x)$  qualifies as a probability distribution over the integers.
- (c) Define  $P(x|y) = 2^{-K(x|y)}$ . Show that  $\sum_x P(x|y) \leq 1$ , for each fixed  $y$ , so  $P(x|y)$  qualifies as a conditional probability distribution over the integers.
- (d) Define  $P(x) = 2^{-K(x|l(x))}$ . Show that  $\sum_x P(x) = \infty$ , so  $P(x)$  does not qualify as a probability distribution.

*Comments.* Hint for Item (a): use the Kraft Inequality on page 74. Hint for Item (d): if  $C(x) \geq l(x) - 2l(l(x))$ , then  $K(x|l(x)) \leq l^*(x) - l(l(x)) + O(l(l(l(x))))$ .

## 4.3 Discrete Sample Space

We first develop the theory in the discrete domain. This is in a sense a first approximation to the theory in the continuous domain. One interpretation is to set  $\mathcal{B} = \mathcal{N}$  and consider the sample space  $S = \mathcal{N}$ . (In terms of classical measure theory our sample space is  $S = \{\Gamma_x : x \in \mathcal{N}^*, l(x) = 1\}$ .) Since all elements of  $S$  are one-letter strings, Item 2 in Definition 4.2.1 is not applicable. Since the elements of  $\mathcal{N}$  are considered as one-letter strings, none of them is a prefix of any other. All elements of  $\mathcal{N}$  have prefix  $\epsilon$ . Definition 4.2.1 requires such a measure  $\mu$  to satisfy  $\mu(\epsilon) = \sum_{x \in \mathcal{N}} \mu(x) = 1$ . Except for the interpretation, there is no difference between a discrete measure and a probability distribution over sample space  $\mathcal{N}$ . We use the same font (capital italics) to denote them.

**Definition 4.3.1** A *discrete semimeasure* is a function  $P$  from  $\mathcal{N}$  into  $\mathcal{R}$  that satisfies  $\sum_{x \in \mathcal{N}} P(x) \leq 1$ . It is a *probability measure* if equality holds.

**Example 4.3.1** (RELATION DISCRETE AND CONTINUOUS MEASURE)

The discrete Lebesgue measure  $L$  on the set of basic elements  $\mathcal{B} = \mathcal{N}$  is a function  $L : \mathcal{N} \rightarrow \mathcal{R}$  defined by  $L(x) = 2^{-2l(x)-1}$ . We verify that  $L$  is a probability measure:

$$\sum_{x \in \mathcal{N}} L(x) = \sum_{n \in \mathcal{N}} 2^{-n-1} \sum_{l(x)=n} 2^{-l(x)} = \sum_{n \in \mathcal{N}} 2^{-n-1} = 1.$$

Here we use  $l(x)$  not as the number of occurrences of basic symbols in  $x$  (there is only 1 occurrence of an element in  $\mathcal{B}$ ) but rather as just a function.

The continuous Lebesgue measure  $\lambda$  on the set of basic elements  $\mathcal{B} = \{0, 1\}$  is a function  $\lambda : \{0, 1\}^* \rightarrow \mathcal{R}$  defined by  $\lambda(x) = 2^{-l(x)}$  (from Example 4.2.1). In the discrete measure  $L(x)$  the argument  $x$  is an element of  $\mathcal{N}$  with the interpretation that no two different arguments are prefixes of each other. In the continuous measure  $\lambda(x)$  the argument  $x$

is an element of  $\{0, 1\}^*$  with the usual interpretation that arguments can be prefixes of other ones. An incorrect interpretation is seductive by confusing  $\mathcal{B} = \mathcal{N}$ , which is both the basic set and the domain in the discrete case, with the basic set  $\mathcal{B} = \{0, 1\}$  and/or the domain  $\mathcal{B}^*$  in the continuous case.

We give a numerical example. According to our standard correspondence Equation 1.1, we have  $1, 5, 6 \in \mathcal{N}$  correspond to  $0, 10, 11 \in \{0, 1\}^*$ . But  $L(1) > L(5) + L(6)$  since  $L(1) = \frac{1}{4}$  and  $L(5) = L(6) = \frac{1}{16}$ . The interpretation in terms of cylinder sets for  $L$  is that  $\Gamma_1, \Gamma_5, \Gamma_6$  are pairwise disjoint. As a comparison, for the continuous measure  $\lambda$  we have  $\lambda(1) = \lambda(10) + \lambda(11)$ , and the interpretation in cylinder sets is that  $\Gamma_1 = \Gamma_{10} \cup \Gamma_{11}$ .

For  $\lambda$  we have  $\bigcup_{l(x)=n} \Gamma_x = S$  for each  $n$ . Thus,  $\sum_{l(x)=n} \lambda(x) = 1$  for each  $n$  and

$$\sum_{x \in \{0, 1\}^*} \lambda(x) = \infty.$$

In contrast, for the discrete measure  $L$  we have  $\bigcup_{l(x)=n} \Gamma_x \subset S$ , and we have  $\sum_{l(x)=n} L(x) = 2^{-n-1}$  for each  $n$ , and hence  $\sum_{x \in \mathcal{N}} L(x) = 1$ .  $\diamond$

**Example 4.3.2** If an enumerable semimeasure  $P$  is a probability measure, then it must be recursive. By Definition 1.7.3, there exists a recursive function  $g(x, k)$ , nondecreasing in  $k$ , with  $P(x) = \lim_{k \rightarrow \infty} g(x, k)$ . We can compute an approximation  $P^k$  of the function  $P$  from below for which  $\sum_x P^k(x) > 1 - \epsilon$ . This means that  $|P(x) - P^k(x)| < \epsilon$ , for all  $x$ .  $\diamond$

### 4.3.1 Universal Enumerable Semimeasure

In Section 4.1 we defined the notion of a universal two-argument function as being in an appropriate sense able to simulate each element in the class of one-argument enumerable functions. This was similar to the universality notion in Turing machines. We now look at a slightly different notion of “universality” meaning that some one-argument function is the “largest” in a class of one-argument functions.

**Definition 4.3.2** Let  $\mathcal{M}$  be a class of discrete semimeasures. A semimeasure  $P_0$  is *universal* (or maximal) for  $\mathcal{M}$  if  $P_0 \in \mathcal{M}$ , and for all  $P \in \mathcal{M}$ , there exists a constant  $c_P$  such that for all  $x \in \mathcal{N}$ , we have  $c_P P_0(x) \geq P(x)$ , where  $c_P$  possibly depends on  $P$  but not on  $x$ .

We say that  $P_0$  (multiplicatively) *dominates* each  $P \in \mathcal{M}$ . It is easy to prove that the class of all semimeasures has no universal semimeasure. This is also the case for the class of recursive semimeasures (Lemma 4.3.1 below).

**Theorem 4.3.1** *There is a universal enumerable discrete semimeasure. We denote it by  $\mathbf{m}$ .*

**Proof.** We prove the theorem in two steps. In Stage 1 we show that the enumerable discrete semimeasures can be effectively enumerated as

$$P_1, P_2, \dots$$

In Stage 2 we show that  $P_0$  as defined below is universal:

$$P_0(x) = \sum_{n \geq 1} \alpha(n) P_n(x),$$

with  $\sum \alpha(n) \leq 1$  where  $\alpha(n) > 0$ , for all  $n$ . Stage 1 consists of two parts. In the first part, we enumerate all enumerable functions; and in the second part we effectively change the enumerable functions to enumerable discrete semimeasures, leaving the functions that were already discrete semimeasures unchanged.

**STAGE 1** Let  $\psi_1, \psi_2, \dots$  be an effective enumeration of all real-valued partial recursive functions. Consider a function  $\psi$  from this enumeration (where we drop the subscript for notational convenience). Without loss of generality, assume that each  $\psi$  is approximated by a rational-valued two-argument partial recursive function  $\phi'(x, k) = p/q$  (this is the interpretation of the literal  $\phi'(\langle x, k \rangle) = \langle p, q \rangle$ ). Without loss of generality, each such  $\phi'(x, k)$  is modified to a rational-valued two-argument partial recursive function  $\phi(x, k)$  so as to satisfy the properties below. For all  $x \in \mathcal{N}$ , for all  $k > 0$ ,

- if  $\phi(x, k) < \infty$ , then also  $\phi(x, 1), \phi(x, 2), \dots, \phi(x, k - 1) < \infty$  (this can be achieved by the trick of dovetailing the computation of  $\phi'(x, 1), \phi'(x, 2), \dots$  and assigning computed values in enumeration order to  $\phi(x, 1), \phi(x, 2), \dots$ );
- $\phi(x, k + 1) \geq \phi(x, k)$  (dovetail the computation of  $\phi'(x, 1), \phi'(x, 2), \dots$  and assign the enumerated values to  $\phi(x, 1), \phi(x, 2), \dots$  satisfying this requirement and ignoring the other computed values); and
- $\lim_{k \rightarrow \infty} \phi(x, k) = \psi(x)$  (as does  $\phi'(\cdot)$ ).

The resulting  $\psi$ -list contains all enumerable real-valued functions, and is actually represented by the approximators in the  $\phi$ -list. Each enumerable function  $\psi$  (rather, the approximating function  $\phi$ ) defines a discrete semimeasure  $P$ . In the algorithm below, the local variable array  $P$  contains the current approximation to the values of  $P$  at each stage of the computation. This is doable because the nonzero part of the approximation is always finite.

**Step 1** Initialize by setting  $P(x) := 0$  for all  $x \in \mathcal{N}$ ; and set  $k := 0$ .

**Step 2** Set  $k := k + 1$ , and compute  $\phi(1, k), \dots, \phi(k, k)$ . {If any  $\phi(i, k)$ ,  $1 \leq i \leq k$  is undefined, then  $P$  will not change any more and it is trivially a discrete semimeasure.}

**Step 3** If  $\phi(1, k) + \dots + \phi(k, k) \leq 1$  then set  $P(i) := \phi(i, k)$  for all  $i = 1, 2, \dots, k$  else terminate.

{Step 3 is a test of whether the new assignment of  $P$ -values satisfies the discrete semimeasure requirements.}

**Step 4** Go to Step 2.

If  $\psi$  is already a discrete semimeasure, then  $P = \psi$  and the algorithm never finishes but approximates  $P$  from below. If for some  $x$  and  $k$  with  $x \leq k$  the value  $\phi(x, k)$  is undefined, then the last assigned values of  $P$  do not change any more even though the computation goes on forever. Because the condition in Step 3 is satisfied by the values of  $P$ , it is a discrete semimeasure. If the condition in Step 3 gets violated, then the computation terminates and the  $P$ -approximation to  $P$  is a discrete semimeasure—even a recursive one.

Executing this procedure on all functions in the list  $\phi_1, \phi_2, \dots$  yields an effective enumeration  $P_1, P_2, \dots$  of all enumerable discrete semimeasures (and only enumerable discrete semimeasures).

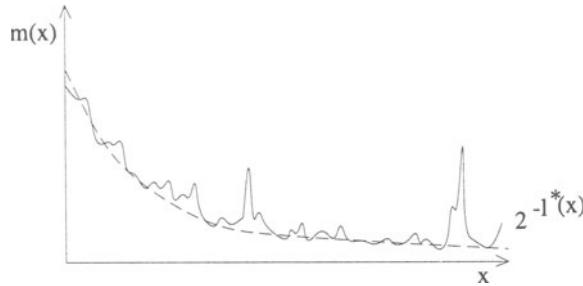
STAGE 2 Define the function  $P_0$  as

$$P_0(x) = \sum_{n \geq 1} \alpha(n) P_n(x),$$

with  $\sum_n \alpha(n) \leq 1$  and  $\alpha(n) > 0$  for all  $n$ . Then  $P_0$  is a discrete semimeasure since

$$\begin{aligned} \sum_{x \geq 0} P_0(x) &= \sum_{n \geq 1} \alpha(n) \sum_{x \geq 0} P_n(x) \\ &\leq \sum_{n \geq 1} \alpha(n) \\ &\leq 1. \end{aligned}$$

The function  $P_0$  is also enumerable, since  $P_n(x)$  is enumerable in  $n$  and  $x$ . (Use the universal partial recursive function  $\phi_0$  and the construction above.) Finally,  $P_0$  dominates each  $P_n$  since  $P_0(x) \geq \alpha(n) P_n(x)$ . Therefore,  $P_0$  is a universal enumerable discrete semimeasure. Obviously, there are countably infinitely many universal enumerable semimeasures. We now fix a *reference* universal enumerable discrete semimeasure and denote it by  $\mathbf{m}$ , depicted in Figure 4.1.  $\square$



**FIGURE 4.1.** Graph of  $m(x)$  with lower bound  $1/x \cdot \log x \cdot \log \log x \cdots$

**Example 4.3.3** In the definition of  $\mathbf{m}(x)$  as  $\sum_n \alpha(n)P_n(x)$  we can choose  $\alpha(n) = 2^{-n}$  or  $\alpha(n) = 6/(\pi n)^2$ . In choosing  $\alpha$  we must take care that the resulting  $\mathbf{m}$  is enumerable. In particular, we can define

$$\mathbf{m}(x) = \sum_{n \geq 1} 2^{-K(n)} P_n(x),$$

the form of which will turn out to be convenient later. Namely, this assignment yields the domination relation  $\mathbf{m}(x) \geq 2^{-K(n)} P_n(x)$ . The domination constant  $2^{-K(n)}$  is for simple  $n$  much greater than the domination constant  $2^{-n}$ . In short, we write

$$\mathbf{m}(x) \geq 2^{-K(P)} P(x), \quad (4.2)$$

for all enumerable discrete semimeasures  $P = P_n$ , where  $K(P) = K(n)$ .  $\diamond$

Let  $f(x, y)$  be a function such that for each fixed  $y$  we have  $\sum_x f(x, y) \leq 1$ . Such functions  $f$  define conditional probabilities  $P(x|y) := f(x, y)$ .

**Theorem 4.3.2** *If  $P(x|y)$  is enumerable, then  $2^{K(P)} \mathbf{m}(x|y) \geq P(x|y)$ , for all  $x, y$ .*

**Lemma 4.3.1** *The class of recursive semimeasures has no universal element.*

**Proof.** Suppose a recursive semimeasure  $P$  is universal for the class of recursive semimeasures. Since for each  $x \in \mathcal{N}$  there is a recursive semimeasure  $P_x$  such that  $P_x(x) > 0$ , it follows from the universality of  $P$  that  $P(x) > 0$  for all  $x$ . Since also  $\sum_x P(x) \leq 1$ , the function  $P$  also converges to zero. We can compute a sequence  $x_0 < x_1 < \dots$  such that

$$\sum_{x_i < x < x_{i+1}} P(x) > x_i P(x_i),$$

for all  $i$ . Therefore, the function  $Q$ , defined by  $Q(x) := xP(x)$  for  $x = x_1, x_2, \dots$ , and zero otherwise, is recursive. Moreover,  $\sum_x Q(x) \leq 1$ , so  $Q$  is a semimeasure. However, for each constant  $c$  there is an  $x$  such that  $Q(x) > cP(x)$ , which contradicts the universality of  $P$ .  $\square$

This lemma is one of the reasons for introducing the notion of enumerable semimeasures, rather than sticking to recursive ones. Compare this to the introduction of enumerable Martin-Löf tests in Sections 2.4, 2.5: among the recursive Martin-Löf tests there is no universal one.

**Lemma 4.3.2** *The function  $\mathbf{m}$  is not recursive and  $\sum_x \mathbf{m}(x) < 1$ .*

**Proof.** If  $\mathbf{m}$  is recursive, then it is universal for the class of recursive semimeasures, by Theorem 4.3.1. But there is no such universal element by Lemma 4.3.1.

We know from Example 4.3.2 that if an enumerable semimeasure is also a probability distribution, then it is recursive. Since  $\mathbf{m}$  is enumerable but not recursive this implies  $\sum_x \mathbf{m}(x) < 1$ .  $\square$

Let us look at the dependency between recursiveness and measurehood of enumerable semimeasures. One reason for introducing semimeasures, instead of just restricting consideration to probability distributions summing to one, is due to the fact that on some input the reference prefix machine  $U$  runs forever. Normalizing  $\mathbf{m}$  by dividing each value of it by  $\Omega = \sum_x \mathbf{m}(x)$  yields a proper probability distribution, say  $\mathbf{P}$ , such that  $\sum_x \mathbf{P}(x) = 1$ . We know from Example 4.3.2 that if an enumerable semimeasure is also a probability distribution, then it is recursive. So either  $\mathbf{P}$  is recursive or it is not enumerable. Measure  $\mathbf{P}$  is not recursive by the same argument that works for  $\mathbf{m}$  in Lemma 4.3.2. Hence,  $\mathbf{P}$  is not even enumerable.

**Example 4.3.4** We look at the behavior of  $\mathbf{m}(x)$  as  $x$  runs over the natural numbers. Define two enumerable measures:  $v(x) = 6/(\pi x)^2$ , and  $w$  defined as

$$w(x) = \begin{cases} 1/x & \text{for } x = 2^k \text{ and } k \in \mathbb{N}^+, \\ 0 & \text{otherwise.} \end{cases}$$

Prove that  $\sum_x w(x) = \sum_x v(x) = 1$ . The function  $\mathbf{m}(x)$  dominates  $v(x)$ . It is less obvious, but nonetheless true, that even though the series  $\sum_x 1/x$  associated with the upper bound  $1/x$  on  $w(x)$  diverges,  $\mathbf{m}(x)$  also dominates  $w(x)$ . From the Kraft Inequality, page 74, we know that the series  $\sum_x 1/x \log^2 x$  converges. The function  $\mathbf{m}(x)$  dominates  $1/x \log^2 x$ , but jumps at many places higher than that as shown in Figure 4.1. This is witnessed by the domination of  $\mathbf{m}(x)$  over  $w(x)$ .  $\diamond$

### 4.3.2 A Priori Probability

Let  $P_1, P_2, \dots$  be the effective enumeration of all enumerable semimeasures constructed in Theorem 4.3.1. There is another way to effectively enumerate the enumerable semimeasures that predate it. Think of the input to a prefix machine  $T$  as being provided by an indefinitely long sequence of fair coin flips. The probability of generating an initial input segment  $p$  is  $2^{-l(p)}$ . If  $T(p) < \infty$ , that is,  $T$ 's computation on  $p$  terminates, then presented with any infinitely long sequence starting with  $p$ , the machine  $T$ , being a prefix machine, will read exactly  $p$  and no further.

Let  $T_1, T_2, \dots$  be the standard enumeration of prefix machines of Theorem 3.1.1, page 193. For each prefix machine  $T$ , define

$$Q_T(x) = \sum_{T(p)=x} 2^{-l(p)}. \quad (4.3)$$

In other words,  $Q_T(x)$  is the probability that  $T$  computes output  $x$  if its input is provided by successive tosses of a fair coin. This means that  $Q_T$  satisfies

$$\sum_{x \in \mathcal{N}} Q_T(x) \leq 1.$$

Equality holds exactly for those  $T$  for which each one-way infinite input contains a finite initial segment constituting a halting program.

We can approximate  $Q_T$  as follows. (The algorithm uses the local variable  $Q(x)$  to store the current approximation to  $Q_T(x)$ .)

**Step 1** Initialize  $Q(x) := 0$  for all  $x$ .

**Step 2** Dovetail the running of all programs on  $T$  so that in stage  $k$  step  $k - j$  of program  $j$  is executed. If the computation of some program  $p$  halts with, say,  $T(p) = x$ , then go to Step 3.

**Step 3** Set  $Q(x) := Q(x) + 2^{-l(p)}$  and go to Step 2.

The algorithm approximates the displayed sum in Equation 4.3 for each  $x$  by the contents of  $Q(x)$ . This shows that  $Q_T$  is enumerable. Starting from a standard enumeration of prefix machines  $T_1, T_2, \dots$ , this construction gives an enumeration of only enumerable semimeasures

$$Q_1, Q_2, \dots$$

The  $P$ -enumeration of Theorem 4.3.1 contains all elements enumerated by this  $Q$ -enumeration. We only need to prove that the  $Q$ -enumeration contains all enumerable measures (Lemma 4.3.4).

**Definition 4.3.3** The *universal a priori probability* on the positive integers is defined as

$$Q_U(x) = \sum_{U(p)=x} 2^{-l(p)},$$

where  $U$  is the reference prefix machine of Theorem 3.1.1.

The use of prefix machines in the present discussion rather than plain Turing machines is necessary. By Kraft's Inequality on page 74, the series  $\sum_p 2^{-l(p)}$  converges (to  $\leq 1$ ) if the summation is taken over all halting programs  $p$  of any fixed prefix machine. In contrast, if the summation is taken over all halting programs  $p$  of an universal plain Turing machine, then the series  $\sum_p 2^{-l(p)}$  diverges.

In Section 3.6.2 we have studied the real number  $\Omega = \sum_x Q_U(x)$  and called it the “halting probability.” No matter how we choose reference  $U$ , the halting probability is less than 1. Namely,  $U$  does not halt for some finite input  $q$  (the *halting problem* in Section 1.7). That is,  $\sum_x Q_U(x) \leq 1 - 2^{-l(q)}$ . If we normalize  $Q_U(x)$  by  $\mathbf{P}(x) = Q_U(x)/\Omega$ , then the resulting function  $\mathbf{P}$  is not enumerable. Namely, if it were enumerable, it would also be recursive by Example 4.3.2. By Theorem 4.3.3 the function  $\mathbf{P}$  would also be universal, and by Lemma 4.3.1 this is impossible.

### 4.3.3 Algorithmic Probability

It is common to conceive of an object as being simpler if it can be briefly described. But the shortness of description of an object depends on the description methods we allow, the “admissible description syntax,” so to speak. We want to effectively reconstruct an object from its description. The shortest self-delimiting effective description of an object  $x$  is quantified by  $K(x)$ .

This leads to a recursively invariant notion of algorithmic probability, which can be interpreted as a form of “Occam’s razor”: the statement that one object is simpler than the other is equivalent to saying that the former object has higher probability than the latter.

**Definition 4.3.4** The *algorithmic probability*  $R(x)$  of  $x$  is defined as

$$R(x) = 2^{-K(x)}.$$

Let us see what this means. Firstly, consider a simple object. If  $x$  consists of a string of  $n$  zeros, then  $K(x) \leq \log n + 2 \log \log n + c$ , where  $c$  is a constant independent of  $n$ . Hence,

$$R(x) \geq \frac{1}{cn \log^2 n}.$$

Generate a binary sequence  $y$  by  $n$  tosses of a fair coin. With overwhelming probability,  $K(y) \geq n$ . For such complex objects  $y$

$$R(y) \leq 1/2^n.$$

### 4.3.4 The Coding Theorem

Now we are ready to state the remarkable and powerful fact that the universal enumerable discrete semimeasure  $\mathbf{m}(x)$ , the universal a priori probability  $Q_U(x)$ , and the algorithmic probability  $R(x) = 2^{-K(x)}$  all coincide up to an independent fixed multiplicative constant. In mathematics the fact that quite different formalizations of concepts turn out to be equivalent is often interpreted as saying that the captured notion has an inherent relevance that transcends the realm of pure mathematical abstraction. We call the generic distribution involved the *universal distribution*. The following is called the *Coding Theorem*.

**Theorem 4.3.3** *There is a constant  $c$  such that for every  $x$ ,*

$$-\log \mathbf{m}(x) = -\log Q_U(x) = K(x),$$

*with equality up to an additive constant  $c$ .*

**Proof.** Since  $2^{-K(x)}$  represents the contribution to  $Q_U(x)$  by a shortest program for  $x$ , we have  $2^{-K(x)} \leq Q_U(x)$ , for all  $x$ .

Clearly,  $Q_U(x)$  is enumerable. Namely, enumerate all programs for  $x$ , by running reference machine  $U$  on all programs at once in dovetail fashion: in the first phase execute step 1 of program 1; in the second phase execute step 2 of program 1 and step 1 of program 2; in the  $i$ th phase execute step  $j$  of program  $k$  for all positive  $j$  and  $k$  such that  $j+k=i$ . By the universality of  $\mathbf{m}(x)$  in the class of enumerable discrete semimeasures,  $Q_U(x) = O(\mathbf{m}(x))$ .

It remains to show that  $\mathbf{m}(x) = O(2^{-K(x)})$ . This is equivalent to showing that  $K(x) \leq -\log \mathbf{m}(x) + O(1)$ . This is shown as follows: Exhibit a prefix-code  $E$  encoding each source word  $x$  as a code word  $E(x) = p$ , satisfying

$$l(p) \leq -\log \mathbf{m}(x) + O(1),$$

together with a decoding prefix machine  $T$  such that  $T(p) = x$ . Then,

$$K_T(x) \leq l(p).$$

Then, by the Invariance Theorem 3.1.1, page 193,

$$K(x) \leq K_T(x) + O(1).$$

On the way to constructing  $E$  as required, we recall a construction for the Shannon-Fano code:

**Lemma 4.3.3** *If  $P$  is a semimeasure on the integers,  $\sum_x P(x) \leq 1$ , then there is a binary prefix-code  $E$  such that the code words  $E(1), E(2), \dots$  are in lexicographical order and  $l(E(x)) \leq -\log P(x) + 2$ . This is the Shannon-Fano code.*

**Proof.** Let  $[0, 1)$  be the half-open real unit interval, corresponding to the sample space  $S = \{0, 1\}^\infty$ . Each element  $\omega$  of  $S$  corresponds to a real number  $0.\omega$ . Let  $x \in \{0, 1\}^*$ . The half-open interval  $[0.x, 0.x + 2^{-l(x)})$  corresponding to the cylinder (set) of reals  $\Gamma_x = \{0.\omega : \omega = x\dots \in S\}$  is called a *binary interval*. We cut off disjoint, consecutive, adjacent (not necessarily binary) intervals  $I_x$  of length  $P(x)$  from the left end of  $[0, 1)$ ,  $x = 1, 2, \dots$ . Let  $i_x$  be the length of the longest binary interval contained in  $I_x$ . Set  $E(x)$  equal to the binary word corresponding to the leftmost such interval. Then  $l(E(x)) = -\log i_x$ . It is easy to see that  $I_x$  is covered by at most four binary intervals of length  $i_x$ , from which the claim follows.  $\square$

In contrast with the proof of Theorem 1.11.1, the Kraft Inequality, we cannot assume here that the sequence  $-\log P(1), -\log P(2), \dots$  is non-decreasing. This causes a loss of two bits in the upper bound.

We use this construction to find a prefix machine  $T$  such that  $K_T(x) \leq -\log \mathbf{m}(x) + c$ . That  $\mathbf{m}(x)$  is not recursive but only enumerable results in  $c = 4$ .

Since  $\mathbf{m}(x)$  is enumerable, there is a partial recursive function  $\phi(t, x)$  with  $\phi(t, x) \leq \mathbf{m}(x)$  and  $\phi(t + 1, x) \geq \phi(t, x)$ , for all  $t$ . Moreover,  $\lim_{t \rightarrow \infty} \phi(t, x) = \mathbf{m}(x)$ . Let  $\psi(t, x)$  be the greatest partial recursive lower bound of special form on  $\phi(t, x)$  defined by

$$\psi(t, x) = \max\{2^{-k} : 2^{-k} \leq \phi(t, x) \text{ and } k \in \mathcal{N}\}.$$

We can assume that  $\psi$  enumerates its range without repetition. Then,

$$\begin{aligned} \sum_{x,t} \psi(t, x) &= \sum_x \sum_t \psi(t, x) \\ &\leq \sum_x 2\mathbf{m}(x) \leq 2. \end{aligned}$$

The series  $\sum_t \psi(t, x)$  can only converge to precisely  $2\mathbf{m}(x)$  in case there is a positive integer  $k$  such that  $\mathbf{m}(x) = 2^{-k}$ .

In a manner similar to the proof of Lemma 4.3.3 on page 253, we chop off consecutive, adjacent, disjoint half open intervals  $I_{t,x}$  of length  $\psi(t, x)/2$ , in enumeration order of a dovetailed computation of all  $\psi(t, x)$ , starting from the left-hand side of  $[0, 1)$ . We have already shown that this is possible. It is easy to see that we can construct a prefix machine  $T$  as follows: If  $\Gamma_p$  is the leftmost largest binary interval of  $I_{t,x}$ , then  $T(p) = x$ . Otherwise,  $T(p) = \infty$  ( $T$  does not halt).

By construction of  $\psi$ , for each  $x$  there is a  $t$  such that  $\psi(t, x) > \mathbf{m}(x)/2$ . Each interval  $I_{t,x}$  has length  $\psi(t, x)/2$ . Each  $I$ -interval contains a binary interval  $\Gamma_p$  of length at least one-quarter of that of  $I$ . Therefore, there

is a  $p$  with  $T(p) = x$  such that  $2^{-l(p)} \geq \mathbf{m}(x)/16$ . This implies  $K_T(x) \leq -\log \mathbf{m}(x) + 4$ , which was what we had to prove.  $\square$

**Corollary 4.3.1** If  $P$  is an enumerable discrete semimeasure, then there is a constant  $c_P = K(P) + O(1)$  such that  $K(x) \leq -\log P(x) + c_P$ , for all  $x$ .

The *conditional universal distribution* is  $Q_U(x|y) = \sum_{U(p,y)=x} 2^{-l(p)}$ . Let  $f(x, y)$  be nonnegative function such that for each fixed  $y$  we have  $\sum_x f(x, y) \leq 1$ . Such functions define conditional probabilities  $P(x|y) := f(x, y)$ . Let  $P(\cdot|y)$  be an enumerable function. By Theorem 4.3.2, we have  $K(x|y) \leq -\log P(x|y) + K(P) + O(1)$ . For fixed  $y$ , define the *universal conditional enumerable semimeasure*  $\mathbf{m}(x|y)$  as the largest (within a constant factor) nonnegative enumerable function  $f(x, y)$ . The proof that  $\mathbf{m}(x|y)$  exists is essentially the same as the proof of the existence of the unconditional universal enumerable semimeasure  $\mathbf{m}(x)$  in Section 4.3.1. Then as a corollary of Theorem 4.3.3 (rather, of its proof), we have the *Conditional Coding Theorem*:

**Theorem 4.3.4** For all  $x, y$ ,

$$-\log \mathbf{m}(x|y) = -\log Q_U(x|y) + O(1) = K(x|y) + O(1).$$

Together, Theorems 4.3.4 and Theorem 4.3.2 on page 249 show the following:

**Corollary 4.3.2**  $K(x|y) \leq -\log P(x|y) + K(P) + O(1)$ .

Theorem 4.3.3 shows that the universal enumerable discrete semimeasure  $\mathbf{m}$  in the  $P$ -enumeration, defined in terms of a function computed by a prefix machine, and the universal a priori probability distribution  $Q_U$  in the  $Q$ -enumeration, defined as the distribution to which the universal prefix machine transforms the uniform distribution, are equal up to a multiplicative constant. This is a particular case of the more general fact that both sequences enumerate the same functions and there are recursive isomorphisms between the two.

**Lemma 4.3.4** There are recursive functions  $f, g$  such that  $Q_n = \Theta(P_{f(n)})$  and  $P_n = \Theta(Q_{g(n)})$ .

**Proof.** First, we construct  $f$ . Let  $Q = Q_n$  be the probability distribution induced by prefix machine  $T_n$  if its programs are generated by fair coin flips. We compute  $Q(x)$  from below by a recursive function  $\phi(x, t)$ , such that  $\phi(t+1, x) \geq \phi(t, x)$  and  $\lim_{t \rightarrow \infty} \phi(t, x) = Q(x)$ . The function  $\phi$  is defined as follows:

**Step 1** For all  $x$ , set  $\phi(0, x) := 0$ . Dovetail the computation of  $T$  on all of its programs  $x$ . Let variable  $t$  count the steps of the resulting computation. Set  $t := 0$ .

**Step 2** Set  $t := t + 1$ .

**Step 3** For all  $x$  with  $x \leq t$  do:

if the computation  $T(p)$  terminates in the  $t$ th step with  $T(p) = x$   
**then** set  $\phi(t+1, x) := \phi(t, x) + 2^{-l(p)}$  **else** set  $\phi(t+1, x) := \phi(t, x)$ .

**Step 4** Go to Step 2.

We can make the described procedure rigorous in the form of a prefix machine  $T$ . Let  $T = T_m$  in the standard enumeration of prefix machines. The construction in Theorem 4.3.1 that transforms every prefix machine into a prefix machine computing a semimeasure leaves  $P_m$  invariant. Machine  $T_m$  computes  $P_m$  in the  $P$ -enumeration. Hence,  $Q = P_m$ . It suffices to set  $f(n) = m$ . Clearly,  $f$  is recursive: we have just outlined the algorithm to compute it.

Second, we construct  $g$ . Let  $P = P_n$  be the  $n$ th element in the effective enumeration constructed in Theorem 4.3.1. Similar to the construction in the proof of the Coding Theorem 4.3.3, since  $P$  is enumerable, we can find a prefix machine  $T$  such that for each  $x$ ,

1.  $P(x)/16 \leq 2^{-K_T(x)}$ , which means that  $K_T(P) \leq -\log P(x) + 4$ .
2. In the proof of Theorem 4.3.3, we had  $\sum_t \psi(t, x) \leq 2m(x)$ , and hence  $\sum_{T(p)=x} 2^{-l(p)} \leq m(x)$ , where  $T$  is as in that proof. Similarly, for  $P$  with the current prefix machine  $T$ ,  $\sum_{T(p)=x} 2^{-l(p)} \leq P(x)$ .

Therefore, if  $Q = Q_m$  in the  $Q$ -enumeration, then the probability distribution  $Q_m$  satisfies

$$2^{-K_T(x)} \leq \sum_{T(p)=x} 2^{-l(p)} = Q_m(x),$$

which by Items 1 and 2 says that  $Q_m(x) = \Theta(P(x))$ . Obviously, the function  $g$  defined by  $g(n) = m$  is recursive.  $\square$

**Example 4.3.5** A priori an outcome  $x$  may have high probability because it has many long descriptions. The Coding Theorem 4.3.3 tells us that in that case it must have a short description too. In other words, the a priori probability of  $x$  is dominated by the shortest program for  $x$ .

Just as we have derived the probability distribution  $Q_U$  from  $U$ , we can derive the probability distribution (actually semimeasure)  $Q_T$  from the

prefix machine  $T$  constructed in the proof of Theorem 4.3.3. Since  $Q_T$  is enumerable, we have  $Q_T(x) = O(\mathbf{m}(x))$ . By definition,  $2^{-K_T(x)} = O(Q_T(x))$ . In the proof of Theorem 4.3.3 it was shown that  $K_T(x) \leq -\log \mathbf{m}(x) + 4$ . Using Theorem 4.3.3 once more, we have  $K(x) \leq K_T(x) + O(1)$ . Altogether, this gives

$$-\log Q_T(x) = -\log \mathbf{m}(x) = K_T(x)$$

up to additive constants. Again therefore, if  $x$  has many long programs with respect to  $T$ , then it also has a short program with respect to  $T$ . In general, we can ask the question of how many descriptions of what length does a finite object have? This leads us to the statistics of description length. For instance, it turns out that there is a universal constant limiting the number of shortest descriptions of *any* finite object, Exercise 4.3.5 on page 266.  $\diamond$

We compare the Coding Theorem 4.3.3 with Shannon's Noiseless Coding Theorem 1.11.2, page 75, in classical information theory. The latter states that given any probability distribution  $P$  on the positive integers, we can construct a binary prefix-code  $E$  in such a way that on the average,  $l(E(x)) \leq -\log P(x) + 1$ . Recall from Section 1.4 the notion of "universal code" as a code that gives near-optimal encodings for any probability distribution on the source alphabet. The Coding Theorem 4.3.3 shows that there is a *single fixed* enumerable universal code  $E'$  for *any* enumerable probability distribution. The code  $E'(x)$  is the shortest program to compute  $x$  by the reference prefix machine  $U$ .

The code-word lengths satisfy  $l(E'(x)) = K(x)$  up to a fixed additive constant independent of  $x$ . We have shown that for any enumerable discrete semimeasure  $P$ ,  $P(x) \leq c_P \mathbf{m}(x)$ , where  $c_P$  is a constant depending on the distribution  $P$  but not on  $x$ . By Theorem 4.3.3, the code-word lengths satisfy  $l(E'(x)) \leq -\log P(x) + c_P$ . This is far better than the performance of any universal code we have met in Section 1.11.1.

### 4.3.5 Randomness by Sum Tests

In Theorem 2.4.1, page 131, we have exhibited a universal  $P$ -test for randomness of a string  $x$  of length  $n$  with respect to an arbitrary recursive distribution  $P$  over the sample set  $S = \mathcal{B}^n$  with  $\mathcal{B} = \{0, 1\}$ .

The universal  $P$ -test measures how justified is the assumption that  $x$  is the outcome of an experiment with distribution  $P$ . We now use  $\mathbf{m}$  to investigate alternative characterizations of random elements of the sample set  $S = \mathcal{B}^*$  (equivalently,  $S = \mathcal{N}$ ).

**Definition 4.3.5** Let  $P$  be a recursive probability distribution on  $\mathcal{N}$ . A *sum  $P$ -test* is a nonnegative enumerable function  $\delta$  satisfying

$$\sum_x P(x) 2^{\delta(x)} \leq 1. \tag{4.4}$$

A *universal sum  $P$ -test* is a test that additively dominates each sum  $P$ -test.

The sum tests of Definition 4.3.5 are slightly stronger than the tests according to Martin-Löf's original Definition 2.4.1, page 129.

**Lemma 4.3.5** *Each sum P-test is a P-test. If  $\delta(x)$  is a P-test, then there is a constant  $c$  such that  $\delta'(x) = \delta(x) - 2 \log \delta(x) - c$  is a sum P-test.*

**Proof.** It follows immediately from the new definition that for all  $n$

$$\sum \{P(x|l(x) = n) : \delta(x) > k\} \leq 2^{-k}. \quad (4.5)$$

Namely, if Equation 4.5 is false, then we contradict Equation 4.4 by

$$\sum_{x \in \mathcal{N}} P(x) 2^{\delta(x)} > \sum_{l(x)=n} P(x|l(x) = n) 2^k \geq 1.$$

Conversely, if  $\delta(x)$  satisfies Equation 4.5 for all  $n$ , then for some constant  $c$ , the function  $\delta'(x) = \delta(x) - 2 \log \delta(x) - c$  satisfies Equation 4.4. Namely,

$$\begin{aligned} \sum_{l(x)=n} P(x|l(x) = n) 2^{\delta'(x)} &= \sum_k \sum \{P(x|l(x) = n) : \delta(x) > k\} 2^{\delta'(x)} \\ &\leq \sum_k 1/(2^k k^2). \end{aligned} \quad (4.6)$$

Choose constant  $c$  such that the last sum converges to at most 1. Since  $P(x|l(x) = n) = P(x)/\sum_{l(x)=n} P(x)$  we have

$$\sum_n \left( \sum_{l(x)=n} P(x) \right) \sum_{l(x)=n} P(x|l(x) = n) 2^{\delta'(x)} = \sum_x P(x) 2^{\delta'(x)}.$$

In the left-hand side of this equality, the expression corresponding to the left-hand side of Equation 4.6 is bounded from above by 1. Therefore, the right-hand side of the equality is at most 1, as desired.  $\square$

this shows that the sum test is not much stronger than the original test. One advantage of Equation 4.4 is that it is just one inequality instead of infinitely many, one for each  $n$ . We give an exact expression for a universal sum P-test in terms of complexity.

**Theorem 4.3.5** *Let  $P$  be a recursive probability distribution. The function  $\kappa_0(x|P) = \log(\mathbf{m}(x)/P(x))$  is a universal sum P-test.*

**Proof.** Since  $\mathbf{m}$  is enumerable and  $P$  is recursive,  $\kappa_0(x|P)$  is enumerable. We first show that  $\kappa_0(x|P)$  is a sum P-test:

$$\sum_x P(x) 2^{\kappa_0(x|P)} = \sum_x \mathbf{m}(x) \leq 1.$$

It only remains to show that  $\kappa_0(x|P)$  additively dominates all sum  $P$ -tests. For each sum  $P$ -test  $\delta$ , the function  $P(x)2^{\delta(x)}$  is a semimeasure that is enumerable. By Theorem 4.3.1, there is a positive constant  $c$  such that  $c \cdot \mathbf{m}(x) \geq P(x)2^{\delta(x)}$ . Hence, there is another constant  $c$  such that  $c + \kappa_0(x|P) \geq \delta(x)$ , for all  $x$ .  $\square$

**Example 4.3.6** An important case is as follows: If we consider a distribution  $P$  restricted to a domain  $A \subset \mathcal{N}$ , then the universal sum  $P$ -test becomes  $\log(\mathbf{m}(x|A)/P(x|A))$ . For example, if  $L_n$  is the uniform distribution on  $A = \{0, 1\}^n$ , then the universal sum  $L_n$ -test for  $x \in A$  becomes

$$\kappa_0(x|L_n) = \log(\mathbf{m}(x|A)/L_n(x)) = n - K(x|n) - O(1).$$

Namely,  $L_n(x) = 1/2^n$  and  $\log \mathbf{m}(x|A) = -K(x|A)$  by Theorem 4.3.4, where we can describe  $A$  by giving  $n$ .  $\diamond$

**Example 4.3.7** The Noiseless Coding Theorem 1.11.2 on page 75 says that the Shannon-Fano code, which codes a source word  $x$  straightforwardly as a word of about  $-\log P(x)$  bits (Example 1.11.2 on page 67, Lemma 4.3.3 on page 253), nearly achieves the optimal expected code-word length. This code is uniform in the sense that it does not use any characteristics of  $x$  itself to associate a code word with a source word  $x$ . The code that codes each source word  $x$  as a code word of length  $K(x)$  also achieves the optimal expected codeword length. This code is nonuniform in that it uses characteristics of individual  $x$ 's to obtain shorter code words. Any difference in code word length between these two encodings for a particular object  $x$  can only be due to exploitation of the individual regularities in  $x$ .

Following Section 2.2.1, define the *randomness deficiency* of a finite object  $x$  with respect to  $P$  as

$$-\lfloor \log P(x) \rfloor - K(x) = -\log P(x) + \log \mathbf{m}(x) + O(1) = \kappa_0(x|P) + O(1),$$

by Theorems 4.3.3, 4.3.5. That is, the randomness deficiency is the outcome of the universal sum  $P$ -test of Theorem 4.3.5.  $\diamond$

**Example 4.3.8** Let us compare the randomness deficiency as measured by  $\kappa_0(x|P)$  with that measured by the universal test  $\delta_0(x|L)$ , for the uniform distribution  $L$ , in Section 2.4. That test consisted actually of tests for a whole family  $L_n$  of distributions, where  $L_n$  is the uniform distribution such that each  $L_n(x) = 2^{-n}$  for  $l(x) = n$ , and zero otherwise. Rewrite  $\delta_0(x|L)$  as

$$\delta_0(x|L_n) = n - C(x|n),$$

for  $l(x) = n$ , and  $\infty$  otherwise. This is close to the expression for  $\kappa_0(x|L_n)$  obtained in Example 4.3.6. From the relations between  $C$  and  $K$  we have established in Chapter 3, it follows that

$$|\delta_0(x|L_n) - \kappa_0(x|L_n)| \leq 2 \log C(x) + O(1).$$

The formulation of the universal sum test in Theorem 4.3.5 can be interpreted as follows: An element  $x$  is random with respect to distribution  $P$ , that is,  $\kappa_0(x|P) = O(1)$ , if  $P(x)$  is large enough, not in absolute value but relative to  $\mathbf{m}(x)$ . If we did not have this relativization, then we would not be able to distinguish between random and nonrandom outcomes for the uniform distribution  $L_n(x)$  above.

Let us look at an example. Let  $x = 00\dots 0$  of length  $n$ . Then  $\kappa_0(x|L_n) = n - K(x|n) + O(1) = n + O(1)$ . If we flip a coin  $n$  times to generate  $y$ , then with overwhelming probability  $K(y|n) \geq n$  and  $\kappa_0(y|L_n) = O(1)$ .

◇

**Example 4.3.9** According to modern physics, electrons, neutrons, and protons satisfy the Fermi-Dirac distribution (Exercise 1.3.7, page 11). We distribute  $n$  particles among  $k$  cells, for  $n \leq k$ , such that each cell is occupied by at most one particle; and all distinguished arrangements satisfying this have the same probability.

We can treat each arrangement as a binary string: an empty cell is a zero and a cell with a particle is a one. Since there are  $\binom{k}{n}$  possible arrangements, the probability for each arrangement  $x$  to happen, under the Fermi-Dirac distribution, is  $FD_{n,k}(x) = 1/\binom{k}{n}$ . According to Theorem 4.3.5,

$$\begin{aligned} \kappa_0(x|FD_{n,k}) &= \log(\mathbf{m}(x|k, n)/FD_{n,k}(x)) \\ &= -K(x|n, k) + \log \binom{k}{n} + O(1) \end{aligned}$$

is a universal sum test with respect to the Fermi-Dirac distribution. It is easy to see that a binary string  $x$  of length  $k$  with  $n$  ones has complexity  $K(x|n, k) \leq \log \binom{k}{n}$ , and  $K(x|n, k) \geq \log \binom{k}{n} - O(1)$  for most such  $x$ . Hence, a string  $x$  with maximal  $K(x|n, k)$  will pass this universal sum test. Each individual such string possesses all effectively testable properties of typical strings under the Fermi-Dirac distribution. Hence, in the limit for  $n$  and  $k$  growing unboundedly, we cannot effectively distinguish one such a string from other such strings.

It is known that photons, nuclei, and some other elementary particles behave according to the Bose-Einstein distribution. Here, we distribute  $n$  particles in  $k$  cells, where each cell may contain many particles. All

possible arrangements are equally likely. By Exercise 1.3.7, the probability of each arrangement  $x$  under the Bose-Einstein distribution is  $BE_{n,k}(x) = 1/A_{n,k}$ , where

$$A_{n,k} = \binom{k+n-1}{n} = \binom{k+n-1}{k-1}.$$

Similar to Example 4.3.9, use Theorem 4.3.5 to obtain a universal sum test with respect to the Bose-Einstein distribution:

$$\kappa_0(x|BE_{n,k}) = \log(\mathbf{m}(x)/BE_{n,k}(x)) = -K(x|n, k) + \log A_{n,k} + O(1).$$

◇

**Example 4.3.10** *Markov's Inequality* says the following: Let  $P$  be any probability distribution; let  $f$  be any nonnegative function with  $P$ -expected value  $\mathbf{E} = \sum_x P(x)f(x) < \infty$ . For  $\mathbf{E} \geq 0$  we have  $\sum\{P(x) : f(x)/\mathbf{E} > k\} < 1/k$ .

Let  $P$  be any probability distribution (not necessarily recursive). The  $P$ -expected value of  $\mathbf{m}(x)/P(x)$  is

$$\sum_x P(x) \frac{\mathbf{m}(x)}{P(x)} \leq 1.$$

Then, by Markov's Inequality

$$\sum_x \{P(x) : \mathbf{m}(x) \leq kP(x)\} \geq 1 - \frac{1}{k}. \quad (4.7)$$

Since  $\mathbf{m}$  dominates all enumerable semimeasures multiplicatively, we have for all  $x$ ,

$$P(x) \leq c_P \mathbf{m}(x), \text{ with } c_P = 2^{K(P)}. \quad (4.8)$$

Equations 4.7 and 4.8 have the following consequences:

1. If  $x$  is a random sample from a simple recursive distribution  $P$ , where “simple” means that  $K(P)$  is small, then  $\mathbf{m}$  is a good estimate for  $P$ . For instance, if  $x$  is randomly drawn from distribution  $P$ , then the probability that

$$c_P^{-1} \mathbf{m}(x) \leq P(x) \leq c_P \mathbf{m}(x)$$

is at least  $1 - 1/c_P$ .

2. If we know or believe that  $x$  is random with respect to  $P$ , and we know  $P(x)$ , then we can use  $P(x)$  as an estimate of  $\mathbf{m}(x)$ .

In both cases the degree of approximation depends on the index of  $P$  and the randomness of  $x$  with respect to  $P$ , as measured by the randomness deficiency  $\kappa_0(x|P) = \log(\mathbf{m}(x)/P(x))$ . For example, the uniform discrete distribution on  $\mathcal{B}^*$  can be defined by  $L(x) = 2^{-2l(x)-1}$ . Then for each  $n$  we have  $L_n(x) = L(x|l(x) = n)$ . To describe  $L$  takes  $O(1)$  bits, and therefore

$$\kappa_0(x|L) = l(x) - K(x) + O(1).$$

The randomness deficiency  $\kappa_0(x|L) = O(1)$  iff  $K(x) \geq l(x) - O(1)$ , that is, iff  $x$  is random.

The nonrecursive “distribution”  $\mathbf{m}(x) = 2^{-K(x)}$  has the remarkable property that the test  $\kappa_0(x|\mathbf{m}) = O(1)$  for all  $x$ : the test shows all outcomes  $x$  random with respect to it. We can interpret Equations 4.7, 4.8 as saying that if the real distribution is  $P$ , then  $P(x)$  and  $\mathbf{m}(x)$  are close to each other with large  $P$ -probability. Therefore, if  $x$  comes from some unknown recursive distribution  $P$ , then we can use  $\mathbf{m}(x)$  as an estimate for  $P(x)$ . In other words,  $\mathbf{m}(x)$  can be viewed as the universal a priori probability’ of  $x$ .

The universal sum  $P$ -test  $\kappa_0(x|P)$  can be interpreted in the framework of hypothesis testing as the likelihood ratio between hypothesis  $P$  and the fixed alternative hypothesis  $\mathbf{m}$ . In ordinary statistical hypothesis testing, some properties of an unknown distribution  $P$  are taken for granted, and the role of the universal test can probably be reduced to some tests that are used in statistical practice.

◇

### 4.3.6 Randomness by Universal Gambling

We toss a fair coin a hundred times and it shows heads every time. The argument that a hundred heads in a row are just as probable as any other outcome convinces us only that the axioms of probability theory do not solve all mysteries as they are supposed to. We feel that a sequence consisting of a hundred heads is not due to pure chance, while some other sequences with the same probability are.

#### Example 4.3.11

In some innominate country with a ruling party and free elections, the share of votes for the ruling party is  $x_i.y_i$  % in thirty successive elections, with  $x_i \geq 50$  and  $y_i$  is the  $i$ th digit in the decimal expansion of  $\pi = 3.1415\dots$ ,  $i = 0, 1, \dots, 29$ . However, if we complain about this, the election organizers tell us that some sequence has to come up, and the actual outcome is as likely as any other. We cannot criticize a regularity we discover *after the fact*, but only those regularities we have excluded in advance. ◇

In probability theory one starts with the assumption that we have a sample space  $S$  of outcomes, with a probability distribution  $P$ . This  $P$  is either discovered empirically or simply hypothesized, for instance by analogy to similar

processes or considerations of symmetry. It is customary to call properties that hold with  $P$ -probability one “laws of probability.”

Consider a Bernoulli process  $(\frac{1}{2}, \frac{1}{2})$  like the repeated tossing of a fair coin. Each outcome  $x$  is an infinite sequence of zeros and ones. It is customary to predict that a random  $x$  will have each property that holds with probability one. But  $x$  cannot be predicted to have all such properties. To see this, consider the property of belonging to the complement of a given singleton set. Each such property has probability one, but jointly they have probability zero. That is, a random outcome  $x$  cannot be expected to withstand *all* statistical tests chosen afterwards together. But we can expect  $x$  to satisfy a few standard laws, like the Law of Large Numbers, and presume them always chosen. However, the classical theory of probability gives us no criteria for selection of such standard laws.

Kolmogorov’s solution is to select those randomness properties with probability close to one that are “simply expressible.” The objects that do not satisfy such a property have a corresponding regularity and form a simply described set of small measure and correspondingly small cardinality. Then each such object is simply described by the set it is an element of and its position in that set. This allows substitution of the multiple requirement of “satisfying all regularities involved” by a single requirement of “not being a simple object.” In the betting approach we place a single bet that gives a huge payoff in case the outcome is not complex, and which thereby safeguards us against all simple ways of cheating.

**Definition 4.3.6** A nonnegative function  $t : \mathcal{N} \rightarrow \mathcal{R}$  is a *P-payoff function* if

$$\sum_{x \in \mathcal{N}} P(x)t(x) \leq 1.$$

The definition says that the logarithm of an enumerable  $P$ -payoff function is a sum  $P$ -test as in Definition 4.3.5. Among the enumerable payoff functions there is a universal payoff function that incorporates all particular payoff functions: a universal betting strategy.

**Definition 4.3.7** An enumerable  $P$ -payoff function  $t_0 : \mathcal{N} \rightarrow \mathcal{R}$  is *universal* if it multiplicatively dominates each enumerable  $P$ -payoff function  $t$  (that is,  $t(x) = O(t_0(x))$ ).

**Lemma 4.3.6** Let  $P$  be a recursive probability distribution. The function  $t_0(x|P) = m(x)/P(x)$  is a universal enumerable  $P$ -payoff function.

**Proof.** Each enumerable  $P$ -payoff function  $t$  can be expressed as  $t(x) = 2^{\delta(x)}$  with  $\delta$  a sum  $P$ -test, Definitions 4.3.5, 4.3.6. Since  $t_0(x|P) = 2^{\kappa_0(x|P)}$  with  $\kappa_0$  the universal sum  $P$ -test of Theorem 4.3.5, and  $\kappa_0$  dominates each  $\delta$  additively, it follows that  $t_0$  dominates each  $t$  multiplicatively.  $\square$

**Betting Against a Crooked Player**

Suppose you meet a street gambler tossing a coin and offering odds to all passers-by on whether the next toss will be heads “1” or tails “0.” He offers to pay you two dollars if the next toss is heads; you pay him one dollar if the next toss is tails. Should you take the bet? If the gambler is tossing a fair coin, it is a great bet. Probably you will win money in the long run. After all, you can expect that half of the tosses come up heads and half tails. Losing up only one dollar on each heads toss and getting two for each tails makes you rich fast. After some observation you notice that the sample sequence of outcomes looks like 01010101010... . Perhaps the gambler manipulates the outcomes. Expecting foul play, you make the following offer as a bet for 1,000 coin tosses.

You pay \$1 first and propose that your opponent pays you  $2^{1000-K(x)}$  dollars with  $x$  the binary sequence of outcomes of the 1,000 coin flip results. This is better than fair since the gambler is only expected to pay

$$\sum_{l(x)=1000} 2^{-1000} 2^{1000-K(x)} < \$1,$$

by Kraft’s inequality. So he should be happy to accept the proposal. But if the gambler cheats, then, for example, you receive  $2^{1000-\log 1000}$  dollars for a sequence like 01010101010... !

In the 1 versus 2 dollars scheme, you can also propose to add this as an extra bonus pay. This way, you are guaranteed to win big: either polynomially increase your money (when the gambler does not cheat) or exponentially increase your money (when the gambler cheats).

**Example 4.3.12**

Suppose a gambler proposes the following wager to the election organizers in Example 4.3.11. He will bet one dollar in each election. The organizers claim that each outcome  $x$  associated with  $n$  elections has probability  $L_n(x)$ , where  $L_n(x) = 10^{-n}$  is the uniform distribution on decimal strings of length  $n$  and zero otherwise. We formulate a payoff function that is a winning strategy against all simply describable malversations. To back up their claim, the organizers ought to agree to pay  $t(x)$  dollars on outcome  $x$  on any payoff function  $t$  we propose. Namely, the expected amount of payoff is at most the gambler’s total original wager of  $n$  dollars. Accordingly, we propose as payoff function  $t = t_0(\cdot|L_n)$ , the universal payoff function with respect to  $L_n$  defined by

$$\begin{aligned} t_0(x|L_n) &= 2^{-\log L_n(x)-K(x|n)} \\ &= 2^{n \log 10 - K(x|n)}. \end{aligned}$$

If the  $x$  consists of the first  $n$  digits of the decimal expansion of  $\pi$ , the election organizers have to pay the gambler the staggering amount of

$$2^{n \log 10 - K(\pi_{1:n}|n)} \geq c10^n$$

dollars for some fixed constant  $c$  independent of  $n$ , even though the bet did not refer to  $\pi$ . Even if the organizers are smart and switch to some pseudorandom sequence algorithmically generated by their computer, they will have to pay such an amount. In other words, since we propose the payoff function beforehand, it is unlikely that we define precisely the one that detects a particular fraud. However, fraud implies regularity, and the number of regularities is so small that we can afford to make a combined wager on all of them in advance.  $\diamond$

The fact that  $t_0(x|P)$  is a payoff function implies by Markov's Inequality, Equation 4.7, that for each  $k > 0$ ,

$$\sum_x \{P(x) : K(x) \geq -\log P(x) - k\} \geq 1 - 1/2^k. \quad (4.9)$$

By Equation 4.2 and Theorem 4.3.3, for all  $x$ ,

$$K(x) \leq -\log P(x) + K(P) + O(1). \quad (4.10)$$

Setting  $k := K(P)$  we find that with large probability, the complexity  $K(x)$  of a random outcome  $x$  is close to its upper bound  $-\log P(x) + O(K(P))$ . If an outcome  $x$  violates any “laws of probability,” then the complexity  $K(x)$  falls far below the upper bound. Indeed, a proof of some law of probability like the Law of Large Numbers or the Law of the Iterated Logarithm always gives rise to some simple recursive payoff function  $t(x)$  taking large values on the outcomes violating the law.

We can phrase the relations as follows: Since the payoff function  $t_0(\cdot|P)$  dominates all  $P$ -payoff functions that are enumerable,  $\kappa_0(\cdot|P)$  is a universal test of randomness—it measures the deficiency of randomness in the outcome  $x$  with respect to distribution  $P$ , or the extent of justified suspicion against hypothesis  $P$  given the outcome  $x$ .

## Exercises

**4.3.1.** [21] Show that the universal distribution  $\mathbf{m}$  (universal enumerable discrete semimeasure) converges to zero more slowly than any positive recursive function that converges to zero.

**4.3.2.** [15] Show that the class of recursive measures does not contain a universal element.

**4.3.3.** [25] Show that  $\sum_y 2^{-K(x|y)} \leq 1$ .

*Comments.* Hint: use Theorem 4.3.1.

**4.3.4.** [23] Show that if  $K(x) < \log x$  then  $\sum_{i=1}^x 2^{-K(i)} \leq x 2^{-K(x)}$ .

*Comments.* Hint: use Theorem 4.3.3.

**4.3.5.** [32] We study the statistics of description length. By the Coding Theorem 4.3.3,  $K(x) = -\log Q_U(x)$  up to an additive constant. Informally, if an object has many long descriptions, then it also has a short one.

(a) Let  $f(x, n)$  be the number of binary strings  $p$  of length  $n$  with  $U(p) = x$ , where  $U$  is the reference prefix machine of Theorem 3.1.1, page 193. Show that for all  $n \geq K(x)$ , we have  $\log f(x, n) = n - K(x, n) + O(1)$ .

(b) Use Item (a) to show that  $\log f(x, K(x)) = K(x) - K(x, K(x)) + O(1) = O(1)$ . The number of shortest descriptions of any object is bounded by a universal constant.

*Comments.* Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

**4.3.6.** [19] Show that up to a fixed additive constant  $\sum_{l(x)=n} \mathbf{m}(x) = \mathbf{m}(n)$ .

*Comments.* Source: P. Gács, *Ibid.*

**4.3.7.** [17] Give an example of a recursive sequence of rational numbers  $a_n > 0$  such that the series  $\sum_n a_n < \infty$ , but for each other sequence  $b_n > 0$  if  $b_n/a_n \rightarrow \infty$  then  $\sum_n b_n = \infty$ .

*Comments.* Hint: Let  $r_n$  be a recursive increasing sequence of rational numbers with  $\lim_n r_n = \sum_x \mathbf{m}(x)$  and let  $a_n = r_{n+1} - r_n$ . Source: P. Gács, *Ibid.*

**4.3.8.** [29] We can also express statistics of description length with respect to  $C$ . For an enumerable function  $f$  with  $\{f(k) : k \geq 1\}$  satisfying the Kraft Inequality, there exist fewer than  $2^{k+f(k)+O(1)}$  programs of length  $C(x) + k$  for  $x$ .

*Comments.* Hint: consider a machine that assigns a code of length  $m$  to  $x$  iff  $x$  has at least  $2^{k+f(k)}$  programs of length  $m+k$ . Then the number of strings that are assigned a code of length  $m$  is at most  $\sum_k (2^{m+k}/2^{k+f(k)}) = \sum_k 2^{m-f(k)}$ , which by Kraft's Inequality, page 74, is at most  $2^m$ . Hence, this is valid recursively enumerable code. Since  $x$  has no program of length less than  $C(x)$ ,  $x$  has fewer than  $2^{k+f(k)+O(1)}$  programs of length  $C(x)+k$ . Source: J. Tromp, personal communication, March 13, 1991.

**4.3.9.** [39] How many objects are there of a given complexity  $n$ ? Let  $g(n)$  be the number of objects  $x$  with  $K(x) = n$ , and let  $D_n$  be the set of binary strings  $p$  of length  $n$  such that  $U(p)$  is defined. Define the moving average  $h(n, c) = 1/(2c+1) \sum_{i=-c}^c g(n+i) + O(1)$ .

(a) First show  $\sum_y \mathbf{m}(x, y) = \mathbf{m}(x) + O(1)$ .

(b) Show that there is a natural number  $c$  such that  $\log d(D_n) = n - K(n) + O(1)$  and also  $\log h(n, c) = n - K(n)$ .

*Comments.* Hint for Item (b): use Exercise 4.3.5 and Item (a). Since we are interested only in equality up to an additive constant we can omit the normalizing factor  $1/(2c+1)$  from the definition of  $h$ . But we do not know whether we can replace  $h$  by  $g$ . Namely, one can choose a reference prefix machine  $U'$  such that  $g'(n) = 0$  for all odd  $n$ . For instance,  $U'(00p) = U(p)$  if  $l(p)$  is even,  $U'(1p) = U(p)$  for  $l(p)$  is odd, and  $U'(p)$  is undefined otherwise. Then  $U'$  is only defined for inputs of even length, and for all  $x$  we have  $K_{U'}(x) \leq K(x) + 2$ . Source: R.M. Solovay, *Lecture Notes*, 1975, unpublished; and P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

**4.3.10.** [32] Suppose we want to obtain information about a certain object  $x$ . It is not a good policy to guess blindly. The mutual information of two objects  $x$  and  $y$  was given in Example 3.9.1 on page 235 as  $I(x; y) = K(y) - K(y|x, K(x))$ . Show that  $\sum_y \mathbf{m}(y)2^{I(x; y)} = O(1)$ .

*Comments.* In words, the expected value of  $2^{I(x; y)}$  is small, even with respect to the universal distribution  $\mathbf{m}(x)$ . Hint: by the Coding Theorem 4.3.3,  $2^{I(x; y)} = 2^{-K(y|x, K(x))}/\mathbf{m}(x) + O(1)$ . Source: P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987.

**4.3.11.** [37] The *randomness deficiency* of an element  $x$  in a finite set  $A$  is the quantity  $\delta_0(x|A) = \log d(A) - C(x|A)$ . The quantities  $\delta_0(x|A)$  and  $\delta_0(x|\mu)$ , the randomness of  $x$  with respect to distribution  $\mu$ , are related as follows. (The  $c$ 's are fixed constants independent of  $A$  and  $x$ .)

(a) For an arbitrary finite set  $A \subseteq \mathcal{B}^*$ , let  $\mu(x) = 1/d(A)$  if  $x \in A$ , and zero otherwise. Show that

$$\begin{aligned}|C(\mu) - C(A)| &\leq c_1, \\ |\delta_0(x|\mu) - \delta_0(x|A)| &\leq 2 \log C(x|A) + c_2.\end{aligned}$$

(b) For an arbitrary measure  $\mu$  on  $\mathcal{B}^*$  and  $x \in \mathcal{B}^*$  define  $A = \{y | l(y) = l(x), \mu(y) \geq 2^{-(k+1)}\}$ , where  $k$  is an integer such that  $k \leq -\log \mu(x) \leq k+1$ . Then  $x \in A$ , and

$$\begin{aligned}C(A) &\leq C(\mu) + 2 \log(\delta_0(x|\mu) + C(x)) + c_3, \\ \delta_0(x|A) &\leq \delta_0(x|\mu) + 2 \log C(x) + 2 \log(\delta_0(x|\mu) + C(x)) + c_4.\end{aligned}$$

*Comments.* Source: V.V. V'yugin, *SIAM Theory Probab. Appl.*, 32(1987), 508–512; the notion  $\delta_0(x|A)$  was introduced by A.N. Kolmogorov, *Lecture Notes in Math.*, vol. 1021, Springer-Verlag, 1983, pp. 1–5.

**4.3.12.** [34] Let  $X = x_1, x_2, \dots$  be a recursive sequence of natural numbers (in  $\mathcal{N}$  or the corresponding binary strings). The *lower frequency* of some element  $x$  in the sequence is defined as

$$q_X(x) = \liminf_{n \rightarrow \infty} \frac{1}{n} d(\{i : i < n \text{ and } x_i = x\}).$$

- (a) Show that there is a *universal recursive sequence*  $U = u_1, u_2, \dots$ , such that for any recursive sequence  $X = x_1, x_2, \dots$  there is a constant  $c > 0$  such that  $c q_U(x) \geq q_X(x)$ , for all  $x$  in  $\mathcal{N}$ .
- (b) Show that if  $U$  and  $V$  are universal recursive sequences, then  $q_U(x) = \Theta(q_V(x))$ . Fix a *reference universal recursive sequence*  $U$ , and define the *a priori frequency* of  $x$  as  $\mathbf{q}(x) = q_U(x)$ .
- (c) Show that  $\mathbf{q}(x) \neq \Theta(\mathbf{m}(x))$ . ( $\mathbf{m}(x)$  is the *a priori probability* of  $x$ .)
- (d) A set is *enumerable relative to  $0'$* , if it is the range of a function computable by an algorithm with an oracle for some recursively enumerable set. An “algorithm with an oracle for set  $A$ ” is an algorithm that (apart from the usual things) at each step can ask a question of the form “is  $a$  in  $A$ ?” and get the true answer “yes/no” for free. The recursively enumerable sets correspond to the  $0'$ -enumerable sets where the oracle is restricted to recursive sets. Define the notion of  *$0'$ -enumerable semimeasures*, and show that there is a *universal  $0'$ -enumerable semimeasure*  $\mathbf{p}$  such that for each  $0'$ -enumerable semimeasure  $\nu$  there is a constant  $c > 0$  such that  $\mathbf{p}(x) \geq c\nu(x)$ . We call  $\mathbf{p}$  the *a priori probability relative to  $0'$* .
- (e) Show that  $\mathbf{p}(x) = \Theta(\mathbf{q}(x))$ . Compare this with Item (c).

*Comments.* Source: An.A. Muchnik, *SIAM Theory Probab. Appl.*, 32 (1987), 513–514; A.N. Kolmogorov, V.A. Uspensky, *SIAM Theory Probab. Appl.*, 32(1987), 389–412.

## 4.4

### Universal Average-Case Complexity

---

The universal distribution  $\mathbf{m}$  is one of the foremost notions in the theory of Kolmogorov complexity. It multiplicatively dominates all enumerable distributions (and therefore also all computable ones). Therefore, a priori it maximizes ignorance by assigning maximal probability to all objects. It has many remarkable properties and applications. Here we observe that the average-case computational complexity of *any algorithm whatsoever* under the universal distribution turns out to be of the same order of magnitude as the worst-case complexity. This holds both for time complexity and for space complexity.

For many algorithms the average-case running time under some distributions on the inputs is less than the worst-case running time. For instance, using (nonrandomized) Quicksort on a list of  $n$  items to be sorted gives under the uniform distribution on the inputs an average running time of  $O(n \log n)$  while the worst-case running time is  $\Omega(n^2)$ . The worst-case running time of Quicksort is typically reached if the list is already sorted or almost sorted, that is, exactly in cases where we actually should not have to do much work at all. Since in practice the lists to be sorted occurring in computer computations are often sorted or almost sorted, programmers often prefer other sorting algorithms that might run faster

with almost sorted lists. Without loss of generality we identify inputs of length  $n$  with the natural numbers corresponding to binary strings of length  $n$ .

**Definition 4.4.1** Consider a discrete sample space  $\mathcal{N}$  with probability density function  $P$ . Let  $t(x)$  be the running time of algorithm  $A$  on problem instance  $x$ . Define the *worst-case time complexity* of  $A$  as  $T(n) = \max\{t(x) : l(x) = n\}$ . Define the  *$P$ -average time complexity* of  $A$

$$T(n|P) = \frac{\sum_{l(x)=n} P(x)t(x)}{\sum_{l(x)=n} P(x)}.$$

**Example 4.4.1** (**Quicksort**) We compare the average time complexity for Quicksort under the Uniform Distribution  $L(x)$  and under the Universal distribution  $\mathbf{m}(x)$ . Define  $L(x) = 2^{-2l(x)-1}$ , such that the conditional probability  $L(x|l(x) = n) = 2^{-n}$ . We encode the list of elements to be sorted as nonnegative integers in some standard way.

For Quicksort,  $T(n|L) = \Theta(n \log n)$ . We may expect the same complexity under  $\mathbf{m}$ , that is,  $T(n|\mathbf{m}) = \Omega(n \log n)$ . But Theorem 4.4.1 will tell us much more, namely,  $T(n|\mathbf{m}) = \Omega(n^2)$ . Let us give some insight why this is the case.

With the low average time complexity under the Uniform Distribution, there can only be  $o((\log n)2^n/n)$  strings  $x$  of length  $n$  with  $t(x) = \Omega(n^2)$ . Therefore, given  $n$ , each such string can be described by its sequence number in this small set, and hence for each such  $x$  we find  $K(x|n) \leq n - \log n + 3 \log \log n$ . (Since  $n$  is known, we can find each  $n - k$  by coding  $k$  self-delimiting in  $2 \log k$  bits. The inequality follows by setting  $k \geq \log n - \log \log n$ .)

Therefore, no really random  $x$ 's, with  $K(x|n) \geq n$ , can achieve the worst-case running time  $\Omega(n^2)$ . Only strings  $x$  that are nonrandom, with  $K(x|n) < n$ , among which are the sorted or almost sorted lists, and lists exhibiting other regularities, can have  $\Omega(n^2)$  running time. Such lists  $x$  have relatively low Kolmogorov complexity  $K(x)$  since they are regular (can be compactly described), and therefore  $\mathbf{m}(x) = 2^{-K(x)+O(1)}$  is very high. Therefore, the contribution of these strings to the average running time is weighted very heavily.  $\diamond$

**Theorem 4.4.1** (**m-Average Complexity**) *Let  $A$  be an algorithm with inputs in  $\mathcal{N}$ . Let the inputs to  $A$  be distributed according to the universal distribution  $\mathbf{m}$ . Then the average-case time complexity is of the same order of magnitude as the corresponding worst-case time complexity.*

**Proof.** We define a probability distribution  $P(x)$  on the inputs that assigns high probability to the inputs for which the worst-case complexity is reached, and zero probability for other cases.

Let  $A$  be the algorithm involved. Let  $T(n)$  be the worst-case time complexity of  $A$ . Clearly,  $T(n)$  is recursive (for instance by running  $A$  on all  $x$ 's of length  $n$ ). Define the probability distribution  $P(x)$  as follows:

**Step 1** For each  $n = 0, 1, \dots$ , set  $a_n := \sum_{l(x)=n} \mathbf{m}(x)$ .

**Step 2** If  $l(x) = n$  and  $x$  is lexicographically least with  $t(x) = T(n)$  then  $P(x) := a_n$  else  $P(x) := 0$ .

It is easy to see that  $a_n$  is enumerable since  $\mathbf{m}(x)$  is enumerable. Therefore,  $P(x)$  is enumerable. Below we use  $c_P \mathbf{m}(x) \geq P(x)$ , where  $c_P = K(P) + O(1)$  is a constant depending on  $P$  but not on  $x$ , Theorem 4.3.1 on page 247 and Example 4.3.3 on page 249. We have defined  $P(x)$  such that  $\sum_{x \in N} P(x) \geq \sum_{x \in N} \mathbf{m}(x)$ , and  $P(x)$  is an enumerable probability distribution. The average-case time complexity  $T(n|\mathbf{m})$  with respect to the  $\mathbf{m}$  distribution on the inputs, is now obtained by

$$\begin{aligned} T(n|\mathbf{m}) &= \sum_{l(x)=n} \frac{\mathbf{m}(x)t(x)}{\sum_{l(x)=n} \mathbf{m}(x)} \\ &\geq \frac{1}{c_P} \sum_{l(x)=n} \frac{P(x)}{\sum_{l(x)=n} \mathbf{m}(x)} T(n) \\ &= \frac{1}{c_P} \sum_{l(x)=n} \frac{P(x)}{\sum_{l(x)=n} P(x)} T(n) = \frac{1}{c_P} T(n). \end{aligned}$$

The inequality  $T(n) \geq T(n|\mathbf{m})$  holds vacuously.  $\square$

**Corollary 4.4.1** The analogue of the theorem holds for other complexity measures (like space complexity) by about the same proof.

If the algorithm to approximate  $P(x)$  from below is the  $k$ th algorithm in the standard effective enumeration of all algorithms, then  $c_P = K(P) + O(1) < k \log^2 k$ . To approximate the optimal value we must code the algorithm to compute  $P$  as compactly as possible. The ease with which we can describe (algorithmically) the strings that produce a worst-case running time determines the closeness of the average time complexity to the worst-case time complexity. Let  $S \subseteq \{0, 1\}^n$ . Denote by  $T(n|P, S)$  the  $P$ -average computation time as in Definition 4.4.1 but with the average taken over  $S$  instead of  $\{0, 1\}^n$  as with  $T(n|P)$ .

**Lemma 4.4.1** Let  $Q$  be a recursive probability distribution. There is a set  $S$  of inputs with  $Q(S) \geq 1 - 2^{-k}$  such that  $T(n|Q, S) \geq T_S(n|\mathbf{m}, S)/2^{K(Q)+k+O(1)}$ .

**Proof.** If the probability distribution  $Q$  is enumerable (which by Example 4.3.2 on page 246 means it is recursive), then by Markov's Inequality, Equation 4.7 on page 261, and  $2^{K(Q)+O(1)}\mathbf{m}(x) \geq Q(x)$ , substitution in Definition 4.4.1 restricted to  $S$  shows the lemma.  $\square$

**Example 4.4.2 (Quicksort continued)** The average time complexity of Quicksort with the inputs distributed according to the uniform distribution is order  $n \log n$ . By Lemma 4.4.1 and Theorem 4.4.1 we find (for simplicity ignoring  $O(1)$  factors) that the same holds in the computational reality for a set of inputs of combined  $L$ -probability at least  $1 - 2^k$  as long as

$$n \log n \geq n^2 / 2^{K(L)+K(P)+k}.$$

Here,  $L$  is the uniform distribution substituted for  $Q$  in Lemma 4.4.1 (generated by a program of at least length  $K(L)$ ) and  $P$  is the particular distribution used in the proof of Theorem 4.4.1. For  $k$  large and  $n$  so large that  $K(L) + K(P) + k < \log n - \log \log n$  the square average-case running time must take over. Clearly, increasing  $K(L)$  (more complex algorithmic random number generator) increases the size of  $n$  at which the square running time starts to take over.

Frequently, algorithmically generated random numbers are used in order to reduce the average-case computation time to below the worst-case computation time. The above example gives evidence that for every input length only sufficiently complex algorithmic random number generators can achieve reduced average-case computation time.  $\diamond$

**Example 4.4.3** In learning applications in Section 5.4.3 we want to draw elements from the  $\mathbf{m}$  distribution. Since  $\mathbf{m}$  is not computable, we can't have a program for it. Suppose some powerful source deems it fit to give us a table with sufficiently many  $\mathbf{m}$  values. We use this table to randomly draw according to  $\mathbf{m}$  as follows:

Our prospective algorithm has access to an  $\mathbf{m}$  table in the form of a division of the real open interval  $[0, 1)$  into nonintersecting half-open subintervals  $I_x$  such that  $\bigcup I_x = [0, 1)$ . For each  $x$ , the length of interval  $I_x$  is  $\mathbf{m}(x) / \sum_y \mathbf{m}(y)$ . For each finite binary string  $r$ , the *cylinder*  $\Gamma_r$  is the set of all infinite binary strings starting with  $r$ . That is,  $\Gamma_r$  is a half open interval  $[0.r, 0.r + 2^{-l(r)})$  in  $[0, 1)$ . To draw a random example from  $\mathbf{m}$ , the algorithm uses a sequence  $r_1 r_2 \dots$  of outcomes of fair coin flips until the cylinder  $\Gamma_r$ ,  $r = r_1 r_2 \dots r_k$ , is contained in some interval  $I_x$ . It is easy to see that this procedure of selecting  $x$ , using a table for  $\mathbf{m}$  and fair coin flips, is equivalent to drawing an  $x$  randomly according to distribution  $\mathbf{m}$ .

We are often interested in drawing an element of a subset  $D$  of  $\mathcal{N}$ . For instance, we want to draw an  $n$ -length binary vector ( $D = \{0, 1\}^n$ )

when learning Boolean functions. To draw from  $\mathbf{m}(\cdot|D)$ , we simply draw examples from  $\mathbf{m}(\cdot)$  and discard the ones not in  $D$ . If we need to draw  $m$  examples according to  $\mathbf{m}(\cdot|D)$ , then it suffices to draw  $\Omega(2^{K(D)}m)$  examples under  $\mathbf{m}(\cdot)$ . Namely, for each  $x \in D$ ,

$$\begin{aligned}\mathbf{m}(x|D) &= \mathbf{m}(x) \sum_{y \in \mathcal{N}} \mathbf{m}(y) \Bigg/ \sum_{y \in D} \mathbf{m}(y) \\ &= \Theta(2^{K(D)} \mathbf{m}(x)).\end{aligned}$$

◊

Computable versions of  $\mathbf{m}(\cdot)$  and their applications are treated in Section 7.6.

## Exercises

**4.4.1.** [12] Show that the  $\mathbf{m}$ -average time complexity of Quicksort is  $\Omega(n^2)$ .

*Comments.* Source: M. Li and P.M.B. Vitányi, *Inform. Process. Lett.*, 42(1992), 145–149.

**4.4.2.** [12] Show that for each NP-complete problem, if the problem instances are distributed according to  $\mathbf{m}$ , then the average running time of any algorithm that solves it is superpolynomial unless  $P = NP$ .

*Comments.* Source: M. Li and P.M.B. Vitányi, *Ibid.*

## 4.5 Continuous Sample Space

### 4.5.1 Universal Enumerable Semimeasure

Is there a universal enumerable semimeasure in the continuous setting? In the discrete version we had only to satisfy that the probabilities summed to less than or equal one. Here we have to deal with the additional subadditive property. Let  $\mathcal{B} = \{0, 1\}$  be the set of basic elements.

The development of the theory for continuous semimeasures is quite similar to that for discrete semimeasures, except that the analogue of the Coding Theorem 4.3.3 does not hold. Let  $\mathcal{M}$  be a class of continuous semimeasures as in Definition 4.2.1. The definition of universal continuous semimeasure is analogous to Definition 4.3.2 for the discrete case.

**Definition 4.5.1** A semimeasure  $\mu_0$  is *universal* (or maximal) for  $\mathcal{M}$  if  $\mu_0 \in \mathcal{M}$ , and for all  $\mu \in \mathcal{M}$ , there exists a constant  $c > 0$  such that for all  $x \in \{0, 1\}^*$ , we have  $\mu_0(x) \geq c\mu(x)$ .

**Theorem 4.5.1** *There is a universal enumerable continuous semimeasure. We denote it by  $\mathbf{M}$ .*

**Proof.** We prove the theorem in two stages. In Stage 1 we show that the enumerable semimeasures can be effectively enumerated as

$$\mu_1, \mu_2, \dots$$

In Stage 2 we show that

$$\mu_0(x) = \sum_{n \geq 1} \alpha(n) \mu_n(x), \text{ with } \sum \alpha(n) \leq 1,$$

is a universal semimeasure. Stage 1 is broken up into two parts. In the first part we enumerate all enumerable functions; and in the second part we effectively change the enumerable functions to enumerable semimeasures, leaving the functions that were already semimeasures unchanged.

**STAGE 1** Let  $\psi_1, \psi_2, \dots$  be the effective enumeration of all enumerable continuous functions. Fix any  $\psi$  (we drop the subscript for notational convenience). Without loss of generality we can assume (as we have done before in similar cases) that we are actually dealing with rational-valued two argument partial recursive functions  $\phi(x, k) = p/q$  (rather  $\phi(\langle x, k \rangle) = \langle p, q \rangle$ ) such that for all  $x \in \mathcal{B}^*$ , for all  $k > 0$ ,

- if  $\phi(x, k)$  is defined, then for all  $y \leq x$  ( $\leq$  in the sense of the natural lexicographical length-increasing order on  $\mathcal{B}^*$ ),  $\phi(y, 1), \dots, \phi(y, k - 1)$  are all defined;
- $\phi(x, k + 1) \geq \phi(x, k)$ ;
- $\lim_{k \rightarrow \infty} \phi(x, k) = \psi(x)$ ;
- $\psi(x) > \phi(x, k)$  for every  $k$ . (This is achieved by replacing  $\phi(x, k)$  by  $\phi(x, k) := \phi(x, k)/(1 + 1/k)$ . This replacement affects neither the monotonicity of  $\phi$  nor the represented semimeasure—if any.)

Next we use each  $\phi$  associated with  $\psi$  to compute a semimeasure  $\mu$  by approximation from below. In the algorithm, at each stage of the computation the local variable  $\mu$  contains the current approximation to the function  $\mu$ . This is doable because the nonzero part of the approximation is always finite.

We describe a sequence of enumerable semimeasures  $\psi_k(x)$  computed from  $\phi(x, k)$  such that if  $\phi(x, k)$  represents an enumerable semimeasure  $\psi(x)$ , then  $\lim_{k \rightarrow \infty} \psi_k(x) = \psi(x)$ .

**Step 1** Initialize by setting  $\mu(x) := \psi_k(x) := 0$ , for all  $x$  in  $\mathcal{B}^*$  and  $k \in \mathcal{N}$ ; and set  $k := 0$ .

**Step 2** Set  $k := k + 1$ . Compute  $\phi(x, k)$  and set  $\psi_k(x) := \phi(x, k)$  for all  $x \in \mathcal{B}^k$ . {If the computation does not terminate, then  $\mu$  will not change any more and is trivially a semimeasure.}

**Step 3 For**  $i := k - 1, k - 2, \dots, 0$  **do**

**for** each  $x$  of length  $i$  **do**

search for the least  $K > k$  such that  $\phi(x, K) > \sum_{b \in \mathcal{B}} \psi_k(xb)$ ;

set  $\psi_k(x) := \phi(x, K)$ ;

**if**  $\psi_k(\epsilon) \leq 1$  **then**  $\mu := \psi_k$  **else** terminate.

{Step 3 tests whether the new values in Step 2 satisfy the semimeasure requirements. Note that if  $\psi$  is an enumerable semimeasure, then the  $K$ 's always exist. Namely, for each  $x$  we have  $\sum_{b \in \mathcal{B}} \psi_k(xb) < \sum_{b \in \mathcal{B}} \psi(xb) \leq \psi(x)$  and  $\lim_{k \rightarrow \infty} \phi(x, k) = \psi(x)$ .}

**Step 4 Go to Step 2.**

Since  $\phi$  represents  $\psi(x)$ , by monotonicity of  $\phi$  we have  $\psi_k(x) \geq \phi(x, k)$  for all  $x$  of length at most  $k$ , which implies  $\lim_{k \rightarrow \infty} \psi_k(x) = \psi(x)$ . If  $\psi$  is already a semimeasure, then  $\mu := \psi$  and the algorithm never finishes but continues to approximate  $\mu$  from below. If for some  $k$  and  $x \in \mathcal{B}^k$  the value of  $\phi(x, k)$  is undefined, then the values of  $\mu$  do not change any more even though the computation of  $\mu$  goes on forever. If the condition in Step 3 is violated, then the algorithm terminates, and the constructed  $\mu$  is a semimeasure—even a recursive one. Clearly, in all cases,  $\mu$  is an enumerable semimeasure.

The current construction was suggested by J. Tyszkiewicz [personal communication of April 1996]. It can be made to handle  $\mathcal{B} = \mathcal{N}$  if in the construction of  $\psi_k$  one considers and gives possibly nonzero measures only to sequences of length at most  $k$  and consisting of natural numbers  $\leq k$ .

Executing the above procedure on all functions in the list  $\phi_1, \phi_2, \dots$  yields an effective enumeration  $\mu_1, \mu_2, \dots$  of all enumerable semimeasures.

**STAGE 2** Let  $\alpha : \mathcal{N} \rightarrow \mathcal{R}$  be any enumerable function satisfying  $\alpha(n) > 0$  for all  $n$  and  $\sum_n \alpha(n) \leq 1$ . Define the function  $\mu_0$  from  $\mathcal{B}^*$  into  $[0, 1]$  as

$$\mu_0(x) = \sum_n \alpha(n) \mu_n(x).$$

We show that  $\mu_0$  is a universal enumerable semimeasure. The first condition in Definition 4.2.1 of being a semimeasure is satisfied since

$$\mu_0(\epsilon) = \sum_n \alpha(n) \mu_n(\epsilon) \leq \sum_n \alpha(n) \leq 1.$$

The second condition in Definition 4.2.1 of being a semimeasure is satisfied since, for all  $x$  in  $\mathcal{B}^*$

$$\mu_0(x) = \sum_n \alpha(n) \mu_n(x) \geq \sum_n \alpha(n) \sum_{b \in \mathcal{B}} \mu_n(xb) = \sum_{b \in \mathcal{B}} \mu_0(xb).$$

The function  $\mu_0$  is enumerable, since  $\mu_n(x)$  is enumerable in  $n$  and  $x$ . (Use the universal partial recursive function  $\phi_0$  and the construction above.)

Finally,  $\mu_0$  multiplicatively dominates each  $\mu_n$ , since  $\mu_0(x) \geq \alpha(n) \mu_n(x)$ . Therefore,  $\mu_0$  is a universal enumerable semimeasure. There are more than one such universal enumerable semimeasures. We fix a *reference* universal enumerable semimeasure  $\mu_0$  and denote it by  $\mathbf{M}$ .  $\square$

The universal enumerable semimeasure  $\mathbf{M}(x)$  captures the notion of a universal a priori probability needed for application to inductive reasoning (Chapter 5).

Above we can set  $\alpha(n) = 2^{-n}$ . But we can also choose  $\alpha(n) = 2^{-K(n)}$ . For  $\mu = \mu_n$  we can define  $K(\mu) = K(n)$ . Therefore,

$$\mathbf{M}(x) \geq 2^{-K(\mu)} \mu(x), \quad (4.11)$$

for all  $x \in \mathcal{B}^*$ .

At the risk of beating a dead horse (Example 4.3.1), we belabor the distinction between “continuous semimeasure” and “discrete semimeasure,” and the relation between  $\mathbf{M}$  and  $\mathbf{m}$ .

The discrete sample space theory is simply a restriction of the more sophisticated continuous approach we take in this section. Theorem 4.5.1 is a lifted version of Theorem 4.3.1. Namely, if we set  $\mathcal{B} = \mathcal{N}$  and restrict the arguments of the measure functions to sequences of natural numbers of length one, and incorporate the resulting simplifications in the proof of Theorem 4.5.1, then we obtain the proof of Theorem 4.3.1, and instead of  $\mathbf{M} : \mathcal{B}^* \rightarrow \mathcal{R}$ , we obtain its discrete version  $\mathbf{m} : \mathcal{B} \rightarrow \mathcal{R}$ .

**Lemma 4.5.1** *If a continuous enumerable semimeasure is a measure, it is recursive.*

**Proof.** Let  $\mu$  be an enumerable semimeasure with  $\sum_{b \in \{0,1\}} \mu(xb) = \mu(x)$  for all  $x \in \{0,1\}^*$  and  $\mu(\epsilon) = 1$ . Then, we can approximate all  $\mu(x)$  to any degree of precision starting with  $\mu(0), \mu(1)$  and determining  $\mu(x)$  for all  $x$  of length  $n$ , for consecutive  $n = 1, 2, \dots$ .  $\square$

**Lemma 4.5.2** *The set of recursive continuous semimeasures has no universal element.*

**Proof.** Set  $\mathcal{B} = \mathcal{N}$ . If there is a universal recursive continuous semimeasure, then its restriction to domain  $\mathcal{B}$  would by definition be a universal recursive discrete semimeasure, contradicting Lemma 4.3.1.  $\square$

**Lemma 4.5.3** *The function  $\mathbf{M}$  is not recursive and  $\mathbf{M}$  is not a probability measure.*

**Proof.** If  $\mathbf{M}$  is recursive, then it is universal for the class of recursive continuous semimeasures, by Theorem 4.5.1. This contradicts Lemma 4.5.2.

For  $\mathbf{M} : \mathcal{B}^* \rightarrow \mathcal{R}$  we prove  $\sum_{b \in \mathcal{B}} \mathbf{M}(b) < 1$  by the same proof of Lemma 4.3.2 (with  $\mathbf{M}$  instead of  $\mathbf{m}$ ).  $\square$

#### 4.5.2

#### A Priori Probability

As in the discrete case, we can interpret the enumerable semimeasures in a different way. To do so we require an appropriate new Turing machine variant.

**Definition 4.5.2**

*Monotone machines* are Turing machines with a one-way read-only input tape, some work tapes, and a one-way write-only output tape. The input tape contains a one-way infinite sequence of 0's and 1's and initially the input head scans the leftmost bit. The output tape is written one symbol in  $\mathcal{B}$  at the time, and the output is defined as the finite binary sequence on the output tape if the machine halts, and the possibly infinite sequence appearing on the output tape in a never-ending process if the machine does not halt at all. For a (possibly infinite) sequence  $x$  we write  $M(p) = x$  if  $M$  outputs  $x$  after reading  $p$  and no more. (Machine  $M$  either halts or computes forever without reading additional input.)

We define a sample space  $S_{\mathcal{B}}$  consisting of all finite and infinite sequences over  $\mathcal{B}$ :

$$S_{\mathcal{B}} = \mathcal{B}^* \bigcup \mathcal{B}^\infty.$$

**Definition 4.5.3**

Monotone machines compute partial functions  $\psi : \{0, 1\}^* \rightarrow S_{\mathcal{B}}$  such that for all  $p, q \in \{0, 1\}^*$  we have  $\psi(p)$  is a prefix of  $\psi(pq)$ . The function  $\psi$  induces a mapping  $\psi' : \{0, 1\}^\infty \rightarrow S_{\mathcal{B}}$  as follows: Let  $\omega = \omega_1\omega_2\dots \in \{0, 1\}^\infty$ .

**Case 1**  $\psi'(\omega)$  is the infinite sequence  $\zeta = \zeta_1\zeta_2\dots$  if for each  $n$ , there is an  $m$  such that  $\psi(\omega_{1:n}) = \zeta_{1:m}$ , and  $m$  goes to infinity with  $n$ .

**Case 2** If for some  $n$ ,  $\psi(\omega_{1:n})$  is defined, and for all  $m > n$  we have  $\psi(\omega_{1:m})$  equals  $\psi(\omega_{1:n})$ , then  $\psi'(\omega) = \psi(\omega_{1:n})$ . If for all  $n$ ,  $\psi(\omega_{1:n})$  is the empty word  $\epsilon$ , then  $\psi'(\omega) = \epsilon$ .

**Case 3** If there is an  $n$  such that  $\psi(\omega_{1:n})$  is undefined, then  $\psi'(\omega)$  is undefined.

We call such functions  $\psi'$  *monotone functions*. For convenience we drop the prime on the extension  $\psi'$  from now on.

**Definition 4.5.4** A monotone machine  $M$  maps subsets of  $\{0, 1\}^\infty$  to subsets of  $S_B$ . If  $\psi$  is the function computed by  $M$ , and  $A \subseteq \{0, 1\}^\infty$ , then define

$$\psi(A) = \{x \in S_B : \psi(\omega) = x, \omega \in \{0, 1\}^\infty\}.$$

A *cylinder* set  $\Gamma_x$  of  $S_B$  is defined as

$$\Gamma_x = \{x\omega : x \in B^*, \omega \in S_B\}.$$

Each semimeasure  $\mu$  is transformed by a monotone machine  $M$ , computing a monotone function  $\psi$ , to  $\mu_\psi$  (also a semimeasure) as follows: For each  $x \in B^*$ , let  $X \subseteq \{0, 1\}^*$  be the set of  $y$ 's such that  $\psi(y) = xz$  for some  $z \in S_B$ . Then  $M$  maps  $\bigcup_{y \in X} \Gamma_y \in \{0, 1\}^\infty$  to  $\Gamma_x \subseteq S_B$ . It is possible that for some  $y, z \in X$  we have  $\Gamma_y \cap \Gamma_z \neq \emptyset$ . This is the case precisely if  $y$  is a proper prefix of  $z$  or conversely. That is, either  $\Gamma_y$  is contained in  $\Gamma_z$  or vice versa. To obtain the total  $\mu$ -measure of  $\bigcup_{y \in X} \Gamma_y \in \{0, 1\}^\infty$  we sum of the  $\mu$ -measures of all the constituent cylinders that are not contained in other cosntituency cylinders. This is done by restricting  $X$  to a subset  $Y$  obtained from  $X$  by eliminating all strings that have a proper prefix in  $X$  and summing over the cylinders associated with elements in  $Y$ . The probability  $\mu_\psi(x)$  that  $M$  computes a sequence starting with  $x$  on  $\mu$ -random input from  $\{0, 1\}^\infty$  is given by

$$\mu_\psi(x) = \sum_{y \in Y} \mu(y). \quad (4.12)$$

Clearly, one can effectively enumerate all monotone machines  $M_1, M_2, \dots$  and therefore the associated monotone functions  $\psi_1, \psi_2, \dots$  they compute. We show that the corresponding enumeration  $\mu_{\psi_1}, \mu_{\psi_2}, \dots$ , with the  $\mu_\psi$ 's defined as in Equation 4.12, is an enumeration of all and only enumerable measures.

**Definition 4.5.5** A monotone function  $\psi$  is  $\mu$ -regular if the set of sequences  $\omega \in \{0, 1\}^\infty$  for which  $\psi'(\omega)$  is defined by Case 1 of Definition 4.5.3 has  $\mu$ -measure one. In other words, except for a set of negligible probability,  $\psi(\omega)$  is defined by Case 1.

**Lemma 4.5.4** (i) *For each recursive measure  $\mu$  and each  $\mu$ -regular monotone function  $\psi$ , we have  $\mu_\psi$  is a recursive measure as well.*

(ii) For each recursive measure  $\mu$  there exists a  $\lambda$ -regular  $\psi$ , with  $\lambda$  the uniform measure, such that  $\lambda_\psi = \mu$ ; moreover, there is a  $\mu$ -regular  $\phi$  such that  $\mu_\phi = \lambda$  and  $\phi$  is the inverse of  $\psi$  in the domain of definition of  $\phi\psi$ , and  $\phi(\omega) \in \mathcal{B}^\infty$  except for possibly the recursive  $\omega$ 's and  $\omega$ 's lying in countably many intervals of  $\mu$ -measure 0.

**Proof.** (i) For each  $x \in \mathcal{B}^*$  and  $n$  we must be able to approximate  $\mu_\psi(x)$  within accuracy  $2^{-n}$ . Choose  $m$  such that

$$\sum \{\mu(y) : l(y) = m, l(\psi(y)) > l(x)\} > 1 - 2^{-(n+1)}. \quad (4.13)$$

Such an  $m$  exists since  $\psi$  is  $\mu$ -regular, and it is easy to find such an  $m$  effectively.

Let  $Z \subseteq Y$  (with  $Y$  as in Equation 4.12) consisting of the  $y \in Y$  that additionally satisfy the condition in Equation 4.13. Since by assumption  $\mu$  is a recursive measure, the  $\mu(y)$ 's can be computed to within an accuracy of  $2^{-(m+n+1)}$ . Let  $\hat{\mu}(y)$  be such an approximation, and let  $\alpha(x, n) = \sum_{y \in Z} \hat{\mu}(y)$ . There are at most  $2^m$   $y$ 's satisfying the condition in Equation 4.13, and therefore

$$\begin{aligned} |\mu_\psi(x) - \alpha(x, n)| &< 2^{-(n+1)} + 2^m \cdot 2^{-(m+n+1)} \\ &= 2^{-n}. \end{aligned}$$

(ii) Omitted. □

For the less nice cases, namely, the enumerable semimeasures, the following property is easy to prove:

**Lemma 4.5.5** *If  $\mu$  is an enumerable semimeasure and  $\psi$  a monotone function, then  $\mu_\psi$  is again an enumerable semimeasure.*

**Proof.** By enumerating the semimeasure  $\mu$  (from below) for all its arguments  $x$ , and dovetailing the computations of  $\psi$  for all these arguments as well, we approximate  $\mu_\psi$  from below using Equation 4.12. Therefore,  $\mu_\psi$  is enumerable. It is straightforward to verify that  $\mu_\psi$  is also a semimeasure. □

The theorem below shows that  $\mu$  is an enumerable semimeasure iff it can be obtained as the semimeasure on the output sequences of some monotone machine where the input is supplied by independent tosses of a fair coin. In other words, each enumerable semimeasure  $\mu : \mathcal{B}^* \rightarrow \mathcal{R}$  can be obtained from some monotone function  $\psi$  and the uniform measure  $\lambda$ , in the sense of Equation 4.12:

$$\mu(x) = \lambda_\psi(x) = \sum_{y \in Y} 2^{-l(y)},$$

for the prefix-free set of programs  $y$  such that  $\psi(y)$  starts with  $x$ .

**Theorem 4.5.2** A semimeasure  $\mu$  is enumerable if and only if there is a monotone function  $\psi$  such that  $\mu = \lambda_\psi$ , where  $\lambda$  is the uniform measure.

**Proof. (IF)** By Lemma 4.5.5.

(ONLY IF) We construct a monotone function  $\psi$  such that  $\mu = \lambda_\psi$ , where  $\lambda(x) = 2^{-l(x)}$  is the uniform measure (which is recursive). For this construction, we have to decompose the interval  $[0, 1]$  into nonintersecting sets of measure  $\mu(x)$ . We will represent  $x \in \mathcal{B}^*$  by a set of intervals  $\Phi(x)$  of  $[0, 1]$  with the property that when  $x$  is a prefix of  $y$  then

$$\begin{aligned}\bigcup_{z \in \Phi(x)} \Gamma_z &\supseteq \bigcup_{z \in \Phi(y)} \Gamma_z, \\ \mu(x) &= \lambda \left( \bigcup_{z \in \Phi(x)} \Gamma_z \right).\end{aligned}$$

This can be achieved incrementally. Analogous to the construction in the proof of the Coding Theorem 4.3.3, since  $\mu$  is enumerable, there is a rational-valued recursive function  $\phi(x, k)$ , nondecreasing in  $k$  for fixed  $x$ , such that  $\lim_{k \rightarrow \infty} \phi(x, k) = \mu(x)$ . By definition,  $\mu(x) \geq \sum_{b \in \mathcal{B}} \mu(xb)$ . Without loss of generality, we can assume that  $\phi(x, k) \geq \sum_{b \in \mathcal{B}} \phi(xb, k)$ , for all  $t$ . (Whenever this inequality is not satisfied, we can decrease the  $\phi(xb, k)$ 's proportionally to the extent that the inequality becomes valid, without compromising the limiting behavior of  $\phi$  for fixed  $x$  and  $t$  going to infinity.) To obtain  $\mu(x)$  we approximate it by successive representations  $\Phi_1(x), \Phi_2(x), \dots$ , such that  $\Phi_k(x)$  is a prefix-free set with

$$\lambda \left( \bigcup_{z \in \Phi_k(x)} \Gamma_z \right) = \sum_{y \in \Phi_k(x)} 2^{-l(y)} = \phi(x, k),$$

satisfying the following: If  $x$  is a prefix of  $y$ , then

$$\bigcup_{z \in \Phi_k(x)} \Gamma_z \supseteq \bigcup_{z \in \Phi_k(y)} \Gamma_z.$$

If  $x$  and  $y$  are incomparable, then

$$\left( \bigcup_{z \in \Phi_k(x)} \Gamma_z \right) \cap \left( \bigcup_{z \in \Phi_k(y)} \Gamma_z \right) = \emptyset.$$

If  $k < k'$ , then

$$\bigcup_{z \in \Phi_{k'}(x)} \Gamma_z \supseteq \bigcup_{z \in \Phi_k(y)} \Gamma_z.$$

The construction of the  $\Phi_k$ 's is straightforward. Hence,  $\lambda_\psi(x) = \mu(x)$ .  $\square$

**Definition 4.5.6** Let  $U$  be the reference monotone machine. Denote the function computed by  $U$  by  $\psi$ . The *universal a priori probability* that a binary sequence starts with  $x$  is  $\lambda_\psi(x)$ , with  $\lambda_\psi$  in the sense of Equation 4.12 ( $\mu$  replaced by the uniform measure  $\lambda$ ).

It turns out that the universal a priori probability and the universal enumerable semimeasure are *equal*, just as in the discrete case (Theorem 4.3.3). If we provide  $U$ 's input by tosses of a fair coin, then the probability that  $U$ 's output starts with  $x$  is given by  $\mathbf{M}(x)$ . This  $\mathbf{M}$  is a central concept in this area.

**Theorem 4.5.3**  $-\log \lambda_U(x) = -\log \mathbf{M}(x) + O(1)$ .

**Proof.** This follows from the fact that for each monotone machine  $M$  in the effective enumeration  $M_1, M_2, \dots$ , we have that  $U(1^{n(M)}0p) = M(p)$ . Namely, this shows that  $\lambda_U(x) = 2^{-(n(M)+1)}\lambda_M(x)$ . The  $\lambda_M$ 's contain all enumerable semimeasures by Theorem 4.5.2. Since  $\lambda_U$  multiplicatively dominates each  $\lambda_M$ , it qualifies as the universal enumerable semimeasure  $\mathbf{M}$ .  $\square$

**Corollary 4.5.1** If the monotone machine  $T$  defines the enumerable semimeasure  $\mu$ , then  $\mathbf{M}(x) \geq 2^{-K(\mu)}\mu(x)$  with  $K(\mu)$  the shortest self-delimiting description of  $T$ .

### 4.5.3 \*Solomonoff Normalization

We have chosen to develop the theory essentially along lines reminiscent of Martin-Löf's theory of tests for randomness: we view it as mathematically elegant that the universal element, dominating all elements in a given class, belongs to that class. While our basic goal was to obtain a universal measure in the class of recursive measures, satisfying the above criteria turned out to be possible only by weakening both the notion of measure and the notion of effective computability. Another path was taken by R.J. Solomonoff. He viewed the notion of measure as sacrosanct. He normalized each semimeasure  $\mu$  to a measure using a particular transformation.

**Definition 4.5.7** The *Solomonoff normalization* of a semimeasure  $\mu$  on the sample space  $\mathcal{B}^\infty$  is defined by

$$\begin{aligned}\mu_{norm}(\epsilon) &= 1, \\ \mu_{norm}(\omega_{1:n}b) &= \mu(\omega_{1:n}b) \frac{\mu_{norm}(\omega_{1:n})}{\sum_{a \in \mathcal{B}} \mu(\omega_{1:n}a)},\end{aligned}$$

for all  $n \in \mathcal{N}$  and  $b \in \mathcal{B}$ .

This makes

$$\mu_{norm}(\omega_{1:n}) = \mu(\omega_{1:n}) \prod_{i=0}^{n-1} \frac{\mu(\omega_{1:i})}{\sum_{a \in \mathcal{B}} \mu(\omega_{1:i}a)}.$$

It is straightforward to verify that if  $\mu$  is a semimeasure, then  $\mu_{norm}$  is a measure. We call  $\mathbf{M}_{norm}$  the *Solomonoff measure*. It is at once clear that  $\mathbf{M}_{norm}$  dominates all enumerable semimeasures as well. This comes at the price that  $\mathbf{M}_{norm}$  is *not* enumerable itself by Lemmas 4.5.1, 4.5.2. Nonetheless, for many applications (as in Chapter 5) it may be important to have a universal *measure* available with most of the right properties, and the fact that it is not enumerable may not really bother us. Another possible objection is that there are more ways to normalize a semimeasure to a measure. So why choose this one? The choice for  $\mathbf{M}_{norm}$  is not unique. Solomonoff justifies this choice by his particular interest in the interpretation of  $\mathbf{M}$  as a priori probability: the measure  $\mathbf{M}(x)$  is the probability that the monotone reference machine outputs a sequence starting with  $x$  if its input is supplied by fair coin flips.

Suppose the output of the machine is known thus far, say  $x$ . One wants to know the relative probability that 1 rather than 0 will be the next symbol when we know there is a next symbol. This relative probability is  $P = \mathbf{M}(x1)/\mathbf{M}(x0)$ . Many, if not most, applications of probability involve ratios of this sort. For example, this is the case in computing conditional probabilities, relative probabilities of various scientific hypotheses, and mathematical decision theory. If one wants to divide up  $\mathbf{M}(xu)$  (the unnormalized probability that the machine does *not* print another symbol after emitting  $x$ ), then only the normalization above leaves the ratio  $P$  invariant. The ratio  $P$  is sacred since it happens to be given by the reference universal monotone machine. The justification for  $\mathbf{M}_{norm}$  lies in the fact that it is the unique measure such that  $\mathbf{M}_{norm}(x1)/\mathbf{M}_{norm}(x0) = \mathbf{M}(x1)/\mathbf{M}(x0)$ .

This normalization eliminates all probability concentrated on the finite sequences and divides all  $\mu(x)$ 's by the corresponding remaining sums. This way we obtain an object related to the greatest measure contained in the semimeasure. Another possibility would be to use the unnormalized conditional probability  $P_u = \mathbf{M}(x1)/\mathbf{M}(x)$  instead of  $P_{norm} = \mathbf{M}(x1)/(\mathbf{M}(x0) + \mathbf{M}(x1))$ , the normalized conditional probability. Using  $P_u$  is equivalent to considering  $xu$  to be a real possibility. In most (if not all) applications we have the *auxiliary information* that either  $x0$  or  $x1$  *did* occur, so the probability of  $xu$  is zero. In cases of this sort,  $P_{norm}$  is correct and  $P_u$  is incorrect as values for the conditional probability of 1, given that  $x$  has occurred. Solomonoff argues (somewhat weakly) that since such probability ratios *need* to be used in applications, we somehow have to find ways to deal with it, and we cannot just refuse to accept normalized probability because it is not enumerable. This tentative section is based on R.J. Solomonoff [*IEEE Trans. Inform. Theory*, IT-24(1978), 422–432], discussions with L.A. Levin, P. Gács, and Solomonoff;

and in good part on Solomonoff's arguments [Letter, October 4, 1991]. See also Exercise 4.5.6, page 301, and the discussion in the History and References Section 4.7.

#### 4.5.4 \*Monotone Complexity and a Coding Theorem

There are two possibilities to associate complexities with machines. The first possibility is to take the length of the shortest program, while the second possibility is to take the negative logarithm of the universal probability. In the discrete case, using prefix machines, these turned out to be the same by the Coding Theorem 4.3.3. In the continuous case, using monotone machines, it turns out they are different.

**Definition 4.5.8** The complexity  $KM$  is defined as

$$KM(x) = -\log \mathbf{M}(x).$$

In contrast with  $C$  and  $K$  complexities, in the above definition the greatest prefix-free subset of *all* programs that produce output starting with  $x$  on the reference monotone machine  $U$  are weighed.

**Definition 4.5.9** Let  $U$  be the reference monotone machine. The complexity  $Km$ , called *monotone complexity*, is defined as

$$Km(x) = \min\{l(p) : U(p) = x\omega, \omega \in S_B\}.$$

We omit the invariance theorems for  $KM$  complexity and  $Km$  complexity, stated and proven completely analogously to Theorems 2.1.1, 3.1.1 on pages 97, 193. By definition,  $KM(x) \leq Km(x)$ . In fact, all complexities coincide up to a logarithmic additive term, Section 4.5.5.

It follows directly from Definition 4.5.8 and Equation 4.11 that for each enumerable semimeasure  $\mu$  we have

$$KM(x) \leq -\log \mu(x) + K(\mu).$$

The following theorem can be viewed as a coding theorem, continuous version, for recursive measures. It states that for recursive measures  $\mu$ , each sample has a minimal description bounded by the logarithm of its probability plus a constant.

**Theorem 4.5.4** Let  $\mu : \{0, 1\}^* \rightarrow \mathcal{R}$  be a recursive measure. Then

$$Km(x) \leq -\log \mu(x) + c_\mu,$$

where  $c_\mu$  is a constant depending on  $\mu$  but not on  $x$ .

**Proof.** Recall the terminology in the proof of Lemma 4.3.3 on page 253. We give an inductive procedure for assigning code words to samples starting with  $x$  (rather, the cylinder  $\Gamma_x$  consisting of finite and infinite sequences starting with  $x$ ). Start by setting interval  $I_\epsilon := [0, 1)$  (of length  $\mu(\epsilon) = 1$ ).

Inductively, divide interval  $I_x$  into two half-open intervals  $I_{x0}$  (the left one) and interval  $I_{x1}$  (the right one), of lengths  $\mu(x0)$  and  $\mu(x1)$ , respectively. In each interval  $I_x$  determine the length of the largest binary interval. Since  $\mu$  is a real-valued recursive function, we cannot achieve that the lengths of the intervals  $I_{x0}$  and  $I_{x1}$  are exactly  $\mu(x0)$  and  $\mu(x1)$ , but we can achieve a good rational approximation (good enough that even the product of the error factors is bounded).

Let the binary string representing the leftmost such interval be  $r$ . Select as code for a sample starting with  $x$  the code word  $r$ . Note that samples starting with  $x$  and  $x0$  may be represented by the same code word. This is fine. According to this procedure, each sample of strings starting with  $x$  gets assigned a code word  $r$  with  $l(r) \leq -\log \mu(x) + 2$  (analogous to the proof of Lemma 4.3.3).

By construction the code is monotone: if  $r$  is a code word for a sample starting with  $x$ , then so is each code word starting with  $r$ . Since also  $\mu$  is recursive, there is a monotone machine  $M$  as follows: If  $r$  is a code word for a sample starting with  $x$ , then  $M$  recovers  $x$  as prefix of an output string  $x'$  (possibly longer than  $x$ ) by following the above procedure. For each prefix  $q$  of  $r$  it computes the longest  $y$  such that  $q$  is a code for a sample starting with  $y$ . Since the code is monotone,  $M$  can operate monotonically as well, by outputting subsequent digits of  $x'$  for each additional digit read from the one-way input  $r$ .

If  $n(M)$  is the index of  $M$  in the standard enumeration of monotone machines, then input  $\overline{n(M)}r$  is a program for the reference monotone machine  $U$  to output a string starting with  $x$ . Hence, the length  $Km(x)$  of the shortest program for  $U$  for outputting a string with prefix  $x$  satisfies  $Km(x) \leq -\log \mu(x) + c_\mu$  with  $c_\mu = 2n(M) + 2$ . Refinement of the argument lets us set  $c_\mu = K(M) + 2$ .  $\square$

Theorem 4.5.4 is the continuous analogue for the Coding Theorem 4.3.3 for discrete enumerable semimeasures. We know that  $KM(x) \leq Km(x) + O(1)$ . It has been shown that equality does not hold: the difference between  $KM(x)$  ( $= -\log M(x)$ ) and  $Km(x)$  is very small, but still rises unboundedly. This contrasts with the equality between  $-\log m(x)$  and  $K(x)$  in Theorem 4.3.3. Intuitively, this phenomenon is justified by exposing the relation between  $M$  and  $m$ .

The Coding Theorem 4.3.3 states that  $K(x) = -\log m(x) + O(1)$ . L.A. Levin [*Soviet Math. Dokl.*, 14(1973), 1413–1416] conjectured that the analogue would

hold for the unrestricted continuous version. But it has been shown that

$$\sup_{x \in \mathcal{B}^*} |KM(x) - Km(x)| = \infty,$$

[P. Gács, *Theoret. Comput. Sci.*, 22(1983), 71–93].

The exact relation is (for each particular choice of basis  $\mathcal{B}$  such as  $\mathcal{B} = \mathcal{N}$ , the natural numbers, or  $\mathcal{B} = \{0, 1\}$ )

$$KM(x) \leq Km(x) \leq KM(x) + Km(l(x)) + O(1). \quad (4.14)$$

The left-hand inequality follows from the definitions. The right-hand inequality can be understood as follows: First look at  $\mathcal{B} = \mathcal{N}$ . The Coding Theorem 4.3.3 can be proved for every recursively enumerable prefix-free subset of  $\mathcal{N}^*$ , and not only for  $\mathcal{N}$ . (It is enough to prove it for  $\mathcal{N}$ ; the extension to prefix-free sets follows by encoding.) That is, if  $F$  is a recursively enumerable prefix-free subset of  $\mathcal{N}^*$  (like  $\mathcal{N}^n$  or another “cut” through the prefix tree representation of  $\mathcal{N}^*$ ), then  $Km_F(x) = -\log \mathbf{m}_F(x)$ . (Here, the subscript  $F$  means that both  $Km$  and  $\mathbf{m}$  are defined with respect to an effective enumeration of monotone machines with programs in  $F$ .)

Now,  $Km(x) \leq Km_F(x) + Km(n(F)) + O(1)$ , with  $n(F)$  the index of the partial recursive function describing  $F$ . Since obviously,  $\mathbf{m}_F(x) \leq \mathbf{M}(x)$ , we have proved the right-hand inequality of Equation 4.14. Similar reasoning can improve the estimate to

$$KM(x) \leq Km(x) \leq KM(x) + Km(KM(x)) + O(1).$$

These results show that differences between  $Km(x)$  and  $KM(x)$  cannot exceed those accounted for by the tree structure of  $\mathcal{N}^*$ . The problem is equivalent for binary basis  $\mathcal{B} = \{0, 1\}$ . Indeed, the estimate is the best possible one, since the following can be shown:

**Claim 4.5.1** For each co-enumerable function  $g : \mathcal{N} \rightarrow \mathcal{N}$  for which

$$Km(x) - KM(x) \leq g(l(x)),$$

we have  $Km(n) \leq g(n) + O(1)$ .

**Proof.** It follows from  $KM(x) \leq Km(x) + O(1)$  that

$$\sum_n 2^{-Km(x)} \leq \sum_n \mathbf{M}(n) \leq \mathbf{M}(\epsilon) \leq 1.$$

Therefore,  $Km(n) \leq g(n) + O(1)$  is equivalent to saying that

$$\sum_n 2^{-g(n)} < \infty.$$

Conversely, if  $\sum_n 2^{-g(n)} < \infty$  then on the domain  $\mathcal{N}$  we have that  $2^{-g(n)}$  is multiplicatively dominated by  $\mathbf{M}(n)$ . Since on the domain  $\mathcal{N}$  we have  $KM(n) = Km(n) + O(1)$ , it follows that  $Km(n) \leq KM(n) + O(1) \leq g(n) + O(1)$  for  $n \in \mathcal{N}$ .  $\square$

This shows that the differences between  $Km(x)$  and  $KM(x)$  must in some sense be very small. The next question to ask is whether the quantities involved are usually different, or whether this is a rare occurrence. In other words, whether for almost all infinite sequences  $\omega$ , the difference between  $Km$  and  $KM$  is bounded by a constant. The following facts have been proven [P. Gács, *Theoret. Comput. Sci.*, 22(1983), 71–93].

- Lemma 4.5.6**
- (i) For random strings  $x \in \mathcal{B}^*$  we have  $Km(x) - KM(x) = O(1)$ .
  - (ii) There exists a function  $f(n)$  that goes to infinity as  $n \rightarrow \infty$  such that  $Km(x) - KM(x) \geq f(l(x))$ , for infinitely many  $x$ .
    - (a) If  $x$  is a finite string of natural numbers ( $\mathcal{B} = \mathbb{N}$ ), then we can choose  $f(n) = \log n$ .
    - (b) If  $x$  is a finite binary string ( $\mathcal{B} = \{0, 1\}$ ), then we can choose  $f(n)$  as the inverse of some version of Ackermann's function (Exercise 1.7.18, page 45).

That is, for almost all infinite  $\omega$  the difference  $Km(\omega_{1:n}) - KM(\omega_{1:n})$  has an upper bound that is smaller than any unbounded recursive function.

#### 4.5.5 \*Relation Between Complexities

There are five complexity variants discussed in this book: the plain complexity  $C(\cdot)$ , the complexity  $KM(\cdot) = -\log \mathbf{M}(\cdot)$  associated with the universal measure  $\mathbf{M}$ , the monotone complexity  $Km(\cdot)$ , Loveland's uniform complexity  $C(\cdot; \cdot)$  of Exercise 2.3.3 on page 124, and the prefix complexity  $K(\cdot)$ . What is the quantitative difference between them? Partially, such as between  $C(\cdot)$ ,  $l(\cdot)$ , and  $K(\cdot)$ , these relations follow from Chapters 2, 3. It is easy to see that  $Km(x) \leq l(x) + O(1)$ , and by definition  $KM(x) \leq Km(x)$ . Moreover,  $C(x; l(x)) \leq C(x) + O(1)$  by Exercise 2.3.3.

We include the table of relations, Figure 4.2, and omit the proofs. The table gives the differences resulting from subtraction of the complexities naming the rows from the complexities naming the columns. The upper entry in an intersection gives an upper bound on the differences of the

	$C(x)$	$KM(x)$	$Km(x)$	$K(x)$
$l(x)$	$\leq O(1)$	$\leq O(1)$	$\leq O(1)$	$\ell^k(\log x, \epsilon)$
				$\ell^k(\log x, 0)$
$C(x)$		$\ell^k(x, \epsilon)$	$\ell^k(x, \epsilon)$	$\ell^k(x, \epsilon)$
		$\ell^k(x, 0)$	$\ell^k(x, 0)$	$\ell^k(x, 0)$
$KM(x)$	$\log l(x) + O(1)$		$\ell^k(x, \epsilon)$	$\ell^k(x, \epsilon)$
	$\log l(x) - O(1)$		$A^{-1}(l(x))$	$\ell^k(x, 0)$
$Km(x)$	$\log l(x) + O(1)$	$< 0$		$\ell^k(x, \epsilon)$
	$\log l(x) - O(1)$	$< 0$		$\ell^k(x, 0)$
$C(x; l(x))$	$\ell^k(x, \epsilon)$	$\ell^k(x, \epsilon)$	$\ell^k(x, \epsilon)$	$\log l(x) + \ell^k(x, \epsilon)$
	$\ell^k(x, 0)$	$\ell^k(x, 0)$	$\ell^k(x, 0)$	$\log l(x) + \ell^k(x, 0)$

FIGURE 4.2. Relations between five complexities

complexities for *all*  $x$ . The lower entry gives a lower bound on the differences between the complexities for *infinitely many*  $x$ . These bounds hold for any fixed  $k$ . The function  $A^{-1}$  is a version of the very slow-growing inverse of the Ackermann function in Exercise 1.7.18. (Interestingly, while  $A$  grows faster than any primitive recursive function, its inverse  $A^{-1}$  is primitive recursive.) We denote  $\log x + \log \log x + \dots + (1 + \epsilon) \log^k x$  by  $\ell^k(x, \epsilon)$ , where  $\log^1 = \log$  and  $\log^k = \log \log^{k-1}$  for  $k > 1$ . For background material, see the History and References Section 4.7.

#### 4.5.6

#### \*Randomness by Integral Tests

The randomness of infinite sequences with respect to the uniform measure has been treated in Sections 2.5 and 3.6. In the latter section, Corollary 3.6.1, page 214, gives the following exact expression characterizing the infinite sequences that are random with respect to the uniform measure  $\lambda$ :

$$\rho_0(\omega|\lambda) = \sup_{n \in \mathcal{N}} \{n - K(\omega_{1:n})\}. \quad (4.15)$$

That is,  $\rho_0(\omega|\lambda) < \infty$  iff  $\omega$  is random with respect to  $\lambda$ . We generalize this result to exact expressions testing the randomness of infinite sequences for arbitrary recursive measures  $\mu$ .

A universal sequential  $\mu$ -test  $\delta_0(\cdot|\mu)$  of Section 2.5 additively majorizes all other sequential  $\mu$ -tests. It distinguishes the random infinite sequences from the nonrandom ones. However, we did not find an exact expression of the universal sequential  $\mu$ -test in terms of complexity. But we can develop other types of tests that do have exact expressions in terms of Kolmogorov complexity for a universal test separating the random infinite sequences from the nonrandom ones.

Sequential  $\mu$ -tests were essentially functions of a continuous variable in  $S = \mathcal{B}^\infty$ . Constructivity of such functions of infinite sequences was ensured by having a sequential test approximate its value by taking the supremum of all tests of finite initial segments of the infinite sequence.

We start by considering the notion of enumerable unit integrable functions of a continuous variable. It is easy to show that the logarithms of such functions are slightly stronger versions of sequential tests. This parallels the development of the sum tests of Definition 4.3.5, page 257. Subsequently, we give the exact expression of a universal enumerable unit integrable function with respect to a recursive measure  $\mu$ . The logarithm of this quantity is the exact expression for a universal test. This approach requires the application of a few elementary properties of integration.

Let  $\mathcal{B} = \{0, 1, \dots, k-1\}$  be a finite nonempty alphabet with  $k \geq 2$ . The continuous variable ranges over the set of one-way infinite sequences  $S = \mathcal{B}^\infty$ . The set  $S$  has the power of the continuum since it can be

mapped onto the set  $\mathcal{R}$  of real numbers, Example 1.4.1. The set  $X$  below is a subset of  $S$ .

**Definition 4.5.10** A nonnegative function  $f : S \rightarrow \mathcal{R}$  is *unit integrable* (over a set  $X$  with respect to measure  $\mu$ ) if

$$\int_X f(\omega) \mu(d\omega) \leq 1.$$

Consider recursive functions of the form  $g(x, k) = \langle r, s \rangle$ . The interpretation is that  $g$  has an argument  $x \in \mathcal{B}^*$  and a rational value  $r/s$ . The cylinder set  $\Gamma_x$  is the set of infinite sequences over  $\mathcal{B}$  starting with  $x$ .

**Definition 4.5.11** A nonnegative function  $f : S \rightarrow \mathcal{R}$  is *enumerable* if there exists a recursive function  $g(x, k)$  satisfying  $g(x, k+1) \geq g(x, k)$  and  $g(xy, k) \geq g(x, k)$  such that

$$f(\omega) = \sup_{\substack{\omega \in \Gamma_x \\ k \in \mathcal{N}}} \{g(x, k)\}.$$

(Furthermore,  $f$  is co-enumerable iff  $-f$  is enumerable, and  $f$  is recursive iff it is both enumerable and co-enumerable.) We determine the relation between an enumerable unit integrable function and a sequential test.

**Lemma 4.5.7** *If  $f$  is an enumerable unit integrable function over  $S$  with respect to a recursive measure  $\mu$ , then  $\log f(\cdot)$  is a sequential  $\mu$ -test. Moreover, if  $\delta$  is a sequential  $\mu$ -test, then the function  $f$  defined by  $\log f(\omega) = \delta(\omega) - 2 \log \delta(\omega) - c$  is an enumerable unit integrable function.*

**Proof.** Let  $f$  be an enumerable unit integrable function. Without loss of generality we can assume that the values  $f(\omega)$  are of the form  $2^m$ . Define

$$\gamma(x) = \inf_{\omega \in \Gamma_x} [\log f(\omega)].$$

Since  $f$  is enumerable, so is  $\gamma$ , even though its definition uses “inf,” by the way we defined enumerable functions. Also,  $\gamma(\omega_{1:n})$  is monotonic in  $n$ . It follows that  $\log f(\omega) = \sup_{n \in \mathcal{N}} \{\gamma(\omega_{1:n})\}$ . Moreover, for each  $n$ ,

$$\sum \{\mu(x) 2^{\gamma(x)} : l(x) = n, \gamma(x) > k\} \leq 2^{-k}.$$

Otherwise,

$$\int_S f(\omega) \mu(d\omega) > \sum_{l(x)=n} \mu(x) 2^{\gamma(x)} \geq \sum_{l(x)=n} \mu(x) 2^k \geq 1.$$

Hence,  $f$  is a sequential  $\mu$ -test according to Definition 2.5.1, page 141.

Conversely, assume that  $\delta$  is a sequential  $\mu$ -test. By Definition 2.5.1 the function  $\delta$  is enumerable. Therefore,  $f$  defined as in the lemma is enumerable. By Definition 2.5.1 we also have  $\mu\{\omega : \delta(\omega) \geq m\} \leq 2^{-m}$ , for each  $m$ . Define  $f(\omega) = c2^{\delta(\omega)} / \delta(\omega)^2$  for a constant  $c = (\pi^2/3)^{-1}$  and  $f(x) = \inf_{\omega \in \Gamma_x} \{f(\omega)\}$ . Since  $2^x/x^2$  is monotonic, these functions are also enumerable. Define

$$E_m = \{\omega : m \leq \delta(\omega) < m+1\}, \text{ and } c_m = \sup_{\omega \in E_m} \{f(\omega)\}.$$

Then (with  $m \geq 1$ ),

$$\begin{aligned} \int_S f(\omega) \mu(d\omega) &\leq \sum_m \mu(E_m) \\ &\leq c \sum_m \frac{2^{m+1}}{m^2} 2^{-m} \\ &\leq 2c \sum_m 1/m^2 \leq 1. \end{aligned}$$

□

**Definition 4.5.12** Let  $f$  be a unit integrable function over  $S$  with respect to  $\mu$ . A function  $\delta$  is an *integral  $\mu$ -test* iff  $\delta(\omega) = \log f(\omega)$ . It is a *universal integral  $\mu$ -test* if it additively dominates all integral  $\mu$ -tests.

Lemma 4.5.7 states that sequential  $\mu$ -tests and integral  $\mu$ -tests correspond up to a logarithmic additive term. It remains to show that there exists a universal integral  $\mu$ -test. We proceed by demonstrating the existence of a universal unit integrable function.

**Definition 4.5.13** A function  $f_0$  is *universal* for a class of unit integrable functions over domain  $X$  if  $f_0$  belongs to the class and for each  $f$  in the class there is a constant  $c$  such that for all  $\omega \in X$  we have  $cf_0(\omega) \geq f(\omega)$ .

**Theorem 4.5.5** *The class of enumerable unit integrable functions (over  $X$  with respect to a recursive measure  $\mu$ ) has a universal element, the function  $f_0$  defined by*

$$\log f_0(\omega) = \sup_{\omega \in \Gamma_x} \{-K(x|\mu) - \log \mu(x)\}.$$

**Proof.** We need to show that  $f_0$  is enumerable, unit integrable, and that it multiplicatively dominates all enumerable unit integrable functions. Since  $K(\cdot)$  is co-enumerable and  $\mu(\cdot)$  is recursive, it follows that  $f_0$  is enumerable.

**Claim 4.5.2** The function  $f_0$  is unit integrable.

**Proof.** First, write

$$f_0(\omega) = \sup_{\omega \in \Gamma_x} \left\{ 2^{-K(x|\mu) - \log \mu(x)} \right\}.$$

We have defined  $f_0$  only on elements of  $\mathcal{B}^\infty$ . If we want to define  $f_0$  on finite sequences  $x \in \mathcal{B}^*$ , then the natural way is

$$f_0(x) \stackrel{\text{def}}{=} \inf_{\omega \in \Gamma_x} \{f_0(\omega)\}.$$

This makes  $f_0(\omega_{1:n})$  monotonic nondecreasing in  $n$ . Let

$$g(x) = 2^{-K(x|\mu) - \log \mu(x)}.$$

Using Theorems 4.3.1 and 4.3.3, pages 247 and 253,

$$\begin{aligned} \int_X f_0(\omega) \mu(d\omega) &= \int_X \sup_n \{g(\omega_{1:n})\} \mu(d\omega) \\ &\leq \sum_n \int_X g(\omega_{1:n}) \mu(d\omega) \\ &= \sum_n \sum_{l(x)=n} g(x) \mu(x) \\ &= \sum_x g(x) \mu(x) = \sum_x 2^{-K(x|\mu)} \leq 1. \end{aligned}$$

□

**Claim 4.5.3** The function  $f_0$  multiplicatively dominates all enumerable unit integrable functions.

**Proof.** Let  $g$  be an enumerable unit integrable function. Without loss of generality, assume that  $g$  has values of the form  $2^m$  only. Enumerability of  $g$  is equivalent to saying that there exists a recursively enumerable set  $T$  of pairs  $(x, m)$  such that

$$g(\omega) = \sup \{2^m : \omega \in \Gamma_x \text{ and } (x, m) \in T\}.$$

For any set  $T$  of pairs  $(x, m)$  define the subset of elements “stabbed” by  $\omega$  as

$$T(\omega) \stackrel{\text{def}}{=} \{(x, m) \in T : \omega \in \Gamma_x\}.$$

The proof goes by replacing  $T$  by a set  $T'$  such that  $T'(\omega)$  contains at most one element  $(x, m)$  for each  $m$ . For this purpose, proceed as follows: Use  $T$  to define the recursively enumerable set  $T_m$  by

$$T_m = \{x : (x, m) \in T\}.$$

Each element  $x \in T_m$  is associated with a cylinder  $\Gamma_x \subseteq S$ . Similarly,  $T_m$  is associated with the union of all these cylinders:  $R = \bigcup_{x \in T_m} \Gamma_x$ . However, the constituent cylinders may overlap: some infinite sequences  $\omega$  have two different prefixes in  $T_m$ . By replacing  $T_m$  by a set  $T'_m$  in which no element is a prefix of another element, and such that  $\omega$  has a prefix in  $T_m$  iff it has a prefix in  $T'_m$ , we achieve that  $R$  equals the union of the nonoverlapping cylinders of elements in  $T'_m$ .

From the enumeration of  $T_m$  we obtain an enumeration of such a prefix-free set  $T'_m$  by the following “processing” step. Starting with an empty set  $T'_m$ , we put each enumerated element  $x$  from  $T_m$  into  $T'_m$  as long as  $x$  is not a prefix of an element already in  $T'_m$  and there is no element in  $T'_m$  that is a prefix of  $x$ . If  $T'_m$  contains a prefix of  $x$ , then we simply discard  $x$ . If  $x$  is the prefix of one or more elements in  $T'_m$ , then we replace  $x$  by the smallest prefix-free set  $A$  of elements of the form  $xy$  such that each infinite sequence starting with  $x$  has a prefix in  $T'_m \cup A$ . Then the latter set is prefix-free and we set  $T'_m := T'_m \cup A$ . Since  $T'_m$  is finite at each stage, each such set  $A$  is finite and can be effectively determined. We give a formal description of this “processing step”:

**Initialize:**  $T'_m := \emptyset$ .

**For** each enumerated  $x \in T_m$  **do**

**if**  $\Gamma_x \cap \Gamma_y = \emptyset$  for all  $y \in T'_m$  **then**  $T'_m := T'_m \cup \{x\}$

**else if**  $\Gamma_x \subseteq \Gamma_y$  for some  $y \in T'_m$  **then** discard  $x$

**else**  $T'_m := T'_m \cup A$

{where we choose  $A := \{xy_1, \dots, xy_k\}$  the smallest set such that:

1.  $\Gamma_z \cap \Gamma_{z'} = \emptyset$  for all  $z, z' \in T'_m \cup A$ ; and

2.  $\bigcup_{z \in T'_m \cup A} \Gamma_z = \bigcup_{z \in T'_m \cup \{x\}} \Gamma_z$ .

By construction,  $T'_m$  is recursively enumerable and prefix-free. Because it is prefix-free, the subset  $T'_m(\omega)$  stabbed by  $\omega$  is a singleton set or  $\emptyset$ . We are now ready to define the promised  $T'$  as  $T' = \bigcup_m \{(x, m) : x \in T'_m\}$ . The subset  $T'(\omega)$  stabbed by  $\omega$  contains at most one element  $(x, m)$  for each  $m$ .

By the construction above,

$$g(\omega) = \sup\{2^m : T'_m(\omega) \neq \emptyset\}.$$

Define further

$$g(x) \stackrel{\text{def}}{=} \sup_{x \in T'_m} \{2^m\}, \text{ and } g'(x) \stackrel{\text{def}}{=} \sum_{\{m : x \in T'_m\}} 2^m,$$

where  $g(x) = g'(x) = 0$  if there is no  $m$  with  $x \in T'_m$ . This way,

$$g(\omega) = \sup_{\omega \in \Gamma_x} \{g(x)\}, \text{ and } g(x) \leq g'(x). \quad (4.16)$$

Since each  $m$  occurs at most once in an element in  $T'(\omega)$ , we have

$$2g(\omega) \geq \sum_{T'_m(\omega) \neq \emptyset} 2^m = \sum_{\{x: \omega \in \Gamma_x\}} \sum_{\{m: x \in T'_m\}} 2^m = \sum_{\{x: \omega \in \Gamma_x\}} g'(x).$$

Therefore,

$$\begin{aligned} \sum_x g'(x) \mu(x) &= \sum_x \int_{\Gamma_x} g'(x) \mu(d\omega) \\ &= \int_X \sum_{\{x: \omega \in \Gamma_x\}} g'(x) \mu(d\omega) \\ &\leq 2 \int_X g(\omega) \mu(d\omega) \leq 2. \end{aligned}$$

Since  $g'(x) \mu(x)/2$  is enumerable and by the above sums to at most 1, it follows from Theorems 4.3.1 and 4.3.3, pages 247, 253, that there is a constant  $c$  such that

$$g'(x) \leq c \frac{2^{-K(x|\mu)}}{\mu(x)}.$$

Since  $g(x) \leq g'(x)$ , using Equation 4.16,

$$g(\omega) \leq \sup_{\omega \in \Gamma_x} \left\{ c 2^{-K(x|\mu) - \log \mu(x)} \right\}.$$

□

**Corollary 4.5.2** Let  $\mu$  be a recursive measure. The function

$$\rho_0(\omega|\mu) = \sup_{\omega \in \Gamma_x} \{-K(x|\mu) - \log \mu(x)\}$$

is a universal integral  $\mu$ -test.

**Example 4.5.1** With respect to the special case of the uniform distribution  $\lambda$ , Corollary 4.5.2 sets  $\rho_0(\omega|\lambda) = \sup_{n \in \mathcal{N}} \{n - K(\omega_{1:n})\}$  up to a constant additional term. This is the expression we found already in Corollary 3.6.1, page 214. ◇

**Example 4.5.2** To quantify the domination constants between  $f_0$  and the enumerable unit integrable functions we can proceed by an argument we have met several times before (for example, Theorem 4.5.1, page 272). First one shows that the enumerable unit integrable functions can be effectively enumerated as

$$f_1, f_2, \dots$$

Second, one shows that

$$f'_0(\omega) = \sum_{n \geq 1} \alpha(n) f_n(\omega), \text{ with } \sum \alpha(n) \leq 1,$$

is unit integrable. We can choose  $\alpha(n) = 2^{-K(n)}$ . Since the individual  $f_n$ 's are enumerable,  $f'_0$  is also enumerable. Since the individual  $f_n$ 's are unit integrable,  $f'_0$  is also unit integrable. Since  $f'_0(\omega) \geq 2^{-K(n)} f_n(\omega)$ , for all  $n$  and  $\omega$ , the function  $f'_0$  is a universal enumerable unit integrable function. By the above theorem,  $f_0$  is also a universal enumerable unit integrable function. Therefore, there is a constant  $c$  such that

$$c f_0(\omega) \geq f'_0(\omega) \geq 2^{-K(n)} f_n(\omega).$$

◇

**Example 4.5.3** Consider a discrete sample space  $S$  and the class of enumerable unit integrable functions over  $S$  with a recursive measure  $\mu$ . For convenience we set  $S \subseteq \mathcal{N}$ . Then  $\sum_{x \in S} f(x)\mu(x) \leq 1$ . For each such function  $f$  we have that  $f(x)\mu(x)$  is enumerable and sums up to at most 1. By Theorem 4.3.1 on page 247, there is a constant  $c$  such that  $f(x)\mu(x) \leq c \cdot \mathbf{m}(x|\mu)$  for all  $x$ . We have  $K(x|\mu) = -\log \mathbf{m}(x|\mu) + O(1)$ , by Theorem 4.3.3, page 253, and therefore

$$f(x) \leq c 2^{-K(x|\mu) - \log \mu(x)}.$$

Since  $\mathbf{m}$  is enumerable unit integrable, we find that the universal enumerable unit integrable function over  $S$  with respect to  $\mu$  is

$$f_0(x) = 2^{-K(x|\mu) - \log \mu(x)}.$$

The discrete sample space approach is connected to the continuous sample space treatment as follows: If we induce from  $f_0$  a function  $f'_0$  over  $\mathcal{B}^\infty$  by defining  $f'_0(\omega) = \sup_n \{f_0(\omega_{1:n})\}$ , then we obtain the function of Theorem 4.5.5 again. ◇

#### 4.5.7

#### \*Randomness by Martingale Tests

This section parallels the discussion in Sections 4.3.5, 4.3.6. We are presented with an infinite sequence of outcomes by an adversary who claims

that the distribution of outcomes is the recursive measure  $\mu$  on  $S = \mathcal{B}^\infty$ , where  $\mathcal{B}$  is a finite nonempty set of basic symbols. We would like to verify this claim under the assumption (guaranteed by the adversary) that the game is fair in the following sense: Given any function  $f$  such that for all  $n$ ,

$$\sum_{l(x)=n} f(x)\mu(x) \leq 1,$$

the adversary accepts odds of  $f(x)$  to one against outcome  $x$ .

We start out with the definition of a type of test embodying a betting strategy. Let  $x \in \mathcal{B}^*$  and  $b \in \mathcal{B}$ . You bet  $\gamma(x)$  units that the outcome  $x$  does not occur. The second inequality of Definition 4.5.14 says that if you continue playing another step after history  $x$ , then on average (using the conditional probability  $\mu(xb)/\mu(x)$ ), your payoff will not increase.

**Definition 4.5.14** Let  $\mu : \mathcal{B}^* \rightarrow \mathcal{R}$  be a recursive measure on the sample space  $\mathcal{B}^\infty$ . Let  $\gamma : \mathcal{B}^* \rightarrow \mathcal{R}$  be a nonnegative, enumerable function with the property

$$\begin{aligned}\mu(\epsilon)2^{\gamma(\epsilon)} &\leq 1, \\ \mu(x)2^{\gamma(x)} &\geq \sum_{b \in \mathcal{B}} \mu(xb)2^{\gamma(xb)}.\end{aligned}$$

Then  $\gamma$  is a *martingale  $\mu$ -test*. A martingale  $\mu$ -test is *universal* for a class of martingale  $\mu$ -tests if it additively dominates all martingale  $\mu$ -tests.

This test is slightly stronger than the original Martin-Löf test of Definition 2.4.1, page 129.

**Lemma 4.5.8** *Each martingale  $\mu$ -test is a  $\mu$ -test. If  $\delta(x)$  is a  $\mu$ -test, then there exists a constant  $c$  such that  $\delta(x) - 2 \log \delta(x) - c$  is a martingale  $\mu$ -test.*

**Proof.** It follows immediately from the new definition that for all  $n$ ,

$$\sum_{l(x)=n} \{\mu(x) : \gamma(x) > k\} \leq 2^{-k}. \quad (4.17)$$

Conversely, if  $\gamma(x)$  satisfies Equation 4.17, then for some constant  $c$  the function  $\gamma(x) - 2 \log \gamma(x) - c$  satisfies Definition 4.5.14.  $\square$

Thus, the universal (Martin-Löf)  $\mu$ -test of Definition 2.4.2 on page 130, the universal sum  $\mu$ -test of Definition 4.3.5 on page 257, and the universal martingale  $\mu$ -test all yield the same values up to a logarithmic additive term.

**Definition 4.5.15** A *sequential martingale  $\mu$ -test*  $\delta$  for  $\omega \in \mathcal{B}^\infty$  is obtained from a martingale  $\mu$ -test  $\gamma$  by defining

$$\delta(\omega) = \sup_{n \in \mathcal{N}} \{\gamma(\omega_{1:n})\}.$$

A sequential martingale  $\mu$ -test  $\sigma_0(\cdot|\mu)$  is *universal* if it additively dominates all sequential martingale  $\mu$ -tests  $\delta$ : there is a constant  $c$  such that for all  $\omega \in \mathcal{B}^\infty$  we have  $c \cdot \sigma_0(\omega|\mu) \geq \delta(\omega)$ .

**Theorem 4.5.6** Let  $\mu : \mathcal{B}^* \rightarrow \mathcal{R}$  be a recursive measure on  $\mathcal{B}^\infty$ .

- (i) The function  $\gamma_0 : \mathcal{B}^* \rightarrow \mathcal{R}$  defined by  $\gamma_0(x|\mu) = \log(\mathbf{M}(x)/\mu(x))$  is a universal martingale  $\mu$ -test.
- (ii) The function  $\sigma_0 : \mathcal{B}^\infty \rightarrow \mathcal{R}$  defined by

$$\sigma_0(\omega|\mu) = \sup_n \{\log(\mathbf{M}(\omega_{1:n})/\mu(\omega_{1:n}))\}$$

is a universal sequential martingale  $\mu$ -test .

**Proof.** Let  $V_k$  be the set of infinite sequences  $\omega$  such that there is an index  $n$  with  $k < -\log \mu(\omega_{1:n}) - Km(\omega_{1:n})$ . Using  $-\log \mathbf{M}(\omega_{1:n}) \leq Km(\omega_{1:n})$ , this implies that

$$\gamma_0(\omega_{1:n}|\mu) = \log \frac{\mathbf{M}(\omega_{1:n})}{\mu(\omega_{1:n})} > k,$$

for all  $\omega$  in  $V_k$ . The left-hand side of the inequality satisfies Definition 4.5.14. By Equation 4.17, therefore,  $\mu(V_k) \leq 2^{-k}$ . Since  $\mathbf{M}$  is enumerable and  $\mu$  is recursive,  $\gamma_0$  is enumerable. Hence it is a martingale  $\mu$ -test, and  $\sigma_0$  is a sequential martingale  $\mu$ -test. By Definition 4.5.14 the function  $\mu(x)2^{\gamma(x)}$  is an enumerable semimeasure. Thus, for each martingale  $\mu$ -test  $\gamma$ , there is a positive constant  $c_\gamma$  such that

$$\mathbf{M}(x) \geq c_\gamma \mu(x)2^{\gamma(x)}.$$

Therefore,  $\gamma_0$  and  $\sigma_0$  additively dominate all martingale  $\mu$ -tests and sequential martingale  $\mu$ -tests, respectively.  $\square$

As before, we call  $\gamma_0(x|\mu)$  the *randomness deficiency* of  $x$ . Theorem 4.5.6 yields yet another characterization of an infinite random sequence, equivalent to Martin-Löf's characterization of randomness in Section 2.5. Namely, for each such  $\mu$ , an infinite binary sequence  $\omega$  is  $\mu$ -random iff

$$\sigma_0(\omega|\mu) < \infty.$$

By the proof of Theorem 4.5.6 this is equivalent to

$$\sup_n \{-\log \mu(\omega_{1:n}) - Km(\omega_{1:n})\} < \infty.$$

**Corollary 4.5.3** An infinite sequence  $\omega$  is random iff  $\sup_n \{-\log \mu(\omega_{1:n}) - Km(\omega_{1:n})\} < \infty$ . (The expression on the left side is yet another  $\mu$ -randomness test.)

Altogether, we have found the following characterizations in this section. An infinite sequence  $\omega$  is  $\mu$ -random iff

$$KM(\omega_{1:n}) = Km(\omega_{1:n}) + O(1) = -\log \mu(\omega_{1:n}) + O(1).$$

(The  $O(1)$  term is independent of  $n$  but may depend on  $\omega$ .) The  $\mu$ -random infinite sequences by definition have  $\mu$ -measure one in the set of all infinite sequences. Such sequences, and only such sequences, satisfy all “effective” laws of probability theory (that is, withstand any sequential  $\mu$ -test in Martin-Löf’s sense). There is also a conditional version of Theorem 4.5.6:

**Corollary 4.5.4** Let  $f(\omega, \zeta)$  be a function with  $\int_S 2^{-f(\omega, \zeta)} \lambda(d\omega) \leq 1$ , where  $\lambda$  is the uniform measure. Such functions  $f$  define conditional measures as follows: Define  $f(x, \zeta) = \sup_{\omega \in \Gamma_x} \{f(\omega, \zeta)\}$ . Then,

$$\mu(x|\zeta) = \int_{\Gamma_x} 2^{-f(\omega, \zeta)} \lambda(d\omega).$$

If  $f$  is co-enumerable (the corresponding  $\mu(\omega|\zeta)$  is enumerable), then there is a constant  $c$  such that  $c \cdot M(\omega|\zeta) \geq 2^{-f(x, \zeta)}$ , for all  $\zeta$  and all  $x$  with  $\omega \in \Gamma_x$ .

**Example 4.5.4** Let  $\mu = \lambda$ , the uniform measure. The above discussion says that an infinite binary sequence  $\omega$  is random with respect to the uniform measure (that is, in the original sense of Martin-Löf) iff

$$KM(\omega_{1:n}) = Km(\omega_{1:n}) + O(1) = n + O(1).$$

◇

#### 4.5.8

#### \*Randomness by Martingales

In gambling, the “martingale” is a well-known system of play. With this system one bets in a casino that the outcome of a roulette run will be red each time, and after each loss the new bet is double the old one. No matter how long the run of losses, once red comes up the gambler wins. The current gain minus past loss in this run equals the amount of the first bet. But even though red is bound to come up sometime, the gambler may go broke and lose all before this happens. In his 1853 book *The Newcomes*, W.M. Thackeray admonishes: “You have not played as yet? Do not do so; above all avoid the martingale if you do.”

According to von Mises’ Definition 1.9.1 on page 51, a random sequence in the sense of collective defies a player betting at fixed odds, and in

fixed amounts, on the tosses of a fair coin, to make unbounded gain in the long run. Von Mises does not require that the player can have no debt. Obviously, if the player cannot go broke, and he bets according to a martingale, his gain will *eventually* exceed each bound.

We can test for randomness by betting. Suppose a casino claims that the distribution of outcomes  $\omega$  in sample space  $S = \mathcal{B}^\infty$  is the measure  $\mu$ . Then given any function  $f(\omega)$ , with  $\int_S f(\omega) \mu(d\omega) < 1$ , the casino should accept 1 unit for an obligation to pay  $f(\omega)$  on outcome  $\omega$ . (We have called such functions “unit integrable” in Definition 4.5.10 on page 287.) A martingale is a payoff function that leads to such a global payoff.

**Definition 4.5.16** Let  $\mu : \mathcal{B}^* \rightarrow [0, 1)$  be a measure on  $S = \mathcal{B}^\infty$ . Let  $t$  be a nonnegative function from  $\mathcal{B}^*$  into  $\mathcal{R}$  such that, with  $x \in \mathcal{B}^*$ ,

$$\begin{aligned}\mu(\epsilon)t(\epsilon) &\leq 1, \\ \mu(x)t(x) &\geq \sum_{b \in \mathcal{B}} \mu(xb)t(xb).\end{aligned}$$

Let  $\omega$  be an element of  $S$ . We call  $t(\omega_{1:n})$  a  $\mu$ -submartingale. We call  $t(\omega_{1:n})$  a  $\mu$ -martingale if the equations hold with equality.

**Definition 4.5.17** An enumerable  $\mu$ -submartingale  $t_0$  is *universal* for the class of enumerable  $\mu$ -submartingales if it multiplicatively dominates all enumerable  $\mu$ -submartingales  $t'$ . (That is, there is a constant  $c$  such that  $c \cdot t_0(x) \geq t'(x)$  for all  $x \in \mathcal{B}^*$ .)

**Theorem 4.5.7** *There is a universal enumerable  $\mu$ -submartingale. We denote it by  $t_0(\cdot|\mu)$ .*

**Proof.** Define  $t_0(x|\mu) = \mathbf{M}(x)/\mu(x)$ . Then  $t_0(x|\mu) = 2^{\gamma_0(x|\mu)}$  with  $\gamma_0(\cdot|\mu)$  as in Theorem 4.5.6. From Definition 4.5.15 and Definition 4.5.16, it follows that  $\gamma(x)$  is a martingale  $\mu$ -test iff  $t(x) = 2^{\gamma(x)}$  is an enumerable  $\mu$ -submartingale. Hence, the theorem follows from Theorem 4.5.6.  $\square$

**Example 4.5.5** We give an alternative proof of Theorem 4.5.7. Comparison with Definition 4.2.1 shows that  $t(x)$  is an enumerable  $\mu$ -submartingale iff  $t(x)\mu(x)$  is an enumerable semimeasure. By Theorem 4.5.1 on page 272, the semimeasure  $\mathbf{M}$  is a universal enumerable semimeasure. Hence, for each enumerable  $\mu$ -submartingale  $t$  there is a constant  $c$  such that  $c \cdot \mathbf{M}(x) \geq t(x)\mu(x)$ , for all  $x$ . Therefore,

$$t_0(x|\mu) = \frac{\mathbf{M}(x)}{\mu(x)}$$

dominates all enumerable  $\mu$ -submartingales within a multiplicative constant. Thus  $t_0(x|\mu)$  is a universal enumerable  $\mu$ -submartingale.  $\diamond$

Intuition tells us to call an outcome “nonrandom” if it allows us to win against the adversary by choosing an appropriate payoff function. That is, if there is an enumerable submartingale  $t$  such that  $t(\omega_{1:n})$  grows unboundedly. From the definitions it follows immediately that this characterization of random infinite sequences is exactly the same as the one using sequential  $\mu$ -tests. Since  $\gamma_0(x|\mu)$  dominates all  $\mu$ -tests, within additive constants, it follows that  $t_0(x|\mu) = \mathbf{M}(x)/\mu(x)$  dominates all enumerable  $\mu$ -submartingales within a multiplicative constant. Since  $\gamma_0$  is a universal martingale  $\mu$ -test, we have the following corollary to Theorem 4.5.7:

**Corollary 4.5.5** An infinite sequence  $\omega$  is  $\mu$ -random iff

$$t_0(\omega|\mu) = \sup_{n \in \mathcal{N}} \{t_0(\omega_{1:n}|\mu)\} < \infty.$$

We recall that the set of such  $\omega$  has  $\mu$ -measure one.

#### 4.5.9 \*Relations Between Tests

There is a simple relation between submartingales and unit integrable functions. The submartingale condition of Definition 4.5.16 implies that for a  $\mu$ -submartingale  $t$ , for all prefix-free sets  $I \subseteq \mathcal{B}^*$  we have

$$\sum_{x \in I} \mu(x)t(x) \leq 1.$$

This implies that  $t$  is a unit integrable function in the sense of Definition 4.5.10.

Unit integrability of a function  $f$  does not imply that it is a submartingale. A submartingale  $t$  is not necessarily monotonic nondecreasing but instead satisfies the restrictive submartingale condition. We can interpret the difference as saying that the unit integrable functions represent betting strategies in a fair game where the expectation of profit in the overall infinite game is at most 1. The submartingale condition says that the expectation of profit for a fair game of some finite length cannot be increased by playing longer. The condition on prefix-free sets states that the expectation with an arbitrary set of stopping times is still at most 1.

Assume a function  $f$  is unit integrable. The integral can be expressed as a sum as follows: Taking the supremum over all prefix-free sets  $I \subseteq Y$  we have

$$\int_X f_0(\omega)\mu(d\omega) = \sup_I \left\{ \sum_{x \in I} \mu(x)f_0(x) \right\}.$$

Then we can increase  $f$  to a  $\mu$ -submartingale  $t$  by defining

$$t(x) = \sup_I \left\{ \frac{1}{\mu(x)} \sum_{y \in I} f(xy) \mu(xy) \right\},$$

where  $I$  ranges over all prefix-free subsets of  $\mathcal{B}^*$ . The proof that  $t$  is a  $\mu$ -submartingale is by simply writing out both sides of the submartingale condition of Definition 4.5.16 in terms of suprema and verifying the required relations. The fact that we increase  $f$  to obtain  $t$  accounts for the fact that the universal enumerable unit integrable function  $f_0$  is slightly smaller than the universal  $\mu$ -submartingale  $t_0(\cdot|\mu)$ . For fixed  $\omega$ ,

$$\sup_{\omega \in \Gamma_x} \left\{ 2^{-K(x|\mu) - \log \mu(x)} \right\} \leq \sup_{\omega \in \Gamma_x} \{ \mathbf{M}(x)/\mu(x) \}.$$

To show that a submartingale can be used to define a function over  $\mathcal{B}^\infty$  we use the so-called Submartingale Convergence Theorem:

- Claim 4.5.4** Consider a sequence of random variables  $\omega_1, \omega_2, \dots$ . If  $t(\omega_{1:n})$  is a  $\mu$ -submartingale and the  $\mu$ -expectation  $\mathbf{E}|t(\omega_{1:n})| < \infty$ , then it follows that  $\lim_{n \rightarrow \infty} t(\omega_{1:n})$  exists with  $\mu$ -probability one.

**Proof.** See J.L. Doob, *Stochastic Processes*, Wiley, 1953, pp. 324–325.  $\square$

This implies that each  $\mu$ -submartingale  $t(\omega_{1:n})$  converges  $\mu$ -almost everywhere to a function  $t(\omega)$ . However, this function may not be enumerable. We obtain enumerability if the function  $t(\omega_{1:n})$  is monotonic in  $n$ .

The relation between the universal enumerable submartingale and the universal enumerable unit integrable function is as follows. The function

$$\rho_0(\omega|\mu) = \sup_n \{-K(\omega_{1:n}|\mu) - \log \mu(\omega_{1:n})\}$$

is a universal integral  $\mu$ -test, and

$$\sigma_0(\omega|\mu) = \sup_n \{\log \mathbf{M}(\omega_{1:n})/\mu(\omega_{1:n})\}$$

is a universal martingale  $\mu$ -test. Hence they are either both finite (for  $\mu$ -random  $\omega$ ) or both infinite (otherwise).

- Lemma 4.5.9** For each infinite sequence  $\omega$ , we have up to fixed additive constants

$$\sigma_0(\omega|\mu) - 2 \log \sigma_0(\omega|\mu) \leq \rho_0(\omega|\mu) \leq \sigma_0(\omega|\mu).$$

**Proof.** The right inequality follows from the fact that we have to increase the universal unit integrable function to make it a universal submartingale.

The left inequality is proved as follows: Let  $\gamma^m(\omega) = m$  if  $\omega$  has a prefix  $x$  with  $\log(\mathbf{M}(x)/\mu(x)) \geq m$  and 0 otherwise. The function  $\sigma(\omega) = \sup_m \{\gamma^m(\omega) - 2 \log m\}$  satisfies  $\sigma(\omega) \leq \rho_0(\omega|\mu) + O(1)$  since

$$2^{\sigma(\omega)} = \max_m \{2^{\gamma^m(\omega)}/m^2\} = \sum_m 2^{\gamma^m(\omega)}/m^2$$

is  $(\pi^2/6)$  integrable (instead of unit integrable) over  $X$  with respect to  $\mu$ , as is easy to verify. Clearly,  $\sigma_0(\omega|\mu) - 2 \log \sigma_0(\omega|\mu) \leq \sigma(\omega)$ , from which the inequality follows.  $\square$

We can similarly analyze the relation between the universal sequential  $\mu$ -test  $\delta_0(\cdot|\mu)$  and the other tests.

**Lemma 4.5.10** *For each infinite sequence  $\omega$ , we have up to fixed additive constants*

$$\delta_0(\omega|\mu) - 2 \log \delta_0(\omega|\mu) \leq \sigma_0(\omega|\mu) \leq \delta_0(\omega|\mu).$$

**Proof.** This follows from the definitions in Theorem 4.5.6, page 294, and Lemma 4.5.8, page 293.  $\square$

## Exercises

**4.5.1.** [09] Show that with basis  $\mathcal{B} = \{0, 1\}$ , we have  $\mathbf{M}(\epsilon) < 1$  and  $\mathbf{M}(x) > \mathbf{M}(x0) + \mathbf{M}(x1)$ , for all  $x$  in  $\mathcal{B}^*$  (strict inequality). Generalize this to arbitrary  $\mathcal{B}$ .

**4.5.2.** [31] Let the probability that an initial segment  $x$  of a binary sequence is followed by  $a \in \{0, 1\}^*$  be  $\mathbf{M}(a|x) = \mathbf{M}(xa)/\mathbf{M}(x)$ .

- (a) Show that there is a constant  $c > 0$  such that the probability (with respect to  $\mathbf{M}$ ) of the next bit being 0 after  $0^n$  is at least  $c$ .
- (b) Show that there exists a constant  $c$  such that the probability (with respect to  $\mathbf{M}$ ) of the next bit being 1 after  $0^n$  is at least  $1/(cn \log^2 n)$ .
- (c) Show that for every constant  $c$  and sufficiently large  $N$ , there are at most  $N/c$  initial segments  $0^n$  ( $1 \leq n \leq N$ ) such that the probability (with respect to  $\mathbf{M}$ ) of the next bit being 1 is larger than  $(c \log^2 n)/n$ .
- (d) Conclude that the probability (with respect to  $\mathbf{M}$ ) of the next bit following  $0^n$  being 1 is approximately  $\Theta(1/n)$ .

*Comments.* This exercise is a simple form of Occam's razor: the conditional  $\mathbf{M}$  probability assigns high probability to the simple explanations and low probability to the complex explanations. The assertion Item (d)

is an  $\mathbf{M}$  prior probability variant of P.S. Laplace's well-known exercise to compute the expectation of a successful trial following  $n$  successful trials in a Bernoulli process  $(p, 1 - p)$  with unknown  $p$ . Using Bayes's Rule with uniform prior probability, this expectation is  $(n + 1)/(n + 2)$  and is a special case of Laplace's Law of Succession (Exercise 1.10.1, page 63). Application of this law sets the probability of a 1 following  $n$  initial 0's at  $1/(n+2)$ . This is pretty close to the order of magnitude  $1/n$  approximation we found in Item (d) using conditional  $\mathbf{M}$  probability. In *A Philosophical Essay on Probabilities*, Laplace uses this rule to compute the probability that the sun will not rise tomorrow, given that it has been rising every morning since the creation of the world 10,000 years ago. Using the above it follows that the probability that the sun will not rise tomorrow is approximately  $1/3,650,000$ . This is correct, in case our information about the sun were exhausted by the fact stated. Hint for Items (a) through (c): use  $|K(x) - (-\log \mathbf{M}(x))| \leq 2 \log K(x) + O(1)$ . For all  $n$  we have  $K(0^n 1) \leq \log n + 2 \log \log n + O(1)$ , and for the majority of  $n$  we have  $K(0^n 1) \geq \log n$ . Source: L.A. Levin and A.K. Zvonkin, *Russ. Math. Surveys*, 25:4(1970), 83–124; see also T.M. Cover and J.A. Thomas, *Elements of Information Theory*, Wiley & Sons, New York, 1991.

**4.5.3.** [27] Even the most common measures can be nonenumerable if the parameters are nonenumerable. Consider a  $(p, 1 - p)$  Bernoulli process and the measure it induces on the sample space  $\{0, 1\}^\infty$ .

- (a) Show that if  $p$  is a recursive real number such as  $\frac{1}{2}$  or  $\frac{1}{3}$  or  $\frac{1}{2}\sqrt{2}$  or  $\pi/4 = \frac{1}{4}3.1415\dots$ , then the measure is recursive.
- (b) Show that if  $p$  is enumerable then the measure can be nonenumerable.
- (c) Show that if  $p$  is a real that is not enumerable, then the measure is not enumerable either.
- (d) Show that if  $p$  is a random real whose successive digits in its binary expansion are obtained by tosses of a fair coin, then with probability one the measure is not enumerable.

*Comments* Hint for Item (b): if  $p$  is enumerable but not recursive, then  $1 - p$  is not enumerable and it is the probability that the first element in the sequence is 0 [P. Gács, personal communication].

**4.5.4.** • [25] The *Solomonoff normalization* of a semimeasure  $\mu$ , with  $\mathcal{B} = \{0, 1\}$ , is  $\mu_{norm}(x) = a(x)\mu(x)$ , with  $a(x)$  as defined in Definition 4.5.7 in Section 4.3.1. We call  $\mathbf{M}_{norm}$ , the *normalized* version of the universal enumerable semimeasure  $\mathbf{M}$ , the *Solomonoff measure*. This leads to an important parallel development of the theory, which may be mathematically less elegant, but is possibly preferable in some applications.

(a) Show that for each semimeasure  $\mu$ , the function  $\mu_{norm}$  is a measure. Conclude that in particular,  $\mathbf{M}_{norm}$  is a measure.

(b) Show that the normalization factor  $a(x) \geq 1$  for all  $x$  in  $\mathcal{B}^*$ .

(c) Show that  $\mathbf{M}_{norm}$  dominates all enumerable semimeasures  $\mu$ . (That is, for each such  $\mu$ , there is a positive constant  $c$ , for all  $x \in \mathcal{B}^*$ , we have  $\mathbf{M}_{norm}(x) \geq c\mu(x)$ .) Conclude that  $\mathbf{M}_{norm}$  dominates  $\mathbf{M}$ .

(d) Show that Solomonoff normalization is not the only normalization.

(e) Let  $\mu$  be an enumerable measure. Does  $\mathbf{M}_{norm}$  dominate all  $\mu_{norm}$ 's?

*Comments.* Hint for Item (e): not all  $\mu_{norm}$ 's are enumerable, so we cannot use Item (c). This elaborates the discussion in Section 4.5.3. Source: Items (a) through (d) R.J. Solomonoff [*IEEE Trans. Inform. Theory*, IT-24(1978), 422–432]; see also History and References Section 4.7.

**4.5.5.** [32] By Exercise 4.5.4,  $\mathbf{M}_{norm}$  dominates  $\mathbf{M}$ .

(a) Show that  $\mathbf{M}$  does not multiplicatively dominate  $\mathbf{M}_{norm}$ .

(b) Show that for each normalizer  $a$  defining measure  $\mathbf{M}'(x) = a(x)\mathbf{M}(x)$  we have  $\mathbf{M}(\omega_{1:n}) = o(\mathbf{M}'(\omega_{1:n}))$ , for some infinite  $\omega$ .

(c) (Open) Item (b) with “all” substituted for “some.”

*Comments.* Item (b) implies that  $\mathbf{M}$  does not dominate any of its normalized versions  $\mathbf{M}'$ . This is a special case of Exercise 4.5.6.

**4.5.6.** • [37] What is the difference between semimeasure  $\mathbf{M}$  and *any* (not necessarily enumerable) measure  $\mu$ ?

(a) Show that  $1/\mathbf{M}(x)$  differs from  $1/\mu(x)$ , for some  $x$ , as the Busy Beaver function  $BB(n)$  from  $n$  (Exercise 1.7.19) for any measure  $\mu$ .

(b) Show the same about  $\mathbf{M}(x)$  versus  $\mathbf{M}(x_0) + \mathbf{M}(x_1)$ .

*Comments.* Hint: take the maximal running time  $T(k)$  of those among the first  $k$  programs that halt. Note that  $T(k)$  is the longest running time of a halting program among the first  $k$  programs and is related to  $BB(k)$ . The following procedure effectively generates the desired  $x$  from  $k$  (even though  $T$  is only enumerable) such that  $\mathbf{M}(x_0) + \mathbf{M}(x_1) < 1/T(k)$ . Start with the empty prefix  $y = \epsilon$  of  $x$ . Find its extension  $y' = y0$  or  $y' = y1$  (whichever you find first), such that  $\mathbf{M}(y') > 1/(2T(k))$ . Extend  $y$  to differ from  $y'$ . Repeat as long as we can (up to  $2T(k)$  times). The resulting string  $x$  (of length  $< 2T(k)$ ) will have both extensions of  $\mathbf{M}$ -semimeasure  $< 1/(2T(k))$ . Since  $x$  is described by  $k$  and an  $O(1)$  program, for large enough random  $k$  we have  $1/k^2 \leq \mathbf{M}(x)$ . See also History and References Section 4.7. Source: one of us (PV) asked L.A. Levin; Levin asked R.M. Solovay, and returned Solovay's solution on September 12, 1989, by e-mail.

**4.5.7.** [17] Let  $\nu_1, \nu_2, \dots$  be an effective enumeration of all discrete enumerable semimeasures. For each  $i$  define  $\gamma_i(x) = \sum_y \{\nu_i(xy) : y \in \{0, 1\}^*\}$ .

(a) Show that  $\gamma_1, \gamma_2, \dots$  is an effective enumeration of the set of enumerable semimeasures: the *extension* semimeasures.

(b) Define  $\mathbf{Mc}(x) = \sum_y \{\mathbf{m}(xy) : y \in \{0, 1\}^*\}$ . We call  $\mathbf{Mc}$  the *extension* semimeasure. Show that  $\mathbf{Mc}$  is universal in the  $\gamma$ -enumeration, that is, for all  $k$ , there is a positive constant  $c_k$  such that for all  $x$  we have  $\mathbf{Mc}(x) \geq c_k \gamma_k(x)$ .

(c) Show that  $\lim_{n \rightarrow \infty} \sum_{l(x)=n} \mathbf{Mc}(x) = 0$ .

(d) Show that the  $\gamma$ -enumeration doesn't contain all enumerable semimeasures.

(e) Show that there is no positive constant  $c$  such that for all  $x$  we have  $\mathbf{Mc}(x) \geq c \mathbf{M}(x)$ . Conclude that  $\mathbf{M}$  dominates  $\mathbf{Mc}$ , but  $\mathbf{Mc}$  does not dominate  $\mathbf{M}$ .

(f) Investigate  $\mathbf{Mc}$  and the normed measure  $\mathbf{Mc}_{norm}$ . (With  $(>) \geq$  denoting (strict) domination we have found  $\mathbf{M}_{norm} \geq \mathbf{M} > \mathbf{Mc}$ .)

*Comments.* Hint: for Item (d), in particular,  $\mathbf{M}$  is not in there. Consider the limit in Item (c) with  $\mathbf{M}$  substituted for  $\mathbf{Mc}$ . The relation between this class of (continuous) measures and the discrete measures has the following extra property: while  $\gamma(\epsilon) \leq 1$ , we have  $\gamma(x) = \gamma(x0) + \gamma(x1) + \nu(x)$  instead of just  $\gamma(x) \geq \gamma(x0) + \gamma(x1)$ . Moreover, in the  $\gamma$  semimeasures all probability is concentrated on the finite sequences and none on the infinite sequences. Source of definition of  $\mathbf{Mc}$ : T.M. Cover, Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin, Tech. Rept. 12, Statistics Dept., Stanford University, Stanford, CA, 1974. Source of Item (c): R.J. Solomonoff, *IEEE Trans. Inform. Theory*, IT-24(1978), 422–432.

**4.5.8.** [15] We compare  $\mathbf{M}(x)$  and  $\mathbf{Mc}(x)$ . Show that, for some infinite  $\omega$  (like recursive reals) we have  $\lim_{n \rightarrow \infty} \mathbf{M}(\omega_{1:n})/\mathbf{Mc}(\omega_{1:n}) = \infty$ .

*Comments.* Hint: since  $0.\omega$  is a recursive real,  $\lim_{n \rightarrow \infty} \mathbf{M}(\omega_{1:n}) > 0$ .

**4.5.9.** [31] (a) Let  $\nu$  be any probability measure and  $G(n) = \mathbf{Ef}(\omega_{1:n})$  with  $f(\omega_{1:n}) = \log \nu(\omega_{1:n}) - \log \mathbf{Mc}(\omega_{1:n})$ . ( $\mathbf{Ef}(\omega_{1:n})$  denotes the  $\nu$ -expectation  $\sum_{l(x)=n} \nu(x)f(x)$ .) Show that  $\lim_{n \rightarrow \infty} G(n) = \infty$ .

(b) Let  $\nu$  be a recursive probability measure, and let  $F$  be a recursive function such that  $\lim_{n \rightarrow \infty} F(n) = \infty$ . Show that there is a constant  $c$  such that for all binary  $x$  of length  $n$  we have  $\log \nu(x) - \log \mathbf{Mc}(x) < c + F(n)$ .

*Comments.* Hint: use Exercises 4.5.7, 4.5.8 to solve Item (a). We can interpret  $-\log \mathbf{M}c(x)$  as the length of Shannon-Fano code using distribution  $\mathbf{M}c$ , and  $-\log \nu(x)$  as the length of Shannon-Fano code using the actually reigning  $\nu$ . Clearly, although  $\log(\nu/\mathbf{M}c)$  approaches infinity with  $n$ , it does so more slowly than any recursive function. In contrast,  $\log(\nu/\mathbf{M})$ , or  $\log(\nu/\mathbf{M}_{norm})$ , is bounded by a constant. Source: R.J. Solomonoff, *IEEE Trans. Inform. Theory*, IT-24(1978), 422–432.

**4.5.10.** [22] Let  $\mu(x)$  be an enumerable probability measure. Suppose we define the co-occurrence of events and conditional events anew as follows: The probability of co-occurrence  $\mu(x, y)$  is  $\mu(x, y) = \mu(x)$  if  $y$  is a prefix of  $x$ ; it is  $\mu(x, y) = \mu(y)$  if  $x$  is a prefix of  $y$ , and equals zero otherwise. The conditional probability is defined as  $\mu(y|x) = \mu(x, y)/\mu(x)$ . Complexities are defined as unconditional  $K\mu(x) = -\log \mu(x)$ ; co-occurrence  $K\mu(x, y) = -\log \mu(x, y)$ ; and conditional  $K\mu(x|y) = -\log \mu(x|y)$ .

Show that these complexities satisfy exactly the information-theoretic equality of symmetry of information:  $K\mu(x, y) = K\mu(x) + K\mu(y|x)$  (Sections 1.11, 2.8, 3.9.1).

*Comments.* Similar (probability) definitions were used by D.G. Willis [*J. ACM*, 17(1970), 241–259]. Source: R.J. Solomonoff, *IEEE Trans. Inform. Theory*, IT-24(1978), 422–432. Solomonoff used  $\mathbf{M}_{norm}$  instead of an arbitrary measure  $\mu$ .

**4.5.11.** [48] The analogue of Theorem 4.3.3 does not hold for continuous semimeasures. Therefore, the inequality in Theorem 4.5.4 cannot be improved to equality.

- (a) Show that for any co-enumerable function  $g : \mathcal{N} \rightarrow \mathcal{N}$  for which  $Km(x) - KM(x) \leq g(l(x))$ , we have  $Km(n) \leq g(n) + O(1)$ .
- (b) Show that for almost all infinite  $\omega$ ,  $Km(\omega_{1:n}) - KM(\omega_{1:n})$  has an upper bound that is smaller than any unbounded recursive function.

*Comments.* See discussion in Section 4.5.4. Source: P. Gács, *Theoret. Comput. Sci.*, 22(1983), 71–93.

**4.5.12.** [15] Show that an infinite sequence  $\omega$  is random with respect to a recursive measure  $\mu$  iff the probability ratio  $\mu(\omega_{1:n})/\mathbf{M}(\omega_{1:n})$  is bounded below.

*Comments.* Hint: see Theorem 4.5.6. M). This ratio can be viewed as the likelihood ratio of hypothesis  $\mu(x)$  and the fixed alternative hypothesis  $\mathbf{M}(x)$ . Source: L.A. Levin, *Soviet Math. Dokl.*, 14(1973), 1413–1416.

**4.5.13.** [42] Consider a finite or countably infinite basis  $\mathcal{B}$ , and define a *probability function*  $p : \mathcal{B} \rightarrow [0, 1]$  such that  $\sum_{b \in \mathcal{B}} p(b) \leq 1$ . If equality holds, we call the probability function *proper*. The *Hellinger distance*  $\rho(q, p)$  between two probability functions  $q$  and  $p$  is defined as

$\sum_{b \in \mathcal{B}} (\sqrt{q(b)} - \sqrt{p(b)})^2$ . The  $\chi^2$  distance, denoted by  $\rho_2(q, p)$ , is defined as  $\sum_{b \in \mathcal{B}} (q(b) - p(b))^2 / q(b)$ . (Use  $0/0 = 0$  and  $\infty/\infty = 0$ .) If  $\mu$  is a semimeasure over  $\mathcal{B}^*$ , and  $\omega \in \mathcal{B}^\infty$ , then the probability function  $\mu(\cdot | \omega_{1:n}) : \mathcal{B} \rightarrow [0, 1]$  is defined by  $\mu(b | \omega_{1:n}) = \mu(\omega_{1:n}b) / \mu(\omega_{1:n})$ . If  $\mu$  is a measure and  $\mu(\omega_{1:n}) \neq 0$ , then  $\mu(\cdot | \omega_{1:n})$  is proper. The randomness deficiency of  $\omega_{1:n}$  with respect to  $\mu$  is  $\gamma_0(\omega_{1:n} | \mu) = \log(\mathbf{M}(\omega_{1:n}) / \mu(\omega_{1:n}))$ . An infinite sequence  $\omega$  is  $\mu$ -random iff  $\gamma_0(\omega_{1:n} | \mu) = O(1)$ .

- (a) Suppose  $\mu$  is a recursive measure,  $\sigma$  is a recursive semimeasure over  $\mathcal{B}^*$ , and  $\omega \in \mathcal{B}^\infty$ . Show that

$$\begin{aligned} & \sum_{i=0}^{n-1} \rho(\sigma(\cdot | \omega_{1:i}), \mu(\cdot | \omega_{1:i})) - \gamma_0(\omega_{1:n} | \mu) - O(1) \leq \gamma_0(\omega_{1:n} | \sigma) \\ & \leq \sum_{i=0}^{n-1} \rho_2(\sigma(\cdot | \omega_{1:i}), \mu(\cdot | \omega_{1:i})) + 2\gamma_0(\omega_{1:n} | \mu) + O(1). \end{aligned}$$

- (b) Suppose that  $\mu$  and  $\sigma$  are recursive semimeasures over  $\mathcal{B}^*$ , and  $\omega \in \mathcal{B}^\infty$  is both  $\mu$ -random and  $\sigma$ -random. Show that

$$\sum_{i=0}^{n-1} \left( \frac{\mu(\omega_{i+1} | \omega_{1:i})}{\sigma(\omega_{i+1} | \omega_{1:i})} - 1 \right)^2 < \infty.$$

- (c) Suppose  $\mu$  is a recursive measure,  $\sigma$  is a recursive semimeasure over  $\mathcal{B}^*$ ,  $\omega \in \mathcal{B}^\infty$  is  $\sigma$ -random, and  $\mu(\omega_{1:n}) > 0$  for all  $n$ . Show that  $\omega$  is  $\mu$ -random iff  $\sum_{i=0}^{\infty} \rho(\sigma(\cdot | \omega_{1:i}), \mu(\cdot | \omega_{1:i})) < \infty$ .

- (d) Show that if  $\omega$  is both  $\mu$ -random and  $\sigma$ -random, then  $\mu(0 | \omega_{1:n}) - \sigma(0 | \omega_{1:n}) \rightarrow 0$ , as  $n \rightarrow \infty$ .

*Comments.* Hint for Item (c): use Items (a) and (b). Conclude from Item (a) that if  $\omega$  is random relative to a recursive measure  $\mu$ , and the recursive measure  $\sigma$  is chosen so that  $\gamma_0(\omega_{1:n} | \sigma) = o(n)$ , then the “mean Hellinger distance”  $n^{-1} \sum_{i=0}^{n-1} \rho(\sigma(\cdot | \omega_{1:i}), \mu(\cdot | \omega_{1:i})) \rightarrow 0$ . Item (c) gives a randomness criterion for objects with respect to a recursive measure, in terms of Hellinger distance with a recursive semimeasure with respect to which the object is known to be random. Source: V.G. Vovk, *Soviet Math. Dokl.*, 35(1987), 656–660. See also the estimate of prediction errors as in Theorem 5.2.1.

**4.5.14.** [28] Let  $\mathcal{B}$  be a finite nonempty set of basic symbols. Let  $\delta_0(\omega | \mu)$  be a universal sequential  $\mu$ -test for sequences in  $\mathcal{B}^\infty$  distributed according to a recursive measure  $\mu$ . Let  $\phi$  be a monotone function as in Definition 4.5.3, page 276, that is  $\mu$ -regular as in Definition 4.5.5.

- (a) Show that  $\delta_0(\omega | \mu) \geq \delta_0(\phi(\omega) | \mu_\phi) + K(\phi)$ .

- (b) Show that if  $\mu_\phi$  is a recursive measure, then there exists a  $\mu_\phi$ -regular monotone function  $\psi$  such that  $\mu_{\psi\phi} = \mu$  and  $\delta_0(\phi(\omega)|\mu_\phi) \geq \delta_0(\omega|\mu) + K(\psi)$ .

*Comments.* This generalizes Exercise 3.6.2, page 219. In particular it shows that a real number has a random binary representation iff it has a random representation in any base  $r \geq 2$ . Note that a sequence is not random in itself, but random with respect to a particular measure. Thus, if we recursively transform a sequence, then its randomness properties and the complexities of its initial segments are preserved up to an additive constant with respect to the transformed measure. Source: L.A. Levin, [Problems Inform. Transmission, 10:3(1974), 206–210; Inform. Contr., 61(1984), 15–37].

- 4.5.15.** [39] Let  $\mu_1, \mu_2, \dots$  be the standard enumeration of enumerable semimeasures, so that it makes sense to talk about  $K(\mu)$  for  $\mu$  in this enumeration. If  $\gamma_0(x|\mu) = \log \mathbf{M}(x)/\mu(x)$  is the universal martingale  $\mu$ -test test for finite  $x$ , then the greater this quantity is, the “less random” is  $x$  with respect to  $\mu$ . Consider the minimal randomness deficiency of  $x$  with respect to each “sufficiently simple” semimeasure  $\mu$  with  $K(\mu) \leq \alpha$ :

$$\beta_x(\alpha) = \inf\{\gamma_0(x|\mu) : C(\mu) \leq \alpha\}.$$

Using randomness deficiency with respect to finite sets  $A$ , we can define (similar to Exercise 4.3.11)

$$\hat{\beta}_x(\alpha) = \inf\{\gamma_0(x|A) : C(A) \leq \alpha, x \in A\}.$$

- (a) Show that

$$\begin{aligned}\beta_x(\alpha + c_1) &\leq \tilde{\beta}_x(\alpha) + 2 \log C(x) + c_5, \\ \beta_x(\alpha + 2 \log(\beta_x(\alpha) + C(x)) + c_3) &\leq \beta_x(\alpha) + 2 \log(\beta_x(\alpha) + C(x)) + c_4, \\ \tilde{\beta}_x(\alpha + 2 \log l(x) + c_7) &\leq \beta_x(\alpha) + 4 \log l(x) + c_6 \text{ for } \alpha \geq c_8.\end{aligned}$$

- (b) For a wide class of sequences to which probability laws are applicable,  $\beta_x(\alpha)$  is “small” (with respect to  $l(x)$ ) for “small” values of  $\alpha$ . We give an example of this: Let  $\{\mu_r\}$  be a family of measures recursive and continuous in a real parameter  $r$ . Let  $\omega$  be an infinite  $\mu_r$ -random sequence for some  $\mu_r$  in this family ( $\sigma_0(\omega|\mu_r) < \infty$ ). Let  $\mu_r(\omega_{1:n}) > 0$ , for all  $n$ , and let the *Fisher Information*

$$I(r|\omega_{1:n}) = \mathbf{E}_r \left( \frac{d \log \mu_r(\omega_{1+n:\infty}|\omega_{1:n})}{d r} \right)$$

exist and be finite. (Recall that  $\mu(y|x)$  is defined as  $\mu(xy)/\mu(x)$ .) Show that for any  $\epsilon > 0$  and  $y = \omega_{1:n}$ , we have  $\beta_y(((1+\epsilon)/2) \log n) < c$ , for all  $n$ , where  $c$  is a constant.

*Comments.* Source: V.V. V'yugin, *SIAM Theory Probab. Appl.*, 32(1987), 508–512. The notion  $\beta_x(\alpha)$  was proposed by A.N. Kolmogorov in a talk at the *Information Theory Symposium*, Tallin, Estonia, 1974.

**4.5.16.** [43] Sequences with maximal Kolmogorov complexity of the initial segments are random in Martin-Löf's sense of passing all effective statistical tests for randomness. Hence, they must satisfy laws like the Law of the Iterated Logarithm.

- (a) Show that if  $\omega$  is an infinite sequence such that  $Km(\omega_{1:n}) \geq n - o(\log \log n)$  for a fixed constant  $c$  and infinitely many  $n$ , then with  $f_n = \omega_1 + \dots + \omega_n$

$$\limsup_{n \rightarrow \infty} |f_n - n/2| \leq \sqrt{(n \log \log n)/2} \leq 1.$$

- (b) Show that if  $Km(\omega_{1:n}) \geq n - o(n)$ , then  $\lim_{n \rightarrow \infty} f_n/n = \frac{1}{2}$ . Conversely, for any  $\epsilon > 0$  there is an  $\omega$  with  $Km(\omega_{1:n}) \geq (1 - \epsilon)n$  for which the above limit is not equal to  $\frac{1}{2}$ .

- (c) We say that an infinite binary sequence  $\omega$  satisfies the *Infinite Recurrence Law* if  $f_n = n/2$  infinitely often. Let  $\epsilon > 0$ . Show the following.

- (i) If  $Km(\omega_{1:n}) > n - (1/2 - \epsilon) \log n$ , for all  $n$ , then  $\omega$  is recurrent.
- (ii) There is a nonrecurrent  $\omega$  ( $f_n/n \geq \frac{1}{2}$  for all  $n$ ) such that we have  $Km(\omega_{1:n}) > n - (2 + \epsilon) \log n$ , for all  $n$ .

*Comments.* Item (a) is one side of the Law of the Iterated Logarithm (which actually states equality instead of inequality). By Theorem 2.5.5, Item (a) holds for almost all infinite binary sequences  $\omega$ . In other words, the (one-half) Law of the Iterated Logarithm holds for all infinite  $\omega$  in a set that (obviously strictly) contains the Martin-Löf random sequences. This is a special case of Equation 2.3 in Section 2.6. There it was shown that for  $C(x) \geq n - \delta(n)$ ,  $n = l(x)$ , we have  $|f_n - n/2| \leq \sqrt{n\delta(n)\ln 2}$ . Item (b) is a form of the Strong Law of Large Numbers (Section 1.10). By Theorem 2.5.5 this gives an alternative proof that this law holds for almost all infinite binary sequences. (Hint for the second part of Item (b): insert ones in an incompressible sequence at  $1/\epsilon$  length intervals.) Source: V.G. Vovk, *SIAM Theory Probab. Appl.*, 32(1987), 413–425.

## 4.6 Universal Average-Case Complexity, Continued

The discrete universal distribution has the remarkable property that the average computational complexity is of the same order of magnitude as the worst-case complexity, Section 4.4. What about the continuous version? This relates to algorithms for online computations that in principle never end. Such processes are abundant: sequence of key strokes from a computer key board; data-base insertions, deletions, and searches; network browser requests; and so on.

Formally, assume the input sequence is infinite over the set  $\mathcal{B}$  of basic symbols, say  $\mathcal{B} := \{0, 1\}$ . Let the probability distribution of the inputs be an enumerable semimeasure  $\mu$  according to Definition 4.2.1 on page 243. Let  $A$  be an algorithm processing inputs  $\omega = \omega_1\omega_2\dots \in \{0, 1\}^\infty$ . The *computation time* of processing input  $\omega_{1:n}$  up to inputting of symbol  $\omega_{n+1}$  is  $t(\omega_{1:n})$ .

**Definition 4.6.1** Consider a continuous sample space  $\{0, 1\}^\infty$  with semimeasure  $\mu$ . Let  $t(\omega_{1:n})$  be the running time of algorithm  $A$  on initial segment instance  $\omega_{1:n}$ . Define the *worst-case time complexity* of  $A$  as  $T(n) = \max\{t(\omega_{1:n}) : \omega_{1:n} \in \{0, 1\}^n\}$ . Define the  $\mu$ -*average time complexity* of  $A$

$$T(n|\mu) = \frac{\sum_{\omega_{1:n}} \mu(\omega_{1:n}) t(\omega_{1:n})}{\sum_{\omega_{1:n}} \mu(\omega_{1:n})}.$$

**Theorem 4.6.1 (M-Average Complexity)** *Let  $A$  be an online algorithm with inputs in  $\{0, 1\}^\infty$ . Let the inputs to  $A$  be distributed according to the universal semimeasure  $\mathbf{M}$ . Then, the average-case time complexity is of the same order of magnitude as the corresponding worst-case time complexity.*

**Proof.** Substitute  $\mu$  for  $P$  and  $\mathbf{M}$  for  $\mathbf{m}$  in Theorem 4.4.1's proof on page 269.  $\square$

## 4.7 History and References

Napoleon's contemporary and friend, Pierre-Simon Laplace, later Marquis de Laplace, may be regarded as the founder of the modern phase in the theory of probability. His anticipation of Kolmogorov complexity reasoning we quoted at the beginning of this chapter occurs in the "sixth principle: the reason why we attribute regular events to a particular cause" of his *Essai philosophique sur les probabilités*. The remarks of Dr. Samuel Johnson (1709–1784) are taken from the monumental J. Boswell, *Life of Johnson*, 1791, possibly the most famous biography written in the English language. For the basics of probability theory (as in Section 1.4) we have primarily used [A.N. Kolmogorov, *Grundbegriffe der Wahrscheinlichkeitsrechnung*, Springer-Verlag, 1933; English translation published by Chelsea, 1956].

The notion of recursively enumerable functions, or enumerable functions, originates, perhaps, from L.A. Levin and A.K. Zvonkin [*Russ. Math. Surveys*, 25(1970), 83–124], but is so natural that it might have been used earlier (unknown to either of the authors of this book, or Levin). Levin and Zvonkin and other authors use “semicomputable function” for “enumerable function” and “computable function” for “recursive function.” The notion of “semimeasure” is most similar to “lower measure” for non-measurable sets. They are also called “defective measure” in W. Feller, *An Introduction to Probability Theory and Its Applications, vol. II*, Wiley, 1970. In the current setting, semimeasures were apparently used first by Levin and Zvonkin [*Russ. Math. Surveys*, 25(1970), 83–124], not by this name but in terms of an equivalent (but awkward) measure on a non-Hausdorff space denoted there as  $\Omega^*$ . The name “semimeasure,” together with the explicit definition as we used it, may originate from [L.A. Levin and V.V. V'yugin, *Lecture Notes in Computer Science*, vol. 53, Springer-Verlag, 1977, pp. 359–364]. The combination “enumerable probability distribution” (measure) as a framework for Solomonoff’s approach below is due to Levin and Zvonkin [*Russ. Math. Surveys*, 25(1970), 83–124], but a related approach using recursive probability distributions (measures) was given by D.G. Willis [*J. ACM*, 17(1970), 241–259]. R.J. Solomonoff used the notion of recursive measures and informally noticed enumerable semimeasures previously in [*Inform. Contr.*, 7(1964), 1–22].

Kolmogorov’s introduction of complexity was motivated by information theory and problems of randomness. Solomonoff introduced algorithmic complexity independently and for a different reason: inductive reasoning. Universal a priori probability, in the sense of a single prior probability that can be substituted for each actual prior probability in Bayes’s Rule, was invented by Solomonoff, with Kolmogorov complexity as a side product, several years before anybody else did. R.J. Solomonoff obtained a Ph.B. (Bachelor of Philosophy) and M.Sc. in Physics at the University of Chicago. He was already interested in problems of inductive inference and exchanged viewpoints with the resident philosopher of science Rudolf Carnap, who taught an influential course in probability theory [*Logical Foundations of Probability*, Univ. Chicago Press, 1950].

In 1956, Solomonoff attended the Dartmouth Summer Study Group on Artificial Intelligence, at Dartmouth College in Hanover, New Hampshire, organized by M. Minsky, J. McCarthy, and C.E. Shannon, and in fact stayed on to spend the whole summer there. (This meeting gave AI its name.) There Solomonoff wrote a memo on inductive inference. McCarthy had the idea that given any mathematical problem, it could be brought into the form of “given a machine and a desired output, find an input from which the machine computes that output.” Solomonoff suggested that there was a class of problems that was not of that form: “given an initial segment of a sequence, predict its continuation.” Mc-

Carthy then thought that if one saw a machine producing the initial segment, and then continuing past that point, would you not think that the continuation was a reasonable extrapolation? With that the idea got stuck, and the participants left it at that.

Later, Solomonoff presented a paper “An Inductive Inference Machine,” at the IEEE Symposium on Information Theory, 1956, describing a program to unsupervisedly learn arithmetic formulas from examples. At the same meeting, there was a talk by N. Chomsky, based on his paper ‘Three models for the description of language’ [*IRE Trans. Inform. Theory*, IT-2(1956), 113–126]. The latter talk started Solomonoff thinking anew about formal machines in induction. In about 1958 he left his half-time position in industry and full-time joined Zator Company, a small research outfit located in some rooms at Mount Auburn Street 140½, Cambridge, Massachusetts, which had been founded by Calvin Mooers sometime around 1954 for the purpose of developing information retrieval technology. Floating mainly on military funding, Zator Co. was a research front organization, employing Mooers, Solomonoff, Mooers’s wife, and a secretary, as well as at various times visitors like Marvin Minsky. It changed its name to the more martial sounding Rockford Research (Rockford, Illinois was the place where Mooers had lived) sometime around 1962. In 1968 Solomonoff left and founded his own (one-man) company, Oxbridge Research in Cambridge in 1970, and has been there ever since, apart from spending nine months as research associate at MIT’s Artificial Intelligence Laboratory, and 1990–1991 at the University of Saarland, Saarbrücken, Germany.

In 1960 Solomonoff published a report “A preliminary report on a general theory of inductive inference” [Tech. Rept. ZTB-138, Zator Company, Cambridge, Mass.] in which he gave an outline of the notion of universal a priori probability and how to use it in inductive reasoning (rather, prediction) according to Bayes’s Rule (Chapter 5). This was sent out to all contractors of the Air Force who were even vaguely interested in this subject. In his paper of 1964 [A formal theory of inductive inference, Part 1, *Inform. Contr.*, 7(1964), 1–22], Solomonoff developed these ideas further and defined the notion of enumeration of monotone machines and a notion of universal a priori probability based on the universal monotone machine.

This way it came about that the original incentive to develop a theory of algorithmic information content of individual objects was Ray Solomonoff’s invention of a *universal* a priori probability that can be used instead of the actual a priori probability in applying Bayes’s Rule. His original suggestion in 1960 was to set the universal a priori probability  $P(x)$  of a finite binary string  $x$  as  $\sum 2^{-l(p)}$ , the sum taken over all programs  $p$  with  $U(p) = x$  where  $U$  is the reference Turing machine of Theorem 2.1.1 on page 97. However, using plain Turing machines this is

improper, since not only does  $\sum_x P(x)$  diverge, but  $P(x)$  itself diverges for each  $x$ . To counteract this defect, Solomonoff in 1964 introduced a machine model tantamount to prefix machines/monotone machines. This left the problem of the corresponding  $P(x)$  not being a probability measure. For this Solomonoff in 1964 suggested, and in 1978 exhibited, a normalization. However, the resulting probability measure is not even enumerable. According to Solomonoff this is a small price to pay. In fact, in some applications we may like the probability measure property and do not care about enumerability (Section 4.5.3 and Chapter 5). The universal distribution has remarkable properties and applications. Such applications are “simple pac-learning” in Section 5.4; the MDL principle in statistical inference and learning, Section 5.5; and the notion and development of “logical depth,” Section 7.7.

The remarkable property of  $\mathbf{m}(\cdot)$ , that the average computational complexity of any algorithm is always of the order of magnitude of the worst-case complexity, Sections 4.4 and 4.6, is from [M. Li and P.M.B. Vitányi, *Inform. Process. Lett.*, 42(1992), 145–149]. For computable versions of  $\mathbf{m}(\cdot)$  this phenomenon is treated in Section 7.6. These considerations involve the maximal gain of average-case complexity over worst-case complexity for (algorithm, distribution) pairs and associated families like polynomial time algorithms and polynomial time computable distributions, [P.B. Miltersen, *SIAM J. Comput.*, 22:1(1993), 147–156; K. Kobayashi, *IEICE Trans. Inform. Systems*, E76-D:6(1993), 634–640; K. Kobayashi, Transformations that preserve malignness of universal distributions, *Theoret. Comput. Sci.*, to appear; A. Jakoby and R. Reischuk, and C. Schindelhauer, *Proc. 12th Symp. Theoret. Aspects Comput. Sci.*, 1995, pp. 628–639]. A.K. Jagota and K.W. Regan [‘Performance of MAX-CLIQUE approximation heuristics under description-length weighed distributions,’ UB-CS-TR 92-24, SUNY at Buffalo, 1992] have extended Theorem 4.4.1 to approximation ratios for approximation algorithms. In this paper, they have also used a distribution  $q(x)$  to approximate  $\mathbf{m}(x)$  and performed extensive experiments for the MAX-CLIQUE problem. In their experiments, it was found that three out of nine algorithms perform much worse under  $q(x)$  than under the uniform distribution, confirming the theory; six other algorithms showed not much difference.

Leonid A. Levin in 1970 gave a mathematical expression of a priori probability as a universal (that is, maximal) enumerable semimeasure, Theorem 4.3.1, and showed that  $-\log \mathbf{m}(x)$  coincides with  $C(x)$  to within an additive term of  $2 \log C(x)$ . In 1974 he explicitly introduced the notion of prefix machines and prefix complexity, and showed the remarkable Theorem 4.3.3, which can be designated as the *Coding Theorem*. To be able to formulate this theorem properly, we first recall the discrete form of the universal a priori probability, using the prefix complexity

$K(x)$ . In their 1970 paper [A.K. Zvonkin and L.A. Levin, *Russ. Math. Surveys*, 25(1970), 83–124], L.A. Levin analyzes the case of continuous semimeasures related to monotone machines and shows the construction of the universal enumerable semimeasure, its equality with the universal a priori probability, and the universal randomness  $\mu$ -test for arbitrary measures  $\mu$ . See also [L.A. Levin, *Soviet Math. Dokl.*, 14(1973), 1477–1480]. The interpretation in terms of discrete semimeasures and the restriction of monotone machines to prefix machines are also due to L.A. Levin [L.A. Levin, *Problems Inform. Transmission*, 10:3(1974), 206–210; P. Gács, *Soviet Math. Dokl.*, 15(1974), 1477–1480]. Prefix complexity was also introduced, independently, by G.J. Chaitin [*J. ACM*, 22(1975), 329–340], including Theorem 4.3.3. The whole theory is brought to majestic (and hard to understand) heights in [L.A. Levin, *Inform. Contr.*, 61(1984), 15–37]. The theory of universal continuous semimeasure  $M(\cdot)$  is used for induction and prediction in Section 5.2 and in the development of the notion of “algorithmic entropy” in Section 8.5.

Solomonoff insists on the use of traditional probabilistic measures (such that  $\mu(x0) + \mu(x1) = \mu(x)$ ). This led to some difficulties with his 1964 paper.  $M$  cannot be uniquely increased to a measure, and it is hard to choose a natural one among possible extensions (so optimality is lost). In the relevant (commissioned) Exercise 4.5.6, R.M. Solovay has shown that *any* such extension would both change  $M$  by more than a constant factor and destroy its algorithmic properties. Solomonoff considers the increase by  $1/o(1)$  to be a merit, and loss of enumerability a small price to pay for it and for avoiding the heresy of redefining the notion of probability. It is not clear from his 1964 paper which extension he wanted (though it is clear that he meant to consider only ordinary probability measures), but in his 1978 paper he rigorously specifies an extension motivated as discussed in Section 4.5.3.

C.P. Schnorr [*Lecture Notes in Mathematics*, vol. 218, Springer-Verlag, 1971] introduced the use of martingales, due to P. Levy and used advantageously by J. Ville in 1939 [*Étude Critique de la Notion de Collectif*, Gauthier-Villars, 1939] in the study of Martin-Löf tests. Independently, C.P. Schnorr [*J. Comput. System Sci.*, 7(1973), 376–388] for the uniform distribution, and L.A. Levin [*Sov. Math. Dokl.*, 14(1973), 1413–1416] for arbitrary computable distributions, introduced the monotone variant of complexity  $Km(x)$  in about 1973. The monotone complexity  $Km$  smoothes out the oscillations in order to characterize randomness. Monotone complexity obliterates all quantitative differences among Martin-Löf random sequences, and hence does not allow us to make distinctions in randomness properties, in contrast to  $K$  complexity [M. van Lambalgen, *Random Sequences*, Ph.D. Thesis, University of Amsterdam, 1987; *J. Symb. Logic*, 54(1989), 1389–1400]. L.A. Levin [*Soviet Math. Dokl.*, 14(1973), 1413–1416] introduced monotone complexity, and C.P. Schnorr

[*J. Comput. System Sci.*, 7(1973), 376–388] introduced another complexity, which he called “process complexity.” The difference between those two complexities is not bounded by any constant [V.V. V'yugin, *Semiotika i Informatika*, 16(1981), 14–43 (p. 35); English translation: *Selecta Mathematica formerly Sovietica*, 13:4(1994), 357–389]. C.P. Schnorr in [*Basic Problems in Methodology and Linguistics*, R.E. Butts and J. Hintikka, Eds., Reidel, 1977, pp. 193–210], introduced a variant of monotone complexity coinciding up to an additive constant with Levin's variant. For further historical notes see [A.N. Kolmogorov and V.A. Uspensky, *SIAM J. Theory Probab. Appl.*, 32(1987), 387–412].

C.P. Schnorr's later definition is as follows: A partial recursive function  $\phi : \{0, 1\}^* \times \mathcal{N} \rightarrow \{0, 1\}^*$  is called a *monotone* interpreter iff

- (i) for all  $(p, n)$  in the domain of  $\phi$  we have that  $l(\phi(p, n)) = n$ , and
- (ii) for all  $(p, n), (pq, n+k)$  in the domain of  $\phi$  we have that  $\phi(p, n)$  is a prefix of  $\phi(pq, n+k)$ .

We can think of monotone interpreters as being computed by *monotone machines* according to Schnorr, which are Turing machines with two one-way read-only input tapes containing  $p$  and  $n$ , respectively; some work tapes; and a one-way write-only output tape. The output tape is written in binary, and the machine halts after it outputs  $n$  bits.

Define  $Km_\phi(x) = \min\{l(p) : \phi(p, l(x)) = x\}$ , and  $Km_\phi(x) = \infty$  if such  $p$  does not exist. There is an additively optimal monotone interpreter  $\phi_0$  such that for any other monotone interpreter  $\phi$  there is a constant  $c$  such that  $Km_{\phi_0}(x) \leq Km_\phi(x) + c$  for all  $x$ . Select one such  $\phi_0$  as reference and set the *monotone complexity* according to Schnorr  $Km(x) = Km_{\phi_0}(x)$ . Similarly, we can define the conditional monotone complexity  $Km(x|y)$ .

L.A. Levin used another definition. Instead of a function, the definition uses a recursively enumerable relation  $A(p, x)$  with the property that if  $p$  is a prefix of  $q$  and  $A(p, x)$ ,  $A(q, y)$  hold, then  $x, y$  must be compatible (one is a prefix of the other). The meaning is that our machine on input  $p$  outputs a (possibly infinite) string with prefix  $x$ . The minimum length of such an input  $p$  is the monotone complexity  $Km_A(x)$  with respect to  $A$  according to Levin. Among all such recursively enumerable relations there is a universal one, say  $U$ , such that for each  $A$  above there is a constant  $c$  such that for all  $x$ , we have  $Km_U(x) \leq Km_A(x) + c$ . We fix such a  $U$  and define the *monotone complexity according to Levin* as  $Km(x) = Km_U(x)$ . Let us call this a type 1 monotone machine.

The following definition of monotone machines is not equivalent but also appropriate. We require that if  $p$  is a prefix of  $q$  and  $A(p, x)$  and  $A(q, y)$  hold then  $x$  is a prefix of  $y$ . Let us call this a type 2 monotone machine.

There is yet another definition, the apparently most obvious one. The machine  $T$  has a one-way input tape and a one-way output tape. It keeps

reading input symbols and emitting output symbols. For a (possibly infinite) string  $x$  we write  $T(p) = x$  if  $T$  outputs  $x$  after reading  $p$  and no more. Let us call this a type 3 monotone machine. This type is used in the main text.

The monotone complexities arising from the latter three different kinds of machine are not necessarily the same. But all interesting upper bounds work for type 3 (the largest), and P. Gács's theorem distinguishing  $Km$  from  $KM$  works for type 1 (the smallest). See [C.P. Schnorr, *J. Comput. Syst. Sci.*, 7(1973), 376–388; A.K. Zvonkin and L.A. Levin, *Russ. Math. Surveys*, 25:6(1970), 83–124, attributed to L.A. Levin; L.A. Levin, *Sov. Math. Dokl.*, 14(1973), 1413–1416; P. Gács, *Theoret. Comput. Sci.*, 22(1983), 71–93] and on generalization of monotone complexity [A.Kh. Shen', *Sov. Math. Dokl.*, 29:3(1984), 569–573].

The relation between different complexities in the table of Figure 4.2 is asserted (many relations without proofs) in V.A. Uspensky's survey [*Kolmogorov Complexity and Computational Complexity*, O. Watanabe (Ed.), Springer-Verlag, 1992, pp. 85–101]. Many of the missing proofs or references are provided in [V.A. Uspensky and A.Kh. Shen', *Math. Systems Theory*, 29(1996), 271–292]. The most difficult relation, the lower bound on  $Km(x) - KM(x)$ , which infinitely often exceeds a variant of the non-primitive recursively slow-growing inverse of the Ackermann function (Exercise 1.7.18, page 45), is due to P. Gács [*Theoret. Comput. Sci.*, 22(1983), 71–93].

In the dictionary the word “martingale” is defined as (a) a betting system; (b) part of a horse’s harness; (c) part of a sailing rig. The delightful remark of Thackeray was quoted secondhand from [J. Laurie Snell, *Mathematical Intelligencer*, 4:3(1982)]. The mathematical study of martingales was started by P. Levy and continued by J. Ville [*Etude Critique de la Concept du Collectif*, Gauthier-Villars, 1939] in connection with von Mises’s notion of a random sequence. Ville showed that von Mises-Wald-Church random sequences defined in Section 1.9 do not satisfy all randomness properties; see the exercises in Section 1.9. But he also developed martingale theory as a tool in probability theory. A successful application of martingales was by J.L. Doob in the theory of stochastic processes in probability theory. In connection with random sequences in the sense of Martin-Löf, the martingale approach was first advocated and developed by C.P. Schnorr [*Lecture Notes in Mathematics*, vol. 218, Springer-Verlag, 1971]. Schnorr gives an overview of his work in [pp. 193–210 in: *Basic Problems in Methodology and Linguistics*, R.E. Butts, J. Hintikka (Eds.), D. Reidel, 1977]. See also [R. Heim, *IEEE Trans. Inform. Theory*, IT-25(1979), 558–566] for relations between recursive payoff functions, martingales, algorithmic information content, effective random tests, and coding theorems. The material used here is gleaned from [P. Gács, *Lecture Notes on Descriptive Complexity and*

*Randomness*, Manuscript, Boston University, 1987; T.M. Cover, P. Gács and R.M. Gray, *Ann. Probab.*, 17:3(1989), 840–865]. A survey of the basics of algorithmic probability and its relation universal betting and to prefix complexity is given in [T.M. Cover and J.A. Thomas, *Elements of Information Theory*, Wiley, 1991]. The election example is perhaps due to L.A. Levin. In general, the material on exact expressions of universal randomness  $\mu$ -tests is partially due to unpublished work of L.A. Levin, and is based on [P. Gács, *Z. Math. Logik Grundl. Math.*, 28(1980), 385–394; *Theoret. Comput. Sci.*, 22(1983), 71–93] and personal suggestions of P. Gács.

# 5

---

## Inductive Reasoning

### 5.1 Introduction

The *Oxford English Dictionary* defines **induction** as “the process of inferring a general law or principle from the observations of particular instances.” This defines precisely what we would like to call *inductive inference*. On the other hand, we regard *inductive reasoning* as a more general concept than inductive inference, namely, as a process of reassigning a probability (or credibility) to a law or proposition from the observation of particular instances.

In other words, inductive inference draws conclusions that *accept or reject* a proposition, possibly without total justification, while inductive reasoning only changes the degree of our belief in a proposition. In *deductive reasoning* one derives the absolute truth or falsehood of a proposition, such as when a mathematical proposition is proved from axioms. In deduction one often discards information: from the conclusion one cannot necessarily deduce the assumptions. In induction one generally increases information but does not discard information: the observed data follow from the induced law. In this view, deduction may be considered a special form of induction.

#### 5.1.1 Epicurus's Principle

Inductive reasoning dates back at least to the Greek philosopher of science Epicurus (342?–270? B.C.), who proposed the following approach:

**Principle of Multiple Explanations.** If more than one theory is consistent with the observations, keep all theories.

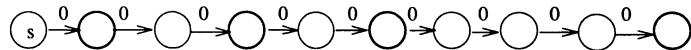
In his *Letter to Pythocles*, Epicurus motivates this as follows. There are cases, especially of events in the heavens such as the risings and settings of heavenly bodies and eclipses, where it is sufficient for our happiness that several explanations be discovered. In these cases, the events “have multiple causes of coming into being and a multiple predication of what exists, in agreement with the perceptions.”

When several explanations are in agreement with the (heavenly) phenomena, we must keep all of them for two reasons. Firstly, the degree of precision achieved by multiple explanations is sufficient for human happiness. Secondly, it would be unscientific to prefer one explanation to another when both are equally in agreement with the phenomena. This, he claims, would be to “abandon physical inquiry and resort to myth.” His follower Lucretius (95–55 B.C.) considered multiple explanations as a stage in scientific progress. According to him, to select one explanation from several equally good ones is not appropriate for the person who would “proceed step by step.”

“There are also some things for which it is not enough to state a single cause, but several, of which one, however, is the case. Just as if you were to see the lifeless corpse of a man lying far away, it would be fitting to state all the causes of death in order that the single cause of this death may be stated. For you would not be able to establish conclusively that he died by the sword or of cold or of illness or perhaps by poison, but we know that there is something of this kind that happened to him.” [Lucretius]

In the calculus of probabilities it has been customary to postulate the “principle of indifference” or the “principle of insufficient reason.” The principle of indifference considers events to be equally probable if we have not the slightest knowledge of the conditions under which each of them is going to occur. When there is an absolute lack of knowledge concerning the conditions under which a die falls, we have no reason to assume that a certain face has a higher probability of coming up than another. Hence, we assume that each outcome of a throw of the die has probability  $\frac{1}{6}$ .

[Bertrand’s Paradox] The principle of indifference is not without difficulties. Consider the following elegant paradox. We are given a glass containing a mixture of water and wine. All that is known is that  $1 \leq \text{water/wine} \leq 2$ . The principle of indifference then tells us that we should assume the probability that the ratio lies between 1 and  $\frac{3}{2}$  is 0.5 and the probability that the ratio lies between  $\frac{3}{2}$  and 2 is also 0.5. Let us take a different approach. We know  $\frac{1}{2} \leq \text{wine/water} \leq 1$ . Hence by the same principle the probabilities that this new ratio lies in the intervals of  $\frac{1}{2}$  to  $\frac{3}{4}$  and  $\frac{3}{4}$  to 1 should each be 0.5. Thus, according to the second calculation, there is 0.5 probability such that the water/wine ratio lies between 1 to  $\frac{4}{3}$  and 0.5 probability such that the water/wine ratio lies between  $\frac{4}{3}$  to 2. But the two hypotheses are incompatible.



**FIGURE 5.1.** Trivial consistent automaton

### 5.1.2 Occam's Razor

The second and more sophisticated principle is the celebrated Occam's razor principle commonly attributed to William of Ockham (1290?–1349?). This was formulated about fifteen hundred years after Epicurus. In sharp contrast to the principle of multiple explanations, it states:

**Occam's Razor Principle.** Entities should not be multiplied beyond necessity.

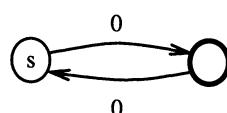
According to Bertrand Russell, the actual phrase used by William of Ockham was, “It is vain to do with more what can be done with fewer.” This is generally interpreted as, “among the theories that are consistent with the observed phenomena, one should select the simplest theory.” Isaac Newton (1642–1727) states the principle as Rule I for natural philosophy in the *Principia*, “We are to admit no more causes of natural things than such as are both true and sufficient to explain the appearances. To this purpose the philosophers say that Nature does nothing in vain, and more is in vain when less will serve; for Nature is pleased with simplicity, and affects not the pomp of superfluous causes.”

**Example 5.1.1** A *deterministic finite automaton* (DFA)  $A$  has a finite number of states, including a starting state and some accepting states. At each step,  $A$  reads the next input symbol and changes its state according to current state and the input symbol. Let us measure “simplicity” by the number of states in the automaton. The sample data are

Accepted inputs: 0, 000, 00000, 000000000;

Rejected inputs:  $\epsilon$ , 00, 0000, 000000.

There are infinitely many finite automata that are consistent with these data. Figure 5.1 shows the trivial automaton, which simply encodes the data. Figure 5.2 shows the simplest automaton. The marker  $S$  indicates the starting state and the bold circles are the accepting states.



**FIGURE 5.2.** Smallest consistent automaton

Since the automaton in Figure 5.1 simply literally encodes the data, we do not expect that the machine anticipates future data. On the other hand, the automaton in Figure 5.2 makes the plausible *inference* that the language accepted consists of strings of an odd number of 0's. It selects the “simplest” described division of the positive and negative data. It therefore also anticipates data it has not yet seen and that do not logically follow from the observed data. The latter appeals to our intuition as a reasonable inference. ◇

A too simplistic application of Occam's razor may also lead to nonsense, as the following story illustrates. Once upon a time, there was a little girl named Emma. Emma had never eaten a banana, nor had she ever been on a train. One day she had to journey from New York to Pittsburgh by train. To relieve Emma's anxiety, her mother gave her a large bag of bananas. At Emma's first bite of her banana, the train plunged into a tunnel. At the second bite, the train broke into daylight again. At the third bite, Lo! into a tunnel; the fourth bite, La! into daylight again. And so on all the way to Pittsburgh. Emma, being a bright little girl, told her grandpa at the station, “Every odd bite of a banana makes you blind; every even bite puts things right again.” Freely adapted from [N.R. Hanson, *Perception and Discovery*, 1969, Freeman and Cooper, p. 359].

In the learning automaton example, it turns out that one can *prove* the following: If sufficient data are drawn randomly from any fixed distribution, then with high probability the smallest consistent automaton (or a reasonably small automaton) will with high probability correctly predict acceptance or rejection of most data that are drawn afterwards from this distribution.

**Example 5.1.2** In spite of common intuitive acceptance of Occam's razor, the notion of simplicity remains a highly controversial and elusive idea. Things are subtler than they seem. For example, consider the following seemingly innocent rule:

Select a hypothesis that is as well in agreement with the observed value as possible; if there is any choice left, choose the simplest possible hypothesis.

Let there be an unknown number of white balls and black balls in a sealed urn. We randomly draw one ball at a time, note its color and replace it, and shake the urn thoroughly. After  $n$  draws we must decide what fraction of the balls in the urn is white. The possible hypotheses state that some rational fraction  $r$  of balls in the urn is white, where  $0 < r < 1$ . By the above rule, if in  $n$  draws  $m$  white balls are selected, then we should formulate the hypothesis  $r = m/n$ . Let there be  $\frac{1}{3}$  white and  $\frac{2}{3}$  black balls. Then the probability of getting the true hypothesis

$r = \frac{1}{3}$  is zero if  $n$  is not divisible by 3, and it tends to 0, even under the assumption that  $n$  is divisible by 3. Even for a sequence of draws where the process does converge, convergence may be too slow for practical use.  $\diamond$

We still have not defined “simplicity.” How does one define it? Is  $\frac{1}{4}$  simpler than  $\frac{1}{10}$ ? Is  $\frac{1}{3}$  simpler than  $\frac{2}{3}$ ? Note that saying that there are  $\frac{1}{3}$  white balls in the urn is the same as that there are  $\frac{2}{3}$  black balls. If one wants to infer polynomials, is  $x^{100} + 1$  more complicated than  $13x^{17} + 5x^3 + 7x + 11$ ?

Can a thing be simple under one definition of simplicity and not simple under another? The contemporary philosopher Karl Popper has said that Occam’s razor is without sense since there is no objective criterion for simplicity. Popper states that every such proposed criterion will necessarily be biased and subjective.

It is widely believed that the better a theory compresses the data concerning some phenomenon under investigation, the better we learn, generalize, and the better the theory predicts unknown data. This is the basis of the “Occam’s razor” paradigm about “simplicity.” Making these ideas rigorous involves the length of the shortest effective description of the theory: its Kolmogorov complexity; which is the size in bits of the shortest binary program to compute a description of the theory on a universal computer. This complexity, although defined in terms of a particular machine model, is independent up to an additive constant and acquires an asymptotically universal and absolute character through Church’s thesis, from the ability of universal machines to simulate one another and execute any effective process. This train of thought will lead us to a rigorous mathematical relation between data compression and learning.

### 5.1.3 Bayes’s Rule

In contrast to Epicurus and Ockham, Thomas Bayes took a probabilistic view of nature. Assume that we have observational data  $D$ .

**Bayes’s Rule.** The probability of hypothesis  $H$  being true is proportional to the learner’s initial belief in  $H$  (the prior probability) multiplied by the conditional probability of  $D$  given  $H$ .

The two fundamental components in the general inductive reasoning theory we are developing are Bayes’s formula and the universal prior probability. They both bear the same characteristics: superficially trivial but philosophically deep. We have studied the mathematical theory of the universal distribution in Chapter 4. In Section 5.2 we will develop the underlying mathematics and validation of Solomonoff’s predictive theory. But first we develop Bayesian

theory starting from the motivation in Section 1.10 and the formal definition in Section 1.6.

Consider a situation in which one has a set of observations of some phenomenon and also a (possibly countably infinite) set of hypotheses that are candidates to explain the phenomenon. For example, we are given a coin and we flip it 100 times. We want to identify the probability that the coin has outcome “heads” in a single coin flip. That is, we want to find the bias of the coin. The set of possible hypotheses is uncountably infinite if we allow each real bias in  $[0, 1]$ , and countably infinite if we allow each rational bias in  $[0, 1]$ .

For each hypothesis  $H$  we would like to assess the probability that  $H$  is the “true” hypothesis, given the observation of  $D$ . This quantity,  $\Pr(H|D)$ , can be described and manipulated formally in the following way:

Consider a discrete sample space  $S$ . Let  $D$  denote a sample of outcomes, say experimental data concerning a phenomenon under investigation. Let  $H_1, H_2, \dots$  be an enumeration of countably many hypotheses concerning this phenomenon, say each  $H_i$  is a probability distribution over  $S$ . The list  $\mathcal{H} = \{H_1, H_2, \dots\}$  is called the *hypothesis space*. The hypotheses  $H_i$  are exhaustive (at least one of them is true) and mutually exclusive (at most one of them is true).

For example, say the hypotheses enumerate the possible rational (or computable) biases of the coin. As another possibility there may be only two possible hypotheses: hypothesis  $H_1$ , which says the coin has bias 0.2, and hypothesis  $H_2$ , which puts the bias at 0.8.

Suppose we have a priori a distribution of the probabilities  $P(H)$  of the various possible hypotheses in  $\mathcal{H}$ . Because the list of hypotheses is exhaustive and mutually exclusive we have  $\sum_i P(H_i) = 1$ . Assume furthermore that for all  $H \in \mathcal{H}$  we can compute the probability  $\Pr(D|H)$  that sample  $D$  arises if  $H$  is the case. Then we can also compute (or approximate in case the number of hypotheses with nonzero probability is infinite) the probability  $\Pr(D)$  that sample  $D$  arises at all

$$\Pr(D) = \sum_i \Pr(D|H_i)P(H_i).$$

From the definition of conditional probability it is easy to derive **Bayes’s Rule** (Example 1.6.3, page 19),

$$\Pr(H_i|D) = \frac{\Pr(D|H_i)P(H_i)}{\Pr(D)}, \quad (5.1)$$

and substitution yields

$$\Pr(H_i|D) = \frac{\Pr(D|H_i)P(H_i)}{\sum_i \Pr(D|H_i)P(H_i)}. \quad (5.2)$$

Despite the fact that Bayes's Rule essentially rewrites the definition of conditional probability, and nothing more, it is its interpretation and application that are profound and controversial. The different  $H$ 's represent the possible alternative hypotheses concerning the phenomenon we wish to discover. The term  $D$  represents the empirically or otherwise known data concerning this phenomenon. The term  $\Pr(D)$ , the probability of data  $D$ , is considered as a normalizing factor so that  $\sum_i \Pr(H_i|D) = 1$ .

The term  $P(H_i)$  is called the *a priori*, *initial*, or *prior* probability of hypothesis  $H_i$ . It represents the probability of  $H_i$  being true before we have obtained any data. The prior probability  $P(H)$  is often considered as the learner's *initial degree of belief* in hypothesis  $H$ .

The term  $\Pr(H_i|D)$  is called the *final*, *inferred*, or *posterior* probability, which represents the adapted probability of  $H_i$  after seeing the data  $D$ . In essence Bayes's Rule is a mapping from prior probability  $P(H)$  to posterior probability  $\Pr(H|D)$  determined by data  $D$ .

Continuing to obtain more and more data, and repeatedly applying Bayes's Rule using the previously obtained inferred probability as the current prior, eventually the inferred probability will concentrate more and more on the "true" hypothesis. It is important to understand that one can find the true hypothesis also, using many examples, by the law of large numbers. In general, the problem is not so much that in the limit the inferred probability would not concentrate on the true hypothesis, but that the inferred probability give as much information as possible about the possible hypotheses from only a limited number of data. Given the prior probability of the hypotheses, it is easy to obtain the inferred probability, and therefore to make informed decisions.

In many learning situations, if the data are consistent with the hypothesis  $H_i$ , in the strict sense of being forced by it, then  $P(D|H_i) = 1$ . Example: outcome of a throw with a die is "even" while the hypothesis says "six." If the data are inconsistent with the hypothesis, then  $P(D|H_i) = 0$ . We assume that there is no noise that distorts the data.

**Example 5.1.3** We reconsider the example given in Section 1.9. An urn contains many dice with different biases of outcome 6 in a random throw. The set of biases is  $A$ . Assume that the difference between each pair of biases is greater than  $2\epsilon$  with  $\epsilon > 0$  and the biases are properly in between 0 and 1. Randomly drawing a die from the urn, our task is to determine its bias. This is done by experimenting. We throw the die  $n$  times, independently. If 6 shows up  $m$  times, then our learning algorithm chooses the least bias in  $A$  that is nearest to  $m/n$ .

Let  $H_p$  be the event of drawing a die with bias  $p$  for outcome "6" from an urn. For  $q \in A$ , let  $D_q$  be the experimental data such that  $m$  successes

(6's) were observed out of  $n$  throws and  $|m/n - q| \leq \epsilon$ . Then,

$$P(H_p|D_q) = \frac{P(D_q|H_p)P(H_p)}{P(D_q)},$$

where  $P(D_q) = \sum_p P(D_q|H_p)P(H_p)$ . With  $H_p$  being true, the probability of  $m$  successes out of  $n$  throws is given by the binomial distribution, page 61, as

$$\binom{n}{m} p^m (1-p)^{n-m}.$$

The deviation  $\epsilon$  (where  $0 \leq \epsilon \leq 1$ ) from the average number of successes  $pn$  in a series of  $n$  experiments is analyzed by estimating the combined tail probability

$$P(|m - pn| > \epsilon n) = \sum_{|m - pn| > \epsilon n} \binom{n}{m} p^m (1-p)^{n-m}$$

of the binomial distribution. The estimates are given by Chernoff's bound of Lemma 1.10.1 on page 61,

$$P(|m - pn| > \epsilon pn) \leq 2e^{-\epsilon^2 pn/3}.$$

If  $p$  is the true attribute of the die we have drawn and  $q \neq p$ , then  $P(D_q|H_p) < 2^{-\Omega(n)}$ . (We have assumed that  $|p - q| \geq 2\epsilon$  and  $0 < p, q < 1$ .) The probability  $P(D_p|H_p) \geq 1 - 1/2^{\Omega(n)}$ . Altogether this means that the posterior probability  $P(H_p|D_p)$  goes to 1 exponentially fast with  $n$  (and  $P(H_q|D_p)$  goes to 0 as fast, for all  $q \neq p$ ).

From the Bayesian formula one can see that if the number of trials is small, then the posterior probability  $P(H_p|D_q)$  may strongly depend on the prior probability  $P(H_p)$ . When  $n$  grows large, the posterior probability condenses more and more around  $m/n$ .  $\diamond$

What we want from Bayes's Rule is that the posterior probability give as much information as possible about the possible hypothesis from only a limited number of data. This appears to require that the prior probability be correctly chosen.

In real-world problems, the prior probabilities may be unknown, uncomputable, or conceivably may not exist. (What is the prior probability of use of words in written English? There are many different sources of different social backgrounds living in different ages.) This problem would be solved if we could find a *single* probability distribution to use as the prior distribution in each different case, with approximately the same result as if we had used the real distribution. Surprisingly, this turns

out to be possible up to some mild restrictions on the class of prior distributions being taken into account.

#### 5.1.4 Hume on Induction

The philosopher D. Hume (1711–1776) argued that true induction is impossible because we can only reach conclusions by using known data and methods. Therefore, the conclusion is logically already contained in the start configuration. Consequently, the only form of induction possible is deduction. Philosophers have tried to find a way out of this deterministic conundrum by appealing to probabilistic reasoning such as using Bayes's Rule. One problem with this is where the "prior probability" one uses has to come from. Unsatisfactory solutions have been proposed by philosophers like R. Carnap and K. Popper.

However, R.J. Solomonoff's inductive method of Section 5.2, of which we have already seen a glimpse in Section 1.10, may give a rigorous and satisfactory solution to this old problem in philosophy.

Essentially, combining the ideas of Epicurus, Ockham, Bayes, and modern computability theory, Solomonoff has successfully invented a "perfect" theory of induction. It incorporates Epicurus's multiple explanations idea, since no hypothesis that is still consistent with the data will be eliminated. It incorporates Ockham's simplest explanation idea since the hypotheses with low Kolmogorov complexity are more probable. The inductive reasoning is performed by means of the mathematically sound rule of Bayes.

#### 5.1.5 Hypothesis Identification and Prediction by Compression

Our aim is to demonstrate that data compression is the answer to many questions about how to proceed in inductive reasoning. Given a body of data concerning some phenomenon under investigation, we want to select the most plausible hypothesis from among all appropriate hypotheses or predict future data. It is widely believed that the better a theory compresses the data concerning some phenomenon under investigation, the better we have learned, generalized, and the better the theory predicts unknown data, following the "Occam's razor" paradigm about "simplicity." This belief is vindicated in practice but apparently has not been rigorously proved before. Making these ideas rigorous involves the length of the shortest effective description of some object: its Kolmogorov complexity. We treat the relation between data compression and learning and show that compression is almost always the best strategy, both in hypotheses identification by using the minimum description length (MDL) principle and in prediction methods in the style of R. Solomonoff. The argument says that among all hypotheses consistent with the data the one with least Kolmogorov complexity is the most likely one. *Prediction* in Solomonoff's manner uses a complexity-weighted combination of all hypotheses in the form of the universal prior  $M(\cdot)$  (Section 5.2). *Hypothesis identification* by minimum description length (Section 5.5) balances

the complexity of the model—and its tendency for overfitting—against the preciseness of fitting the data—the error of the hypothesis.

## 5.2 Solomonoff's Theory of Prediction

Let us consider theory formation in science as the process of obtaining a compact description of past observations together with predictions of future ones. Ray Solomonoff argues that the preliminary data of the investigator, the hypotheses he proposes, the experimental setup he designs, the trials he performs, the outcomes he obtains, the new hypotheses he formulates, and so on, can all be encoded as the initial segment of a potentially infinite binary sequence. The investigator obtains increasingly longer initial segments of an infinite binary sequence  $\omega$  by performing more and more experiments on some aspect of nature. To describe the underlying regularity of  $\omega$ , the investigator tries to formulate a theory that governs  $\omega$  on the basis of the outcome of past experiments. Candidate theories (hypotheses) are identified with computer programs that compute binary sequences starting with the observed initial segment.

There are many different possible infinite sequences (histories) on which the investigator can embark. The phenomenon he wants to understand or the strategy he uses can be stochastic. In this view each phenomenon can be identified with a measure on the continuous sample space. (For the definition of “measure” see Section 1.6.)

We express the task of learning a certain concept in terms of sequences over a basic alphabet  $\mathcal{B}$ . We express what we know as a finite sequence over  $\mathcal{B}$ ; an experiment to acquire more knowledge is encoded as a sequence over  $\mathcal{B}$ ; the outcome is encoded over  $\mathcal{B}$ ; new experiments are encoded over  $\mathcal{B}$ ; and so on. This way, a concept can be viewed as a probability distribution (rather, measure)  $\mu$  over a sample space  $S = \mathcal{B}^\infty$  of all one-way infinite binary sequences. Each such sequence corresponds to one never-ending sequential history of conjectures and refutations and confirmations. The distribution  $\mu$  can be said to be the concept or phenomenon involved.

Clearly, if we know  $\mu$ , then we can best predict which element  $a \in \mathcal{B}$  is likely to turn up after a finite initial segment  $x$ . That is, we want to predict or extrapolate according to the conditional distribution  $\mu(a|x)$ . Below we show (by Bayes's Rule) that  $\mu(a|x)$  satisfies Equation 5.3.

**Example 5.2.1** Let  $\mathcal{B} = \{0, \dots, 9\}$ . The phenomenon consists of those sequences  $\omega = \omega_1\omega_2\dots$ , where  $\omega_{2i+1}$  is the  $i$ th decimal in the decimal expansion of  $\pi = 3.1415\dots$ , and  $\omega_{2i} = a$ , with  $a \in \{0, \dots, 9\}$ , with equal probabilities  $\frac{1}{10}$ . Some example values of the measure  $\mu$  describing the phenomenon are  $\mu(3) = 1$ ,  $\mu(4) = 0$ ,  $\mu(31) = \frac{1}{10}$ ,  $\mu(41) = 0$ ,  $\mu(311) = \frac{1}{10}$ ,

$\mu(314) = 0$ ,  $\mu(3114) = \frac{1}{100}$ . Prediction should follow conditional probabilities:  $\mu(0|3) = \frac{1}{10}$ ,  $\mu(0|31) = 0$ ,  $\mu(1|31) = 1$ , and  $\mu(4|311) = \frac{1}{10}$  while  $\mu(4|3114) = 1$ . But  $\mu(3|3114) = 0$  and  $\mu(3|311) = \frac{1}{10}$ .  $\diamond$

**Example 5.2.2** Let our basic alphabet be  $\mathcal{B} = \{0, \dots, 9, E, O, H, .\}$ . Numbers can be denoted as decimals like “5.1” because we have a decimal dot in our system. The phenomenon  $\mu$  that we would like to know about is the value of the gravitational constant  $g$  in the formula  $h = gt^2/2$ , where  $h$  is the height from which we drop an object in a vacuum (on Earth at about sea level), and  $t$  is the time it takes to hit the ground. An experiment  $E10$  means that we drop the object from 10 meters. An outcome  $O2$  means that the object takes 2 seconds to hit the ground. A hypothesis  $H7.8$  means that we hypothesize  $g = 7.8$  meters/seconds squared. We take  $\mu$  to be positive:  $\mu(\cdot|\cdot) > 0$  for all arguments. Then,  $\mu(O2|H9.8E19.6) = 1$ ,  $\mu(O2|H1E2) \ll 1$ , and  $\mu(O1|H9.8E4.9) = 1$ , because  $g = 9.8$ . Experimenting, the investigator will eventually discover the law of  $\mu$  and be able to confidently predict  $\mu$  after every initial sequence  $w$ , by  $\mu(Ox|wHyEz) = 1$  iff  $z = yx^2/2$  and  $\ll 1$  otherwise.  $\diamond$

### 5.2.1 Universal Prediction

Following Examples 5.2.1 and 5.2.2, the aim is to *predict* outcomes concerning a phenomenon  $\mu$  under investigation. In this case we have some prior evidence (prior distribution over the hypotheses, experimental data) and we want to predict future events. This situation can be modelled by considering a sample space  $S$  of one-way infinite sequences of basic elements  $\mathcal{B}$  defined by  $S = \mathcal{B}^\infty$ . We assume a prior distribution  $\mu$  over  $S$  with  $\mu(x)$  denoting the probability of a sequence starting with  $x$ . Here  $\mu(\cdot)$  is a *semimeasure* satisfying

$$\begin{aligned}\mu(\epsilon) &\leq 1 \\ \mu(x) &\geq \sum_{a \in \mathcal{B}} \mu(xa).\end{aligned}$$

Given a previously observed data string  $x$ , the inference problem is to predict the next symbol in the output sequence, that is, to extrapolate the sequence  $x$ . In terms of the variables in formula 5.1,  $H_{xy}$  is the hypothesis that the sequence starts with initial segment  $xy$ .

Traditional notation is “ $\mu(\Gamma_x)$ ” instead of “ $\mu(x)$ ” where *cylinder*  $\Gamma_x = \{\omega \in S : \omega \text{ starts with } x\}$ . We use “ $\mu(x)$ ” for convenience.  $\mu$  is a *measure* if equalities hold.

Given a previously observed data string  $x$ , the inference problem is to predict the next symbol in the output sequence, that is, to extrapolate the sequence  $x$ . In terms of the variables in formula 5.1,  $H_{xy}$  is the hypothesis that the sequence starts with initial segment  $xy$ . Data

$D_x$  consists of the fact that the sequence starts with initial segment  $x$ . Then,  $\Pr(D_x|H_{xy}) = 1$ , that is, the data are forced by the hypothesis, or  $\Pr(D_z|H_{xy}) = 0$  for  $z$  is not a prefix of  $xy$ , that is, the hypothesis contradicts the data. For  $P(H_{xy})$  and  $\Pr(D_x)$  in formula 5.1 we substitute  $\mu(xy)$  and  $\mu(x)$ , respectively. For  $P(H_{xy}|D_x)$  we substitute  $\mu(y|x)$ . This way the formula is rewritten as

$$\mu(y|x) = \frac{\mu(xy)}{\mu(x)}. \quad (5.3)$$

The final probability  $\mu(y|x)$  is the probability of the next symbol string being  $y$ , given the initial string  $x$ . Obviously we now only need the prior probability  $\mu$  to evaluate  $\mu(y|x)$ . The goal of inductive inference in general is to be able to either (i) predict, or extrapolate, the next element after  $x$  or (ii) to infer an underlying effective process that generated  $x$ , and hence to be able to predict the next symbol. In the most general deterministic case such an effective process is a Turing machine, but it can also be a probabilistic Turing machine or, say, a Markov process (which makes its brief and single appearance here). The central task of inductive inference is to find a universally valid approximation to  $\mu$  which is good at estimating the conditional probability that a given segment  $x$  will be followed by a segment  $y$ .

In general this is impossible. But suppose we restrict the class of priors to the *recursive* semimeasures as defined in Chapter 4, and restrict the set of basic elements  $\mathcal{B}$  to  $\{0, 1\}$ . Under this relatively mild restriction, it turns out that we can use the single universal semimeasure  $\mathbf{M}$  (Theorem 4.5.1, page 272) as a “universal” prior (replacing the real prior) for prediction. The remainder of this section is devoted to the formal expression and proof of this loosely stated fact. Subsequently we demonstrate that ultimate compression and shortest programs almost always lead to optimal prediction.

**Definition 5.2.1** Let  $\mu : \mathcal{B}^* \rightarrow [0, 1]$  be a semimeasure and let  $x, y \in \mathcal{B}$ . The conditional semimeasure  $\mu(y|x)$  is defined as in Equation 5.3 above.

For each  $x$ , the function  $\mu(y|x)$  is a semimeasure on the sample space (or the cylinder)

$$\Gamma_x = \{x\omega : x \in \mathcal{B}^*, \omega \in \mathcal{B}^\infty\}.$$

If the unconditional  $\mu$  is a measure, then so is  $\mu(\cdot|x)$ . This is easily verified from the definitions. We can interpret  $\mu(y|x)$  as the  $\mu$ -probability that a string starting with  $x$  continues with  $y$ . Hence, if we have observed that a given sequence from a  $\mu$ -distribution starts with  $x$ , then  $\mu(y|x)$  allows us to predict whether it will continue with  $y$ . This formula solves the problem of extrapolation of sequences once we know  $\mu$ .

However, in many cases  $\mu$  is unknown or unknowable. We would like a standard procedure to estimate  $\mu$ . It turns out that  $\mathbf{M}(y|x)$  is a good estimator for *all recursive*  $\mu(y|x)$ .

When we predict the continuation of a sequence over  $\mathcal{B}^*$  randomly drawn from a distribution  $\mu$ , can we estimate the error we make using  $\mathbf{M}$  for prediction instead of the actually recursive distribution  $\mu$ ?

**Example 5.2.3** Let  $\mu$  be a recursive semimeasure on  $\mathcal{B}^\infty$ , and  $\omega \in \mathcal{B}^\infty$ . If for all  $n$ , we have  $\mu(\omega_{1:n}) > 0$ , then

$$\frac{\mathbf{M}(\omega_{1:n})}{\mu(\omega_{1:n})} = \prod_{i=1}^n \frac{\mathbf{M}(\omega_i|\omega_{1:i-1})}{\mu(\omega_i|\omega_{1:i-1})} \geq 2^{-K(\mu)}.$$

This expression is the product of the ratios of the conditional probabilities for the successive symbols of  $\omega$ . The geometric mean of the factors of this  $n$ -fold product is

$$r_n = \left( \prod_{i=1}^n \frac{\mathbf{M}(\omega_i|\omega_{1:i-1})}{\mu(\omega_i|\omega_{1:i-1})} \right)^{1/n}. \quad (5.4)$$

By Theorem 4.5.1 on page 272, for each  $\mu$ -random infinite  $\omega$  there is a constant  $c$  such that

$$\sup_{n \in \mathbb{N}} \left\{ \frac{\mathbf{M}(\omega_{1:n})}{\mu(\omega_{1:n})} \right\} \leq c.$$

Thus, for each element of a subset of  $\mathcal{B}^\infty$  of  $\mu$ -measure one (the  $\mu$ -random  $\omega$ 's), we can bound  $r_n$  from two sides:

$$2^{-K(\mu)/n} \leq r_n \leq c^{1/n}.$$

That is, for  $\mu$ -random  $\omega$ 's the associated  $r_n \rightarrow 1$  for  $n \rightarrow \infty$ . But, one may ask, would it not be possible that some factors in Equation 5.4 are much greater than 1 while other factors are much less than 1? Is it possible that the geometric mean  $r_n$  goes to 1, but the ratio

$$\frac{\mathbf{M}(\omega_{n+1}|\omega_{1:n})}{\mu(\omega_{n+1}|\omega_{1:n})}$$

fluctuates forever? The amplitude of the fluctuations might even increase? This could mean that predictions according to  $\mathbf{M}$  could be far off the mark from predictions according to  $\mu$  infinitely often. However, Theorem 5.2.1 will show that if there are fluctuations of this sort, then they must damp out quickly.  $\diamond$

In the following analysis we assume that the one-way infinite binary sequences in  $S = \{0, 1\}^\infty$  are distributed according to a recursive measure  $\mu$ . If we want to predict the next element (0 or 1) that will appear after an initial sequence  $x \in \{0, 1\}^*$ , then we can do no better than predicting according to  $\mu(0|x)$ . If instead we predict according to another distribution  $\rho$ , then the prediction error can be measured as the difference of the probabilities involved:

$$|\rho(0|x) - \mu(0|x)|.$$

Because of the binary choice between 0 and 1 we only have to consider the conditional probability of predicting 0. Instead of the absolute value of the difference we take the square of the difference in the sequel. Instead of looking at the error in the  $n$ th prediction after a *particular* initial segment of length  $n - 1$  we consider the  $\mu$ -average of the  $n$ th prediction error over all initial segments of length  $n - 1$ . Of course, we are interested in the prediction error we make if  $\rho$  is set to  $\mathbf{M}$ .

**Definition 5.2.2** Let  $S_n$  be the  $\mu$ -expected value of the square of the difference in  $\mu$ -probability and  $\mathbf{M}$ -probability of 0 occurring at the  $n$ th prediction

$$S_n = \sum_{l(x)=n-1} \mu(x)(\mathbf{M}(0|x) - \mu(0|x))^2.$$

We may call  $S_n$  the *expected squared error at the nth prediction*.

**Theorem 5.2.1** Let  $\mu$  be a recursive semimeasure. Using the notation above,  $\sum_n S_n \leq k/2$  with  $k = K(\mu) \ln 2$ . (Hence,  $S_n$  converges to 0 faster than  $1/n$ .)

**Proof.** Let  $\mu$  be a recursive semimeasure and  $\mathbf{M}$  the universal enumerable semimeasure, both over the set of infinite binary sequences. We will make  $\mathbf{M}$  a measure not as Solomonoff did, by normalization, but rather by introducing an extra value  $u$ , meaning “undefined,” for the variables  $\omega_i$ . The  $\mu$ -probability of any sequence  $\omega_{1:n}$ , where  $\omega_i = u$  for some  $i$  with  $1 \leq i \leq n$ , is zero, but the  $\mathbf{M}$ -probability is positive for this. With this convention, we do not have to change  $\mathbf{M}$  by “normalization.”

Let us introduce some notation. Suppose that  $P, Q$  are probability distributions over some discrete set. Define

$$D(P \parallel Q) = \sum_x P(x) \ln(P(x)/Q(x)), \quad (5.5)$$

with “ln” the natural logarithm. This quantity 5.5 is called the *Kullback divergence* of the two distributions. By the concavity of the logarithm, it is always nonnegative and is 0 only if  $P = Q$ . It is a measure of the difference of the two distributions. It can be related to the more

familiar distance measure of squared difference in the case when  $x$  takes the values 0,1:

$$D(P \parallel Q) \geq 2(P(0) - Q(0))^2. \quad (5.6)$$

Namely, consider the function  $f(Q(0)) = D(P \parallel Q) - 2(P(0) - Q(0))^2$ , for some fixed  $P(0)$ . Firstly,  $f(P(0)) = 0$ . Secondly, for  $P(0) > Q(0)$  the function  $f(Q(0))$  increases since  $df(Q(0))/dQ(0) > 0$ , and for  $P(0) < Q(0)$  the function  $f(Q(0))$  decreases since  $df(Q(0))/dQ(0) < 0$ .

The discrete distributions above can be considered as distributions over the initial segment of length 1 of an infinite sequence. We now consider initial segments of length 2 of infinite sequences. Let  $\mu$  be a semimeasure over  $\mathcal{B}^\infty$ , for some nonempty basic set  $\mathcal{B}$ . Define the discrete distributions  $\mu_1, \mu_2$  over  $\mathcal{B}$  and  $\mu^2$  over  $\mathcal{B}^2$  by

$$\begin{aligned}\mu_1(x) &= \mu(x), \\ \mu_2(y|x) &= \mu(xy)/\mu(x) \\ \mu^2(xy) &= \mu(xy).\end{aligned}$$

That is,  $\mu_1(x)$  is the  $\mu$ -probability that the initial element is  $x$ ; the  $\mu$ -probability that the second element is  $y$ , given the first one is  $x$ , is  $\mu_2(y|x)$ ; and the  $\mu$ -probability that the first two elements are  $xy$  is  $\mu^2(xy)$ . For a pair of semimeasures  $\mu$  and  $\rho$  over  $\mathcal{B}^\infty$ , define the distances  $D_1, D_2, D^2$  by

$$\begin{aligned}D_1(\mu \parallel \rho) &= D(\mu_1 \parallel \rho_1), \\ D_2(\mu \parallel \rho) &= \sum_x \mu_1(x) D(\mu_2(\cdot|x) \parallel \rho_2(\cdot|x)), \\ D^2(\mu \parallel \rho) &= D(\mu^2 \parallel \rho^2).\end{aligned}$$

**Claim 5.2.1** Let  $\mu$  be a measure and  $\rho$  a semimeasure. Then  $D^2(\mu \parallel \rho) = D_1(\mu \parallel \rho) + D_2(\mu \parallel \rho)$ .

**Proof.** Rewrite  $D(\mu^2 \parallel \rho^2)$  as

$$\begin{aligned}\sum_{xy} \mu(xy) \ln \frac{\mu(xy)}{\rho(xy)} &= \sum_{xy} \mu(x) \mu_2(y|x) \ln \frac{\mu(x) \mu_2(y|x)}{\rho(x) \rho_2(y|x)} \\ &= \sum_x \mu(x) \ln \frac{\mu(x)}{\rho(x)} \sum_y \mu_2(y|x) + \sum_x \mu(x) \sum_y \mu_2(y|x) \ln \frac{\mu_2(y|x)}{\rho_2(y|x)}.\end{aligned}$$

Since  $\sum_y \mu_2(y|x) = 1$ , the first term of the derived sum equals  $D_1(\mu \parallel \rho)$  as defined above. The second term of the derived sum equals  $D_2(\mu \parallel \rho)$  as defined above.  $\square$

An immediate consequence of Claim 5.2.1 is the following inequality. We say that distribution  $P'$  is a *refinement* of distribution  $P$  if each  $P(x)$  is the sum of one or more  $P'(y_i)$ 's. If  $P', Q'$  are refinements of the distributions  $P, Q$ , then we find

$$D(P' \parallel Q') \geq D(P \parallel Q). \quad (5.7)$$

To see this, we can view each  $y_i$  as a pair  $(x, y_i)$ . Then with the appropriate distributions  $\mu, \rho$  we have  $P = \mu^1, Q = \rho^1, P' = \mu^2, Q' = \rho^2$ .

Another consequence comes when we extend the approach to infinite sequences  $x_1 x_2 \dots$ . For  $l(x) = n$ , define the discrete distributions  $\mu^n(x) = \mu(x)$  and  $\mu_{n+1}(y|x) = \mu(xy)/\mu(x)$ . Define furthermore

$$D_i(\mu \parallel \rho) = \sum_{l(x)=i-1} \mu^{i-1}(x) D(\mu_i(\cdot|x) \parallel \rho_i(\cdot|x)).$$

Here  $D_i$  is defined exactly as  $D_2$  (except that  $l(x) = i - 1$  instead of  $l(x) = 1$ ). It is the  $\mu$ -expected value of the divergence of the conditional probability  $\mu_i(\cdot|x_1 \dots x_{i-1})$  from the corresponding conditional  $\rho$ -probability. Similar to Claim 5.2.1 it can now be shown that

$$D(\mu^n \parallel \rho^n) = \sum_{i=1}^n D_i(\mu \parallel \rho).$$

To prove the theorem we set  $\rho := M$  and  $k := K(\mu) \ln 2$ . By Theorem 4.3.5, page 258, for all  $x \in \mathcal{B}^*$ ,

$$\ln \frac{\mu(x)}{M(x)} \leq k.$$

It follows that  $D(\mu^n \parallel M^n) \leq k$ , independently of  $n$ . Therefore,

$$\lim_{n \rightarrow \infty} D(\mu^n \parallel M^n) = \lim_{n \rightarrow \infty} \sum_{i=1}^n D_i(\mu \parallel M) \leq k.$$

This is the summed  $\mu$ -expected value of the divergence of the conditional probability  $\mu_i(\cdot|x_1 \dots x_{i-1})$  from the corresponding  $M$ -conditional probability, the summation taken over all  $i \geq 1$ . Let us derive from it a bound on the more familiar square difference. Let us assume that  $x_i$  takes only the values 0, 1 and the “undefined” value  $u$ . We define new conditional distributions  $\mu'_i, M'_i$  running only over the values 0, 1 by  $\mu'_i = \mu_i$  and

$$\begin{aligned} M'_i(0|x_1 \dots x_{i-1}) &= M_i(0|x_1 \dots x_{i-1}), \\ M'_i(1|x_1 \dots x_{i-1}) &= 1 - M_i(0|x_1 \dots x_{i-1}). \end{aligned}$$

Now,  $D_i(\mu \parallel M)$  is the  $\mu$ -expected value of the divergence of the conditional  $\mu_i$  and  $M_i$  over 0, 1,  $u$ . By an extension of the refinement Equation 5.7 from  $D^2(\cdot)$  to  $D^n(\cdot)$ , this value will only decrease if we replace

these with the conditional probabilities  $\mu'_i$  and  $\mathbf{M}'_i$  over 0, 1. But then, we can apply the analogue of the quadratic estimate of Equation 5.6 for each  $D_i(\cdot)$ , for all  $i \geq 1$  and obtain the statement of the theorem.  $\square$

Take in Theorem 5.2.1 any (possibly nonrecursive) two (semi)measures  $\mu, \rho$  over  $\mathcal{B}^\infty$  such that with  $n = l(x)$

$$\mu(x) \geq 2^{-k(n)} \rho(x),$$

for all  $x \in \mathcal{B}^*$ , and define  $S_n = \sum_{l(x)=n-1} \rho(x)(\mu(0|x) - \rho(0|x))^2$ . The same proof now shows the following:

**Corollary 5.2.1**  $\sum_{i=1}^n S_i \leq k(n) \ln \sqrt{2}$ .

Ordinary statistical analysis of a Bernoulli sequence gives an expected squared error for the probability of the  $n$ th symbol proportional to  $1/n$ , and an expected total squared error for the first  $n$  symbols of  $\ln n$ . This is clearly much larger than the total squared error of  $K(\mu) \ln \sqrt{2}$  we found in Theorem 5.2.1. This discrepancy can be understood by noting that the theorem requires  $\mu$  to be recursive. The set of recursive measures is countable. But the parameter defining a Bernoulli process can be any real number in between 0 and 1. The set of these measures is uncountable. It has been shown that if one restricts the possible hypotheses to countably many, then the statistical error converges much more rapidly than if one considers uncountably many hypotheses [T.M. Cover, *Ann. Stat.*, 1:5(1973), 862–871]. For further discussion see [R.J. Solomonoff, *IEEE Trans. Inform. Theory*, IT-24(1978), 422–432]. This remark may be misleading in this form. Namely, consider prediction of sequences generated by a Bernoulli process with parameter  $p$ . Is it true that one will be able to predict better on a limited initial segment, using Solomonoff's method, if one knows  $p$  is rational? No, unless one also knows it is a *simple* rational since the error has a multiplicative factor proportional to  $K(p)$ . Therefore, for a finite segment of length  $n$  (which is what really counts), the statistical method taking into account all  $p$ 's is inferior to the Solomonoff method only if one knows that there are at most  $n$  possibilities for  $p$ .

Assume that a probabilistic theory concerning some phenomenon is expressible as a computable measure  $\mu$  on  $\{0, 1\}^\infty$ . Then *Solomonoff's Inductive Formula*  $\mathbf{M}(y|x)$ , to estimate the actual probabilities  $\mu(y|x)$  to predict outcomes  $y$  given a sequence of observed outcomes  $x$ , can be viewed as a mathematical form of Occam's razor:

find all rules fitting the data and then predict  $y$  according to the universal distribution on them.

Formalization of this principle for probabilistic theories  $\mu$  encounters difficulties because of a tradeoff between the complexity  $K(\mu)$  and the randomness deficiency  $\log \mathbf{M}(x)/\mu(x)$ . The current approach is accurate for large  $x$  and simple  $\mu$ .

Note that while the following Theorem *does* imply the convergence of the conditional probabilities similarly to Theorem 5.2.1, it *does not* imply the speed of convergence estimate.

**Theorem 5.2.2** *Let  $\mu$  be a positive recursive measure. If the length of  $y$  is fixed and the length of  $x$  grows to infinity, then*

$$\frac{\mathbf{M}(y|x)}{\mu(y|x)} \rightarrow 1,$$

*with  $\mu$ -probability one. The infinite sequences  $\omega$  with prefixes  $x$  satisfying the displayed asymptotics are precisely the  $\mu$ -random sequences.*

**Proof.** Unfortunately, no elementary proof is known to the authors. We use an approach based on the so-called Submartingale Convergence Theorem 4.5.4, page 298, which states that the following property holds for each sequence of random variables  $\omega_1, \omega_2, \dots$ . If  $f(\omega_{1:n})$  is a  $\mu$ -submartingale, and the  $\mu$ -expectation  $\mathbf{E}|f(\omega_{1:n})| < \infty$ , then it follows that  $\lim_{n \rightarrow \infty} f(\omega_{1:n})$  exists with  $\mu$ -probability one.

In our case,

$$t(\omega_{1:n}|\mu) = \frac{\mathbf{M}(\omega_{1:n})}{\mu(\omega_{1:n})}$$

is a  $\mu$ -submartingale, and the  $\mu$ -expectation  $\mathbf{E}t(\omega_{1:n}|\mu) \leq 1$ . Therefore, there is a set  $A \subseteq \mathcal{B}^\infty$  with  $\mu(A) = 1$ , such that for each  $\omega \in A$  the limit  $\lim_{n \rightarrow \infty} t(\omega_{1:n}|\mu) < \infty$ . These are precisely the  $\mu$ -random  $\omega$ 's by Corollary 4.5.5 on page 297. Consequently, for fixed  $m$ , for each  $\omega$  in  $A$ , we have

$$\lim_{n \rightarrow \infty} \frac{\mathbf{M}(\omega_{1:n+m})/\mu(\omega_{1:n+m})}{\mathbf{M}(\omega_{1:n})/\mu(\omega_{1:n})} = 1,$$

provided the limit of the denominator is not zero. The latter fact is guaranteed by the universality of  $\mathbf{M}$ : for each  $x \in \mathcal{B}^*$  we have  $\mathbf{M}(x)/\mu(x) \geq 2^{-K(\mu)}$  by Theorem 4.5.1, page 272.  $\square$

**Example 5.2.4** We continue Example 5.2.1. Suppose we are given an infinite decimal sequence  $\omega$ . The even positions contain the subsequent digits of  $\pi = 3.1415\dots$ , and the odd positions contain uniformly distributed, independently drawn random decimal digits. Then,  $\mathbf{M}(a|\omega_{1:2i}) \rightarrow \frac{1}{10}$  for  $a = 0, 1, \dots, 9$ , while  $\mathbf{M}(a|\omega_{1:2i+1}) \rightarrow 1$  if  $a$  is the  $i$ th digit of  $\pi$ , and to 0 otherwise.  $\diamond$

Instead of using the result in Theorem 5.2.1, it is perhaps conceptually easier to use Theorem 5.2.2 for the estimation required in formula 5.3.

To estimate the conditional probability  $\mu(y|x)$  that the next segment after  $x$  is  $y$  we use the conditional prior probability

$$\mathbf{M}(y|x) = \frac{\mathbf{M}(xy)}{\mathbf{M}(x)}. \quad (5.8)$$

To be able to do so we have to assume that  $\mu$  in formula 5.3 is a recursive measure. If the length of  $y$  is fixed and the length of  $x$  grows to infinity, then Theorem 5.2.2 states that the conditional a priori probability is almost always asymptotically equal to the conditional probability.

The convergence is very fast. With  $S_n$  as in Definition 5.2.2, the expected squared error in the  $n$ th prediction of 0, and with  $\mathcal{B} = \{0, 1\}$ , this figure of merit accounts for all error in the  $n$ th prediction. By Theorem 5.2.1, the summed expected error over all predictions is bounded by a constant,  $\sum_{n=1}^{\infty} S_n < k/2$ , where  $k$  is a constant depending only on  $\mu$ . (We can set  $k = K(\mu) \ln 2$ , where  $K(\mu)$  is the length of the shortest program computing  $\mu$  in a self-delimiting binary programming language.) Using  $\mathbf{M}$ , the expected prediction error  $S_n$  in the  $n$ th prediction goes to 0 faster than  $1/n$ . Used as the prior in Bayes's Rule, this proves mathematically that the inferred probability using prior  $\mathbf{M}$  converges very fast to the inferred probability using the actual prior  $\mu$ . The problem with Bayes's Rule has always been the determination of the prior. Using  $\mathbf{M}$  universally gets rid of that problem and is provably perfect.

### 5.2.2 Prediction by Data Compression

Above it is shown that the universal distribution *itself* is directly suited for prediction. The universal distribution combines a weighted version of the predictions of all enumerable semimeasures, including the prediction of the semimeasure with the shortest program. It is not a priori clear that the shortest program dominates in all cases—and in fact it does not. However, we show that in the overwhelming majority of cases—the typical cases—the shortest program dominates sufficiently to validate the approach that uses only shortest programs for prediction. The properties of  $\mathbf{M}(x)$  allow us to demonstrate that a minimum description length procedure is almost always optimal for prediction.

Given a semimeasure on  $\{0, 1\}^{\infty}$  and an initial binary string  $x$ , our goal is to find the most probable extrapolation of  $x$ . That is, taking the negative logarithm on both sides of Equation 5.3, we want to determine  $y$  with  $l(y) = n$  that minimizes

$$-\log \mu(y|x) = -\log \mu(xy) + \log \mu(x).$$

In the remainder of this section, let  $\mu$  be a *recursive* semimeasure. An infinite binary sequence  $\omega$  is  $\mu$ -random iff

$$\sup_n \mathbf{M}(\omega_1 \dots \omega_n) / \mu(\omega_1 \dots \omega_n) < \infty$$

(see Theorem 4.5.6 and further discussion on page 294). Let  $\omega$  be a  $\mu$ -random infinite binary sequence and  $xy$  a finite prefix of  $\omega$ . For  $l(x)$  that grows unboundedly with  $l(y)$  fixed, we have by Theorem 5.2.2

$$\lim_{l(x) \rightarrow \infty} \log \mu(y|x) - \log M(y|x) = 0. \quad (5.9)$$

Therefore, if  $x$  and  $y$  satisfy the above conditions, then maximizing  $\mu(y|x)$  over  $y$  means minimizing  $-\log M(y|x)$ . It is shown in Lemma 4.5.6 on page 285 that  $-\log M(x)$  is slightly smaller than  $Km(x)$ , the length of the shortest program for  $x$  on the reference universal monotonic machine. For binary programs this difference is very small (Equation 4.14 on page 284) but possibly unbounded in the length of  $x$ .

Together this shows the following. Given  $xy$  that is a prefix of a (possibly not  $\mu$ -random)  $\omega$ , optimal prediction of fixed-length extrapolation  $y$  from an unboundedly growing prefix  $x$  of  $\omega$  need not necessarily be reached by the shortest programs for  $xy$  and  $x$  minimizing  $Km(xy) - Km(x)$ , but is reached by considering the weighted version of all programs for  $xy$  and  $x$ , which is represented by

$$-\log M(xy) + \log M(x) = (Km(xy) - g(xy)) - (Km(x) - g(x)).$$

Here  $g(x)$  is a function that can rise to infinity between the inverse of the Ackermann function and  $Km(l(x)) \leq \log \log x$ —but only in case  $x$  is not  $\mu$ -random.

Therefore, for certain  $x$  and  $y$  that are *not*  $\mu$ -random, optimization using the minimum length programs may result in very incorrect predictions. However, for  $\mu$ -random  $x$  we have that  $-\log M(x)$  and  $Km(x)$  coincide up to an additional constant independent of  $x$ , that is,  $g(xy) = g(x) = O(1)$  (Lemma 4.5.6). Hence, together with Equation 5.9, we find the following.

**Theorem 5.2.3** *Let  $\mu$  be a recursive semimeasure, and let  $\omega$  be a  $\mu$ -random infinite binary sequence and  $xy$  a finite prefix of  $\omega$ . For  $l(x)$  that grows unboundedly and  $l(y)$  fixed,*

$$\lim_{l(x) \rightarrow \infty} -\log \mu(y|x) = Km(xy) - Km(x) + O(1) < \infty,$$

*where  $Km(xy)$  and  $Km(x)$  grow unboundedly.*

By its definition  $Km$  is monotone in the sense that always  $Km(xy) - Km(x) \geq 0$ . If  $y$  makes this difference equal 0, then the shortest effective monotone program for  $x$  is also a shortest effective monotone program for  $xy$  and hence predicts  $y$  given  $x$ . For all large enough  $\mu$ -random  $x$ , predicting by determining  $y$  that minimizes the difference of the minimum program lengths for  $xy$  and  $x$  gives a good prediction. Here  $y$

should be preferably large enough to eliminate the influence of the  $O(1)$  term.

**Corollary 5.2.2 (Prediction by Data Compression)** Assume the conditions of Theorem 5.2.3. With  $\mu$ -probability going to one as  $l(x)$  grows unboundedly, a fixed-length  $y$  extrapolation from  $x$  maximizes  $\mu(y|x)$  iff  $y$  can be maximally compressed with respect to  $x$  in the sense that it minimizes  $Km(xy) - Km(x)$ . That is,  $y$  is the string that minimizes the length difference between the shortest program that outputs  $xy\dots$  and the shortest program that outputs  $x\dots$ .

The universal prior distribution  $M(\cdot)$  is not computable and neither is the complexity  $Km(\cdot)$ . But later we will see that many inference models or principles can be viewed as computable approximations to Solomonoff's inference principles.

## 5.3 Universal Recursion Induction

---

Partial aspects of the inductive inference question in the real world can be formalized in many ways. We consider a particular abstraction associated with the name of E.M. Gold.

We are given an effective enumeration of partial recursive functions  $f_1, f_2, \dots$ . Such an enumeration can be the functions computed by Turing machines, but also the functions computed by finite automata. We want to infer a particular function  $f$ . To do so, we are presented with a sequence of examples,  $D = e_1, e_2, \dots, e_n$ , containing elements (possibly with repetitions) of the form

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

For  $n \rightarrow \infty$  we assume that  $D$  contains all elements of the displayed form.

### 5.3.1 Hypothesis Identification

Let the different hypotheses  $H_k$  be " $f = f_k$ ." Since  $P(D|H_k)$  is 1 or 0 according to whether  $D$  is consistent with  $f_k$  or not, take any positive prior distribution  $P(H_k)$ , say  $P(H_k) = 1/k(k + 1)$ , and apply Bayes's Rule 5.1 to obtain

$$P(H_k|D) = \frac{P(D|H_k)P(H_k)}{\sum\{P(H_j) : f_j \text{ is consistent with } D\}}. \quad (5.10)$$

With increasing  $n$ , the denominator term is monotonically nonincreasing. Since all examples eventually appear, the denominator converges to a limit.

For each  $k$ , the inferred probability of  $f_k$  is monotonically nondecreasing with increasing  $n$ , until  $f_k$  is inconsistent with a new example, in which case it falls to zero and stays there henceforth. Only the  $f_k$ 's that are consistent with the sequence of presented examples have positive inferred probability. Order the  $f_k$ 's by decreasing probability. The ordering between the  $f_k$ 's with positive probability never changes; the only thing that can happen is that  $f_k$ 's are removed from this set because their probability drops to zero. At each step we infer the  $f_k$  with the highest posterior probability. This is always the first element in the current ordering. Assuming that  $f$  occurs in the list, at some step all preceding functions before the first occurrence of  $f$  have been deleted. From that step onwards  $f$  has the highest probability.

Analyzing this process we find that it is equivalent to starting with  $f_k$ 's ordered by decreasing probability according to  $P(H_k) = 1/k(k+1)$ , which induces the enumeration  $f_1, f_2, \dots$ . As in the classical Gold method, after receiving each new example we eliminate all remaining  $f_k$ 's that are inconsistent from the beginning onward up to the position of the first consistent function. We receive a new example  $e$ , set  $D := D, e$ , and repeat this process. Eventually, the new first function in the enumeration is a copy of  $f$  and it doesn't change any more. This algorithm is called *learning by enumeration*. One learns more and more about the unknown target function and approximates it until the correct identification has been achieved. This learning model is called *learning in the limit*. The learner eventually learns the concept exactly but never knows when that has happened. This deceptively simple idea has generated a large body of sophisticated literature.

How do we use the universal prior probability in this setting? Choose  $P(H_k) = m(k)$ , with  $m(\cdot)$  the universal discrete probability. By Theorem 4.3.1, page 247, we have  $m(x) = 2^{-K(x)}$ , with  $K(\cdot)$  the prefix complexity. With this prior, at each stage,  $P(\cdot|D)$  will be largest for the “simplest” consistent hypothesis. That is, for the one with the least prefix complexity. In the limit, this will be the case for  $H_k$  such that  $f_k = f$  with  $K(k)$  is minimal. This process corresponds to enumeration of the  $f_k$ 's as  $f_{\pi(1)}, f_{\pi(2)}, \dots$ , where  $\pi$  is a permutation such that  $K(\pi(i)) \leq K(\pi(i+1))$  for all  $i = 1, 2, \dots$

This way one enumerates the functions by increasing prefix complexity. Assume that the looked-for  $f = f_k$  is a “simple” function, as it in practice always is, with  $K(k) \ll k$ . (If  $f$  is a truly random function then it is highly unlikely that anyone will ever conceive of its existence and would want to learn it.) We will find simple  $f_k$  much faster using  $m(\cdot)$  as prior than using  $1/k(k+1)$ ; see Exercise 5.3.2. In fact, the speed-up can be uncomputably fast! But since the  $m$  is uncomputable, one can use  $1/k(k+1)$  or something better as a (trivially computable) approximation.

### 5.3.2 Prediction

Suppose we want to infer the correct value of  $f(x)$  after having seen data sample  $D$ . We can refer to the analysis above and simply predict by

$$P(e|D) = \frac{\sum P(H_k|D, e)}{\sum P(H_k|D)}. \quad (5.11)$$

But let us use the universal measure  $\mathbf{M}$ . For this analysis, replace the examples by binary self-delimiting codes (see for example Section 1.4);  $e = (x, y, 1)$  by  $\bar{e} = \bar{x}\bar{y}1$  and  $e = (x, y, 0)$  by  $\bar{e} = \bar{x}\bar{y}0$ . Here,  $\bar{x} = 1^{l(x)}0x$ , so the machine can see where the encoding of  $x$  ends without having to look at the next symbol. For convenience, we denote this binary encoding of  $D$  also by “ $D$ .” Let  $\mathcal{D}$  be the largest set of  $D$ ’s (possibly infinite) such that  $D$  is consistent with  $f_k$ . Now set

$$P(H_k) = \mathbf{M} \left( \bigcup_{D \in \mathcal{D}} D \right).$$

If we assume a recursive distribution on the examples, Solomonoff’s maxim says we must predict according to

$$\mathbf{M}(e|D). \quad (5.12)$$

By Theorem 5.2.1, the expected squared error  $S_n$  in the  $n$ th prediction goes to zero faster than  $1/n$ . We hasten to remark that this does not say much about the number of mistakes in a particular single sequence.

---

## Exercises

**5.3.1.** [14] Consider Equation 5.10 on page 335. Show that as the number  $n$  of examples  $D = e_1, \dots, e_n$  grows, the inferred probability  $P(f_k|D)$  monotonically increases either indefinitely or until it suddenly falls to 0 and stays there thereafter, for each  $k$ . Argue why this process is called “learning in the limit” and “learning by enumeration.”

**5.3.2.** [17] We continue Exercise 5.3.1. (a) Give the implicit effective enumeration of the hypotheses for the prior  $P(H_k) = 1/k(k+1)$  and give the implicit uncomputable enumeration for the prior  $P'(H_k) = \mathbf{m}(k)$ .

(b) Show that if  $f = f_k$  is the function to be learned with  $k$  minimal, then we learn  $f$  about  $k/K(k)$  times faster using  $P'(H_k) = \mathbf{m}(k)$  than using  $P(H_k) = 1/k(k+1)$ . (If  $k = 2^r$  then we learn exponentially faster using  $\mathbf{m}$ , and so on.)

*Comments.* Source: M. Li, P.M.B. Vitányi, *J. Comput. System Sci.*, 44:2(1992), 343–384.

**5.3.3.** [25] Consider an effective enumeration  $f_1, f_2, \dots$  of partial recursive functions with values in the set  $\{0, 1\}$  only. Each such function  $f$  defines an infinite binary sequence  $\omega = \omega_1\omega_2\dots$  by  $\omega_i = f(i)$ , for all

*i.* This way, we have an enumeration of infinite sequences  $\omega$ . These sequences form a binary tree with the root labeled  $\epsilon$  and each  $\omega$  an infinite path starting from the root. We are trying to learn a particular function  $f$ , in the form that we predict  $\omega_n$  from the initial sequence  $\omega_1 \dots \omega_{n-1}$  for all  $n \geq 1$ . We want to analyze the number of *mistakes* we make in this process. This is a version of prediction where the probability  $P(e|D)$  of Equation 5.11 is either 0 or 1. If our prediction is wrong (say, we predict a 0 and it should have been a 1), then this counts as one mistake. Show there is an algorithm that infers  $f = f_n$  making fewer than  $2 \log n$  mistakes in all infinitely many predictions.

*Comments.* Hint: define for each  $f_i$  with associated infinite sequence  $\omega^i$  a measure  $\mu_i$  by  $\mu_i(\omega^i) = 1$ . This implies that also  $\mu_i(\omega_1^i \dots \omega_n^i) = 1$  for all  $n$ . Let  $\mu$  be a semimeasure defined by

$$\mu(x) = \sum_i \frac{1}{i(i+1)} \mu_i(x),$$

for each  $x \in \{0,1\}^*$ . (Note that  $\mu$  is a simple computable approximation to  $\mathbf{M}$ .) The prediction algorithm is very simple: “**If**  $\mu(0|x) \geq \frac{1}{2}$  **then** predict 0 **else** predict 1.”

**5.3.4.** [25] Suppose we put weight  $2^{-K(n)}$  on  $\mu_n$  in Exercise 5.3.3 instead of  $1/(n(n+1))$ . (Of course, the prediction algorithm becomes non-effective because we cannot compute these weights since  $K(\cdot)$  is uncomputable).

- (a) Show that then the number of mistakes is at most  $k < \log n + 2 \log \log n$  for all  $n$ .
- (b) Show that for regular  $n$  (say,  $n = 2^k$ ) we have  $k < (1 + \epsilon) \log \log n$ , for all  $\epsilon > 0$ .

*Comments.* Hint: show  $k < K(n)$ .

**5.3.5.** [17] Show that if the target function is  $f$  and we make  $k$  errors in the first  $m$  predictions in Exercise 5.3.3, then  $\log \binom{m}{k} + K(m) + O(1) \geq K(f(1) \dots f(m))$ .

**5.3.6.** [22] Denote  $f(1) \dots f(m)$  by  $x$  in Exercise 5.3.5. Show that if  $k/m$  is small, then  $k \approx K(x)/\log(m/K(x))$ . For instance, with  $K(x) = \sqrt{m}$  we find  $k > 2\sqrt{m}/\log m$ .

**5.3.7.** [27] We ask for the number of examples required to effectively infer a particular function  $f$  by determining the number of examples needed to *describe* or *certify* a particular function  $f$  in any effective way. Let  $D = e_1 e_2 \dots e_n$  be a sequence of examples  $e_i = (x_i, y_i, b_i)$  and let  $x = x_1 x_2 \dots x_n$ ,  $y = y_1 y_2 \dots y_n$ , and  $b = b_1 b_2 \dots b_n$ . We must cope with pathological cases such as that  $x$  simply spells out  $f$  in some programming language. Let  $c$  be an appropriate fixed constant. Show that if  $K(f|x, y) > K(b|x, y) + c$ , then we cannot effectively find  $f$ .

## 5.4 Simple Pac-Learning

The model of pac-learning introduced in 1984 by L.G. Valiant is aimed at describing “feasible” learning in a rigorous manner to remedy lack of precision and objectivity in the ad hoc and vague descriptions of learning used previously. While the pac model is a first step towards a rigorous mathematical formulation of the notion of learning, it has turned out that its requirements are too strong. That is, many tasks that intuitively are learnable, or are clearly learnable in practice, turn out not to be learnable in the pac model. Moreover, there is a certain triviality in the pac learning idea in the sense that pac learning algorithms are just algorithms that may return any hypothesis that is consistent with the examples.

It turns out that using the notion of universal distribution one can both refine the pac learning concept and make it more sensible while also the algorithms required to achieve this new *simple pac-learning* must do something more clever than just be consistent.

### 5.4.1 Pac-Learning

We start with the basic theory of pac-learning. According to commonly accepted views in the theory of computation, “feasibility” means that the learning algorithm should run in polynomial time and use a polynomial number of examples. This requirement implies that not all examples can turn up. Hence, it is impossible to infer a concept precisely. This means that one can only hope to learn the concept approximately. Moreover, one could be presented with unrepresentative examples for the concept to be inferred. It seems reasonable to assume that the examples are drawn randomly from a sample space according to a probability distribution. The approximation between the target concept and the learned concept can be expressed in the probability of the set of examples on which the two concepts disagree.

The sample space  $S$  can be discrete ( $\mathcal{N}$ ) or continuous ( $\mathcal{R}$ ). The elements of  $S$  are called *examples*. A *concept*  $c$  is a subset of  $S$ . Abusing notation, we use  $c$  and the *characteristic function*  $f : S \rightarrow \{0, 1\}$  of  $c$  interchangeably for the concept  $c \subseteq S$  and for the syntactic representation of  $c$ . A *concept class*  $\mathcal{C}$  is a set of concepts.

Consider a concept class  $\mathcal{C}$ . For each concept  $f \in \mathcal{C}$  and example  $v \in S$ ,

$$f(v) = \begin{cases} 1 & \text{if } v \text{ is a positive example of } f, \\ 0 & \text{if } v \text{ is a negative example of } f. \end{cases}$$

The learning algorithm draws examples from the sample space  $S$  according to a fixed but unknown probability distribution  $P$ . Each example in the sample comes with a label “Positive” or “Negative.”

**Definition 5.4.1** A concept class  $\mathcal{C}$  is *pac-learnable* (probably approximately correct learnable) iff there exists a (possibly randomized) learning algorithm  $A$  such

that for each  $f \in \mathcal{C}$  and  $\epsilon (0 < \epsilon < 1)$ , algorithm  $A$  halts in a finite number of steps and examples, and outputs a concept  $h \in \mathcal{C}$  that satisfies the following: With probability at least  $1 - \epsilon$ ,

$$\sum_{f(v) \neq h(v)} P(v) < \epsilon.$$

A concept class is *polynomially pac-learnable* iff it is pac-learnable and the learning algorithm always halts within time and number of examples  $p(l(f), 1/\epsilon)$ , for some polynomial  $p$ .

Examples of concept classes are the set of finite automata; the set of Boolean formulas; the set of  $k$ -DNF formulas (Boolean formulas in disjunctive normal form such that each term has at most  $k$  literals).

#### 5.4.2 Occam's Razor Formalized

Given a set of examples, Occam's razor tells us to choose the simplest concept consistent with the data. However, the simplest concept may be noncomputable (in the sense of shortest effective program) or infeasible to find (in the sense of shortest description in some fixed syntax). For instance, computing a shortest  $k$ -DNF formula consistent with a given set of examples is NP-hard. It turns out that each polynomial time “reasonable approximation” to a shortest rule can be used to achieve polynomial pac-learning.

- **Definition 5.4.2**

Let  $\alpha < 1$  and  $\beta \geq 1$  be constants,  $m$  be the number of examples, and  $s$  be the length (in number of bits) of the smallest concept in  $\mathcal{C}$  consistent with the examples. An *Occam algorithm* is a polynomial time algorithm that finds a hypothesis  $h \in \mathcal{C}$  consistent with the examples and satisfying  $C(h) \leq s^\beta m^\alpha$ .

Intuitively, for a sufficiently large (but still polynomial) set of examples, generated according to the unknown probability distribution  $P$ , any consistent hypothesis that is significantly smaller than the examples will be consistent with most of the unseen examples as well. The commonly used form of Occam algorithm results from replacing  $C(h)$  by the size of  $h$ . Our definition yields a stronger version (which is sometimes required) of *Occam's Razor Theorem* below.

**Theorem 5.4.1** *A concept class  $\mathcal{C}$  is polynomially pac-learnable if there is an Occam algorithm for it.*

**Proof.** Fix an error tolerance  $\epsilon (0 < \epsilon < 1)$ . Choose  $m$  such that

$$m \geq \max \left\{ \left( \frac{2s^\beta}{\epsilon} \right)^{1/(1-\alpha)}, \frac{2}{\epsilon} \log \frac{1}{\epsilon} \right\}. \quad (5.13)$$

This is polynomial in  $s$  and  $1/\epsilon$ .

**Claim 5.4.1** Let  $m$  be as in Equation 5.13. Let  $\mathcal{C}$  be a set of  $r$  concepts, and let  $f$  be one of them. The probability that any concept  $h \in \mathcal{C}$  satisfies  $P(f \neq h) \geq \epsilon$  and is consistent with  $m$  independent examples of  $f$  is less than  $(1 - \epsilon)^m r$ .

**Proof.** Let  $E_h$  be the event that hypothesis  $h$  agrees with all  $m$  examples of  $f$ . If  $P(h \neq f) \geq \epsilon$ , then  $h$  is a *bad hypothesis*. That is,  $h$  and  $f$  disagree with probability at least  $\epsilon$  on a random example. The set of bad hypotheses is denoted by  $B$ . Since the  $m$  examples of  $f$  are independent,

$$P(E_h) \leq (1 - \epsilon)^m.$$

Since there are at most  $r$  bad hypotheses,

$$P\left(\bigcup_{h \in B} E_h\right) \leq (1 - \epsilon)^m r.$$

□

The postulated Occam algorithm finds a hypothesis of Kolmogorov complexity at most  $s^\beta m^\alpha$ . The number  $r$  of hypotheses of this complexity satisfies

$$\log r \leq s^\beta m^\alpha.$$

By assumption on  $m$ ,

$$r \leq (1 - \epsilon)^{-m/2}.$$

(Use  $\epsilon < -\log(1 - \epsilon) < \epsilon/(1 - \epsilon)$  for  $0 < \epsilon < 1$ ). Using the claim, the probability of producing a hypothesis with error larger than  $\epsilon$  is less than

$$(1 - \epsilon)^m r \leq (1 - \epsilon)^{m/2}.$$

Substituting  $m$  we find that the right-hand side is at most  $\epsilon$ . □

**Corollary 5.4.1** According to Definition 5.4.2 with  $\alpha = 0$ , the theorem says that if a learning algorithm compresses the data (of  $m$  examples where  $m$  satisfies Equation 5.13) to a representation of length  $s^\beta$ , that is, length polynomial in the length of the target concept, then that algorithm is a polynomially pac learning algorithm.

Informally, Occam's Razor Theorem says that given a set of positive and negative data, any consistent concept of size "reasonably" shorter than the size of data is an "approximately" correct concept with high probability. That is, if one finds a shorter representation of data, then

one learns. The shorter the conjecture is, the more efficiently it explains the data, hence the more precise the future prediction.

### 5.4.3 Making Pac-Learning Simple

In the pac-learning model we require that the algorithm learn under all distributions. Therefore, pac-learning is also called *distribution-free learning*. It has turned out that many problems are intractable (NP-hard) in this model. Maybe it is only feasible to learn under *some* distributions, like the computable ones. And maybe it is too much to ask to be able to learn all finite automata fast (humans cannot either), but surely we ought to be able to learn a sufficiently *simple* finite automaton fast (as humans can).

Pac learning is not really distribution-free since *only* product distributions are considered: all examples are drawn from the same distribution. This is much different from the general Solomonoff setting we saw before.

Certain concepts that are not known to be polynomially pac-learnable in the distribution-free model can be polynomially pac-learned under the uniform distribution by a specialized learning algorithm. However, learning under one distribution may be too restrictive to be useful. It may be useful to investigate polynomial pac-learnability under various classes of distributions. One would like this class to be wide enough to be interesting, yet restrictive enough to make more concept classes polynomially learnable. An outline of this section is as follows:

**Discrete sample set** Let  $Q$  be a distribution on a discrete sample set.

A concept class is polynomially pac-learnable under all distributions that are multiplicatively dominated by  $Q$  *provided you draw your sample according to  $Q$  in the learning phase* iff the concept class is polynomially pac-learnable under  $Q$ . The caveat is that while pac-learning for distribution  $P$ , we need to draw the sample according to  $Q$ . We shall develop this theory with  $Q = \mathbf{m}$ , the universal enumerable distribution. The remarkable property of this distribution is that in a polynomial sample, with overwhelming probability all examples of logarithmic complexity (the simple examples) will be represented. Hence, a learning algorithm just needs to reconstruct a concept approximately when all simple examples are given. This leads to a new type of learning algorithm.

**Continuous sample set** Let  $\mu$  be a semimeasure on a continuous sample set. A concept class is polynomially pac-learnable under all semimeasures that are multiplicatively dominated by  $\mu$  iff the concept class is pac-learnable under  $\mu$ . Here we have dropped the polynomial time requirement and also eliminated the need for drawing according to  $\mu$  in the learning phase. We will take for  $\mu$  the universal enumerable semimeasure  $\mathbf{M}$  and show that there is a class that

is pac-learnable under  $\mathbf{M}$  but not pac-learnable under all distributions.

#### 5.4.4 Discrete Sample Space

Consider the discrete sample space  $S = \mathcal{N}$ . A distribution  $P$  is *simple* if it is (multiplicatively) dominated by the universal distribution  $\mathbf{m}$ ; see Section 4.3.1. That is, there exists a constant  $c_P$ , such that for all  $x$ ,

$$c_P \mathbf{m}(x) \geq P(x).$$

The first question is how large the class of simple distributions is. It includes the enumerable (and a fortiori the computable) distributions like the uniform distribution, normal distribution, geometric distribution, and Poisson distribution (with computable parameters). It can be shown that there is a distribution that is simple but not enumerable and that there is a distribution that is not simple. See Exercise 5.4.1, page 349.

In the discrete case we needed to modify the standard pac-learning model by the requirement that the sample in the learning phase is drawn according to  $\mathbf{m}$ . A possible justification (which we give for what it is worth) is that in real life the examples are sometimes provided by mechanical or artificial means or good-willed teachers, rather than provided according to the underlying distribution. Naturally, the simpler examples are provided first—that is, more or less according to  $\mathbf{m}(x) = 2^{-K(x)}$ . We prove the following completeness result:

**Theorem 5.4.2** *A concept class  $\mathcal{C}$  is polynomially learnable under the universal distribution  $\mathbf{m}$  iff it is polynomially learnable under each simple distribution  $P$ , provided that in the learning phase the set of examples is drawn according to  $\mathbf{m}$ .*

**Proof.** Since  $P$  is simple, there is a constant  $c_P > 0$  such that for all  $x$ ,

$$c_P \mathbf{m}(x) \geq P(x).$$

Assume  $\mathcal{C}$  is polynomially pac-learnable under distribution  $\mathbf{m}$ . Then one can run the learning algorithms with error parameter  $\epsilon/c_P$  in polynomial time  $t$ . Let  $err$  be the set of strings that are misclassified by the learned concept. So with probability at least  $1 - \epsilon$ ,

$$\sum_{x \in err} \mathbf{m}(x) \leq \epsilon/c_P.$$

Then,

$$\sum_{x \in err} P(x) \leq c_P \sum_{x \in err} \mathbf{m}(x) \leq \epsilon.$$

If the underlying distribution is  $P(\cdot)$  rather than  $\mathbf{m}(\cdot)$ , then we are guaranteed to “pac-learn”  $\mathcal{C}$  (in time  $t$ ), if sampling according to  $\mathbf{m}(\cdot)$ .  $\square$

There are two undesirable aspects of Theorem 5.4.2 in this easy textbook version. The algorithm must know  $c_P = 2^{K(P)}$  depending on the unknown distribution  $P$  in order to determine error parameter  $\epsilon/c_P$ . Secondly, in general we are not dealing with the sample space  $\mathcal{N}$  but with a subset like  $D = \{0, 1\}^n$ . Hence we actually deal with  $\mathbf{m}(\cdot|D)$ . Both problems are technically resolved (essentially by polynomial oversampling) at the cost of somewhat complicating the statement of Theorem 5.4.2 and its proof in [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5 (1991), 915–935].

Since  $\mathbf{m}(\cdot)$  assigns higher probabilities to simpler strings, one could suspect that after polynomially many examples, *all* simple strings are sampled and the probability that is concentrated on the unsampled strings is very low (inverse polynomial). However, this is not the case. Let  $E$  be any set of  $n^c$  examples. Then it can be shown (left to the reader as Exercise 5.4.2, page 349) that

$$\sum_{x \notin E} \mathbf{m}(x) = \Omega\left(\frac{1}{\log^2 n}\right). \quad (5.14)$$

Previous approaches considered only syntactically described classes of concepts. Kolmogorov complexity allows us to introduce the idea of restricting a syntactically described class of concepts to the concepts that are simple in the sense of having short effective descriptions.

**Example 5.4.1** Learning log-DNF. We consider a class not known to be polynomial pac-learnable (under all distributions). We show that this class is polynomially pac-learnable under simple distributions.

**Definition 5.4.3** Consider concepts that are Boolean formulas over variables  $x_1, \dots, x_n$ . A *literal* is either a variable  $x$  or its negation  $\neg x$ . Denote logical disjunction “ $\vee$ ” by “ $+$ ” and logical conjunction “ $\wedge$ ” by concatenation. A *DNF formula* is a Boolean formula  $f$  in disjunctive normal form  $f = m_1 + \dots + m_s$ . The terms  $m_i$  are called *monomials*. Each  $m_i$  is a product of literals. A *CNF formula* is a Boolean formula  $g$  in conjunctive normal form  $g = c_1 \dots c_s$ . The terms  $c_i$  are called *clauses*. Each  $c_i$  is a sum of literals.

The sample space is  $S = \{0, 1\}^n$ . An example vector  $v \in S$  represents a truth assignment to the  $n$  variables: if the  $i$ th element of  $v$  is 1, then  $x_i = \text{true}$ , otherwise  $x_i = \text{false}$ . An example  $v$  is a positive example for  $f$  if its truth assignment makes the whole formula  $f$  true. Denote this by  $f(v) = 1$ . Otherwise  $v$  is a negative example for  $f$ , denoted by  $f(v) = 0$ .

Assume that the positive examples are distributed according to some distribution  $P^+$ , and the negative examples according to some distribution  $P^-$ . The learning algorithm can choose to press a button to obtain an example from  $P^+$ , and it can press a button to obtain an example from  $P^-$ . (This is easily shown to be equivalent to the standard pac model with one distribution of positive and negative examples.) A  $k$ -DNF formula has at most  $k$  literals per monomial. Without loss of generality, let us assume that all  $k$ -DNF (over  $n$  variables) formulas have length at least  $n$ . It is known that  $k$ -DNF is polynomially pac-learnable only for the case when  $k$  is a fixed constant independent of  $n$ .

Let us generalize this a little. Write log-DNF to denote DNF formulas over  $n$  variables where each monomial term  $m$  has complexity  $K(m) = O(\log n)$  and the length of the formula does not exceed a polynomial in  $n$ . This is a nontrivial superset of  $k$ -DNF. It is not known whether log-DNF is polynomially pac-learnable; we will show that it is polynomially pac-learnable under  $\mathbf{m}$  (and hence simple-distribution-free polynomially pac-learnable in the sense of Theorem 5.4.2).

**Lemma 5.4.1** *The class log-DNF is polynomially learnable under  $\mathbf{m}$ .*

**Proof.** Let  $f(x_1, \dots, x_n)$  be a log-DNF with each term of prefix complexity at most  $c \log n$ . If  $m$  is a monomial term in  $f$ , we write  $m \in f$ . Sample  $n^d$  examples, where  $d$  is large enough to satisfy the argument below.

**Claim 5.4.2** With probability greater than  $1 - n^c/e^n$ , all examples of the following form will be drawn:

For each monomial term  $m \in f$ : the example vectors  $0_m$ , defined as the vectors that satisfy  $m$  and have 0 entries for all variables not in  $m$ ; the example vectors  $1_m$ , defined as the vectors that satisfy  $m$  and have 1 entries for all variables not in  $m$ .

**Proof.** Each monomial  $m \in f$  satisfies  $K(m) \leq c \log n$ , and therefore  $K(0_m) \leq c \log n + O(1)$ . Therefore,

$$\mathbf{m}(0_m) \geq 2^{-c \log n - O(1)} \geq n^{-c-1},$$

for large enough  $n$ . This is the probability that  $0_m$  will be sampled in one draw. Let  $E$  be the event that  $0_m$  does not occur in  $n^d$  draws. Then

$$\Pr(E) < (1 - n^{-c-1})^{n^d} \leq e^{-n}/2,$$

for large enough  $d$ . The same estimate holds for the probability that the example  $1_m$  is not sampled in  $n^d$  draws. There are only  $n^c$  possible monomials  $m$  such that  $K(m) \leq c \log n$ . Hence, the probability such that

all vectors  $0_m$  and  $1_m$  associated with such monomials  $m$  are sampled is at least  $1 - n^c/e^n$ .  $\square$

Now we approximate  $f$  by the following learning procedure:

**Step 0** Sample  $n^d$  examples according to  $\mathbf{m}$ . Set POS (NEG) to the set of positive (negative) examples sampled.

**Step 1** For each pair of examples in POS, construct a monomial that contains  $x_i$  if both vectors have “1” in position  $i$ , contains  $\neg x_i$  if both vectors have “0” in position  $i$ , and does not contain variable  $x_i$  otherwise ( $1 \leq i \leq n$ ).

**Step 2** Among the monomials constructed in Step 1, delete the ones that imply examples in NEG. {The remainder forms a set  $M$ .}

**Step 3** Set  $A_m = \{v : m(v) = 1\}$ . {That is,  $A_m$  is the set of positive examples implied by monomial  $m$ .} Use a greedy set-cover algorithm to find a small set, say  $C$ , of monomials  $m \in M$  such that  $\bigcup_{m \in C} A_m$  covers all positive examples in POS.

We need to prove the correctness of the algorithm.

**Claim 5.4.3** With probability greater than  $1 - n^c/e^n$ , the set of monomials  $\{m : m \in f\} \subseteq M$ .

**Proof.** By Claim 5.4.2, with probability at least  $1 - n^c/e^n$ , we draw all vectors  $1_m$  and  $0_m$  such that monomial  $m$  has prefix complexity at most  $c \log n$ . In Step 1 of the algorithm we form the monomial  $m$  from examples  $1_m$  and  $0_m$ . Thus, with probability at least  $1 - n^c/e^n$ , all monomials of  $f$  belong to  $M$ .  $\square$

Of course, many other monomials consistent with the examples may also be in  $M$ . Finding all of the original monomials of  $f$  precisely is NP-hard. For the purpose of learning it is sufficient to approximate  $f$ . We use the following well-known approximation of set-cover result.

**Claim 5.4.4** Let  $A_1, \dots, A_n$  be sets such that  $\bigcup_{i=1}^n A_i = A = \{1, \dots, q\}$ . If there exist  $k$  sets  $A_{i_1}, \dots, A_{i_k}$  such that  $A = \bigcup_{j=1}^k A_{i_j}$ , then it is possible to find in polynomial time  $l = O(k \log q)$  sets  $A_{h_1}, \dots, A_{h_l}$  such that  $A = \bigcup_{j=1}^l A_{h_j}$ .

Let  $f$  have  $k$  monomials. The  $k$  monomials cover the positive examples in the sense that  $\text{POS} \subseteq \bigcup_{m \in f} A_m$ . There are  $2^n$  examples. By Claim 5.4.4, we can find  $O(kn)$  monomials to approximate  $f$  and cover the examples

in POS in polynomial time. Occam's Razor Theorem 5.4.1 implies that our algorithm polynomially learns log-DNF.  $\square$

Step 3 of the above algorithm is needed to compress the data in order to apply Occam's Razor Theorem. One may wonder whether one can encode each monomial as binary vectors efficiently, and hence sample all binary vectors of prefix complexity  $c \log n$ ; then decode these into monomials to get all monomials of prefix complexity  $c \log n$ ; then run the set-cover algorithm to choose a small set of monomials to achieve learning. But this does not work since there are  $2^n$  0-1 vectors and  $3^n$  monomials of  $n$  variables. It can be shown that there is no effective encoding scheme that selectively codes only  $2^n$  monomials including all monomials of prefix complexity  $c \log n$ .

$\diamond$

If  $\mathbf{m}$  is given as a table, then one can randomly sample according to  $\mathbf{m}$  as explained in Example 4.4.3 on page 271. Unfortunately  $\mathbf{m}$  is not computable. This problem is eliminated by using a time-bounded version of  $\mathbf{m}$ . Then the entire theory still holds in a scaled down version as discussed in Example 7.6.2 on page 508.

### 5.4.5 Continuous Sample Space

Consider a continuous sample space  $S = \{0, 1\}^\infty$ , that is, the set of one-way infinite binary sequences or, equivalently, the set of real numbers in the unit interval  $[0, 1]$  (Section 1.6). For example, the uniform distribution now is defined as  $\lambda(\Gamma_x) = 2^{-l(x)}$ , where  $\Gamma_x$  denotes the set of all one-way infinite binary strings starting with  $x$ . This is the uniform measure on the interval  $[0, 1]$ .

We will use  $\mathbf{M}$ , the continuous universal enumerable measure of Theorem 4.3.1, page 247.

**Definition 5.4.4** A measure  $\mu$  over  $S$  is *simple* if it is dominated multiplicatively by  $\mathbf{M}$ , in the sense of  $\mu(x) = O(\mathbf{M}(x))$ .

Similar to the discrete case one can show that there are simple measures that are not enumerable, and there are measures that are not simple.

A concept class is a subset  $\mathcal{C} \subseteq 2^S$  of concepts, each of which is a measurable set. If  $c$  is a concept to be learned, then  $x \in S$  is a positive example if  $x \in c$ , and it is a negative example if  $x \in S - c$ . The remaining definitions of pac-learning can now be rephrased in the continuous setting in the obvious way. If  $c, c'$  are two sets, then  $c\Delta c'$  is the symmetric set difference  $(c \cup c') - (c \cap c')$ .

While for discrete sample spaces all concept classes are pac-learnable (although not all are polynomially pac-learnable), this is not the case for continuous sample spaces. Here we show that all continuous concept

classes are pac-learnable over each simple measure  $\mu$  iff they are pac-learnable under the universal measure  $\mathbf{M}$ . In contrast with polynomial pac-learning of discrete concepts, we do not need to require (but do allow) that the learning algorithm samples according to the universal measure.

**Theorem 5.4.3** *A concept class  $\mathcal{C}$  of concepts in  $S$  is pac-learnable under  $\mathbf{M}$  iff it is learnable under each simple measure.*

**Proof.** (IF ) This holds vacuously.

(ONLY IF) We say that  $\mathcal{C}_\epsilon \subseteq 2^S$  is an  $\epsilon$ -cover of  $\mathcal{C}$ , with respect to distribution  $\mu$ , if for every  $c \in \mathcal{C}$  there is a  $\hat{c} \in \mathcal{C}_\epsilon$  that is  $\epsilon$ -close to  $c$  ( $\mu(c\Delta\hat{c}) < \epsilon$ ). A concept class  $\mathcal{C}$  is *finitely coverable* if for every  $\epsilon > 0$  there is a finite  $\epsilon$ -cover  $\mathcal{C}_\epsilon$  of  $\mathcal{C}$ , everything with respect to a given measure  $\mu$ . A finite  $\epsilon$ -cover  $\mathcal{C}_\epsilon$  has finitely many concepts  $c$ , and each  $c$  is in the closure of the set of cylinders under finite union and complement (and finite intersection).

**Claim 5.4.5** A concept class  $\mathcal{C}$  is finitely coverable with respect to  $\mu$  iff  $\mathcal{C}$  is learnable with respect to  $\mu$ .

**Proof.** The proof doesn't use Kolmogorov complexity, so there is no reason to give it here. It is due to G. Benedek and A. Itai, *Theoret. Comput. Sci.*, 86:2(1991), 377–390; also reproduced in [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20(1991), 915–935].  $\square$

Using Claim 5.4.5, if  $\mathcal{C}$  can be learned under  $\mathbf{M}$ , then  $\mathcal{C}$  can be finitely covered with respect to  $\mathbf{M}$ . Let  $\mu$  be a simple distribution and  $d > 0$  such that  $\mathbf{M}(x) \geq d\mu(x)$  for all  $x$ . Then any finite  $\epsilon d$ -cover of  $\mathcal{C}$  with respect to  $\mathbf{M}$  is also a finite  $\epsilon$ -cover with respect to  $\mu$ . Using Claim 5.4.5 again, it follows that  $\mathcal{C}$  is learnable with respect to  $\mu$ . This finishes the proof of the theorem.  $\square$

Note that this is a strong statement since we are saying that if one can learn under  $\mathbf{M}$ , then one can also learn under any simple measure  $\mu$ , while sampling according to  $\mu$ . In the polynomial learning of discrete concepts we had to require sampling according to  $\mathbf{m}$ . This improvement is made possible by relaxing the “polynomial time learning” requirement to “learning.”

**Example 5.4.2** If all continuous concept classes that can be simple pac-learned can also be pac-learned, then Theorem 5.4.3 would be vacuously true. However, there exist continuous concept classes that can be simple pac-learned but cannot be pac-learned under all distributions. Let the witness class  $\mathcal{C}$  consist of all concepts of the form  $c = \bigcup\{\Gamma_x : x \in I\}$ , where  $I$  satisfies the condition that if  $x, y \in I$  and  $x \neq y$ , then either  $\mathbf{M}(x) \geq 2\mathbf{M}(y)$

or vice versa. (As usual, the cylinder  $\Gamma_x$  is the set of all infinite binary sequences starting with  $x$ , and  $\mathbf{M}(x)$  is the  $\mathbf{M}$ -measure of  $\Gamma_x$ .)

Given a continuous concept class  $\mathcal{C}$  (as defined in Section 5.4.3) and a finite set  $E \subseteq S = \{0, 1\}^\infty$ . If  $\{E \cap c : c \in \mathcal{C}\} = 2^E$ , then we say  $E$  is *shattered* by  $\mathcal{C}$ . The *Vapnik-Chervonenkis (VC) dimension* of  $\mathcal{C}$  is the smallest integer  $d$  such that no  $E \subseteq S$  of cardinality  $d + 1$  is shattered by  $\mathcal{C}$ ; if no such  $d$  exists, then the dimension of  $\mathcal{C}$  is infinite. It is known that a class  $\mathcal{C}$  has finite VC-dimension iff it is pac-learnable.

This was proven by A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth *J. ACM*, 35:4(1989), 929–965. The fundamental Occam’s Razor Theorem also appears there.

It can be easily shown, and it is left to the reader as Exercise 5.4.8 on page 350, that  $\mathcal{C}$  has an infinite VC-dimension and therefore is not pac-learnable under all measures. On the other hand,  $\mathcal{C}$  is finitely coverable under  $\mathbf{M}$ , and therefore pac-learnable with respect to  $\mathbf{M}$  (and hence under all simple measures).  $\diamond$

## Exercises

**5.4.1.** [27] Consider discrete distributions. Show that there is a nonenumerable distribution that is simple and there is a distribution that is not simple.

*Comments.* Therefore, the simple distributions properly contain the enumerable distributions but do not include all distributions. Source: this and the next three exercises are from [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5(1991), 915–935].

**5.4.2.** [24] Prove Equation 5.14. Can you improve this bound?

**5.4.3.** [30] Consider DNF over  $n$  variables. A DNF formula  $f$  is *simple* if, for each term  $m$  of  $f$ , there is a vector  $v_m \in \{0, 1\}^n$  that satisfies  $m$  but does not satisfy any other monomials of  $f$  even by changing one bit and  $K(v_m) = O(\log n)$ . Simple DNFs can contain many high prefix complexity terms as opposed to  $O(\log n)$ -DNF. For example, take  $y \in \{0, 1\}^n$  with  $K(y) \geq n - O(1)$ . Then the number of 1’s in  $y$  is about  $n/2$ . Construct a term  $m$  containing  $x_i$  if  $y_i = 1$  and neither  $x_i$  or  $\neg x_i$  otherwise. Then  $K(m) \geq n - O(1)$  but the vector  $1_m$  consisting of all 1’s satisfies  $m$  and has  $K(1_m) = O(\log n)$ . The class of simple DNF is pretty general. Show that it is polynomially learnable under  $\mathbf{m}$ .

**5.4.4.** [25] Show that log-DNF of Example 5.4.1 does not contain and is not contained in simple DNF of Exercise 5.4.3.

**5.4.5.** [33] A *k*-decision list over  $n$  variables is a list of pairs  $L = (m_1, b_1), \dots, (m_s, b_s)$ , where  $m_i$  is a monomial of at most  $k$  variables and

$b_i \in \{0, 1\}$ , for  $1 \leq i \leq s$ , except that always  $m_s = 1$ . A decision list  $L$  represents a Boolean function  $f_L$  defined as follows: For each example  $v \in \{0, 1\}^n$ , let  $f_L(v) = b_i$  where  $i$  is the least index such that  $v$  satisfies  $m_i$ . Since there are at most  $(2n)^{k+1}$  monomials of  $k$  literals over  $n$  variables, we have  $s \leq (2n)^{k+1}$ .

(a) Using Occam's Razor Theorem 5.4.1, we show that  $k$ -decision lists are polynomially pac-learnable.

(b) Let us define *log-decision* list to be a decision list with each term having prefix complexity  $O(\log n)$ . Show that log-decision list is polynomially learnable under  $\mathbf{m}$ .

*Comments.* Source: Item (a) [R. Rivest, *Machine Learning*, 2:3(1987), 229–246] Item (b) was stated as an open problem in the first edition of this book and was solved by J. Castro and J.L. Balcázar [pp. 239–248 in: *Proc. 6th Int'l Workshop on Algorithmic Learning Theory, Lect. Notes Artific. Intell.*, Vol. 997, Springer-Verlag, Berlin, 1995]. They also show that simple decision lists, a generalization of the simple DNF of Exercise 5.4.3, are polynomial learnable under  $\mathbf{m}$ .

**5.4.6.** [O35] Are any of log-DNF, simple DNF, log-decision list polynomially pac-learnable?

**5.4.7.** [35] A Boolean formula is *monotone* if no literal in it is negated. A  $k$ -term DNF is a DNF consisting of at most  $k$  monomials.

(a) Show that pac-learning monotone  $k$ -term DNF requires more than polynomial time unless RP=NP.

(b) Show that the class of monotone  $k$ -term DNF is polynomially pac-learnable under  $\mathbf{m}$  (even with the algorithm learning the exact concept with high probability).

*Comments.* This shows that more concepts are simple pac-learnable than are pac-learnable unless RP=NP. Source: for Item (a) L. Pitt and L.G. Valiant, *J. ACM*, 35(1989), 965–984; for Item (b) M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5(1991), 915–935.

**5.4.8.** [M36] Show that the continuous concept class  $\mathcal{C}$  defined in Example 5.4.2 is is pac-learnable under all simple measures but not pac-learnable (that is, under all measures).

*Comments.* Source: M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5 (1991), 915–935.

## 5.5 Hypothesis Identification by Minimum Description Length

---

We can formulate scientific theories in two steps. First, we formulate a set of possible alternative hypotheses, based on scientific observations or other data. Second, we select one hypothesis as the most likely one. Statistics is the mathematics of how to do this. A relatively recent paradigm in statistical inference was developed by J.J. Rissanen and by C.S. Wallace and his coauthors. The method can be viewed as a computable approximation to the noncomputable approach in Section 5.2 and was inspired by it. In accordance with Occam’s dictum, it tells us to go for the explanation that compresses the data the most.

**Minimum description length (MDL) principle.** Given a sample of data and an effective enumeration of the appropriate alternative theories to explain the data, the best theory is the one that minimizes the sum of

- the length, in bits, of the description of the theory;
- the length, in bits, of the data when encoded with the help of the theory.

The idea of a two-part code for a body of data  $D$  is natural from the perspective of Kolmogorov complexity. If  $D$  does not contain any regularities at all, then it consists of purely random data, and there is no hypothesis to identify. Assume that the body of data  $D$  contains regularities. With help of a description of those regularities (a model) we can describe the data compactly. Assuming that the regularities can be represented in an effective manner (that is, by a Turing machine), we encode the data as a program for that machine. Squeezing all effective regularity out of the data, we end up with a Turing machine representing the meaningful regular information in the data together with a program for that Turing machine representing the remaining meaningless randomness of the data. This is the intuition.

Formally, assume that our candidate theories are effective computation procedures, that is, Turing machines. While with the Kolmogorov complexity we represent the amount of information in the individual objects, here we are interested in the division of this information into a two-part code, as in Section 2.1.1. First, the “valuable” information representing regularities that are presumably usable (the model part of length  $K(T)$ ), followed by the “useless” random part of length  $C(D|T) = l(p)$  as in Equation 2.1 on page 99. However, it is difficult to find a valid mathematical way to force a sensible division of the information at hand in a meaningful part and a meaningless part. One way to proceed is suggested by the analysis below.

MDL is based on striking a balance between regularity and randomness in the data. All we will ever see are the data at hand (if we know more, then in fact we possess more data, which should be used as well). The best model or explanatory theory is taken to be the one that optimally uses regularity in the data to compress. “Optimally” is used in the sense

of maximal compression of the data in a two-part code, that is, a model and a description from which the data can be reproduced with help of the model. If the data are truly random, no compression is possible and the optimum is reached for the “empty” model. The empty model has description length 0. If the data are regular, then compression is possible, and using the MDL principle (or the Kolmogorov complexity approach) identifies the optimal model.

We call such a “model,” or “theory,” a “hypothesis.” With a more complex description of the hypothesis  $H$ , it may fit the data better and therefore decrease the misclassified data. If  $H$  describes all the data, then it does not allow for measuring errors. A simpler description of  $H$  may be penalized by an increase in misclassified data. If  $H$  is a trivial hypothesis that contains nothing, then all data are described literally and there is no generalization. The rationale of the method is that a balance between these extremes seems to be required.

Ideally, the description lengths involved should be the shortest effective description lengths, the prefix complexities, which however cannot be effectively computed. (This obviously impedes actual use. In practice, one needs to consider computable approximations to shortest descriptions, for example by restricting the allowable approximation time.) The code of the shortest effective self-delimiting descriptions, the prefix complexity code, gives the least expected code-word length—close to the entropy (pages 181, 231 or Section 8.1)—and moreover compresses the regular objects until all regularity is squeezed out. All shortest effective descriptions are completely random themselves, without any regularity whatsoever. Kolmogorov complexity can be used to develop a theory of (idealized) minimum description length reasoning. We rigorously derive and justify this Kolmogorov complexity based form of minimum description length, “ideal MDL,” via the Bayesian approach using the universal distribution  $\mathbf{m}(\cdot)$  of Section 4.3.1 as the particular prior distribution over the hypotheses. This leads to a mathematical explanation of correspondences and differences between ideal MDL and Bayesian reasoning, and in particular it gives some evidence under what conditions the latter is prone to overfitting while the former is not.

The analysis of both hypothesis identification by ideal MDL in this section, and of prediction in Section 5.2.2, shows that maximally compressed descriptions give good results on data samples that are random with respect to probabilistic hypotheses. These data samples form the overwhelming majority and occur with probability going to one when the length of the data sample grows unboundedly. That is, both for hypothesis identification and prediction, data compression is provably optimal but for a subset of (hypothesis, data sample) pairs of vanishing probability.

### 5.5.1 Derivation of MDL from Bayes's Rule

Let us see how we can rigorously derive an ideal version of MDL from first principles, *in casu* Bayes's rule as given by Equation 5.1. While this does not rigorously prove anything about MDL in its applied forms, the relations with Bayesian inference we establish for ideal MDL are corroborated by empirical evidence for applied MDL.

Recall Bayes's Rule, Sections 1.6, 1.10, 5.2:

$$\Pr(H|D) = \frac{\Pr(D|H)P(H)}{\Pr(D)}.$$

Here  $H$  is a hypothesis,  $P$  is the prior probability of the hypotheses and  $D$  is the set of observed data. In this equation we are only concerned with finding the  $H$  that maximizes  $\Pr(H|D)$  with  $D$  and  $P$  fixed. Taking the negative logarithm of both sides of the equation, this is equivalent to *minimizing* the expression  $-\log \Pr(H|D)$  over  $H$ :

$$-\log \Pr(H|D) = -\log \Pr(D|H) - \log P(H) + \log \Pr(D). \quad (5.15)$$

Since the probability  $\Pr(D)$  is constant under varying  $H$ , this means we want to find an hypothesis  $H$  that minimizes

$$-\log \Pr(D|H) - \log P(H). \quad (5.16)$$

In applied forms of MDL one roughly speaking interprets these negative logarithms  $-\log P(x)$  as the corresponding Shannon-Fano (or Huffman) code-word lengths. But why should one use the Shannon-Fano code (or Huffman code) and no other code reaching an expected code word length equal to the entropy? In particular, ignoring feasibility, why not use the objective shortest effective code, the shortest effective descriptions with code-word length set equal to the prefix complexities. This also has an expected code-word length equal to the entropy (pages 181, 231 or Section 8.1), but additionally, the shortest program compresses the object by effectively squeezing out and accounting for all regularities in it. The resulting code word is maximally random, that is, it has maximal prefix complexity.

For now let us assume that  $H$  and  $D$  are expressed as natural numbers or finite binary strings. To obtain the ideal MDL principle it suffices to replace the probabilities involved in Equation 5.16 by the *universal probability*  $\mathbf{m}(\cdot)$  of Theorem 4.3.1 on page 247. The analysis of the *conditions* under which this substitution is *justified*, or conversely, how application of ideal MDL is equivalent to Bayesian inference using *admissible* probabilities, is deferred to the next section. Therefore, under conditions to be established below, we substitute according to

$$\log P(H) = \log \mathbf{m}(H) + O(1), \quad (5.17)$$

$$\log \Pr(D|H) = \log \mathbf{m}(D|H) + O(1).$$

By Theorem 4.3.3 on page 253, we have

$$\begin{aligned}\log \mathbf{m}(H) &= -K(H) \pm O(1), \\ \log \mathbf{m}(D|H) &= -K(D|H) \pm O(1),\end{aligned}$$

where  $K(\cdot)$  is the prefix complexity. Therefore, using the substitution of Equation 5.17 we can replace the sum of Equation 5.16 by the sum of the minimum lengths of effective self-delimiting programs that compute descriptions of  $H$  and  $D|H$ . That is, we look for the  $H$  that minimizes

$$K(D|H) + K(H), \quad (5.18)$$

which is the code-independent, recursively invariant, absolute form of the MDL principle.

The term  $-\log P(D|H)$  is also known as the *self-information* in information theory and the *negative log-likelihood* in statistics. It can be regarded as the number of bits it takes to redescribe or encode  $D$  with an ideal code relative to  $H$ .

**Example 5.5.1** If we replace all  $P$ -probabilities in Equation 5.15 by the corresponding  $\mathbf{m}$ -probabilities, we obtain in the same way by Theorem 4.3.3

$$K(H|D) = K(H) + K(D|H) - K(D) + O(1).$$

In Theorem 3.9.1 on page 232 it is shown that symmetry of information holds for individual objects in the following sense:

$$\begin{aligned}K(H, D) &= K(H) + K(D|H, K(H)) + O(1) \\ &= K(D) + K(H|D, K(D)) + O(1).\end{aligned} \quad (5.19)$$

Substitution gives  $K(H|D) = K(H, D) - K(D)$ , up to an  $O(\log K(H, D))$  additive term. The term  $K(D)$  is fixed and doesn't change for different  $H$ 's. Minimizing the left-hand term,  $K(H|D)$  can then be interpreted as

**Alternative formulation MDL principle.** Given a hypothesis space  $\mathbf{H}$ , we want to select the hypothesis  $H$  such that the length of the shortest encoding of  $D$  together with hypothesis  $H$  is minimal.

◇

The discussion seems to have arrived at its goal, but a danger of triviality lurks nearby. Yet it is exactly the solution how to prevent trivialities, which gives us the key to the very meaning of ideal MDL, and by extension some insight in applied versions.

Since it is not more difficult to describe some object if we get more conditional information, we have  $K(D|H, K(H)) \leq K(D|H) + O(1)$ . Thus, by Equation 5.19 the quantity in Equation 5.18 satisfies

$$K(H) + K(D|H) \geq K(H, D) + O(1) \geq K(D) + O(1),$$

with equalities for the trivial hypothesis  $H_0 = D$  or  $H_0 = \emptyset$ . At first glance this would mean that the hypothesis  $H_{\text{mdl}}$  that minimizes the sum of Equation 5.18 could be set to  $D$  or  $\emptyset$ , which is absurd in general. However, we have only derived the validity of Equation 5.18 under the condition that Equation 5.17 holds. The crucial part of our justification of MDL is to establish precisely when Equation 5.17 is valid.

## 5.5.2 The Role of Universal Probability

It is well known, see Example 1.11.2, that the so-called Shannon-Fano code for an ensemble of source words distributed according to probability  $Q$  is a prefix code  $E_Q$  with  $l(E_Q(x)) = -\log Q(x)$  satisfying

$$\sum_x Q(x)l(E_Q(x)) = \min_{E'} \left\{ \sum_x Q(x)l(E'(x)) : E' \text{ is a prefix code} \right\} + O(1),$$

that is, it realizes the least expected code-word length among all prefix codes (the entropy of  $Q(\cdot)$ ) by Shannon's Noiseless Coding Theorem. Therefore, the  $H$  which minimizes Equation 5.16, that is,

$$l(E_{\Pr(\cdot|H)}(D)) + l(E_P(H))$$

minimizes the sum of two prefix codes which both have shortest *expected* code word lengths.

But there are many prefix codes which have expected code-word length equal to the entropy. Among those prefix codes there is one which gives the *shortest effective code word* to each individual source word: the prefix code with code word length  $K(x)$  for object  $x$ . In ideal MDL we want minimize the sum of the effective description lengths of the *individual* elements  $H, D$  involved. This means using the *shortest effective description lengths*, as in Equation 5.18. However, we cannot simply replace negative logarithms in Equation 5.16 by corresponding  $K(\cdot)$  terms. We can only do so if Equation 5.17 holds.

To satisfy Equation 5.17 we are free to make the new *assumption* that the prior probability  $P(\cdot)$  in Bayes's rule Equation 5.1 is fixed as  $\mathbf{m}(\cdot)$ . Whether this can be justified or not is a question which we address in Section 5.5.7.

However, we *cannot assume* that the probability  $\Pr(\cdot|H)$  equals  $\mathbf{m}(\cdot|H)$ . Namely, as explained at length in Section 5.1.3, probability  $\Pr(\cdot|H)$  may be totally determined by the hypothesis  $H$ . Depending on  $H$  therefore,  $l(E_{\Pr(\cdot|H)}(D))$  may be *very* different from  $K(D|H)$ . This holds especially for "simple" data  $D$  that have low probability under assumption of hypothesis  $H$ .

**Example 5.5.2** Let us look at a simple example evidencing this discrepancy. Suppose we flip a coin of unknown bias  $n$  times. Let hypothesis  $H$  and data  $D$  be defined by:

$$\begin{aligned} H &:= \text{Probability of "heads" is } \frac{1}{2}, \\ D &:= \underbrace{hh\dots h}_{n \text{ times "h"}(\text{eads})} \end{aligned}$$

Then we have  $\Pr(D|H) = 1/2^n$  and

$$l(E_{\Pr(\cdot|H)}(D)) = -\log \Pr(D|H) = n.$$

In contrast,  $K(D|H) \leq \log n + 2 \log \log n + O(1)$ .  $\diamond$

### 5.5.3 Fundamental Inequality

The theory dealing with randomness of individual objects states that under certain conditions  $-\log \Pr(D|H)$  and  $K(D|H)$  are close. The first condition is that  $\Pr(\cdot|\cdot)$  be a recursive function. That is, it can be computed to any required precision for each argument  $D$  and conditional  $H$ . Then we can appeal to the following known facts: Firstly, by Example 4.3.3 on page 249 following Theorem 4.3.1 on page 247,

$$\mathbf{m}(D|H) \geq 2^{-K(\Pr(\cdot|H))} \Pr(D|H).$$

Therefore,

$$\log \frac{\mathbf{m}(D|H)}{\Pr(D|H)} \geq -K(\Pr(\cdot|H)) \geq -K(H) + O(1). \quad (5.20)$$

The last inequality arises since from  $H$  we can compute  $\Pr(\cdot|H)$  by assumption on  $\Pr(\cdot|\cdot)$ .

Secondly, if the data sample  $D$  is sufficiently *random* with respect to the recursive distribution  $\Pr(\cdot|H)$  (with respect to  $H$  therefore) in the sense of Martin-Löf (and only if it is so), we have

$$\log \frac{\mathbf{m}(D|H)}{\Pr(D|H)} \leq 0, \quad (5.21)$$

where  $\kappa_0((D|H)|\Pr(\cdot|H)) = \log(\mathbf{m}(D|H)/\Pr(D|H))$  is a “universal sum *P*-test” as in Theorem 4.3.5, page 258. The overwhelming majority of  $D$ ’s is random in this sense because for each  $H$  we have

$$\sum_D \Pr(D|H) 2^{\kappa_0((D|H)|\Pr(\cdot|H))} = \sum_D \mathbf{m}(D|H) \leq 1,$$

since  $\mathbf{m}(\cdot|H)$  is a probability distribution. For  $D$ ’s that are random in the appropriate sense, Equations 5.20 and 5.21 mean by Equation 5.17 that

$$K(D|H) - K(\Pr(\cdot|H)) \leq -\log \Pr(D|H) \leq K(D|H).$$

Above, by assuming that the a priori probability  $P(H)$  of hypothesis  $H$  is in fact the universal probability, we obtained  $-\log P(H) = \mathbf{m}(x)$ . However, we do not need to make this assumption. For recursive  $P(\cdot)$ , we can analyze the situation when  $H$  is random in the required sense with respect to  $P(\cdot)$ . The first inequality below holds since  $\mathbf{m}(\cdot)$  majorizes  $P(\cdot)$  this way; the second inequality expresses the assumption of randomness of  $H$ ,

$$\begin{aligned}\mathbf{m}(H) &\geq 2^{-K(P)}P(H), \\ \log(\mathbf{m}(H)/P(H)) &\leq 0,\end{aligned}$$

where  $K(P)$  is the length of the shortest self-delimiting program for the reference universal prefix machine to simulate the Turing machine computing the probability density function  $P : \mathcal{N} \rightarrow [0, 1]$ . That is, it is the shortest effective self-delimiting description of  $P$ . Then,

$$K(H) - K(P) \leq -\log P(H) \leq K(H). \quad (5.22)$$

Altogether, we find

$$\begin{aligned}K(D|H) + K(H) - \alpha(P, H) &\leq -\log \Pr(D|H) - \log P(H) \\ &\leq K(D|H) + K(H),\end{aligned} \quad (5.23)$$

with

$$\alpha(P, H) = K(\Pr(\cdot|H)) + K(P),$$

and we note that  $K(\Pr(\cdot|H)) \leq K(H) + O(1)$ . We call Equation 5.23 the **Fundamental Inequality (FI)** because it describes the fundamental relation between Bayes's Rule and MDL in mathematical terms. It is left for us to interpret it.

#### 5.5.4 Validity Range of FI

We begin by stressing that Equation 5.23 holds only in case simultaneously  $H$  is  $P$ -random and  $D$  is  $\Pr(\cdot|H)$ -random. What is the meaning of this?

$H$  is  $P$ -random means that the true hypothesis must be “typical” for the prior distribution  $P$  in the sense that it must belong to all effective majorities (sets on which the majority of  $P$ -probability is concentrated). In Example 4.3.10 on page 261 it is shown that this is the set of  $H$ 's such that  $K(H) \approx -\log P(H)$ . In case  $P(H) = \mathbf{m}(H)$ , that is, the prior distribution equals the universal distribution, then for *all*  $H$  we have  $K(H) = -\log P(H)$ , that is, all hypotheses are random with respect to the universal distribution.

Let us look at an example of a distribution where some hypotheses are random, and some other hypotheses are nonrandom. Let the possible hypotheses correspond to the binary strings of length  $n$ , while  $P$

is the uniform distribution that assigns probability  $P(H) = 1/2^n$  to each hypothesis  $H \in \{0, 1\}^n$ . Then  $H = 00\dots 0$  has low complexity  $K(H) \leq \log n + 2 \log \log n$ . However,  $-\log P(H) = n$ . Therefore, by Equation 5.22,  $H$  is not  $P$ -random. If we obtain  $H$  by  $n$  flips of a fair coin, then with overwhelming probability we will have that  $K(H) = n + O(\log n)$ , and therefore  $-\log P(H) \approx K(H)$  and  $H$  is  $P$ -random.

$D$  is  $\Pr(\cdot|H)$ -random means that the data are random with respect to the probability distribution  $\Pr(\cdot|H)$  induced by the hypothesis  $H$ . Therefore, we require that the sample data  $D$  are “typical,” that is, “randomly distributed” with respect to  $\Pr(\cdot|H)$ . If, for example,  $H = (\mu, \sigma)$  induces the Gaussian distribution  $\Pr(\cdot|(\mu, \sigma)) = N(\mu, \sigma)$  and the data  $D$  are concentrated in a tail of this distribution, like  $D = 00\dots 0$ , then  $D$  is atypical with respect to  $\Pr(\cdot|H)$  in the sense of being nonrandom because it violates the  $\Pr(\cdot|H)$ -randomness test Equation 5.21.

Note that an hypothesis satisfying FI is a prefix complexity version of the Kolmogorov minimal sufficient statistic of Section 2.2.2. This shows the connections between MDL, Bayes’s rule, and the Kolmogorov minimal sufficient statistic.

### 5.5.5 If Optimal Hypothesis Violates FI

The only way to violate the Fundamental Inequality is that either  $D$  is not  $\Pr(\cdot|H)$ -random and by Equation 5.21,  $-\log \Pr(D|H) > K(D|H)$ , or that  $H$  is not  $P$ -random and by Equation 5.22,  $-\log P(H) > K(H)$ . We give an example of the first case:

We sample a polynomial  $H_2 = ax^2 + bx + c$  at  $n$  arguments chosen uniformly at random from the interval  $[0, 1]$ . The sampling process introduces Gaussian errors in the function values obtained. The set of possible hypotheses is the set of polynomials. Assume that all numbers involved are of fixed bounded accuracy.

Because of the Gaussian error in the measuring process, with overwhelming probability the only polynomials  $H_{n-1}$  that fit the sample precisely are of degree  $n - 1$ . Denote the data sample by  $D$ . Now, this hypothesis  $H_{n-1}$  is likely to minimize  $K(D|H)$  since we just have to describe the  $n$  lengths of the intervals between the sample points along the graph of  $H_{n-1}$ . However, for this  $H_{n-1}$  data sample  $D$  is certainly not  $\Pr(\cdot|H_{n-1})$ -random, since it is extremely unlikely, and hence atypical, that  $D$  arises when sampling  $H_{n-1}$  with Gaussian error. Therefore, Equation 5.21 is violated, which means that  $-\log \Pr(D|H_{n-1}) > K(D|H_{n-1})$ , contrary to what we used in deriving the Fundamental Inequality. With prior probability  $P(\cdot) := \mathbf{m}(\cdot)$ , which means  $-\log P(\cdot) = K(\cdot) + O(1)$ , this moreover violates the Fundamental Inequality.

In contrast, with overwhelming likelihood  $H_2$  will show the data sample  $D$  random to it. That being the case, the Fundamental Inequality holds.

Now what happens if  $H_{n-1}$  is the true hypothesis and the data sample  $D$  by chance is as above? In that case the Fundamental Inequality is violated and Bayes's Rule and MDL may each select very different hypotheses as the most likely ones, respectively.

### 5.5.6 If Optimal Hypothesis Satisfies FI

Given data sample  $D$  and prior probability  $P$ , we call a hypothesis  $H$  *admissible* if the Fundamental Inequality Equation 5.23 holds. Restriction to the set of *admissible* hypotheses excludes setting  $K(D|H) + K(H) \approx K(D)$  for trivial hypothesis (like  $H = \emptyset$  or  $H = D$ , which are not admissible).

**Theorem 5.5.1** *Let  $\alpha(P, H)$  as in Equation 5.23 be small. Then Bayes's Rule and ideal MDL are optimized (or almost optimized) by the same hypothesis among the admissible  $H$ 's. That is, there is one admissible  $H$  that simultaneously (almost) minimizes both  $-\log \Pr(D|H) - \log P(H)$  (selection according to Bayes's Rule) and  $K(D|H) + K(H)$  (selection according to MDL).*

**Proof.** The smallness of  $\alpha(P, H)$  means that both the prior distribution  $P$  is simple, and that the probability distribution  $\Pr(\cdot|H)$  over the data samples induced by hypothesis  $H$  simple. In contrast, if  $\alpha(P, H)$  is large, which means that either of the mentioned distributions is not simple, for example when  $K(\Pr(\cdot|H)) = K(H)$  for complex  $H$ , then there may be some discrepancy. Namely, in Bayes's Rule our purpose is to maximize  $\Pr(H|D)$ , and the hypothesis  $H$  that minimizes  $K(D|H) + K(H)$  also maximizes  $\Pr(H|D)$  up to a  $2^{-\alpha(P, H)}$  multiplicative factor. Conversely, the  $H$  that maximizes  $\Pr(H|D)$  also minimizes  $K(D|H) + K(H)$  up to an additive term  $\alpha(P, H)$ . That is, with

$$\begin{aligned} H_{\text{mdl}} &:= \min_{H \in \mathcal{H}} \{K(D|H) + K(H)\}, \\ H_{\text{bayes}} &:= \max_{H \in \mathcal{H}} \{\Pr(H|D)\}, \end{aligned} \quad (5.24)$$

we have

$$\begin{aligned} 2^{-\alpha(P, H)} &\leq \frac{\Pr(H_{\text{mdl}}|D)}{\Pr(H_{\text{bayes}}|D)} \leq 1, \\ \alpha(P, H) &\geq K(D|H_{\text{mdl}}) + K(H_{\text{mdl}}) - K(D|H_{\text{bayes}}) - K(H_{\text{bayes}}) \geq 0. \end{aligned} \quad (5.25)$$

□

Therefore, if  $\alpha(P, H)$  is small enough and Bayes's rule selects an admissible hypothesis, and so does ideal MDL, then both criteria are (almost) optimized by both selected hypotheses.

### 5.5.7 What MDL Does

We can now assess what prior distributions and assumptions about the relation between the data sample and selected hypothesis MDL assumes.

That is, how we can translate MDL in terms of Bayes's Rule. Identifying application of MDL with application of Bayes's rule on some prior distribution  $P$ , we must assume that given  $D$ , the Fundamental Inequality is satisfied for  $H_{\text{mdl}}$  as defined in Equation 5.24. This means that  $H_{\text{mdl}}$  is  $P$ -random for the prior distribution  $P$  used. One choice to guarantee this is to choose

$$P(\cdot) := \mathbf{m}(\cdot) (= 2^{-K(\cdot)}).$$

This is a valid choice even though  $\mathbf{m}$  is not recursive, since the latter requirement arose from the requirement that  $\mathbf{m}(\cdot)/P(\cdot)$  be enumerable, which is certainly guaranteed by choice of  $P(\cdot) := \mathbf{m}(\cdot)$ . This choice of prior distribution over the hypotheses is an objective and recursively invariant quantified form of Occam's razor: simple hypotheses  $H$  (with  $K(H) \ll l(H)$ ) have high probability, and complex or random hypotheses  $H$  (with  $K(H) \approx l(H)$ ) have low probability, namely,  $2^{-l(H)}$ . This choice of prior distribution is most convenient, since the randomness test  $\log \mathbf{m}(H)/P(H) = 0$  for *each* hypothesis  $H$ . This means that all hypotheses  $H$  are random with respect to distribution  $\mathbf{m}(\cdot)$ . It is easy to verify the following.

**Theorem 5.5.2** *Let  $\alpha(P, H)$  in de FI Equation 5.23 be small (for example  $\alpha = O(1)$ ). With prior  $P(\cdot)$  set to  $\mathbf{m}(\cdot)$ , the Fundamental Inequality Equation 5.23 is satisfied iff data sample  $D$  is  $\Pr(\cdot|H_{\text{mdl}})$ -random.*

With the chosen prior and data sample  $D$ , the required  $\Pr(\cdot|H_{\text{mdl}})$ -randomness of  $D$  constrains the domain of hypotheses from which we can choose  $H_{\text{mdl}}$ . Hence we can interpret ideal MDL as an application of Bayes's Rule with as prior distribution the universal distribution  $\mathbf{m}(\cdot)$  and selection of a hypothesis  $H_{\text{mdl}}$  which shows the given data sample random to it in the precise sense of  $\Pr(\cdot|H_{\text{mdl}})$ -randomness of individual objects as developed in Section 2.4.

Since the notion of individual randomness incorporates all effectively testable properties of randomness, application of ideal MDL will select the simplest hypothesis which balances the  $K(D|H)$  and  $K(H)$  and also shows the data sample  $D$  random (as far as we are able to tell) with respect to the selected hypothesis  $H_{\text{mdl}}$ .

This is the “reason” why the hypothesis selected by ideal MDL is not simply the one that perfectly fits the data. With some amount of overstatement one can say that if one obtains perfect data for a true hypothesis, then ideal MDL interprets these data as data obtained from a simpler hypothesis subject to measuring errors. Consequently, in this case ideal MDL is going to give you the *false simple* hypothesis and *not* the *complex true* hypothesis.

- Ideal MDL only gives us the true hypothesis if the data satisfy certain conditions relative to the true hypothesis. Stated differently: there are only data and no true hypothesis for ideal MDL. The principle simply obtains the hypothesis that is suggested by the data and it assumes that the data are random with respect to the hypothesis.

We have now provided an explanation why and when MDL works within the theory of algorithmic information that agrees with the analysis performed in the theory of probability, namely, that the MDL estimates of the data-generating models are consistent, except for data in small probability. And this is the essence of the requirement of the data being random, relative to the best model, in condition FI, Equation 5.23.

It is only to within terms of order  $O(1)$  that the MDL and the Bayesian techniques are equivalent. In modern forms of MDL one departs from the straight correspondence with Bayes's rule and takes  $-\log[m(D|\hat{H})/\sum_{D'} m(D'|\hat{H})]$ , instead of  $-\log[m(D|\hat{H})]$ , where  $\hat{H}(D)$  is the minimizing hypothesis and the summation runs through all data  $D'$  such that  $\hat{H}(D') = \hat{H}(D)$ , [J.J. Rissanen, *IEEE Trans. Inform. Theory*, IT-42:1(1996), 40–47]. The probability in the denominator gets absorbed by the term  $O(1)$ , however, but for smaller amounts of data it does make a difference.

### 5.5.8 Applying Minimum Description Length

Unfortunately, the function  $K(\cdot)$  of the hypotheses  $H$  is not computable (Section 3.4). For practical applications one must settle for easily computable approximations. One way to do this is as follows: First encode both  $H$  and  $D|H$  by a simply computable bijection as a natural number in  $\mathcal{N}$ . Assume we have some standard procedure to do this. Then consider a simple self-delimiting description of  $x$ . For example,  $x$  is encoded by  $x' = 1^{l(x)}0l(x)x$ . This makes  $l(x') = \log x + 2\log\log x + 1$ , which is a simple upper approximation of  $K(x)$ ; see Section 3.2. Since the length of code-word sets of prefix-codes corresponds to a probability distribution by Kraft's Inequality (page 74), this encoding corresponds to assigning probability  $2^{-l(x')}$  to  $x$ . In the MDL approach, this is the specific usable approximation to the universal prior. In the literature we find a more precise approximation that, however, has no practical meaning. For convenience, we smooth our encoding as follows.

**Definition 5.5.1** Let  $x \in \mathcal{N}$ . The *universal MDL prior* over the natural numbers is  $M(x) = 2^{-\log x - 2\log\log x}$ .

In the Bayesian interpretation the prior distribution expresses one's prior knowledge about the “true” value of the parameter. This interpretation may be questionable, since the used prior is usually not generated by repeated random experiments. In Rissanen's view, the parameter is generated by the selection of the class of hypotheses and it has no inherent meaning. It is just

one means to describe the properties of the data. The selection of  $H$  that minimizes  $K(H) + K(D|H)$  (or Rissanen's approximation thereof) allows one to make statements about the data. Since the complexity of the models plays an important part, the parameters must be encoded. To do so, we truncate them to a finite precision and encode them with the prefix-code above. Such a code happens to be equivalent to a distribution on the parameters. This may be called the universal MDL prior, but its genesis shows that it expresses no prior knowledge about the true value of the parameter. See [J.J. Rissanen, *Stochastic Complexity and Statistical Inquiry*, World Scientific, 1989]. Above we have given a validation of MDL from Bayes's Rule, which holds irrespective of the assumed prior, provided it is recursive and the hypotheses and data are random.

**Example 5.5.3** In statistical applications,  $H$  is some statistical distribution (or model)  $H = P(\theta)$  with a list of parameters  $\theta = (\theta_1, \dots, \theta_k)$ , where the number  $k$  may vary and influence the (descriptive) complexity of  $\theta$ . (For example,  $H$  can be a normal distribution  $N(\mu, \sigma)$  described by  $\theta = (\mu, \sigma)$ .) Each parameter  $\theta_i$  is truncated to finite precision and encoded with the prefix-code above.

The data sample consists of  $n$  outcomes  $\mathbf{y} = (y_1, \dots, y_n)$  of  $n$  trials  $\mathbf{x} = (x_1, \dots, x_n)$  for distribution  $P(\theta)$ . The data sample  $D$  in the above formulas is given as  $D = (\mathbf{x}, \mathbf{y})$ . By expansion of conditional probabilities we have therefore

$$\Pr(D|H) = \Pr(\mathbf{x}, \mathbf{y}|H) = \Pr(\mathbf{x}|H) \cdot \Pr(\mathbf{y}|H, \mathbf{x}).$$

In the argument above we take the negative logarithm of  $\Pr(D|H)$ , that is,

$$-\log \Pr(D|H) = -\log \Pr(\mathbf{x}|H) - \log \Pr(\mathbf{y}|H, \mathbf{x}).$$

Taking the negative logarithm in Bayes's rule and the analysis of the previous section now yields that MDL selects the hypothesis with highest inferred probability satisfying  $\mathbf{x}$  is  $\Pr(\cdot|H)$ -random and  $\mathbf{y}$  is  $\Pr(\cdot|H, \mathbf{x})$ -random. Thus, Bayesian reasoning selects the same hypothesis as MDL does, provided the hypothesis with maximal inferred probability causes  $\mathbf{x}, \mathbf{y}$  to satisfy these randomness requirements.

Under certain general conditions, J.J. Rissanen has shown that with  $k$  parameters and  $n$  data (for large  $n$ ) Equation 5.16 is minimized for hypotheses  $H$  with  $\theta$  encoded by

$$-\log P(H) = \frac{k}{2} \log n$$

bits. This is called the *optimum model cost* since it represents the cost of the hypothesis description at the minimum description length of the total.

As an example, consider a Bernoulli process  $(p, 1 - p)$  with  $p$  close to  $\frac{1}{2}$ . For such processes  $k = 1$ . Let the outcome be  $x = x_1 x_2 \dots x_n$ . Set  $f_x = \sum_{i=1}^n x_i$ . For outcome  $x$  with  $C(x) \geq n - \delta(n)$ , the number of 1's will be (by Lemma 2.3 on page 159)

$$f_x = n/2 \pm \sqrt{\frac{3}{2}(\delta(n) + c)n / \log e}.$$

With  $\delta(n) = \log n$ , the fraction of such  $x$ 's in  $\{0, 1\}^n$  is at least  $1 - 1/n$  and goes to 1 as  $n$  rises unboundedly. Hence, for the overwhelming number of  $x$ 's the frequency of 1's will be within

$$2^{-\frac{1}{2} \log n}$$

of the value  $\frac{1}{2}$ . That is, to express an estimate to parameter  $p$  with high probability it suffices to use a precision of  $\frac{1}{2} \log n$  bits. It is easy to generalize this example to arbitrary  $p$ .  $\diamond$

**Example 5.5.4** In biological modeling, we often wish to fit a polynomial  $f$  of unknown degree to a set of data points

$$D = (x_1, y_1), \dots, (x_n, y_n),$$

such that it can predict future data  $y$  given  $x$ . Even if the data did come from a polynomial curve of degree, say, two, because of measurement errors and noise, we still cannot find a polynomial of degree two fitting all  $n$  points exactly. In general, the higher the degree of fitting polynomial, the greater the precision of the fit. For  $n$  data points, a polynomial of degree  $n - 1$  can be made to fit exactly, but probably has no predicting value. Applying ideal MDL we look for  $H_{\text{mdl}} := \min_{\mathcal{H}} \{K(\mathbf{x}, \mathbf{y}|H) + K(H)\}$ .

Let us apply the ideal MDL principle where we describe all  $(k - 1)$ -degree polynomials by a vector of  $k$  entries, each entry with a precision of  $d$  bits. Then the entire polynomial is described by

$$kd + O(\log kd) \text{ bits.} \quad (5.26)$$

(We have to describe  $k$ ,  $d$ , and account for self-delimiting encoding of the separate items.) For example,  $ax^2 + bx + c$  is described by  $(a, b, c)$  and can be encoded by about  $3d$  bits. Each datapoint  $(x_i, y_i)$  that needs to be encoded separately with precision of  $d$  bits per coordinate costs about  $2d$  bits.

For simplicity assume that probability  $\Pr(\mathbf{x}|H) = 1$  (because  $\mathbf{x}$  is prescribed). To apply the ideal MDL principle we must trade the cost of hypothesis  $H$  (Equation 5.26) against the cost of describing  $\mathbf{y}$  with help

of  $H$  and  $\mathbf{x}$ . As a trivial example, suppose that  $n - 1$  out of  $n$  datapoints fit a polynomial of degree 2 exactly, but only 2 points lie on any polynomial of degree 1 (a straight line). Of course, there is a polynomial of degree  $n - 1$  that fits the data precisely (up to precision). Then the Ideal MDL cost is  $3d + 2d$  for the 2nd degree polynomial,  $2d + (n - 2)d$  for the 1st degree polynomial, and  $nd$  for the  $(n - 1)$ th degree polynomial. Given the choice among those three options, we select the 2nd degree polynomial for all  $n > 5$ .

A more sophisticated approach, accounting for the average encoding cost of exceptions, assumes that the data are Gaussian distributed. Consider polynomials  $f$  of degree at most  $n - 1$  that minimize the error

$$\text{error}(f) = \sum_{i=1}^n (f(x_i) - y_i)^2. \quad (5.27)$$

This way we find an optimal set of polynomials for each  $k = 1, 2, \dots, n$ . To apply the MDL principle we must trade the cost of hypothesis  $H$  (Equation 5.26) against the cost of describing  $D|H$ .

To describe measuring errors (noise) in data it is common practice to use the normal distribution. In our case this means that each  $y_i$  is the outcome of an independent random variable distributed according to the normal distribution with mean  $f(x)$  and variance, say, constant. For each of them we have that the probability of obtaining a measurement  $y_i$ , given that  $f(x)$  is the true value, is of the order of  $\exp(-(f(x) - y_i)^2)$ . Considering this as a value of the universal MDL probability, this is encoded in  $s(f(x) - y_i)^2$  bits, where  $s$  is a (computable) scaling constant. For all experiments together we find that the total encoding of  $D|f, \mathbf{x}$  takes  $s \cdot \text{error}(f)$  bits. The MDL principle thus tells us to choose a  $k$ -degree function  $f_k$ ,  $k \in \{0, \dots, n - 1\}$ , that minimizes (ignoring the vanishing  $O(\log kd)$  term)  $kd + s \cdot \text{error}(f_k)$ .  $\diamond$

**Example 5.5.5** In this example we apply the MDL principle to infer decision trees. We are given a set of data, possibly with noise, representing a collection of examples. Each example is represented by a data item in the data set, which consists of a tuple of *attributes* followed by a binary *Class* value indicating whether the example with these attributes is a positive or negative example.

Figure 5.3 gives a small sample set. The columns in the table describe attributes that are weather conditions. The rows are examples that represent weather conditions in relation to some “unspecified occurrences.” The last column classifies the examples as positive or negative, where “P” means that it happened and “N” means that it did not happen. We would like to obtain good predictions for such occurrences by compressing the data. Our task can now be explained as a communication

No.	ATTRIBUTES				CLASS
	<i>Outlook</i>	<i>Temperature</i>	<i>Humidity</i>	<i>Windy</i>	
1	overcast	hot	high	not	N
2	overcast	hot	high	very	N
3	overcast	hot	high	medium	N
4	sunny	hot	high	not	P
5	sunny	hot	high	medium	P
6	rain	mild	high	not	N
7	rain	mild	high	medium	N
8	rain	hot	normal	not	P
9	rain	cool	normal	medium	N
10	rain	hot	normal	very	N
11	sunny	cool	normal	very	P
12	sunny	cool	normal	medium	P
13	overcast	mild	high	not	N
14	overcast	mild	high	medium	N
15	overcast	cool	normal	not	P
16	overcast	cool	normal	medium	P
17	rain	mild	normal	not	N
18	rain	mild	normal	medium	N
19	overcast	mild	normal	medium	P
20	overcast	mild	normal	very	P
21	sunny	mild	high	very	P
22	sunny	mild	high	medium	P
23	sunny	hot	normal	not	P
24	rain	mild	high	very	N

**FIGURE 5.3.** Sample data set

problem between Alice, who observed the data in Figure 5.3, and Bob. Alice and Bob both know the four parameters (outlook, temperature, humidity, windy) and their attributes. Alice wishes to send Bob the information in the table, using as few bits as possible. Alice and Bob have to agree in advance on an encoding technique to be used.

Alice and Bob do not know in advance which table they have to transmit. The simplest strategy for Alice is to transmit the complete table in Figure 5.3 to Bob literally. There are 24 rows. Each row has four attributes and one Class value. Three attributes have three alternative values each; the other attribute and Class have two alternative values each. Then this requires  $24(3 \log_2 3 + 2) = 24(3 \times 1.585 + 2) \approx 162.12$  bits. Or Alice can agree with Bob beforehand about a fixed order of enumerating all  $3 \times 3 \times 2 \times 3 = 54$  combinations of attributes, and then just send the last column of 54 bits, supplying arbitrary Class values for the 30 rows missing from the table in Figure 5.3. These methods use no data compression.

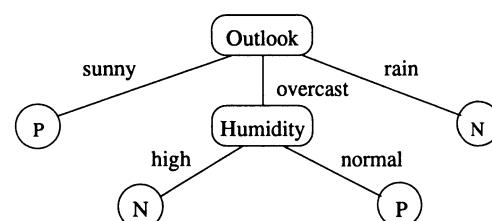
If Alice is clever enough to find some regularity in the data, like “the Class value is ‘N’ iff it rains,” then Alice needs only a few bits to transmit this sentence to Bob, and Bob can use this rule to reconstruct the complete table with all the combinations of attributes with  $3 \times 3 \times 2 \times 3 = 54$  rows.

Let us say that Alice and Bob have agreed to use a decision tree. A decision tree that is consistent with the data can be viewed as a classification procedure. The internal nodes in the tree are *decision nodes*. Each such node specifies a test of a particular attribute; the possible answers are labeled on the arcs leaving the decision node. A leaf of the tree specifies a class to which the object that passed the attribute tests and arrived at this leaf belongs. Given data as in Figure 5.3, we can construct many different decision trees of various sizes that agree with the data. Two such trees are given in Figures 5.4 and 5.5. The tree in Figure 5.4 is imperfect since it makes an error on row 8; the tree in Figure 5.5 classifies all of the sample set correctly. The tree in Figure 5.5 is the smallest perfect decision tree for the data in Figure 5.3.

Some data, for example noise, may not obey the predicting rule defined by the decision tree. One usually has a choice between using a small *imperfect* tree that classifies some data falsely or a big *perfect* tree that correctly classifies all given data. Alice can use a smaller imperfect tree or the bigger perfect tree. The tree in Figure 5.5 grows much bigger just because of a single (perhaps noisy) example (row 8), and Alice may find that it is more economical to code it separately, as an *exception*.

The goal is often to construct a decision tree that has the smallest error rate for classifying unknown future data. Is the *smallest perfect decision tree* really a good predictor? It turns out that in *practice* this is not the case. Due to the presence of noise or inadequacy of the given attributes, selecting a perfect decision tree “overfits” the data and gives generally poor predictions. Many ad hoc rules have been suggested and used for overcoming this problem.

The MDL principle appears to provide a solution and generally works well in practice. Essentially, given the data sample *without* the class val-



**FIGURE 5.4.** Imperfect decision tree

ues, we look for a reasonably small tree such that most data are correctly classified by the tree. We encode the inconsistent data as exceptions. We minimize the sum of

- the number of bits to encode a (not necessarily perfect) decision tree  $T$ , that is, the model that gives the  $-\log P(H)$  term; and
- the number of bits to encode (describe) the examples  $(x_i, y_i)$  in the sample  $(\mathbf{x}, \mathbf{y})$  that are inconsistent with  $T$ , given the entire data sample  $\mathbf{x}$  without the class values  $\mathbf{y}$ . This gives the  $-\log P(\mathbf{y}|H, \mathbf{x})$  term.

We have to provide a coding method. This is important in applications, since it determines where the optimum is. If the encoding of trees is not efficient, then we may end up with a very small tree (with relatively large depth), and too many examples become exceptions. An inefficient encoding of the exceptions would result in overly large trees. In both cases, the prediction of unseen data is affected. The reader should realize that the choice of cost measure and encoding technique cannot be objective. One can encode a tree by making a depth-first traversal. At each internal node, write down the attribute name in some self-delimiting form followed by its edge label. At a leaf write down the class value. If the tree is not perfect, then the data that do not fit in the tree are encoded separately as exceptions (in some economical way using the provided total data sample without the class values).

**Coding the Tree** It is desirable that the smaller trees be represented by shorter encodings. Alice can make a depth-first traversal of the tree in Figure 5.4, and accordingly she writes down

1 Outlook 0 P 1 Humidity 0 N 0 P 0 N.

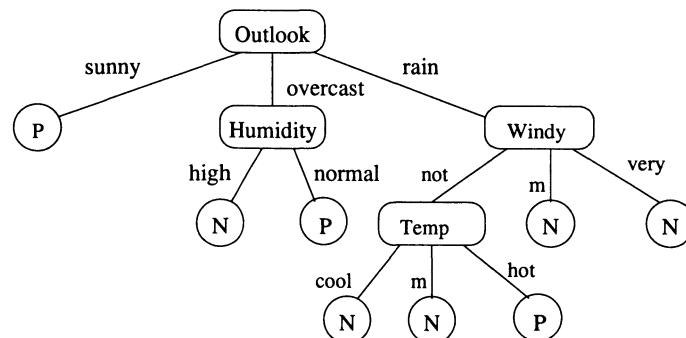


FIGURE 5.5. Perfect decision tree

For the tree in Figure 5.5, she writes down

1 Outlook 0 P 1 Humidity 0 N 0 P 1 Windy 0 N 0 N  
1 Temperature 0 N 0 N 1 P.

Alice uses a “1” to indicate that the next node in the depth-first search is an internal node and then writes the corresponding attribute; Alice writes a 0 followed by “N” (“P”) if she meets a leaf with value “N” (“P”). Representing “N” or “P” requires only one bit. Representing attribute “Outlook” at the root level requires 2 bits since there are 4 possibilities. Encoding the next-level attributes requires  $\log 3$  bits since there are only 3 choices left (“Outlook” is already used). She can use just one bit for “Windy” and one bit for “Temperature” (in fact, this one bit is not needed). Thus, the smaller (but imperfect) tree requires 13.585 bits; the bigger (but perfect) tree requires 25.585 bits.

**Coding the Exceptions** Since the decision tree in Figure 5.4 is not perfect, we need to indicate where the exceptions are. In this case there is a single exception. The most straightforward way is to indicate its *position* among all 54 possible combinations of attributes. This costs  $\log 54 \approx 5.75$  extra bits.

Thus, the encoding using the decision tree in Figure 5.4 uses 19.335 bits; the encoding using the decision tree in Figure 5.5 uses 25.585 bits. The MDL principle tells us to use the former method, which is also much shorter than the 54-bit plain encoding.

Procedures for computing decision trees have been implemented by J.R. Quinlan and R. Rivest [*Inform. Computation*, 80(1989), 227–248]. Computing the absolute minimum decision tree is NP-complete, as shown by T. Hancock, T. Jiang, M. Li, and J. Tromp, *Inform. Comput.*, 126:2(1996), 114–122. They have shown that approximating minimum decision trees is also NP-hard, even approximation to within some polynomial factor. Consequently, approximation heuristics have to be used. See also K. Yamanishi, A Randomized Approximation of the MDL for Stochastic Models with Hidden Variables, Proc. 9th ACM Comput. Learning Conference, ACM Press, 1996; and V. Vovk, Learning about the parameter of the Bernoulli Model, *J. Comput. System Sci.*, to appear.

◇

**Example 5.5.6 (Alternative MDL-Like Principle)** In the above interpretation of MDL we essentially look for a hypothesis  $H$  minimizing  $K(D|H)+K(H)$ . This always satisfies

$$K(D|H) + K(H) \geq K(D).$$

An incorrect interpretation of the way we used MDL in Example 5.5.5 on page 364 is sometimes confused with MDL. In the new approach the idea is that we define  $E := D - D_H$ , where  $D_H$  is the data set classified according to  $H$ . We want to minimize

$$K(H, E) \approx K(H) + K(E|H)$$

over  $H$ . That is,  $E$  denotes the subset of the data sample  $D$  that are *exceptions* to  $H$  in the sense of being “not covered” by  $H$ . We want to find  $H$  such that the description of  $H$  and the exception data  $E$  *not covered* by  $H$  are together minimized. Note that in this case always

$$\min_H \{K(H) + K(E|H)\} \leq K(\emptyset) + K(D) = K(D),$$

in contrast to the standard interpretation of MDL above. This incarnation of MDL is not straightforwardly derived by our approach above. We may interpret it that we look for the shortest description of an accepting program for the data consisting of a classification rule  $H$  and an exception list  $E$ . While this principle sometimes gives good results, application may lead to absurdity as the following shows:

In many problems our data sample  $D$  consists of only positive examples, as when we want to learn (a grammar for) the English language given a corpus of data  $D$  like the *Oxford Dictionary*. Then according to our new MDL rule the best hypothesis is the trivial grammar  $H$  generating *all* sentences over the alphabet. Namely, this grammar gives  $K(H) = O(1)$  independent of  $D$  and also  $E := \emptyset$ . Consequently,

$$\min_H \{K(H) + K(E|H)\} = K(H) = O(1),$$

which is absurd. The principle is vindicated and reduces to the standard one in the context of interpreting  $D = (\mathbf{x}, \mathbf{y})$  as in Example 5.5.3 on page 362, with  $\mathbf{x}$  fixed as in “supervised learning.” This is a correct application as in Example 5.5.5 on page 364. We want to find  $H$  minimizing

$$K(H) + K(\mathbf{y}|H, \mathbf{x}) + K(\mathbf{x}|H),$$

which is the same as minimizing

$$K(H) + K(\mathbf{y}|H, \mathbf{x}),$$

provided we take  $K(\mathbf{x}|H)$  constant. Now,  $K(\mathbf{y}|H, \mathbf{x})$  corresponds to  $K(E|H)$  if we ignore the constant  $\mathbf{x}$  in the conditional.  $\diamond$

**Example 5.5.7 (Maximum Likelihood)** The *maximum likelihood* (ML) principle says that for given data  $D$ , one should select the hypothesis  $H$  that maximizes

$P(D|H)$  or equivalently, minimizes  $-\log P(D|H)$ . In case of finitely many hypotheses, this is a special case of the MDL principle with the hypotheses distributed uniformly (all have equal probability). The principle has many admirers, is supposedly objective, and is due to R.A. Fisher.  $\diamond$

**Example 5.5.8 (Maximum Entropy)** In statistics there are a number of important applications where the ML principle fails, but where the maximum entropy principle has been successful, and conversely.

In order to apply Bayes's Rule, we need to decide what the prior probabilities  $p_i = P(H_i)$  are, subject to the constraint  $\sum_i p_i = 1$  and certain other constraints provided by empirical data or considerations of symmetry, probabilistic laws, and so on. Usually these constraints are not sufficient to determine the  $p_i$ 's uniquely. E.T. Jaynes proposed to select the prior by the *maximum entropy* (ME) principle.

The ME principle selects the estimated values  $\hat{p}_i$  that maximize the entropy function

$$H(p_1, \dots, p_k) = -\sum_{i=1}^k p_i \ln p_i, \quad (5.28)$$

subject to

$$\sum_{i=1}^k p_i = 1 \quad (5.29)$$

and some other constraints. For example, consider a loaded die,  $k = 6$ . If we do not have any information about the die, then using the principle of indifference, we may assume that  $p_i = \frac{1}{6}$  for  $i = 1, \dots, 6$ . This actually coincides with the ME principle, since  $H(p_1, \dots, p_6) = -\sum_{i=1}^6 p_i \ln p_i$ , constrained by Equation 5.29, achieves its maximum  $\ln 6 = 1.7917595$  for  $p_i = \frac{1}{6}$  for all  $i$ .

Now suppose it has been experimentally observed that the die is biased and the average throw gives 4.5, that is,

$$\sum_{i=1}^6 i p_i = 4.5. \quad (5.30)$$

Maximizing the expression in Equation 5.28, subject to the constraints in Equations 5.29 and 5.30, gives the estimates

$$\hat{p}_i = e^{-\lambda i} \left( \sum_j e^{-\lambda j} \right)^{-1}, \quad i = 1, \dots, 6,$$

where  $\lambda = -0.37105$ . Hence,

$$(\hat{p}_1, \dots, \hat{p}_6) = (0.0543, 0.0788, 0, 1142, 0.1654, 0.2398, 0.3475).$$

The maximized entropy  $H(\hat{p}_1, \dots, \hat{p}_6)$  equals 1.61358. How dependable is the ME principle? Jaynes has proven an “entropy concentration theorem” that, for example, implies the following: In an experiment of  $n = 1000$  trials, 99.99% of all  $6^{1000}$  possible outcomes satisfying the constraints of Equations 5.30 and 5.29 have entropy

$$1.602 \leq H\left(\frac{n_1}{n}, \dots, \frac{n_6}{n}\right) \leq 1.614,$$

where  $n_i$  is the number of times the value  $i$  occurs in the experiment. We show that the Maximum Entropy principle can be considered as a special case of the MDL principle, as follows:

Consider the same type of problem. Let  $\theta = (p_1, \dots, p_k)$  be the prior probability distribution of a random variable. We perform a sequence of  $n$  independent trials. Shannon has observed that the real substance of Formula 5.28 is that we need approximately  $nH(\theta)$  bits to record the sequence of  $n$  outcomes. Namely, it suffices to state that each outcome appeared  $n_1, \dots, n_k$  times, respectively, and afterwards give the index of which one of the

$$\binom{n}{n_1, \dots, n_k} = \frac{n!}{n_1! \cdots n_k!} \quad (5.31)$$

possible sequences  $D$  of  $n$  outcomes actually took place. For this no more than

$$k \log n + \log \binom{n}{n_1, \dots, n_k} + O(\log \log n) \quad (5.32)$$

bits are needed. The first term corresponds to  $-\log P(\theta)$ , the second term corresponds to  $-\log P(D|\theta)$ , and the third term represents the cost of encoding separators between the individual items. Using Stirling's approximation of  $n! \sim \sqrt{2\pi n}(n/e)^n$  for the quantity of Equation 5.31, we find that for large  $n$ , Equation 5.32 is approximately

$$n \left( - \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n} \right) = nH\left(\frac{n_1}{n}, \dots, \frac{n_k}{n}\right).$$

Since  $k$  and  $n$  are fixed, the least upper bound on the minimum description length, for an arbitrary sequence of  $n$  outcomes under certain given constraints 5.29 and 5.30, is found by maximizing the term in Equation 5.31 subject to said constraints. This is equivalent to maximizing the entropy function 5.28 under the constraints.

Viewed differently, let  $S_\theta$  be the set of outcomes with values  $(n_1, \dots, n_k)$ , with  $n_i = np_i$ , corresponding to a distribution  $\theta = (p_1, \dots, p_k)$ . Then due to the small number of values ( $k$ ) in  $\theta$  relative to the size of the sets, we have

$$\log \sum_{\theta} d(S_\theta) \approx \max_{\theta} \{\log d(S_\theta)\}. \quad (5.33)$$

The left-hand side of Equation 5.33 is the minimum description; the right-hand side of Equation 5.33 is the maximum entropy.  $\diamond$

## 5.6 History and References

---

The material on Epicurus can be found in E. Asmis [*Epicurus Scientific Method*, Cornell University Press, 1984]. The elegant paper “The use of simplicity in induction,” by J.G. Kemeny [*Phil. Rev.*, 62(1953), 391–408] contains predecessors to the ideas formulated in this chapter. Bayes’s formula originates from Thomas Bayes’s “An essay towards solving a problem in the doctrine of chances” [*Phil. Trans. Roy. Soc.* 25 (1763) 376–398. (*Ibid.* , 54(1764) 298–310, R. Price (Ed.))] posthumously published by his friend Richard Price. Properly speaking, Bayes’s Rule as given in the text is not due to Bayes. P.S. Laplace stated Bayes’s Rule in its proper form and attached Bayes’s name to it in *A philosophical essay on probabilities* (1819). In his original memoir, Bayes assumes the uniform distribution for the prior probability and derives  $P(H_i|D) = P(D|H_i)/\sum_i P(D|H_i)$ . This formula can be derived from Bayes’s Rule in its present form by setting all  $P(H_i)$  equal. Bayes did not state the result in its general form, nor did he derive it through a formula similar to Bayes’s Rule. The books by B. de Finetti [*Probability, Induction, and Statistics*, John Wiley & Sons, 1972], I.J. Good [*Good Thinking*, University of Minnesota Press, 1983], P.S. Laplace [*Ibid.*], R. von Mises [*Probability, Statistics and Truth*, Macmillan, 1939], and T.L. Fine [*Theories of Probability*, Academic Press, 1973] contain excellent discussions on the Bayesian and non-Bayesian views of inductive reasoning.

The idea of using Kolmogorov complexity in inductive inference, in the form of using a universal prior probability, is due to R.J. Solomonoff [*Inform. Contr.*, 7(1964), 1–22, 224–254]. Solomonoff’s original definition of prior probability is problematic through the use of the  $C$ -version of the Kolmogorov complexity instead of the prefix complexity (as used here). Inductive inference, using  $M$  as universal prior, is due to R.J. Solomonoff [*IEEE Trans. Inform. Theory*, IT-24(1978), 422–432]; see also [T.M. Cover, ‘Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin,’ *Tech. Rept.* 12, 1974, Statistics Dept,

Stanford Univ.]. The estimate in Theorem 5.2.1 of the expected prediction error, caused by using  $\mathbf{M}$  instead of the actual distribution to predict, is also due to Solomonoff. The formulation and proof we gave was kindly provided (on commission, so to speak) by P. Gács [private communication of October 5, 1989]. Similarly, Theorem 5.2.2 without speed of convergence estimates, was suggested to us by Gács. This approach seems to encapsulate the substance of the problem of inductive inference in statistics and AI in Chapter 5. The relation between “good prediction” to “better data compression” in Section 5.2.2 and Theorem 5.2.3 is from [P.M.B. Vitányi and M. Li, ‘Minimum Description Length Induction, Bayesianism, and Kolmogorov Complexity,’ Manuscript, CWI, Amsterdam, September 1996]. An good introduction to algorithmic probability, universal betting, and its relations to prefix complexity is given in [T.M. Cover and J.A. Thomas, *Elements of Information Theory*, Wiley, 1991].

Our approach in Section 5.3 is partially based on [M. Li and P.M.B. Vitányi, *J. Comput. System Sci.*, 44:2(1992), 343–384]. It rests on material in Chapter 4 and Section 5.2. The formula in Equation 5.12 and an analysis that  $\mathbf{m}$  has related approximative qualities to the actual recursive prior as  $\mathbf{M}$  has in the continuous setting appeared in [T.M. Cover, ‘Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin,’ Tech. Rept. 12, 1974, Statistics Dept, Stanford Univ.; *The impact of processing techniques on communication*, Nijhoff Publishers, 1985, pp. 23–33]. Related work is that by V.G. Vovk in a time-limited setting [*Problems Inform. Transmission*, 25(1989), 285–292; *Inform. Comput.*, 96(1992), 245–277]. He investigates low-complexity algorithms that give fast probability estimates, or “forecasts,” of the next bit in a sequence (as opposed to noncomputable ones in the Solomonoff procedure). A forecasting algorithm is said to be “simple” if it has a short description that admits fast computation of the forecasts. There is a “universal forecasting algorithm” that for every time bound  $T$  computes forecasts within time  $T$ , and the quality of these forecasts is not much lower than that of any simple forecasting system (of which the descriptive complexity may increase if  $T$  increases). (Solomonoff’s forecasting algorithm is universal in this sense, but noncomputable.) See also the related work [A. DeSantis, G. Markowsky, and M. Wegman, *Proc. 29th IEEE Symp. Found. Comput. Sci.*, 1988, pp. 110–119; N. Littlestone and M. Warmuth, *Proc. 30th IEEE Symp. Found. Comput. Sci.*, 1989, pp. 256–261]. The Gold paradigm was first introduced and studied by E.M. Gold [*Inform. Contr.*, 37(1967), 447–474]. D. Angluin and C. Smith [*Comput. Surveys*, 16(1983), 239–269] and D. Angluin [*Proc. 24th Symp. Theory Comput.*, 1992, pp. 351–369] give the standard survey of this field. The discussion of mistake bounds resulted from conversations with P. Gács, who attributes these and/or related ideas to J.M. Barzdins and R.V. Freivalds.

The pac-learning model was introduced by L.G. Valiant [*Comm. ACM*, 27(1984), 1134–1142] which also contains the learning algorithm for  $k$ -DNF. Although similar approaches were previously studied by V.N. Vapnik and A.Ya. Chervonenkis [*Theory Probab. Appl.*, 16:2(1971), 264–280] and J. Pearl [*Int. J. Gen. Syst.*, 4(1978), 255–264] in somewhat different context, the field of computational learning theory emerged only after the publication of Valiant’s paper. [D. Angluin, *Ibid.*] contains a survey of this field. Textbooks are [B.K. Natarajan, *Machine Learning: A Theoretical Approach*, Morgan Kaufmann, 1991; M. Anthony and N. Biggs, *Computational Learning Theory*, Cambridge University Press, 1992; M. Kearns and U. Vazirani, *Introduction to Computational Learning Theory*, MIT Press, 1994]. The Valiant learning model we have given in the text is slightly simplified for clarity. The important Occam’s Razor Theorem is due to A. Blumer, A. Ehrenfeucht, D. Haussler, M. Warmuth [*Inform. Process. Lett.*, 24(1987), 377–380; *J. ACM*, 35:4(1989), 929–965]. The latter also introduced the VC-dimension into the pac-learning field. In the current (slightly stronger) version of the Occam’s Razor Theorem, Theorem 5.4.1, we have replaced  $\text{size}(h)$  by  $C(h)$ . This subtle change allows one to handle more delicate cases when the structure of a concept requires a concept to have unreasonable size, while the Kolmogorov complexity of the concept is very small. This was useful in the computational molecular biology problem of efficiently learning a string [M. Li, *Proc. 31st IEEE Symp. Found. Comput. Sci.*, 1990, pp. 125–134; revised version: T. Jiang and M. Li, *Math. Systems Theory*, to appear]. The word “pac,” coined by D. Angluin, stands for “probably approximately correct.” The algorithm for  $k$ -decision lists is due to R. Rivest [*Machine Learning*, 2:3(1987), 229–246]. [R. Gavaldà, Ph.D. Thesis, Universitat Politècnica de Catalunya, 1992] contains a study of learning (with queries) using Kolmogorov complexity.

Pac-learning under simple distributions in Section 5.4.3 is from M. Li and P.M.B. Vitányi [*SIAM J. Comput.*, 20:5(1991), 915–935]. Approximation for the set-cover problem is from [V. Chvátal, *Math. Oper. Res.*, 4:3(1979), 233–235; D.S. Johnson, *J. Comput. System Sci.*, 9(1974), 256–276; L. Lovász, *Discrete Math.*, 13(1975), 383–390]. Lemma 5.4.5 and the definition of  $\epsilon$ -cover are due to G. Benedek and A. Itai, *Theoret. Comput. Sci.*, 86:2(1991), 377–390. A drawback of the “simple learning” approach in this section is that  $\mathbf{m}$  is uncomputable. Computable versions of  $\mathbf{m}$  are treated in Section 7.6. Simple pac-learning proceeds the same as before. For example, simple learning of  $\log n$ -DNF with respect to all  $O(n^2)$  time computable distributions  $P^*$  can be achieved using a time-bounded version of  $\mathbf{m}$ . This version of universal distribution requires exponential time to compute—but this needs only to be done once for all applications. Once we have computed an (approximation of) a table for this distribution, we can draw examples (even deterministically) according

to it in the learning phase as explained in Example 7.6.2 on page 508 for all time and for all concepts to be learned. These ideas follow [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5(1991), 915–935].

The minimum description length principle was introduced by J.J. Rissanen in [*Automatica*, 14(1978), 465–471] inspired by R. Solomonoff’s ideas in Section 5.2. A related approach, the minimum message length (MML) approach, is associated with C.S. Wallace and his coauthors. It was introduced earlier and independently (also independent of Solomonoff’s work) by C.S. Wallace and D.M. Boulton [*Computing Journal*, 11(1968), 185–195]. This approach is related to taking the negative logarithm of both sides of Bayes’s Rule, and then replacing the negative logarithms of the probabilities by the associated code-word lengths as justified by the Shannon-Fano Code. In contrast to MDL, the MML method relies on priors. Although different in philosophy and derivation, in practical applications MML and MDL often turn out to be similar. For the fine points in the distinctions we refer to the abundant literature cited below. See also [C.S. Wallace and P.R. Freeman, *J. Royal Stat. Soc. B*, 49:3(1987), 240–252, 252–265; *J. Royal Stat. Soc. B*, 54:1(1992), 195–209] where the MML method is refined and related with (or inspired on) Solomonoff’s work and Kolmogorov complexity, apparently influenced by [G.J. Chaitin, *Scientific American*, 232:5(1975), 47–52].

MDL has been developed by J.J. Rissanen in a series of papers [*Ann. Stat.*, 11(1982), 416–431; *Ann. Stat.*, 14:3(1986), 1080–1100; *Encyclopedia Stat. Sci.*, V, S. Kotz and N.L. Johnson (Eds.), Wiley, 1986; *J. Royal Stat. Soc.*, 49(1987), 223–239, Discussion 252–265; and as a monograph *Stochastic Complexity and Statistical Inquiry*, World Scientific, 1989]. The MDL code consists of a two-part code comprising a model part and a data part conditional on the model, or as a one-part code called the “stochastic complexity” of the data. A recent precise estimate of stochastic complexity is given in [J.J. Rissanen, *IEEE Trans. Inform. Theory*, IT-42:1(1996), 40–47].

Our derivation of MDL from first principles follows [M. Li and P.M.B. Vitányi, *J. Comput. Syst. Sci.*, 44(1992), 343–384; M. Li and P.M.B. Vitányi, *Computer Science Today*, J. van Leeuwen, Ed., Lecture Notes in Computer Science, Vol. 1000, Springer-Verlag, 1995, 518–535; P.M.B. Vitányi and M. Li, *Proc. ISIS: Information, Statistics and Induction in Science*, World Scientific, Singapore, 1996, 282–291; and in particular P.M.B. Vitányi and M. Li, ‘Minimum Description Length Induction, Bayesianism, and Kolmogorov Complexity,’ Manuscript, CWI, Amsterdam, September 1996]. Approximations to the optimum MDL code are considered in [K. Yamanishi, A Randomized Approximation of the MDL for Stochastic Models with Hidden Variables, *Proc. 9th ACM Conf. Comput. Learning Theory*, 1996; and V. Vovk, Learning about the parameter of the Bernoulli Model, *J. Comput. System Sci.*, to appear].

Relations between pac learning and MDL are explored in [K. Yamanishi, *Machine Learning*, 9(1993), 165–203]. The application of the MDL principle to fitting polynomials, as in Example 5.5.4, was originally considered by J.J. Rissanen in [*Ann. Stat.*, 14(1986), 1080–1100] and [‘Stochastic complexity and the maximum entropy principle,’ unpublished]. Example 5.5.5 is based on [J.R. Quinlan and R. Rivest, *Inform. Comput.*, 80(1989), 227–248]. Decision tree algorithms using the MDL principle were also developed by Rissanen and Wax [personal communication with M. Wax, 1988]. Decision trees are coded using the MML principle in [Wallace C.S., Patrick J.D., *Machine Learning*, 11(1993), 7–22]. Using the MDL principle, an application of optimal stochastic complexity in coding decision trees is given in [J.J. Rissanen, ‘Stochastic complexity in learning’ *J. Comput. Syst. Th.*, To appear].

There is a plethora of applications of the MDL principle in many different areas. Some examples are: learning online handwritten characters in [Q. Gao and M. Li, *11th IJCAI*, 1989, pp. 843–848]; surface reconstruction problems in computer vision [E.P.D. Pednault, *11th IJCAI*, 1989, pp. 1603–1609]; and protein structure analysis in [H. Mamitsuka and K. Yamanishi, *Comput. Appl. Biosciences (CABIOS)*, 11:4(1995), 399–411]. Other applications of the MDL principle range from evolutionary tree reconstruction [P. Cheeseman and R. Kanefsky, Working Notes, *AAAI Spring Symposium Series*, Stanford University, 1990]; inference over DNA sequences [L. Allison, C.S. Wallace, and C.N. Yee, *Int. Symp. Artific. Intell. and Math.*, January 1990; pattern recognition; smoothing of planar curves [S. Itoh, *IEEE ISIT*, January 1990]; to neural network computing [A.R. Barron, *Nonparametric Functional Estimation and Related Topics*, G. Roussas, Ed., Kluwer, 1991, pp. 561–576]. See also [A.R. Barron and T.M. Cover, *IEEE Trans. Inform. Theory*, IT-37 (1991), 1034–1054 (Correction Sept. 1991)]. A survey on MDL and its applications to statistical problems is [B. Yu, ‘Minimum description length principle: a review,’ Manuscript, Dept. of Stat., Univ. of Calif. Berkeley, 1996].

The MML principle is applied to the alignment of macro molecules in [L. Allison, C.S. Wallace, C.N. Yee *J. Mol. Evol.*, 35(1992), 77–89]; and to constructing a minimal diagnostic decision trees for medical purposes in [D.P. McKenzie, P.D. McGorry, C.S. Wallace, L.H. Low, D.L. Copolov, B.S. Singh, *Meth. Inform. Medicine*, 32:2(1993), 161–166]. See also a decision graph explanation of protein secondary structure prediction, [D.L. Dowe, J.J. Oliver, T.I. Dix, L. Allison, and C.S. Wallace, *Proc. 26th Hawaii Int. Conf. Syst. Sciences*, 1993, 669–678], and circular clustering of protein dihedral angles by MML [D.L. Dowe, D.L., L. Allison, T.I. Dix, L. Hunter, C.S. Wallace, and T. Edgoose, *Proc. 1st Pacif. Symp. Biocomput. (PSB-1)*, Hawaii, U.S.A., 1996, 242–255]. Applications of the MML mixture modelling program “Snob” by C.S. Wallace and D.L.

Dowe [*Proc. 7th Australian Joint Conf. Artific. Intel.*, Armidale, NSW, Australia, 1994, 37–44]. For inefficiency of financial markets by MML [D.L. Dowe and K.B. Korb, pp. 212–223 in *Proc. Inform. Stat. and Induction in Science (ISIS) Conf.*, World Scientific, Singapore, 1996]. This last ISIS Proc. contains papers on theory and applications of MDL and MML, for example C.S. Wallace’s paper on justification of MML on small samples, and J.J. Rissanen’s paper on improved MDL.

A related method (we have not discussed in the main text) for discovery by minimum length encoding, applied to molecular evolution and sequence similarity, is the “algorithmic significance method” of A. Milosavljević and J. Jurka, [*Machine Learning*, 12(1993), 69–87; *Proc. 1st Int’l Conf. Intelligent Systems for Molecular Biology*, AAAI Press, 1993, 284–291; *CABIOS*, 9:4(1993), 407–411]. This method seems a straight corollary of the fact demonstrated by Lemma 4.3.5 on page 258 that “ $\kappa_0(x|P) = \log(\mathbf{m}(x)/P(x))$ ” is a universal Martin-Löf test for randomness.” This shows that  $\kappa_0(x|P)$  has the associated properties of Definition 2.4.2 on page 130 and in particular the property of Definition 2.4.1 on page 129 that  $\sum\{P(x|l(x) = n) : \kappa_0(x|P) \geq m\} \leq 2^{-m}$ , for all  $n$ . It is interesting that such nonfeasible constructions can yet be directly applied in a practical inference setting. The maximum likelihood (ML) principle was introduced by R.A. Fisher in [*Phil. Trans. Royal Soc. London, Ser. A*, 222(1925), 309–368]. The ML principle has greatly influenced the research, as well as practice, in statistics. Jaynes’ maximum entropy principle was introduced by E.T. Jaynes in [*IEEE Trans. Syst. Sci. Cyb.*, SSC-4(1968), 227–241; *Proc. IEEE*, 70(1982), 939–952]. See also [E.T. Jaynes, *Papers on Probability, Statistics, and Statistical Physics*, 2nd edition, Kluwer Academic Publishers, 1989]. The ME principle has exerted a great influence on statistics and statistical mechanics. The relationship between the maximum entropy principle and the MDL principle was established by J.J. Rissanen [*Ann. Stat.*, 14(1986), 1080–1100] and by M. Feder [*IEEE Trans. Inform. Theory*, IT-32(1986), 847–849]. Rissanen [*Ann. Stat.*, 11(1982), 416–431] demonstrates that the ML principle is a special case of MDL principle.

The role of simplicity, Kolmogorov complexity, and the minimum description length principle in econometric modelling and forecasting is discussed in [H.A. Keuzenkamp and M. McAleer, *The Economic Journal*, 105(1995), 1–21].

Universal probability and Kolmogorov complexity have been applied in cognitive psychology to classical problems of perceptual organization and Gestalt psychology to resolve the central debate between the competing “likelihood principle” and “simplicity principle” [N. Chater, *Psychological Review*, 103(1996), 566–581].

# The Incompressibility Method

The incompressibility of random objects yields a simple but powerful proof technique. The incompressibility method is a general-purpose tool and should be compared with the pigeon-hole principle or the probabilistic method. Whereas the older methods generally show the existence of an object with the required properties, the incompressibility argument shows that almost all objects have the required property. This follows immediately from the fact that the argument is typically used on a Kolmogorov random object. Since such objects are effectively indistinguishable, the proof holds for all such objects. Each class of objects has an abundance of objects that are Kolmogorov random relative to the class.

The method is demonstrated by using it on several examples. We show its versatility and universal applicability by selecting examples from a wide range of applications. This includes combinatorics, random graphs, average-case analysis of Heapsort and other algorithms, routing in communication networks, formal language theory, time bounds on language recognition, string matching, circuit complexity, Turing machine time complexity, and the time complexity of parallel computations.

The method rests on a simple fact: a Kolmogorov random string cannot be compressed. Generally, a proof proceeds by showing that a certain property has to hold for some “typical” instance of a problem. But “typical” instances are difficult to define and impossible to construct. Therefore, a classical proof usually involves all instances of a certain class.

By intention and definition, an individual Kolmogorov random object is a “typical” instance. These are the incompressible objects. Although

individual objects cannot be proved to be incompressible in any given finite axiom system, a simple counting argument shows that almost all objects are incompressible, Theorem 2.2.1 on page 109. In a typical proof using the incompressibility method, one first chooses a random object from the class under discussion. This object is incompressible. Then one proves that the desired property holds for this object. The argument invariably says that if the property does not hold, then the object can be compressed. This yields the required contradiction.

Because we are dealing with only one fixed object, the resulting proofs tend to be simple and natural. They are natural in that they supply rigorous analogues for our intuitive reasoning. In many cases a proof using the incompressibility method implies an average-case result since almost all strings are incompressible.

## 6.1 Three Examples

---

The proposed methodology is best explained by example. The first example contains one of the earliest lower-bound proofs by the incompressibility argument. The second example shows how to use incompressibility to analyze the average-case complexity of an algorithm. The third example was first proved using an incompressibility argument.

### 6.1.1 Computation Time of Turing Machines

Consider the basic Turing machine model in Figure 6.1. This is the model explained in Section 1.7. It has a finite control and a single tape, serving as input tape, output tape, and work tape. The tape is a one-way infinite linear array of squares, each of which can hold a symbol from a finite, nonempty alphabet. The leftmost square is the initial square.

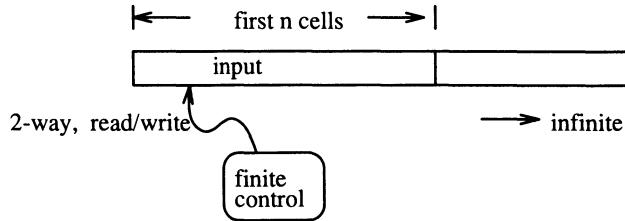
There is a two-way read/write head on the tape. The head movement is governed by the state of the finite control and the symbol in the tape square under scan. In one step, the head may print another symbol in the scanned tape square, move one square left or right (or not move at all), and the state of the finite control may change. At the start of the computation, the input occupies the initial tape segment (one symbol per square) and is delimited by a distinguished endmarker. Initially, the tape head is on the leftmost tape square, and the finite control is in a distinguished initial state. If  $x = x_1 \dots x_n$ , then  $x^R = x_n \dots x_1$ .

#### Definition 6.1.1

Each pair of adjacent squares on the tape is separated by an intersquare boundary. Consider an intersquare boundary  $b$  and the sequence of states of  $T$ 's finite control at the steps when the head crosses  $b$ , first from left to right, and then alternatingly in both directions. This ordered sequence of states is the *crossing sequence* at  $b$ .

#### Lemma 6.1.1

*A Turing machine of the model above requires order  $n^2$  steps to recognize  $L = \{xx^R : x \in \{0, 1\}^*\}$ .*



**FIGURE 6.1.** Single-tape Turing machine

**Proof.** By way of contradiction, assume that there is a Turing machine  $T$  of the above model that recognizes  $L$  in  $o(n^2)$  steps. Without loss of generality assume that the machine halts with the tape head scanning the input endmarker.

Fix a string  $x$  of length  $n$  with  $C(x|T, n) \geq n$ . Such strings exist by Theorem 2.2.1 on page 109. Consider the computation of  $T$  on  $x0^{2n}x^R$ . Let  $l(T)$  and  $l(c.s.)$  denote the lengths of the descriptions of  $T$  and a crossing sequence  $c.s.$ , respectively. If each crossing sequence associated with a square in the  $0^{2n}$  segment in the middle is longer than  $n/2l(T)$ , then  $T$  uses at least  $n^2/l(T)$  steps.

Otherwise there is a crossing sequence of length less than  $n/2l(T)$ . Assume that this is the crossing sequence  $c.s.$  associated with the square in position  $c_0$ . This  $c.s.$  is completely described by at most  $n/2$  bits. Using  $c.s.$ , one can reconstruct  $x$  by exhaustively checking all binary strings of length  $n$ .

For each candidate binary string  $y$  of length  $n$ , put  $y0^{2n}$  on the leftmost  $3n$ -length initial segment of the input tape and simulate  $T$ 's computation from its initial configuration. Each time the head moves from square  $c_0$  to its right neighbor, skip the part of the computation of  $T$  with the head right of  $c_0$ , and resume the computation starting from the next state  $q$  in  $c.s.$  with the head scanning square  $c_0$ .

Suppose that in the computation with  $y$ , each time the head moves from square  $c_0$  to its right neighbor, the current state of  $T$  is the correct next state as specified in  $c.s.$ . Then  $T$  accepts input  $y0^{2n}x^R$ . Namely, the computation right of square  $c_0$  will simply be identical to the computation right of square  $c_0$  on input  $x0^{2n}x^R$ . Since  $T$  halts with its head to the right of square  $c_0$ , it must either accept both  $y0^{2n}x^R$  and  $x0^{2n}x^R$  or reject them both. Since  $T$  recognizes  $L$  we must have  $y = x$ . Therefore, given  $c.s.$ , we can reconstruct  $x$  by a fixed program from the data  $T$ ,  $n$ , and  $c.s.$ . This means that

$$C(x|T, n) \leq l(c.s.) + O(1) \leq n/2 + O(1),$$

which contradicts  $C(w|T, n) \geq n$ , for large  $n$ .  $\square$

### 6.1.2 Adding Fast—On Average

In computer architecture design, efficient design of adders directly affects the length of the CPU clock cycle. Fifty years ago, Burks, Goldstine, and von Neumann obtained a  $\log n$  expected upper bound on the “longest carry sequence” involved in the process of adding two  $n$ -bit binary numbers. This property suggests the design for an efficient adder hardware. We give a simple analysis using the incompressibility method. Let  $x$  and  $y$  be two  $n$ -bit binary numbers and let “ $\oplus$ ” denote the bit-wise “exclusive-or” operator. The following algorithm adds  $x$  and  $y$ .

**Step 1**  $S := x \oplus y$  (add bit-wise, ignoring carries);  $C :=$  carry sequence;

**Step 2** While  $C \neq 0$  do

$$S := S \oplus C;$$

$$C := \text{new carry sequence}.$$

Let us call this the *no-carry adder* algorithm. The expected  $\log n$  upper bound on carry sequence length implies that the algorithm runs in  $1 + \log n$  expected rounds (Step 2). This algorithm is on average the most efficient addition algorithm currently known. But it takes  $n$  steps in the worst case. On average, it is exponentially faster than the trivial linear-time “ripple-carry adder,” and it is two times faster than the well-known “carry-lookahead adder.” In the ripple-carry adder, the carry ripples from right to left, bit by bit, and hence it takes  $\Omega(n)$  steps to compute the sum of two  $n$ -bit numbers. The carry-lookahead adder is used in nearly all modern computers; it is based on a divide-and-conquer algorithm that adds two  $n$ -bit numbers in  $1 + 2 \log n$  steps. We give an easy proof of the  $1 + \log n$  average-case upper bound, using the incompressibility method.

**Lemma 6.1.2** *The “no-carry adder” algorithm has an average running time of at most  $1 + \log n$ .*

**Proof.** Assume both inputs  $x$  and  $y$  have length  $l(x) = l(y) = n$ , with the lower-order bits on the right. If the computation takes precisely  $t$  steps (Step 2 loops  $t$  times), then some thinking shows that there exists a  $u$  such that  $x$  and  $y$  can be written as

$$x = x'bu1x'', y = y'b\neg u1y'',$$

where  $l(u) = t - 1$ ,  $l(x') = l(y')$ ,  $b$  is 0 or 1, and  $\neg u$  is the bitwise complement of  $u$ . Therefore,  $x$  can be described using  $y$ ,  $n$ , and a program  $q$  of  $O(1)$  bits to reconstruct  $x$  from the concatenation of

- the position of  $u$  in  $y$  encoded in exactly  $\log n$  bits (padded with 0’s if needed); and

- the literal representation of  $x'x''$ .

Since the concatenation of the two strings has length  $n - t - 1 + \log n$ , the value  $t$  can be deduced from  $n$  and this length. Therefore,  $t + 1$  bits of  $x$  are saved at the cost of adding  $\log n$  bits. (For  $x' = \epsilon$  bit  $b$  may not exist. But then the algorithm also does not execute the last step because of overflow.)

This shows that  $C(x|n, y, q) \leq n - t - 1 + \log n$ . Hence, for each  $x$  with  $C(x|n, y, q) = n - i$ , the computation must terminate in at most  $\log n + i - 1$  steps. By simple counting as in Theorem 2.2.1 on page 109, there are at most  $2^{n-i}$  strings  $x$  of length  $n$  with Kolmogorov complexity  $C(x|n, y, q) = n - i$ . There are at most  $2^{n-i}$  programs of length  $n - i$ , and hence at most  $2^{n-i}$  strings  $x$  with  $C(x|n, y, q) = n - i$ . Let  $p_i$  denote the fraction of  $x$ 's of length  $l(x) = n$  satisfying  $C(x|n, y, q) = n - i$ . Then,  $p_i \leq 2^{-i}$  and  $\sum_i p_i = 1$ . Hence, averaging over all  $x$ 's (by having  $i$  range from 1 to  $n$ ) with  $y$  fixed, the average computation time for each  $y$  is bounded above by

$$\begin{aligned} \sum_{i=2-\log n}^n p_i(i - 1 + \log n) &= \sum_{i=2-\log n}^n p_i(i - 1) + \sum_{i=2-\log n}^n p_i \log n \\ &\leq \log n + \sum_{i=1}^{\infty} (i - 1)/2^i = 1 + \log n. \end{aligned}$$

Because this holds for every  $y$ , this is also the average running time of the algorithm.  $\square$

### 6.1.3 Boolean Matrix Rank

The rank of a matrix  $R$  is the least integer  $k$  such that each row of  $R$  can be written as a linear sum of  $k$  fixed rows. These  $k$  rows are linearly independent, which means that no row can be written as a linear sum of the others. Our problem is to show the existence of a Boolean matrix with all submatrices of high rank.

Such matrices were used to obtain an optimal lower bound  $TS = \Omega(n^3)$  time-space tradeoff for multiplication of two  $n$  by  $n$  Boolean matrices on random access machines ( $T$  = time and  $S$  = space). Even with integer entries, it is difficult to construct such a matrix. There is no known construction with Boolean values.

Let  $GF(2)$  be the Gaussian field over the two elements 0, 1, with the usual Boolean multiplication and addition:  $0 \times 0 = 0 \times 1 = 1 \times 0 = 0$ ,  $1 \times 1 = 1$ ,  $1 + 1 = 1 + 0 = 0 + 1 = 1$ , and  $0 + 0 = 0$ .

**Lemma 6.1.3** *Let  $n, r, s \in \mathcal{N}$  with  $2 \log n \leq r, s \leq n/4$  and  $s$  even. For each  $n$  there is an  $n \times n$  matrix over  $GF(2)$  such that every submatrix of  $s$  rows and  $n - r$  columns has at least rank  $s/2$ .*

**Proof.** Fix a binary string  $x$  of length  $n^2$ , with  $C(x) \geq n^2$ . This is possible by Theorem 2.2.1 on page 109. Arrange the bits of  $x$  into a square matrix  $R$ , one bit per entry in, say, row-major order. We claim that this matrix  $R$  satisfies the requirement.

Assume by way of contradiction that this were not true. Consider a submatrix of  $R$  of  $s$  rows and  $n - r$  columns, with  $r, s$  as in the condition in the lemma. There are at most  $(s/2) - 1$  linearly independent rows in it. Therefore, each of the remaining  $(s/2) + 1$  rows can be expressed as a linear sum of the other  $(s/2) - 1$  rows. This can be used to describe  $R$  by the following items:

- The characteristic sequence of the  $(s/2) - 1$  independent rows out of the  $s$  rows of  $R$  in  $s$  bits;
- A list of the  $(s/2) - 1$  linearly independent rows in  $((s/2) - 1)(n - r)$  bits;
- List the remainder of  $(s/2) + 1$  rows in order. For each row give only the Boolean coefficients in the assumed linear sum. This requires  $((s/2) - 1)((s/2) + 1)$  bits.

To recover  $x$ , we need only the additional items below:

- A description of this discussion in  $O(1)$  bits;
- The values of  $n, r, s$  in self-delimiting form in  $3 \log n + 6 \log \log n + 3$  bits. For large  $n$ , this is at most  $4 \log n$  bits including the  $O(1)$  bits above;
- $R$  without the bits of the submatrix in row-major order in  $n^2 - (n - r)s$  bits;
- The indices of the columns and rows of the submatrix, in  $(n - r)\log n + s\log n$  bits.

To ensure unique decodability of these binary items, we concatenate them as follows: First list the self-delimiting descriptions of  $n, r, s$ , then append all other items in a fixed order. The length of each of the items can be calculated from  $n, r, s$ . Altogether, this is an effective description of  $x$ , a number of bits of at most

$$\begin{aligned} & n^2 - (n - r)s + (n - r)\log n + s\log n \\ & + \left(\frac{s}{2} - 1\right)(n - r) + \left(\frac{s}{2} - 1\right)\left(\frac{s}{2} + 1\right) + s + 4\log n. \end{aligned}$$

For large  $n$ , this quantity drops below  $n^2$ . But we have assumed that  $C(x) \geq n^2$ , which yields the required contradiction.  $\square$

A proof obtained by the incompressibility method usually implies that the result holds for “almost all strings,” and hence it holds for the “average” case complexity.

## Exercises

---

**6.1.1.** Let the Turing machine in Section 6.1.1 be probabilistic, which means that the machine can flip a fair coin to determine its next move. But the machine is not allowed to err. Prove that such a machine still requires on average order  $n^2$  steps to accept  $L = \{xx^R : x \in \{0, 1\}^*\}$ .

*Comment.* Use Symmetry of Information, Theorem 2.8.2 on page 182. When the machine is allowed to err with probability  $\epsilon$ , then  $L$  can be accepted in time  $O(n \log n)$  by such a machine. It can do so by generating random primes of size  $\log n$  and checking whether both sides are the same modulo these primes [R. Freivalds, *Inform. Process.* 77, 1977, pp. 839-842].

**6.1.2.** [10] (Converting NFA to DFA) A *deterministic finite automaton* (DFA)  $A$  has a finite number of states, including a distinguished start state and some distinguished accepting states. At each step,  $A$  reads the next input symbol and changes its state according to the current state and the input symbol. If  $A$  has more than one alternative at some step, then  $A$  is nondeterministic (NFA). If  $A$  is in an accepting state when it reads a distinguished endmarker, then  $A$  accepts the input. Otherwise  $A$  rejects it. It is well-known that each NFA can be converted to a DFA. Use an incompressibility argument to prove that there exists an NFA with  $n$  states such that the smallest DFA accepting the same language has  $\Omega(2^n)$  states.

*Comment.* Hint: use  $L_k = \{x : \text{the } k\text{th bit of } x \text{ from the right is "1"}\}$ . This problem can also be solved by a simple counting argument.

## 6.2 High- Probability Properties

---

The theory of random individual objects, Sections 2.4 and 2.5, tells us that there is a close relation between high-probability properties and properties of incompressible objects. For infinite binary sequences  $\omega \in \{0, 1\}^\infty$  and  $\lambda$  the uniform (coin-toss) measure, classic probabilistic laws are formulated in global form by

$$\lambda\{\omega : A(\omega)\} = 1, \tag{6.1}$$

where  $A(\omega)$  is some formula expressing some property. In contrast, in the algorithmic theory of random individual objects, the corresponding law is expressed in local form by

$$\text{if } \omega \text{ is random then } A(\omega) \text{ holds.} \tag{6.2}$$

Since the number of classical probabilistic laws as in Equation 6.1 is uncountable, whereas the properties tested by Martin-Löf tests to determine randomness as in Equation 6.2 are only the effectively testable properties which are countable, there are classical probabilistic laws that do not hold in the pointwise sense of Equation 6.2. On the other hand, a classical probabilistic law implies the pointwise algorithmic law: if Equation 6.2 holds for formula  $A$ , then also Equation 6.1 holds for  $A$  (by Theorem 2.5.3 on page 144). It is not immediately obvious how such a relation holds quantitatively for *finite* objects. To fix our thoughts let us look at a simple example.

First we recall the notion of “randomness deficiency” of Section 2.2.1 on page 112. The randomness deficiency of an element in a certain class of objects is the difference between that element’s Kolmogorov complexity and the maximal Kolmogorov complexity of an object in the class (typically the logarithm of its cardinality). Formally, if  $x$  is an element of a finite set of objects  $\mathcal{O}$ , then by Theorem 2.1.3 on page 103 we have  $C(x|\mathcal{O}) \leq l(d(\mathcal{O})) + c$  for some  $c$  independent of  $x$  but possibly dependent on  $\mathcal{O}$ . The *randomness deficiency* of  $x$  relative to  $\mathcal{O}$  is defined as  $\delta(x|\mathcal{O}) = \log(d(\mathcal{O})) - C(x|\mathcal{O})$ .

**Example 6.2.1** Let  $G = (V, E)$  be a graph on  $n$  nodes where each pair of nodes is or is not connected by an edge according to the outcome of a fair coin flip. The probability that a particular node is isolated (has no incident edges) is  $1/2^{n-1}$ . Therefore, the probability that some node is isolated is at most  $n/2^{n-1}$ . Consequently, the probability that the graph is connected is at least  $1 - n/2^{n-1}$ .

Using the *incompressibility method*, the proof that random graphs are connected with high probability is as follows: Each labeled undirected graph  $G = (V, E)$  on  $n$  nodes can be described by giving a characteristic sequence  $\chi$  of the lexicographical enumeration of  $V \times V$  without repetition, namely,  $\chi = \chi_1 \chi_2 \dots \chi_e$  with  $e = \binom{n}{2}$  and  $\chi_i = 1$  if the  $i$ th enumerated edge is in  $E$  and 0 otherwise. There are as many labeled  $n$ -node graphs as there are such characteristic sequences. Therefore, we can consider graphs  $G$  having randomness deficiency at most  $\delta(n)$ ,

$$C(G|n) \geq \binom{n}{2} - \delta(n).$$

Assume by way of contradiction that there is an isolated node  $i$ . Add its identity in  $\log n$  bits to the canonical description of  $G$ , and delete all  $n - 1$  bits indicating presence or absence of edges incident on  $i$ , saving  $n - 1$  bits. From the new description we can reconstruct  $G$  given  $n$ . Then the new description length cannot be smaller than  $C(G|n)$ . Substitution shows that  $\delta(n) \geq n - 1 - \log n$ . Counting the number of programs of length  $C(G|n)$  shows that at most a fraction of  $2^{-(n-1-\log n)}$

of all  $n$ -node graphs contain an isolated node. Consequently, the “connectedness” property for  $n$ -node graphs holds with probability at least  $1 - n/2^{n-1}$ .  $\diamond$

For each class of finite objects there is a close relation between properties that hold with high probability and properties that hold for such objects with small randomness deficiency: the almost incompressible ones. However, the properties and the sets of objects concerned are not identical and should be carefully distinguished. In fact, the following distinctions also indicate in which cases use of which method is preferable:

1. In the probabilistic method, the subset of objects on which the probability of a property is based is the subset of *all* objects satisfying that property. As an example, consider the connectedness property of labeled graphs again. The graphs satisfying this property include the complete graph on  $n$  nodes, the star graph on  $n$  nodes, and the binary hypercube on  $n$  nodes, provided  $n$  is a power of 2. These graphs are certainly not incompressible or random and in fact have complexity  $O(1)$  given  $n$ .
2. If each object with suitable randomness deficiency  $\delta(n)$  has a certain property, then each such object is included in the subset of objects on which the high probability of the property is based.
3. If we prove that properties  $P$  and  $Q$  each hold with probability at least  $1 - \epsilon$  with the probabilistic method, then we can conclude that properties  $P$  and  $Q$  *simultaneously* hold with probability at least  $1 - 2\epsilon$ . In contrast, if both properties  $P$  and  $Q$  hold separately for objects with randomness deficiency  $\delta(n)$ , then they vacuously also hold simultaneously for objects with randomness deficiency  $\delta(n)$ .

In greater generality, suppose that each high-probability property separately holds for an overwhelming majority (say, at least a  $(1 - 1/n)$ th fraction) of all objects. Now consider the situation of  $n$  different properties each of which holds for a  $(1 - 1/n)$ th fraction. Since possibly the subsets on which the different properties fail may be disjoint, possibly their union may constitute the set of all objects. Therefore it is possible that no object at all possesses all the high-probability properties simultaneously.

In contrast, if we prove each out of  $n$  properties separately for objects with randomness deficiency  $\delta(n)$ , then all these properties hold simultaneously for each of these objects.

These considerations show that “high-probability” properties and “incompressibility” properties are not a priori the same. However, we shall prove that they almost coincide under mild conditions on the properties

considered. In fact, the objects with randomness deficiency  $\delta(n)$  satisfy all *simply described* properties that hold with high probability.

If a simple enough property holds with high enough probability then it must necessarily hold for all objects with randomness deficiency  $\delta(n)$ . For example, if  $P$  is a simple property with  $C(P|n) = O(1)$  that holds with probability at least  $1 - 1/n^2$ , then  $P$  holds for all objects with randomness deficiency  $\log n$ .

Consider a class of finite objects  $\mathcal{O}$  parametrized with  $n$  (like  $n$ -node graphs, strings of length  $n$ , groups on  $n$  generators).

**Lemma 6.2.1** *Let  $P$  be a property holding for objects  $O \in \mathcal{O}$  with randomness deficiency  $\delta(n)$ . Then  $P$  holds with probability at least  $1 - 1/2^{\delta(n)-1}$ .*

**Proof.** There are only  $\sum_{i=0}^{\log d(\mathcal{O}) - \delta(n)} 2^i$  programs of length not greater than  $\log d(\mathcal{O}) - \delta(n)$  and there are  $d(\mathcal{O})$  objects.  $\square$

**Lemma 6.2.2** *Let  $f, g$  be two functions such that  $f(n) = o(g(n))$ . Let  $P$  be any property with  $C(P|n) = O(f(n))$  that holds with probability at least  $1 - 1/2^{\delta(n)+g(n)}$  for the class of objects concerned. All such  $P$  hold simultaneously for each object with randomness deficiency  $\delta(n)$  in the class.*

**Proof.** Let  $\mathcal{O}$  be the class of objects under consideration. Suppose  $P$  does not hold for an object  $O \in \mathcal{O}$  and  $O$  has randomness deficiency  $\delta(n)$ . Then we can reconstruct  $O$  from a description of  $P$  and  $O$ 's index  $j$  in an effective enumeration of all objects for which  $P$  doesn't hold. There are at most  $d(\mathcal{O})/2^{\delta(n)+g(n)}$  such objects by assumption, and therefore

$$j \leq \log d(\mathcal{O}) - \delta(n) - g(n).$$

Hence, by the  $\delta(n)$  randomness deficiency of  $O$ , we have  $-\delta(n) - g(n) + C(P|n) + 2 \log C(P|n) \geq -\delta(n)$  which is a contradiction.  $\square$

Typically,  $C(P|n) = O(1)$ , that is, independent of  $n$  for recursive properties. Such properties are bounds on the size of the largest complete subgraph in a  $\log n$ -random graph ( $2 \log n$ ). The quantity  $C(P|n)$  grows unbounded for more “complex” properties that require us to describe a number of parameters that grows unboundedly as  $n$  grows unboundedly.

“ $\delta(n)$  randomness deficiency” is in fact a *universal Martin-Löf test for randomness* incorporating all Martin-Löf tests for randomness with respect to individual properties, Section 2.4.

**Corollary 6.2.1** *The objects of randomness deficiency  $\delta(n)$  possess all properties  $P$  that hold with probability at least  $1 - 2^{-(\delta(n)+C(P|n)+2 \log C(P|n))}$ .*

**Corollary 6.2.2** All recursive properties  $P$  with  $C(P|n) = O(1)$ , each of which holds separately with probability tending to 1 as  $n$  grows unboundedly, hold simultaneously with probability tending to 1 as  $n$  grows unboundedly.

These results mean that if we want to establish that a property holds with *high probability* or for *high-Kolmogorov-complexity objects*, then it suffices to establish either one to prove both. Moreover, the high Kolmogorov complexity objects satisfy all highly probable “simple” properties simultaneously.

## 6.3 Combinatorics

Combinatorial properties are traditionally established by counting arguments or by the probabilistic method. Probabilistic arguments are usually aimed at establishing the existence of an object in a nonconstructive sense. It is ascertained that a certain member of a class has a certain property without actually exhibiting that object. Usually, the method proceeds by exhibiting a random process that produces the object with positive probability. Alternatively, a quantitative property is determined from a bound on its average in a probabilistic situation.

We demonstrate the utility of the incompressibility method in combinatorial theory on several examples. The general pattern is as follows: When we want to prove a certain property of a group of objects (such as graphs), we first fix an incompressible instance of the object, justified by Theorem 2.2.1 on page 109. It is always a matter of using the assumed regularity in this instance to compress the object to reach a contradiction.

### 6.3.1 Transitive Tournament

A *tournament* is defined to be a complete directed graph. That is, for each pair of nodes  $i$  and  $j$ , exactly one of edges  $(i, j)$  or  $(j, i)$  is in  $T$ . The nodes of a tournament can be viewed as players in a game tournament. If  $(i, j)$  is in  $T$ , we say player  $j$  dominates player  $i$ . We call  $T$  *transitive* if  $(i, j), (j, k)$  in  $T$  implies  $(i, k)$  in  $T$ .

Let  $\Gamma = \Gamma_n$  be the set of all tournaments on  $N = \{1, \dots, n\}$ . Given a tournament  $T \in \Gamma$ , fix a standard encoding  $E : T \rightarrow \{0, 1\}^{n(n-1)/2}$ , one bit for each edge. The bit for edge  $(i, j)$  is set to 1 if  $i < j$  and 0 otherwise. There is a one-to-one correspondence between the members of  $\Gamma$  and the binary strings of length  $n(n - 1)/2$ .

Let  $v(n)$  be the largest integer such that every tournament on  $N$  contains a transitive subtournament on  $v(n)$  nodes.

**Theorem 6.3.1**  $v(n) \leq 1 + \lfloor 2 \log n \rfloor$ .

**Proof.** Fix  $T \in \Gamma$  such that

$$C(E(T)|n, p) \geq n(n-1)/2, \quad (6.3)$$

where  $p$  is a fixed program that on input  $n$  and  $E'(T)$  (below) outputs  $E(T)$ . Let  $S$  be the transitive subtournament of  $T$  on  $v(n)$  nodes. We try to compress  $E(T)$ , to an encoding  $E'(T)$ , as follows:

1. Prefix the list of nodes in  $S$  in order of dominance to  $E(T)$ , each node using  $\lceil \log n \rceil$  bits, adding  $v(n)\lceil \log n \rceil$  bits.
2. Delete all redundant bits from the  $E(T)$  part, representing the edges between nodes in  $S$ , saving  $v(n)(v(n)-1)/2$  bits.

Then

$$l(E'(T)) = l(E(T)) - \frac{v(n)}{2}(v(n) - 1 - 2\lceil \log n \rceil). \quad (6.4)$$

Given  $n$ , the program  $p$  reconstructs  $E(T)$  from  $E'(T)$ . Therefore,

$$C(E(T)|n, p) \leq l(E'(T)). \quad (6.5)$$

Equations 6.3, 6.4, and 6.5 are true only when  $v(n) \leq 1 + \lfloor 2 \log n \rfloor$ .  $\square$

The general idea used in the incompressibility proof of Theorem 6.3.1 is the following: If each tournament contains a large transitive subtournament, or any other “regular” property for that matter, then also a tournament  $T$  of maximal complexity contains one. But the regularity induced by too large a transitive subtournament can be used to compress the description of  $T$  to below its complexity, leading to the required contradiction.

P. Stearns showed by induction that  $v(n) \geq 1 + \lfloor \log n \rfloor$ . This is the first problem illustrating the “probabilistic method” in [P. Erdős and J. Spencer, *Probabilistic methods in combinatorics*, Academic Press, 1974]. They collected many combinatorial properties accompanied by elegant proofs using probabilistic arguments. The thrust was to show how to replace counting arguments by pleasant and short probabilistic arguments. To compare the incompressibility method, we include their proofs of Theorem 6.3.1 by counting and probabilistic methods.

**Proof. (Proof by Counting)** Let  $\Gamma = \Gamma_n$  be the class of all tournaments on  $\{1, \dots, n\}$  and  $\Gamma' =$  the class of tournaments on  $\{1, \dots, n\}$  that contain a transitive subtournament on  $v = 2 + \lfloor 2 \log n \rfloor$  players. Then

$$\Gamma' = \bigcup_A \bigcup_{\sigma} \Gamma_{A, \sigma},$$

where  $A \subseteq \{1, \dots, n\}$ ,  $d(A) = v$ ,  $\sigma$  is a permutation on  $A$ , and  $\Gamma_{A,\sigma}$  is the set of  $T$  such that  $T|A$  is generated by  $\sigma$ . If  $T \in \Gamma_{A,\sigma}$ , the  $\binom{v}{2}$  games of  $T|A$  are determined. Thus,

$$d(\Gamma_{A,\sigma}) = 2^{\binom{n}{2} - \binom{v}{2}},$$

and by elementary estimates

$$d(\Gamma') \leq \sum_{A,\sigma} 2^{\binom{n}{2} - \binom{v}{2}} = \binom{n}{v} v! 2^{\binom{n}{2} - \binom{v}{2}} < 2^{\binom{n}{2}} = d(\Gamma).$$

Thus,  $\Gamma - \Gamma' \neq \emptyset$ . That is, there exists  $T \in \Gamma - \Gamma'$  not containing a transitive subtournament on  $v$  players.  $\square$

**Proof.** (Proof by Probabilistic Method) Assume the same notation and suppositions as in the proof by counting. Let  $\mathbf{T} = T_n$  be a random variable. Its values are the members of  $\Gamma$  where for each  $T \in \Gamma$ ,  $\Pr(\mathbf{T} = T) = 2^{-\binom{n}{2}}$ . That is, all members of  $\Gamma$  are equally probable values of  $\mathbf{T}$ . Then the probability that an outcome  $T$  of  $\mathbf{T}$  contains a transitive subtournament on  $v$  players is at most

$$\sum_A \sum_{\sigma} \Pr(\mathbf{T}|A \text{ generated by } \sigma) = \binom{n}{v} v! 2^{-\binom{v}{2}} < 1.$$

Thus, some value  $T$  of  $\mathbf{T}$  does not contain a transitive subtournament on  $v$  players.  $\square$

### 6.3.2 Tournament with $k$ -Dominators

Tournament  $T$  has *property  $S(k)$*  if for every subset  $A$  of  $k$  nodes (players) there is a node (player) in  $N - A$  that dominates (beats) all nodes (players) in  $A$ . Let  $s(k)$  be the minimum number of nodes (players) in a tournament with property  $S(k)$ .

**Theorem 6.3.2**  $s(k) \leq 2^k k^2 (\ln 2 + o(1))$ .

**Proof.** Choose  $n = 2^k k^2 (\ln 2 + O(1))$ . Assume the notation of the previous example. Select  $T$  such that

$$C(E(T)|n, k, p) \geq n(n-1)/2,$$

where  $p$  is a fixed program to compute  $E(T)$  from  $E'(T)$  (given below) and  $n, k$ . By way of contradiction, assume that  $S(k)$  is false for  $T$ . Fix a set  $A$  of  $k$  nodes of  $T$  with no common dominator in  $N - A$ . Describe  $T$  as follows by a compressed description  $E'(T)$ :

- List the nodes in  $A$  first, using  $\lceil \log n \rceil$  bits each.

- List  $E(T)$  with bits representing edges between  $N - A$  and  $A$  deleted (saving  $(n - k)k$  bits).
- Code the edges between  $N - A$  and  $A$ . From each  $i \in N - A$ , there are  $2^k - 1$  possible ways of directing edges to  $A$ , in total  $t = (2^k - 1)^{n-k}$  possibilities. To encode the list of these edges,  $\lceil \log t \rceil$  bits suffice.

This shows that  $C(E(T)|n, k, p) \leq l(E'(T))$ . For large  $k$ ,  $l(E'(T)) < n(n - 1)/2$  bits, which is a contradiction.  $\square$

### 6.3.3 Ramsey Numbers

The previous examples demonstrate a general principle that a random graph (or its complement) cannot contain too large a subgraph that is easily describable. We apply the incompressibility method to obtain a lower bound on Ramsey numbers. A *clique* of a graph is a complete subgraph of that graph. The Ramsey number  $r(k, k)$  is the least integer such that for each graph  $G$  of size  $r(k, k)$ , either  $G$  or  $G$ 's complement contains a clique of size  $k$ . P. Erdős proved in 1947, using the probabilistic method, the following result:

**Theorem 6.3.3**

$$r(k, k) \geq k2^{k/2} \left( \frac{1}{e\sqrt{2}} - o(1) \right). \quad (6.6)$$

**Proof.** To describe a clique (or empty subgraph) of size  $k$  in a graph  $G$  of  $r(k, k)$  vertices we need  $\log \binom{r(k, k)}{k} \leq k \log r(k, k) - \log k!$  bits. Choose  $G$  to be incompressible. Then we must have  $k \log r(k, k) - \log k! \geq k(k - 1)/2$ , since otherwise we can compress  $G$  as in the proof of Theorem 6.3.1. Using Stirling's formula we obtain  $k! \approx k^k e^{-k} \sqrt{2\pi k}$ , and a simple calculation leads to Equation 6.6.  $\square$

### 6.3.4 Coin-Weighing Problem

A family  $\mathcal{D} = \{D_1, D_2, \dots, D_j\}$  of subsets of  $N = \{1, 2, \dots, n\}$  is called a *distinguishing family* for  $N$  if for any two distinct subsets  $M$  and  $M'$  of  $N$  there exists an  $i$  ( $1 \leq i \leq j$ ) such that  $d(D_i \cap M)$  is different from  $d(D_i \cap M')$ . Let  $f(n)$  denote the minimum of  $d(\mathcal{D})$  over all distinguishing families for  $N$ . To determine  $f(n)$  is commonly known as the *coin-weighing problem*. It is known that

$$f(n) \leq \frac{2n}{\log n} \left[ 1 + O\left(\frac{\log \log n}{\log n}\right) \right]. \quad (6.7)$$

Equation 6.7 was independently established by [B. Lindström, *Canad. Math. Bull.*, 8(1965), 477–490] and [D.G. Cantor and W.H. Mills, *Canad. J. Math.*, 18(1966), 42–48]. P. Erdős and A. Rényi [*Publ. Hungar. Acad. Sci.*, 8(1963), 241–254], L. Moser [*Combinatorial structures and their*

applications, Gordon and Breach, 1970, pp. 283–384], and N. Pippenger [*J. Combinat. Theory, Ser. A*, 23(1977), 99–104] have established Theorem 6.3.4 using probabilistic and information theory methods.

Encode each subset  $M$  of  $N$  by  $E(M) \in \{0, 1\}^n$  such that the  $i$ th bit of  $E(M)$  is 1 if  $i$  is in  $M$ , and 0 otherwise.

**Theorem 6.3.4**  $f(n) \geq (2n/\log n)[1 + O(\log \log n/\log n)]$ .

**Proof.** Choose  $M$  such that

$$C(E(M)|\mathcal{D}) \geq n. \quad (6.8)$$

Let  $d_i = d(D_i)$  and  $m_i = d(D_i \cap M)$ . Let  $s_i$  be the subsequence of  $E(M)$  selected from the positions corresponding to 1's in  $E(D_i)$ . Thus,  $l(s_i) = d_i$  and the number of 1's in  $s_i$  is precisely  $m_i$ . Moreover,

$$C(s_i) \geq d_i - O(\log i),$$

since we can use  $\mathcal{D}$ ,  $i$ , the shortest program for  $s_i$ , and  $E(M)$  minus the bits in  $s_i$  to reconstruct  $E(M)$ .

By Equation 2.3, page 159, value  $m_i$  is within range  $d_i/2 \pm O(\sqrt{d_i \log i})$ . Therefore, given  $d_i$ , each  $m_i$  can be described by its discrepancy with  $d_i/2$ , which gives

$$C(m_i|D_i) \leq \frac{1}{2} \log d_i + O(\log \log i).$$

Pad each description of  $m_i$ , given  $D_i$ , to a block of fixed length  $\frac{1}{2} \log n + O(\log \log n)$ . Since  $\mathcal{D}$  is a distinguishing family for  $N$ , given  $\mathcal{D}$ , the values  $m_1, \dots, m_j$  determine  $M$ . Hence, by the established inequalities,

$$C(E(M)|\mathcal{D}) \leq C(m_1, \dots, m_j|\mathcal{D}) \leq \sum_{i=1}^j \left( \frac{1}{2} \log n + O(\log \log n) \right).$$

Together with Equation 6.8 this implies the theorem.  $\square$

### 6.3.5 High-Probability Properties

According to the relations established in Section 6.2, if the property can be proven only with  $O(1)$  randomness deficiency, then this does not by itself imply that the property holds for “almost all” objects in the class. But in many cases a somewhat weaker property can be proven for objects with greater randomness deficiency, as we proceed to show by the following examples.

Almost all strings have high complexity. Therefore, almost all tournaments and almost all undirected graphs have high complexity. Any combinatorial property proven about an arbitrary complex object in such

a class will hold for almost all objects in the class. For example, the proof in Section 6.3.1 can trivially be strengthened as follows: By Theorem 2.2.1, page 109, there are at least  $2^{n(n-1)/2}(1 - 1/n)$  tournaments  $T$  on  $n$  nodes with

$$C(E(T)|n, p) \geq n(n-1)/2 - \log n. \quad (6.9)$$

This is a  $(1 - 1/n)$ th fraction of all tournaments on  $n$  nodes. Using Equation 6.9 in the proof yields the corollary below:

- Corollary 6.3.1** For almost all tournaments on  $n$  nodes (at least a  $(1 - 1/n)$ th fraction), the largest transitive subtournament has at most  $1 + 2\lceil 2 \log n \rceil$  nodes, from some  $n$  onwards.

Similarly, choosing  $C(E(T)|n, k, p) \geq n(n-1)/2 - \log n$  in the proof in Section 6.3.2 yields the following:

- Corollary 6.3.2** For all large enough  $k$ , there is some  $n$  with  $n \leq 2^k k^2(\ln 2 + o(1))$  such that almost all tournaments on  $n$  nodes (at least a  $(1 - 1/n)$ th fraction) have property  $S(k)$ .

The Kolmogorov complexity argument generally yields results on *expected* properties rather than worst-case properties, and is especially suited to obtain results on random structures. Other such applications (like expected maximum vertex degree of randomly generated trees and a related result on random mappings) can be found in the Exercises and in Section 6.4.

---

## Exercises

- 6.3.1.** [17] Let  $w(n)$  be the largest integer such that for each tournament  $T$  on  $N = \{1, \dots, n\}$  there exist disjoint sets  $A$  and  $B$ , each of cardinality  $w(n)$ , in  $N$  such that  $A \times B \subseteq T$ . Prove  $w(n) \leq 2\lceil \log n \rceil$ .

*Comments.* Use Kolmogorov complexity. Hint: add  $2w(n)\lceil \log n \rceil$  bits to describe nodes, and save  $w(n)^2$  bits on edges. Source of the problem: P. Erdős and J. Spencer, *Probabilistic methods in combinatorics*, Academic Press, 1974.

- 6.3.2.** [17] Let  $G = (V, E)$  with  $V = \{1, \dots, n\}$  be an undirected graph on  $n$  nodes with  $C(G|n, p) \geq n(n-1)/2$  where  $p$  is a fixed program to be used to reconstruct  $G$ . A *clique* of a graph is a complete subgraph of that graph. Show that  $G$  does not contain a clique on more than  $1 + \lfloor 2 \log n \rfloor$  nodes.

*Comments.* Hint: use Example 6.3.1 on page 389. N. Alon, J.H. Spencer, and P. Erdős, *The Probabilistic Method*, Wiley, 1992, pp. 86,87 show that

a random graph with edge probability  $\frac{1}{2}$  contains a clique on  $2 \log n$  nodes with probability at least  $1 - 1/e^{n^2}$ .

**6.3.3.** [36] Let  $K(N)$  denote the complete undirected graph of  $n$  nodes  $N = \{1, \dots, n\}$ . If  $A$  and  $B$  are disjoint subsets of  $N$ , then  $K(A, B)$  denotes the complete bipartite graph on sets  $A$  and  $B$ . A set  $\mathbf{C} = (K(A_1, B_1), \dots, K(A_j, B_j))$  is called a covering family of  $K(N)$  if for each edge  $\{u, v\} \in K(N)$  there exists an  $i$  ( $1 \leq i \leq j$ ) such that  $\{u, v\} \in K(A_i, B_i)$ . Let  $g(n)$  denote the minimum of  $\sum_{1 \leq i \leq j} d(A_i \cup B_i)$  over all covering families for  $K(N)$ . Prove by incompressibility that  $g(n)/n \geq \log n + O(\log \log n)$ .

*Comments.* An information-theoretic proof appears in [N. Pippenger, *J. Comb. Theory, Ser. A*, 23(1977), 105–115]. Hint: use Symmetry of Information, Theorem 2.8.2, page 182. Source: M. Li and P.M.B. Vitányi, *J. Comb. Theory Ser. A*, 66:2(1994), 226–236.

**6.3.4.** [25] Consider a random directed graph whose  $n^2$  nodes are on the crosses of a two-dimensional  $n$  by  $n$  grid. All vertical edges (the grid edges) are present and directed upward. For each pair of horizontally neighboring nodes, we flip a three-sided coin; with probability  $p < \frac{1}{2}$  we add an edge from left to right, with probability  $p$  we add an edge from right to left, and with probability  $1 - 2p$  we add no edge. Use incompressibility to prove that the expected maximum path length over all such random graphs is bounded by  $O(n)$ .

*Comments.* Source: T. Jiang and Z.Q. Luo, personal communication, 1992. This problem was studied in connection with communication networks.

**6.3.5.** [35] Given an  $n$ -dimensional cube and a permutation  $\pi$  of its nodes, each node  $v$  wants to send an information packet to node  $\pi(v)$  as fast as possible. Label each edge in the cube with its dimension from  $\{1, \dots, n\}$ . A route  $(v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_k)$  is ascending if  $(v_i, v_{i+1})$  has higher dimension than  $(v_{i-1}, v_i)$  for all  $2 < i < k - 1$ . If two packets use the same edge in the same direction at the same time, then a *collision* occurs, and one packet has to wait. How do we avoid too many collisions on each route? Consider the following probabilistic algorithm  $A_\pi$ :

**Step 1** For each node  $v$ , choose randomly a node  $w$ . Node  $v$  sends its packet over the uniquely determined ascending route to  $w$ .

**Step 2** Send the packet from  $w$  to  $\pi(v)$  through the unique ascending route.

Prove: for any constant  $c$ , algorithm  $A_\pi$  finishes with probability greater than  $1 - 2^{-(c-5n-O(1))/2}$  after at most  $2n + 2c$  steps.

*Comments.* Hint: show that the description of a route on which too many collisions occur can be compressed. Source: L.G. Valiant and G. Brebner, *Proc. 13th ACM Symp. Theory Comput.*, 1981, pp. 263–277; S. Reisch and G. Schnitger give an incompressibility proof in [*Proc. 23rd IEEE Found. Comput. Sci.*, 1982, pp. 45–52].

**6.3.6.** [39] Let  $L \subset \{0,1\}^{2n}$  be a language to be recognized by two parties  $P$  and  $Q$  with unlimited computation power. Party  $P$  knows the first  $n$  bits of the input and party  $Q$  knows the last  $n$  bits.  $P$  and  $Q$  exchange messages to recognize  $L$  according to some bounded-error two-way probabilistic protocol. An input is *accepted* if the probability of acceptance is at least  $1 - \epsilon$  for some fixed  $\epsilon$ ,  $0 \leq \epsilon < \frac{1}{2}$ ; an input is *rejected* if the probability of rejection is at least  $1 - \epsilon$ ; and each input must be either rejected or accepted. The *probabilistic communication complexity* of an input  $(x_1, \dots, x_{2n})$  is the worst case, over all sequences of fair coin tosses, of the number of bits exchanged. The probabilistic communication complexity of the language is the maximum of this over all inputs. The set intersection language SETIN is defined to be the set of all sequences  $a_1 \dots a_n b_1 \dots b_n$  over  $\{0,1\}$  with  $\sum_{i=1}^n a_i b_i \geq 1$ . ( $P$  knows  $a_1, \dots, a_n$  and  $Q$  knows  $b_1, \dots, b_n$ .) Prove: the probabilistic communication complexity of SETIN is  $\Omega(n)$ .

*Comments.* Source: B. Kalyanasundaram and G. Schnitger, *SIAM J. Discrete Math.*, 5:4(1992), 545–557.

## 6.4 Kolmogorov Random Graphs

---

Statistical properties of strings with high Kolmogorov complexity have been studied in Section 2.6. The interpretation of strings as more complex combinatorial objects leads to a completely new set of properties and problems that have no direct counterpart in the “flatter” string world. Here we derive topological, combinatorial, and statistical properties of graphs with high Kolmogorov complexity. Every such graph possesses simultaneously all properties that hold with high probability for randomly generated graphs. They constitute “almost all graphs” and the derived properties a fortiori hold with probability that goes to 1 as the number of nodes grows unboundedly, in the sense of Section 6.2.

### Definition 6.4.1

Each labeled graph  $G = (V, E)$  on  $n$  nodes  $V = \{1, 2, \dots, n\}$  can be coded (up to automorphism) by a binary string  $E(G)$  of length  $n(n - 1)/2$ . We enumerate the  $n(n - 1)/2$  possible edges  $(i, j)$  in a graph on  $n$  nodes in standard lexicographical order without repetitions and set the  $i$ th bit in the string to 1 if the edge is present and to 0 otherwise. Conversely, each binary string of length  $n(n - 1)/2$  encodes a graph on  $n$  nodes. Hence we can identify each such graph with its corresponding binary string.

**Definition 6.4.2** A labeled graph  $G$  on  $n$  nodes has *randomness deficiency* at most  $\delta(n)$ , and is called  $\delta(n)$ -*random*, if it satisfies

$$C(E(G)|n, \delta) \geq n(n-1)/2 - \delta(n). \quad (6.10)$$

**Lemma 6.4.1** A fraction of at least  $1 - 1/2^{\delta(n)}$  of all labeled graphs  $G$  on  $n$  nodes is  $\delta(n)$ -random.

This is a corollary of Lemma 6.2.1. Hence the  $c \log n$ -random labeled graphs constitute a fraction of at least  $(1 - 1/n^c)$  of all graphs on  $n$  nodes, where  $c > 0$  is an arbitrary large but fixed constant, and so on.

High-complexity labeled graphs have many specific topological properties, which seems to contradict their randomness. However, randomness is not “lawlessness” but rather enforces strict statistical regularities. For example, to have diameter exactly two.

**Lemma 6.4.2** The degree  $d$  of each node of a  $\delta(n)$ -random labeled graph satisfies

$$|d - (n-1)/2| = O\left(\sqrt{(\delta(n) + \log n)n}\right).$$

**Proof.** Assume that there is a node such that the deviation of its degree  $d$  from  $(n-1)/2$  is greater than  $k$ . From the lower bound on  $C(E(G)|n, \delta)$  corresponding to the assumption that  $G$  is random, we can estimate an upper bound on  $k$  as follows:

In a description of  $G = (V, E)$  given  $n, \delta$  we can indicate which edges are incident on node  $i$  by giving the index of the interconnection pattern (the characteristic sequence of the set  $V_i = \{j \in V - \{i\} : (i, j) \in E\}$  in  $n-1$  bits where the  $j$ th bit is 1 if  $j \in V_i$  and 0 otherwise) in the ensemble of

$$m = \sum_{|d-(n-1)/2|>k} \binom{n-1}{d} \leq 2^n e^{-2k^2/3(n-1)} \quad (6.11)$$

possibilities. The last inequality follows from a general estimate of the tail probability of the binomial distribution, with  $s_n$  the number of successful outcomes in  $n$  experiments with probability of success  $p = \frac{1}{2}$ . Namely, by Chernoff's bounds, Equation 2.4 on page 159,

$$\Pr(|s_n - pn| > k) \leq 2e^{-k^2/3pn}. \quad (6.12)$$

To describe  $G$  it then suffices to modify the old code of  $G$  by prefixing it with

- A description of this discussion in  $O(1)$  bits;

- the identity of node  $i$  in  $\lceil \log(n + 1) \rceil$  bits;
- the value of  $d$  in  $\lceil \log(n + 1) \rceil$  bits, possibly adding nonsignificant 0's to pad up to this amount;
- the index of the interconnection pattern in  $\log m$  bits (we know  $n, k$  and hence  $\log m$ ); followed by
- the old code for  $G$  with the bits in the code denoting the presence or absence of the possible edges that are incident on node  $i$  deleted.

Clearly, given  $n$  we can reconstruct the graph  $G$  from the new description. The total description we have achieved is an effective program of

$$\log m + \log nk + n(n - 1)/2 - n + O(1)$$

bits. This must be at least the length of the shortest effective binary program, which is  $C(E(G)|n, \delta)$ , satisfying Equation 6.10. Therefore,

$$\log m \geq n - \log nk - O(1) - \delta(n).$$

Since we have estimated in Equation 6.11 that

$$\log m \leq n - (2k^2/3(n - 1)) \log e,$$

it follows that  $k \leq \sqrt{\frac{3}{2}(\delta(n) + \log nk + O(1))(n - 1)/\log e}$ .  $\square$

**Lemma 6.4.3** *All  $o(n)$ -random labeled graphs have diameter 2.*

**Proof.** The only graphs with diameter 1 are the complete graphs that can be described in  $O(1)$  bits, given  $n$ , and hence are not random. It remains to consider  $G = (V, E)$  is an  $o(n)$ -random graph with diameter greater than 2. Let  $i, j$  be a pair of nodes with distance greater than 2. Then we can describe  $G$  by modifying the old code for  $G$  as follows:

- A description of this discussion in  $O(1)$  bits;
- The identities of  $i < j$  in  $2 \log n$  bits;
- The old code  $E(G)$  of  $G$  with all bits representing presence or absence of an edge  $(j, k)$  between  $j$  and each  $k$  with  $(i, k) \in E$  deleted. We know that all the bits representing such edges must be 0 since the existence of any such edge shows that  $(i, k), (k, j)$  is a path of length 2 between  $i$  and  $j$ , contradicting the assumption that  $i$  and  $j$  have distance  $> 2$ . This way we save at least  $n/4$  bits, since we save bits for as many edges  $(j, k)$  as there are edges  $(i, k)$ , that is, the degree of  $i$ , which is  $n/2 \pm o(n)$  by Lemma 6.4.2.

Since we know the identities of  $i$  and  $j$  and the nodes adjacent to  $i$  (they are in the prefix of code  $E(G)$  where no bits have been deleted), we can reconstruct  $G$  from this discussion and the new description, given  $n$ . Since by Lemma 6.4.2 the degree of  $i$  is at least  $n/4$ , the new description of  $G$ , given  $n$ , requires at most

$$n(n-1)/2 - n/4 + O(\log n)$$

bits, which contradicts Equation 6.10 on page 397 for large  $n$ .  $\square$

**Lemma 6.4.4** *Let  $c$  be a fixed constant. If  $G$  is a  $c \log n$ -random labeled graph, then from each node  $i$  all other nodes are either directly connected to  $i$  or are directly connected to one of the least  $(c+3) \log n$  nodes directly adjacent to  $i$ .*

**Proof.** Given  $i$ , let  $A$  be the set of the least  $(c+3) \log n$  nodes directly adjacent to  $i$ . Assume by way of contradiction that there is a node  $k$  of  $G$  that is not directly connected to a node in  $A \cup \{i\}$ . We can describe  $G$  as follows:

- A description of this discussion in  $O(1)$  bits;
- A literal description of  $i$  in  $\log n$  bits;
- A literal description of the presence or absence of edges between  $i$  and the other nodes in  $n-1$  bits;
- A literal description of  $k$  and its incident edges in  $\log n + n-2 - (c+3) \log n$  bits;
- The encoding  $E(G)$  with the edges incident with nodes  $i$  and  $k$  deleted, saving at least  $2n-2$  bits.

Altogether the resultant description has

$$n(n-1)/2 + 2 \log n + 2n - 3 - (c+3) \log n - 2n + 2$$

bits, which contradicts the  $c \log n$ -randomness of  $G$  by Equation 6.10 on page 397. The lemma is proven.  $\square$

In the description we have explicitly added the adjacency pattern of node  $i$ , which we deleted later again. This zero-sum swap is necessary to be able to unambiguously identify the adjacency pattern of  $i$  in order to reconstruct  $G$ . Since we know the identities of  $i$  and the nodes adjacent to  $i$  (they are the prefix where no bits have been deleted), we can reconstruct  $G$  from this discussion and the new description, given  $n$ .

### 6.4.1 Statistics of Subgraphs

It is easy to conclude from the statistics of high-complexity strings in Section 2.6 that the frequency of two-node subgraphs (up to isomorphism) in a  $\delta(n)$ -random graph  $G$  is

$$\frac{n(n-1)}{4} \pm \sqrt{\frac{3}{4}(\delta(n) + O(1))n(n-1)/\log e}.$$

This case is easy since the frequency of such subgraphs corresponds to the frequency of 1's in the  $\binom{n}{2}$ -length standard encoding  $E(G)$  of  $G$ . However, to determine the frequencies of subgraphs on  $k$  nodes (up to isomorphism) for  $k > 2$  is a more complicated matter than the frequencies of substrings of length  $k$ . Clearly, there are  $\binom{n}{k}$  subsets of  $k$  nodes out of  $n$  and hence that many occurrences of subgraphs. Such subgraphs may overlap in more complex ways than substrings of a string.

**Definition 6.4.3** Let  $G = (V, E)$  be a labeled graph on  $n$  nodes. Consider a labeled graph  $H$  on  $k$  nodes  $\{1, 2, \dots, k\}$ . Each subset of  $k$  nodes of  $G$  induces a subgraph  $G_k$  of  $G$ . The subgraph  $G_k$  is an *occurrence* of  $H$  if when we relabel the nodes  $i_1 < i_2 < \dots < i_k$  of  $G_k$  as  $1, 2, \dots, k$ , we obtain  $H$ .

Let  $\#H(G)$  be the *number of times*  $H$  occurs as a subgraph of  $G$  (possibly overlapping). Let  $p$  be the probability that we obtain  $H$  by flipping a fair coin to decide for each pair of nodes whether it is connected by an edge or not,

$$p = 2^{-k(k-1)/2}. \quad (6.13)$$

**Theorem 6.4.1** Assume the terminology above with  $G = (V, E)$  a labeled graph on  $n$  nodes and  $H$  a labeled graph on  $k$  nodes. Let  $C(G|n, \delta) \geq \binom{n}{2} - \delta(n)$ . Then

$$\left| \#H(G) - \binom{n}{k} p \right| \leq \binom{n}{k} \sqrt{\alpha(k/n)p},$$

with  $\alpha = [K(H|n) + \delta(n) + \log \binom{n}{k} - \log(n/k) + O(1)]3/\log e$ .

**Proof.** Assume that  $n$  is divisible by  $k$ . (If it is not, then the following arguments hold with  $n$  replaced by  $n - 1$  until it is.) A *cover* of  $G$  is a set  $C = \{S_1, \dots, S_N\}$ , where the  $S_i$ 's are pairwise disjoint subsets of  $V$  and  $\bigcup_{i=1}^N S_i = V$ . The proof of the following claim is left to the reader as the nontrivial Exercise 6.4.1 on page 403.

**Claim 6.4.1** There is a partition of the  $\binom{n}{k}$  different  $k$ -node subsets into  $h = \binom{n}{k}/N$  distinct covers of  $G$ , each cover consisting of  $N = n/k$  disjoint subsets. That is, each subset of  $k$  nodes of  $V$  belongs to precisely one cover.

Enumerate the covers as  $C_0, C_1, \dots, C_{h-1}$ . For each  $i \in \{0, 1, \dots, h-1\}$  and  $k$ -node labeled graph  $H$ , let  $\#H(G, i)$  be the number of (now nonoverlapping) occurrences of subgraph  $H$  in  $G$  occurring in cover  $C_i$ .

Now consider an experiment of  $N$  trials, each trial with the same set of  $2^{k(k-1)/2}$  outcomes. Let  $U$  be the set of all sequences of  $N$  outcomes. Intuitively, each trial corresponds to an element of a cover, and each outcome corresponds to a  $k$ -node subgraph. Define

$$A_i := \{G \text{ is an } n\text{-node graph} : |\#H(G, i) - Np| > m\},$$

for  $0 \leq i \leq h-1$ . Our proof strategy is as follows: For any  $i$ , if  $m$  is taken too large then we can compress  $G$  below its Kolmogorov complexity  $C(x|n, \delta)$  by first describing  $A_i$  and then giving  $G$ 's index in  $A_i$ , which is a contradiction. One can do this for each  $i$  independently, notwithstanding the dependence between the frequencies of subgraphs in different covers. Namely, the argument depends on the incompressibility of  $x$  alone. If the number of occurrences of a certain subgraph in *any* of the covers is too large or too small then we can compress  $x$ . The dependence of the occurrences of subgraphs between covers is accounted for in the descriptive complexity of the cover index.

We use Chernoff's bounds of Equation 2.4 on page 159 for the tail probability of the binomial distribution in a purely combinatorial fashion of counting occurrences of a particular subgraph in each graph out of the set of all  $n$ -node graphs to determine the cardinality of  $d(A)$ . This yields

$$d(A) \leq 2^{\binom{n}{2}+1} e^{-m^2/3pN},$$

where  $p$  is defined by Equation 6.13. Namely, with  $2^{n(n-1)/2}$  the number of  $n$ -node graphs, the combinatorial quantity  $d(A)/2^{n(n-1)/2}$  is bounded by Chernoff's bound for the probability of the event " $|\#H(G, i) - pN| > m$ " with  $\#H(G, i)$  the number of successful outcomes in  $N$  trials with probability of success  $p$ . Choose  $m$  such that for some constant  $c$  to be determined later,

$$\frac{m^2 \log e}{3pN} = K(\langle H, i \rangle | n) + \delta(n) + c. \quad (6.14)$$

To describe an element  $G$  in  $A_i$  we need only to enumerate  $A_i$  and indicate the index of  $G$  in such an enumeration. The description, given  $n$  for free, contains the following items:

- A description of this discussion in a constant number, say  $c_1$ , bits.
- A description to enumerate  $A_i$ , given  $n$ . To enumerate  $A_i$ , we need to know  $H$  and  $i$ . Therefore, given  $n$ , the required number of bits to enumerate  $A_i$  is at most  $K(\langle H, i \rangle | n) + c_2$ .

- A description of the index  $j$  of  $G$  in  $A_i$  in at most  $\log d(A_i)$  bits. Using Equation 6.14,

$$\begin{aligned}\log d(A_i) &\leq \log \left( 2^{\binom{n}{2}+1} e^{-m^2/3pN} \right) \\ &= \binom{n}{2} + 1 - \frac{m^2 \log e}{3pN} \\ &= \binom{n}{2} + 1 - K(\langle H, i \rangle | n) - \delta(n) - c.\end{aligned}$$

The total description takes at most  $\binom{n}{2} + 1 - \delta(n) - c + c_1 + c_2$  bits. From this description  $G$  can be effectively reconstructed, given  $n$ . Hence the length of the description must be at least  $C(G|n, \delta)$ . Finally, setting  $c := c_1 + c_2 + 2$  contradicts the  $\delta(n)$ -randomness of  $G$  assumed in the theorem. Consequently,

$$|\#H(G, i) - pN| \leq m.$$

Therefore, by Equation 6.14, for all  $i$  ( $0 \leq i \leq h-1$ ) separately and independent of each other,

$$|\#H(G, i) - pN| \leq \sqrt{\frac{K(\langle H, i \rangle | n) + \delta(n) + c}{\log e} 3pN}. \quad (6.15)$$

Estimate the  $K(\cdot)$  terms for all  $i$  by  $K(\langle H, i \rangle | n) \leq K(H | n) + K(i | H, n)$  and  $K(i | H, n) \leq \log \binom{n}{k} - \log(n/k)$ . Substituting the upper bound in the right-hand sides of Equation 6.15 makes them equal for all  $i$ . Summing the left-hand sides of Equation 6.15, and using  $hN = \binom{n}{k}$ , we find

$$\left| \#H(G) - \binom{n}{k} p \right| = \sum_{i=0}^{h-1} |\#H(G, i) - pN|.$$

Hence, summing both sides of Equation 6.15 over  $0 \leq i \leq h-1$  proves the theorem.  $\square$

In Section 2.6 we investigated up to which length  $l$  all blocks of length  $l$  occurred at least once in each  $\delta(n)$ -random string of length  $n$ .

**Theorem 6.4.2** *Let  $\delta(n) = 2\sqrt{2 \log n}/2 / 4 \log n$  and  $G$  be a  $\delta(n)$ -random graph on  $n$  nodes. Then for sufficiently large  $n$ , the graph  $G$  contains all subgraphs on  $\sqrt{2 \log n}$  nodes.*

**Proof.** We are sure that  $H$  on  $k$  nodes occurs at least once in  $G$  if  $\binom{n}{k} \sqrt{\alpha(k/n)p}$  in Theorem 6.4.1 is less than  $\binom{n}{k}p$ . This is the case if

$\alpha < (n/k)p$ , that is,

$$[K(H|n) + \delta(n) + \log \binom{n}{k} - \log \frac{n}{k} + O(1)]3/\log e < (n/k)p.$$

This inequality is satisfied for an overestimate of  $K(H|n)$  by  $\binom{k}{2} + 2\log \binom{k}{2} + O(1)$  (since  $K(H|n) \leq K(H) + O(1)$ ), and  $p = 2^{-k(k-1)/2}$ , and  $k$  set at  $k = \sqrt{2\log n}$ . This proves the theorem.  $\square$

## Exercises

**6.4.1.** [M40] Show the existence of the covers as required in the proof of Theorem 6.4.1 (Claim 6.4.1).

*Comments.* Source: Zs. Baranyai, pp. 91-108 in: A. Hajnal, R. Rado, and V.T. Sós, Eds, *Infinite and Finite Sets, Proc. Coll. Keszthely, Colloq. Math. Soc. János Bolyai*, Vol. 10, North-Holland, Amsterdam, 1975.

**6.4.2.** [26] A labeled graph  $G$  on  $n$  nodes is canonically represented by a binary string  $E(G)$  of length  $n(n - 1)/2$  encoding the presence or absence of possible edges in lexicographical order. An *automorphism* of  $G = (V, E)$  is a permutation  $\pi$  of  $V$  such that  $(\pi(u), \pi(v)) \in E$  iff  $(u, v) \in E$ . Clearly, the set of automorphisms of a graph forms a group with group operation of function composition and one element the identity permutation. It is easy to see that  $G'$  is an automorphism of  $G$  iff  $G'$  and  $G$  have the *same binary string standard encoding*, that is,  $E(G) = E(G')$ . This contrasts with the more general case of permutation relabeling, where the standard encodings may be different. A graph is *rigid* if its only automorphism is the identity automorphism. It turns out that Kolmogorov random graphs are rigid graphs. Show that a labeled graph  $G$ , on  $n$  nodes with  $C(E(G)|n) \geq n(n - 1)/2 - n/4$ , is rigid.

*Comments.* Hint: suppose the contrary and let  $G = (V, E)$ . Let  $\pi$  be a nonidentity permutation of  $V$ . Split  $V$  into subsets  $V_1$  and  $V_2$  where  $V_1$  is the set of nodes that are mapped to a different node by  $\pi$  and  $V_2 = V - V_1$ . Consider two cases:  $d(V_1) \leq \sqrt{n}$  and  $d(V_1) > \sqrt{n}$ . Use Lemma 6.4.2 and Lemma 6.4.3. Source: H. Buhrman, M. Li, and P.M.B. Vitányi, Kolmogorov random graphs, Manuscript, CWI, November 1995.

**6.4.3.** [23] An unlabeled graph is a graph with no labels. For convenience we can define this as follows: Call two labeled graphs *equivalent* (up to relabeling) if there is a relabeling that makes them equal. An *unlabeled graph* is an equivalence class of labeled graphs. It is known that the number of unlabeled graphs is  $\sim 2^{n(n-1)/2}/n!$ . Therefore, by Stirling's formula, if  $G$  is an unlabeled graph then

$$C(E(G)|n) \leq 2^{n(n-1)/2} - n \log n + O(n).$$

(a) Let  $G$  be a labeled graph on  $n$  nodes and let  $G_0$  be the unlabeled version of  $G$ . Show that there exists a graph  $G'$  and a label permutation  $\pi$  such that  $G' = \pi(G)$  and up to additional constant terms  $C(E(G')) = C(E(G_0))$  and  $C(E(G)|n) = C(E(G_0), \pi|n)$ .

(b) Show that there is a set of labeled graphs on  $n$  nodes of cardinality  $(1 - 2^{-n/4})2^{n(n-1)/2}$  such that the set of corresponding unlabeled graphs on  $n$  nodes has cardinality at most  $2^{n(n-1)/2}(1 - 2^{-n/4})/n!$

*Comments.* Hint: for Item (b) use Exercise 6.4.2. Hence, for *every* maximum-complexity graph  $G$  on  $n$  nodes there is a relabeling (permutation) that causes the complexity to drop by as much as  $n \log n$ . Our proofs of topological properties by the incompressibility method required the graph  $G$  to be Kolmogorov random in the sense of  $C(G|n) \geq n(n-1)/2 - O(\log n)$  or for some results  $C(G|n) \geq n(n-1)/2 - o(n)$ . (Rather, we should write here  $C(E(G)|n)$ , where  $E(G)$  is a canonical encoding of  $G$  in a binary string, to be defined below.) Hence by relabeling such a graph we can always obtain a labeled graph that has a complexity too low to use our incompressibility proof. Nonetheless, topological properties do not change under relabeling, and hence the proofs by the incompressibility method for high-complexity graphs can be shown to yield the topological results also for labeled graphs that are not Kolmogorov random and similarly for unlabeled graphs.

**6.4.4.** [26] Show that almost every labeled tree on  $n$  nodes has maximum degree of  $O(\log n / \log \log n)$ .

*Comments.* Hint: represent a labeled tree by a binary sequence of length  $(n - 2) \log n$  (the Prüfer code). Prove a 1-1 correspondence between labeled trees and binary sequences of such length. Then use incompressibility to show that if a tree has larger degree, then one can compress the corresponding binary sequence. Since most binary sequences cannot be compressed, most trees do not have larger degree. Source: W.W. Kirchherr, *Inform. Process. Lett.*, 41(1992), 125–130.

## 6.5 Compact Routing

---

In very large networks like the global telephone network or the Internet, the mass of messages being routed creates major bottlenecks, degrading performance. We analyze a tiny part of this issue by determining the optimal space to represent routing schemes in communication networks on the average for all static networks.

A universal routing strategy for static communication networks will, for every network, generate a *routing scheme* for that particular network. Such a routing scheme comprises a *local routing function* for every node in this network. The routing function of node  $u$  returns for every destination  $v \neq u$  an edge incident to  $u$  on a path from  $u$  to  $v$ . This way,

a routing scheme describes a path, called a *route*, between every pair of nodes  $u, v$  in the network.

It is easy to see that we can do shortest-path routing by entering a *routing table* in each node  $u$  that for each destination node  $v$  indicates to what adjacent node  $w$  a message to  $v$  should be routed first. If  $u$  has degree  $d$ , it requires a table of at most  $n \log d$  bits, and the overall number of bits in all local routing tables never exceeds  $n^2 \log n$ . Several factors may influence the cost of representing a routing scheme for a particular network. We use a basic model and leave variations to the exercises. Here we consider point to point communication networks on  $n$  nodes described by an undirected labeled graph  $G = (V, E)$ , where  $V = \{1, \dots, n\}$ . Assume that nodes know the identities of their neighbors.

In [H. Buhrman, J.H. Hoepman, and P.M.B. Vitányi, *Proc. 15th ACM Symp. Principles Distr. Comput.*, 1996, 134–142], it is shown that, in most models, for almost all graphs (that is, networks),  $\Theta(n^2)$  bits are necessary and sufficient for shortest-path routing. By “almost all graphs” we mean the Kolmogorov random graphs that constitute a fraction of  $1 - 1/n^c$  of all graphs on  $n$  nodes, where  $c \geq 3$  is an arbitrary fixed constant. In contrast, there is a model that causes the average-case lower bound to rise to  $\Omega(n^2 \log n)$  and another model where the average-case upper bound drops to  $O(n \log^2 n)$ . This clearly exposes the sensitivity of such bounds to the model under consideration.

### 6.5.1 Upper Bound

In general (on almost all networks) one can use shortest-path routing schemes occupying at most  $O(n^2)$  bits. Relaxing the requirement of shortest path is expressed in the *stretch factor* of a routing scheme. This equals the maximum ratio between the length of a route it produces and the shortest path between the endpoints of that route. The stretch factor of a routing strategy equals the maximal stretch factor attained by any of the routing schemes it generates. The shortest-path routing strategy has stretch factor equal to 1. Allowing stretch factors larger than 1 reduces the space requirements—to as low as  $O(n)$  bits for stretch factors of  $O(\log n)$ , Exercise 6.5.3.

**Theorem 6.5.1** *For shortest-path routing in  $O(\log n)$ -random graphs, local routing functions can be stored in  $6n$  bits per node. Hence the complete routing scheme is represented by  $6n^2$  bits.*

**Proof.** Let  $G$  be an  $O(\log n)$ -random graph on  $n$  nodes. By Lemma 6.4.4 we know that from each node  $u$  we can route via shortest paths to each node  $v$  through the  $O(\log n)$  directly adjacent nodes of  $u$  that have the least indexes. By Lemma 6.4.3,  $G$  has diameter 2. Once the message has reached node  $v$  its destination is either node  $v$  or a direct neighbor of node  $v$  (which is known in node  $v$  by assumption). Therefore, routing functions of size  $O(n \log \log n)$  can be used to do shortest-path routing. We can do better than this.

Let  $A_0 \subseteq V$  be the set of nodes in  $G$  that are not directly connected to  $u$ . Let  $v_1, \dots, v_m$  be the  $O(\log n)$  least indexed nodes directly adjacent to node  $u$  (Lemma 6.4.4), through which we can route via shortest paths to all nodes in  $A_0$ . For  $t = 1, 2, \dots, l$  define  $A_t = \{w \in A_0 - \bigcup_{s=1}^{t-1} A_s : (v_t, w) \in E\}$ . Let  $m_0 = d(A_0)$  and define  $m_{t+1} = m_t - d(A_{t+1})$ . Let  $l$  be the first  $t$  such that  $m_t < n/\log \log n$ . Then we claim that  $v_t$  is connected by an edge in  $E$  to at least  $\frac{1}{3}$  of the nodes not connected by edges in  $E$  to nodes  $u, v_1, \dots, v_{t-1}$ .

**Claim 6.5.1**  $d(A_t) > m_{t-1}/3$  for  $1 \leq t \leq l$ .

**Proof.** Suppose by way of contradiction that there exists a least  $t \leq l$  such that  $|d(A_t) - m_{t-1}/2| > m_{t-1}/6$ . Then we can describe  $G$ , given  $n$ , as follows:

- This discussion in  $O(1)$  bits;
- Nodes  $u, v_t$  in  $2 \log n$  bits, padded with 0's if need be;
- The presence or absence of edges incident with nodes  $u, v_1, \dots, v_{t-1}$  in  $r = n - 1 + \dots + n - (t - 1)$  bits. This gives us the characteristic sequences of  $A_0, \dots, A_{t-1}$  in  $V$ , where a *characteristic sequence* of  $A$  in  $V$  is a string of  $d(V)$  bits with, for each  $v \in V$ , the  $v$ th bit equals 1 if  $v \in A$  and the  $v$ th bit is 0 otherwise;
- A self-delimiting description of the characteristic sequence of  $A_t$  in  $A_0 - \bigcup_{s=1}^{t-1} A_s$ , using Chernoff's bound, Equation 2.4 on page 159, in at most  $m_{t-1} - \frac{2}{3} \left(\frac{1}{6}\right)^2 m_{t-1} \log e + O(\log m_{t-1})$  bits;
- The description  $E(G)$  with all bits corresponding to the presence or absence of edges between  $v_t$  and the nodes in  $A_0 - \bigcup_{s=1}^{t-1} A_s$  deleted, saving  $m_{t-1}$  bits. Furthermore, we also delete all bits corresponding to presence or absence of edges incident with  $u, v_1, \dots, v_{t-1}$  saving a further  $r$  bits.

This description of  $G$  uses at most

$$n(n-1)/2 + O(\log n) + m_{t-1} - \frac{2}{3} \left(\frac{1}{6}\right)^2 m_{t-1} \log e - m_{t-1}$$

bits, which contradicts the  $O(\log n)$ -randomness of  $G$  by Equation 6.10 on page 397, because  $m_{t-1} > n/\log \log n$ .  $\square$

Recall that  $l$  is the least integer such that  $m_l < n/\log \log n$ . We construct the local routing function  $F(u)$  as follows:

- A table of intermediate routing node entries for all the nodes in  $A_0$  in increasing order. For each node  $w$  in  $\bigcup_{s=1}^l A_s$  we enter in the  $w$ th position in the table the unary representation of the least intermediate node  $v$ , with  $(u, v), (v, w) \in E$ , followed by a 0. For the nodes that are not in  $\bigcup_{s=1}^l A_s$  we enter a 0 in their position in the table indicating that an entry for this node can be found in the second table. By Claim 6.5.1, the size of this table is bounded by

$$n + \sum_{s=1}^l \frac{1}{3} \left(\frac{2}{3}\right)^{s-1} sn \leq n + \sum_{s=1}^{\infty} \frac{1}{3} \left(\frac{2}{3}\right)^{s-1} sn \leq 4n;$$

- A table with explicitly binary coded intermediate nodes on a shortest path for the ordered set of the remaining destination nodes. Those nodes had a 0 entry in the first table and there are at most  $m_l < n / \log \log n$  of them, namely the nodes in  $A_0 - \bigcup_{s=1}^l A_s$ . Each entry consists of the code of length  $\log \log n + O(1)$  for the position in increasing order of a node out of  $v_1, \dots, v_m$  with  $m = O(\log n)$  by Lemma 6.4.4. Hence this second table requires at most  $2n$  bits.

The routing algorithm is as follows: The direct neighbors of  $u$  are known in node  $u$  and are routed without the routing table. If we route from start node  $u$  to target node  $w$  that is not directly adjacent to  $u$ , then we do the following. If node  $w$  has an entry in the first table then route over the edge coded in unary, otherwise find an entry for node  $w$  in the second table.

Altogether, we have  $d(F(u)) \leq 6n$ . Slightly more precise counting and choosing  $l$  such that  $m_l$  is the first such quantity  $< n / \log n$  shows  $d(F(u)) \leq 3n$ .  $\square$

## 6.5.2 Lower Bound

In this section, we show that  $\Omega(n^2)$  bits are required to perform routing on Kolmogorov random graphs. Hence the upper bound in Theorem 6.5.1 is tight up to order of magnitude.

**Theorem 6.5.2** *For shortest-path routing in  $o(n)$ -random graphs each local routing function must be stored in at least  $n/2 - o(n)$  bits per node. Hence the complete routing scheme requires at least  $n^2/2 - o(n^2)$  bits to be stored.*

**Proof.** Let  $G$  be an  $o(n)$ -random graph. Let  $F(u)$  be the local routing function of node  $u$  of  $G$ , and let  $d(F(u))$  be the number of bits used to store  $F(u)$ . Let  $E(G)$  be the standard encoding of  $G$  in  $n(n-1)/2$  bits as in Definition 6.4.1. We now give another way to describe  $G$  using some local routing function  $F(u)$ :

- A description of this discussion in  $O(1)$  bits;

- A description of  $u$  in exactly  $\log n$  bits, padded with 0's if needed;
- A description of the presence or absence of edges between  $u$  and the other nodes in  $V$  in  $n - 1$  bits;
- A self-delimiting description of  $F(u)$  in  $d(F(u)) + 2 \log d(F(u))$  bits;
- The code  $E(G)$  with all bits deleted corresponding to edges  $(v, w) \in E$  for each  $v$  and  $w$  such that  $F(u)$  routes messages to  $w$  through the least intermediary node  $v$ . This saves at least  $n/2 - o(n)$  bits since there are at least  $n/2 - o(n)$  nodes  $w$  such that  $(u, w) \notin E$  by Lemma 6.4.2, and since the diameter of  $G$  is 2 by Lemma 6.4.3 there is a shortest path  $(u, v), (v, w) \in E^2$  for some  $v$ . Furthermore, we delete all bits corresponding to the presence or absence of edges between  $u$  and the other nodes in  $V$ , saving another  $n - 1$  bits. This corresponds to the  $n - 1$  bits for edges connected to  $u$  that we added in one connected block above.

In the description we have explicitly added the adjacency pattern of node  $u$  that we deleted elsewhere. This zero-sum swap is necessary to be able to unambiguously identify the adjacency pattern of  $u$  in order to reconstruct  $G$  given  $n$ , as follows: Reconstruct the bits corresponding to the deleted edges using  $u$  and  $F(u)$  and subsequently inserting them in the appropriate positions of the remnants of  $E(G)$ . We can do so because these positions can be simply reconstructed in increasing order. In total, this new description has

$$n(n - 1)/2 + O(1) + O(\log n) + d(F(u)) - n/2 + o(n)$$

bits, which must be at least  $n(n - 1)/2 - o(n)$  by Equation 6.10. Hence,  $d(F(u)) \geq n/2 - o(n)$ , which proves the theorem.  $\square$

### 6.5.3 Average-Case

Consider the average cost, taken over all labeled graphs of  $n$  nodes, of representing a routing scheme for graphs over  $n$  nodes. For a graph  $G$ , let  $T(G)$  be the number of bits used to store its routing scheme. The *average* total number of bits to store the routing scheme for routing over labeled graphs on  $n$  nodes is  $\sum T(G)/2^{n(n-1)/2}$ , with the sum taken over all graphs  $G$  on nodes  $\{1, 2, \dots, n\}$ , that is, the uniform average over all the labeled graphs on  $n$  nodes.

The results on Kolmogorov random graphs above have the following corollaries. Consider the subset of  $(3 \log n)$ -random graphs within the class of  $O(\log n)$ -random graphs on  $n$  nodes. They constitute a fraction of at least  $(1 - 1/n^3)$  of the class of all graphs on  $n$  nodes. The trivial upper bound on the minimal total number of bits for all routing functions together is  $O(n^2 \log n)$  for shortest-path routing on all graphs on

$n$  nodes (or  $O(n^3)$  for full-information shortest-path routing as in Exercise 6.5.7). Simple computation of the average of the total number of bits used to store the routing scheme over all graphs on  $n$  nodes shows that Theorem 6.5.1, Theorem 6.5.2, and Exercise 6.5.3 all hold for the average case.

## Exercises

---

**6.5.1.** [19] Show that there exist labeled graphs on  $n$  nodes where each local routing function must be stored in at least  $(n/2) \log(n/2) - O(n)$  bits per node (hence the complete routing scheme requires at least  $(n^2/2) \log(n/2) - O(n^2)$  bits to be stored).

*Comments.* Source: [H. Buhrman, J.H. Hoepman, and P.M.B. Vitányi, *Proc. 15th ACM Symp. Principles Distr. Comput.*, 1996, pp. 134–142]. This is also the source for the next six exercises.

**6.5.2.** [23] For routing with stretch factor  $< 2$  there exist graphs on  $n$  nodes (almost  $(n/3)!$  such graphs) where the local routing function must be stored in at least  $(n/3) \log n - O(n)$  bits per node at  $n/3$  nodes (hence the complete routing scheme requires at least  $(n^2/9) \log n - O(n^2)$  bits to be stored).

**6.5.3.** [22] (a) Show that routing with any stretch factor  $> 1$  in  $c \log n$ -random graphs can be done with  $n - 1 - (c + 3) \log n$  nodes with local routing functions stored in at most  $\log(n + 1)$  bits per node, and  $1 + (c + 3) \log n$  nodes with local routing functions stored in  $6n$  bits per node (hence the complete routing scheme is represented by less than  $(6c + 20)n \log n$  bits).

(b) Show that routing with stretch factor 2 in  $c \log n$ -random graphs can be done using  $n - 1$  nodes with local routing functions stored in at most  $\log \log n$  bits per node and 1 node with its local routing function stored in  $6n$  bits (hence the complete routing scheme is represented by  $n \log \log n + 6n$  bits).

(c) Show that routing with stretch factor  $(c + 3) \log n$  in  $c \log n$ -random graphs can be done with local routing functions stored in  $O(1)$  bits per node (hence the complete routing scheme is represented by  $O(n)$  bits).

*Comments.* Hint: use Lemma 6.4.4 on page 399 and restricted use of tables (Items (a) and (b)) as in the proof of Theorem 6.5.1 and no tables in Item (c).

**6.5.4.** [31] Show that for shortest-path routing on  $c \log n$ -random graphs, if nodes know their neighbors and nodes may be relabeled by arbitrary identifiers (which therefore can code information), then with labels of size at most  $(1 + (c + 3) \log n) \log n$  bits the local routing functions can

be stored in  $O(1)$  bits per node. Hence the complete routing scheme including the label information is represented by  $(c+3)n \log^2 n + n \log n + O(n)$  bits.

**6.5.5.** [12] Show that in the model of Exercise 6.5.4 the average number of bits to store the routing scheme for graphs of  $n$  nodes is  $O(n \log^2 n)$  for shortest-path routing.

**6.5.6.** [34] Show that for shortest-path routing in  $o(n)$ -random graphs, if the neighbors are not known then the complete routing scheme requires at least  $n^2/32 - o(n^2)$  bits to be stored.

*Comments.* This holds also under a slightly weaker model.

**6.5.7.** In a *full-information* shortest-path routing scheme, the routing function in  $u$  must, for each destination  $v$ , return all edges incident to  $u$  on shortest paths from  $u$  to  $v$ . These schemes allow alternative, shortest paths to be taken whenever an outgoing link is down. Show that for full-information shortest-path routing on  $o(n)$ -random graphs the local routing function requires  $n^2/4 - o(n^2)$  bits for every node (hence the complete routing scheme requires at least  $n^3/4 - o(n^3)$  bits to be stored). This is also the trivial upper bound.

**6.5.8.** [30] In interval routing on a graph  $G = (V, E)$ ,  $V = \{1, \dots, n\}$  each node  $i$  has for each incident edge  $e$  a (possibly empty) set of pairs of node labels representing disjoint intervals with wrap-around. Each pair indicates the initial edge on a shortest path from  $i$  to any node in the interval, and for each node  $j \neq i$  there is such a pair. We are allowed to permute the labels of graph  $G$  to optimize the interval setting.

(a) Show there are graphs such that for each interval routing scheme some incident edge on each of  $\Omega(n)$  nodes is labeled by  $\Omega(n)$  intervals.

(b) Show that for each  $d \geq 3$  there are graphs of maximal node degree  $d$  such that for each interval routing scheme some incident edge on each of  $\Omega(n)$  nodes is labeled by  $\Omega(n/\log n)$  intervals.

*Comments.* Source: E. Kranakis and D. Krizanc, *Proc. 13th Symp. Theoret. Aspects Comput. Sci., Lecture Notes in Computer Science*, Vol. 1046, Springer-Verlag, 1996, pp. 529–540. Item (b) is improved by C. Gavoille and S. Pérennès [*Proc. 15th ACM Symp. Principles Distr. Comput.*, 1996, pp. 125–133], that is, for each interval routing scheme, each of  $\Omega(n)$  edges are labeled by  $\Omega(n)$  intervals. This shows that interval routing can be worse than straightforward coding of routing tables, which can be done in  $O(n^2 \log d)$  bits total.

**6.5.9.** Consider routing schemes for  $n$ -node graphs  $G = (V, E)$ ,  $V = \{1, \dots, n\}$  with maximal node degree  $d$ . Choose the most convenient labeling to facilitate compact routing schemes.

- (a) Show that for each  $d \geq 3$  there are networks for which any shortest-path routing scheme requires a total of  $\Omega(n^2/\log n)$  bits.
- (b) Same as Item (a) but now with stretch factor  $< 2$  requiring a total of  $\Omega(n^2/\log^2 n)$  bits.

*Comments.* Source: E. Kranakis and D. Krizanc, *Ibid.* Item (a) is improved by C. Gavoile and S. Pérennès [*Ibid.*] for  $3 \leq d \leq \epsilon n$  ( $0 < \epsilon < 1$ ) to  $\Theta(n^2 \log d)$ . This is optimal since straightforward coding of routing tables takes  $O(n^2 \log d)$  bits total.

**6.5.10.** Assume we have a computer network consisting of  $n$  computers connected in a ring by bidirectional communication channels. These channels are asynchronous (message transmission takes unknown time) but “first-in-first-out (FIFO)” (messages do not overtake each other). The computers are *anonymous*, that is, they do not have unique identities. To be able to discuss them individually we number them  $1, \dots, n$ . Let  $x$  be any string in  $\{0, 1\}^n$ . At the start of the computation each computer  $i$  in the ring owns a copy of  $x$  and a bit  $y_i$ . Define  $y \equiv x$  if there is an  $s$  ( $0 \leq s < n$ ) such that  $y_{i+s \bmod n} = x_i$  for all  $i$  ( $1 \leq i \leq n$ ). The problem is to compute a Boolean function  $f_x : \{0, 1\}^n \rightarrow \{0, 1\}$  defined by  $f_x(y) = 1$  if  $y \equiv x$  and 0 otherwise. Each computer executes the same algorithm to compute  $f_x$  and eventually outputs the value  $f_x(y)$ .

Show that there is an algorithm to compute  $f_x(\cdot)$ , with  $C(x) \geq n - O(\log n)$ , on an anonymous ring of  $n$  computers using  $O(n \log n)$  bit exchanges for a fraction of at least  $1 - 1/n$  of all  $2^n$  inputs, and hence  $\Theta(n \log n)$  bit exchanges on average.

*Comments.* S. Moran and M. Warmuth [*SIAM J. Comput.*, 22:2(1993), 379–399] have shown that to compute nonconstant functions  $f$  the computers need to exchange  $\Omega(n \log n)$  bits, and that this bound is tight. This creates a “gap” with the case of computing constant  $f$  requiring zero messages. This model assumes that all the processors in the ring are identical (anonymous), all processors run the same program, and the only parameter of the program is the input to the processor. Here we show that  $f_x$  can be computed and requires  $\Theta(n \log n)$  bits on average. Hint: consider  $x$ ’s with  $O(\log n)$  randomness deficiency, and use Lemma 6.9.1 on page 428 and Lemma 6.9.1 on page 428. Source: E. Kranakis, D. Krizanc, and F.L. Luccio, pp. 392–401 in: *Proc. 13th Symp. Math. Found. Comput. Sci.*, Lect. Notes Comput. Sci., Vol. 969, Springer-Verlag, 1995.

## 6.6 Average-Case Complexity of Heapsort

For many algorithms, it is very difficult to analyze their average-case complexity. In average-case analysis, the incompressibility method has an advantage over a probabilistic approach. In the latter approach, one deals with expectations or variances over some ensemble of objects. Using Kolmogorov complexity, we can reason about an incompressible individual object. Because it is incompressible it has all statistical properties with certainty, rather than having them hold with some (high) probability as in a probabilistic analysis. This fact greatly simplifies the resulting analysis. We demonstrate the utility of the incompressibility method to obtain simple and elegant proofs on some examples.

Heapsort is a widely used sorting algorithm. One reason for its prominence is that its running time is *guaranteed* to be of order  $n \log n$ , and it does not require extra memory space. The method was first discovered by J.W.J. Williams [*Comm. ACM*, 7(1964), 347–348], and subsequently improved by R.W. Floyd. Only recently has one succeeded in giving a precise analysis of its average-case performance. I. Munro has suggested a remarkably simple solution using incompressibility.

A “heap” can be visualized as a complete directed binary tree with possibly some rightmost nodes being removed from the deepest level. The tree has  $n$  nodes, each of which is labeled with a different key, taken from a linearly ordered domain. The largest key  $k_1$  is at the root (on top of the heap), and each other node is labeled with a key that is less than the key of its father.

**Definition 6.6.1** Let  $\text{keys}$  be elements of  $\mathcal{N}$ . An array of keys  $k_1, \dots, k_n$  is a *heap* if they are partially ordered such that

$$k_{\lfloor j/2 \rfloor} \geq k_j \text{ for } 1 \leq \lfloor j/2 \rfloor < j \leq n.$$

Thus,  $k_1 \geq k_2, k_1 \geq k_3, k_2 \geq k_4$ , and so on. We consider “in place” sorting of  $n$  keys in an array  $A[1..n]$  without use of additional memory.

**Heapsort** {Initially,  $A[1..n]$  contains  $n$  keys. After sorting is completed, the keys in  $A$  will be ordered as  $A[1] < A[2] < \dots < A[n]$ .}

**Heapify:** {Regard  $A$  as a tree: the root is in  $A[1]$ ; the two sons of  $A[i]$  are at  $A[2i]$  and  $A[2i + 1]$ , when  $2i, 2i + 1 \leq n$ . We convert the tree in  $A$  to a heap.} **Repeat for**  $i = \lfloor n/2 \rfloor, \lfloor n/2 \rfloor - 1, \dots, 1$ : {the subtree rooted at  $A[i]$  is now almost a heap except for  $A[i]$ } push the key, say  $k$ , at  $A[i]$  down the tree (determine which of the two sons of  $A[i]$  possesses the greatest key, say  $k'$  in son  $A[2i + j]$  with  $j$  equals 0 or 1); **if**  $k' > k$  **then** put  $k$  in  $A[2i + j]$  and **repeat** this process pushing  $k'$  at  $A[2i + j]$  down the tree **until** the process reaches a node that does not have a son whose key is greater than the key now at the father node.

**Sort:** **Repeat** for  $i = n, n - 1, \dots, 2$ :  $\{A[1..i]$  contains the remaining heap and  $A[i+1..n]$  contains the already sorted list  $k_{i+1}, \dots, k_n$  of largest elements. By definition, the element on top of the heap in  $A[1]$  must be  $k_i\}$  switch the key  $k_i$  in  $A[1]$  with the key  $k$  in  $A[i]$ , extending the sorted list to  $A[i..n]$ . Rearrange  $A[1..i-1]$  to a heap with the largest element at  $A[1]$ .

It is well known that the Heapify step can be performed in  $O(n)$  time. It is also known that the Sort step takes no more than  $O(n \log n)$  time. We analyze the precise average-case complexity of the Sort step. There are two ways of rearranging the heap: Williams's method and Floyd's method.

**Williams's Method:** {Initially,  $A[1] = k$ .}

**Repeat** compare the keys of  $k$ 's two direct descendants; **if**  $m$  is the larger of the two **then** compare  $k$  and  $m$ ; **if**  $k < m$  **then** switch  $k$  and  $m$  in  $A[1..i-1]$  **until**  $k \geq m$ .

**Floyd's Method:** {Initially,  $A[1]$  is empty.} Set  $j := 1$ ;  
**while**  $A[j]$  is not a leaf **do**:  
  **if**  $A[2j] > A[2j + 1]$  **then**  $j := 2j$   
  **else**  $j := 2j + 1$ ;  
**while**  $k > A[j]$  **do**:  
  {back up the tree until the correct position for  $k$ }  $j := \lfloor j/2 \rfloor$ ;  
**move** keys of  $A[j]$  and each of its ancestors one node upwards;  
  Set  $A[j] := k$ .

The difference between the two methods is as follows. Williams's method goes from the root at the top down the heap. It makes two comparisons with the son nodes and one data movement at each step until the key  $k$  reaches its final position. Floyd's method first goes from the root at the top down the heap to a leaf, making only one comparison each step. Subsequently, it goes from the bottom of the heap up the tree, making one comparison each step, until it finds the final position for key  $k$ . Then it moves the keys, shifting every ancestor of  $k$  one step up the tree. The final positions in the two methods are the same; therefore both algorithms make the same number of key movements. Note that in the last step of Floyd's algorithm, one needs to move the keys carefully upward the tree, avoiding swaps that would double the number of moves.

The heap is of height  $\log n$ . If Williams's method uses  $2d$  comparisons, then Floyd's method uses  $d + 2\delta$  comparisons, where  $\delta = \log n - d$ . Intuitively,  $\delta$  is generally very small, since most elements tend to be near the bottom of the heap. This makes it likely that Floyd's method performs better than Williams's method. We analyze whether this is the case. Assume a uniform probability distribution over the lists of  $n$  keys, so that all input lists are equally likely.

**Theorem 6.6.1** *On average, Heapsort makes  $n \log n + O(n)$  data movements. Williams's method makes  $2n \log n - O(n)$  comparisons on average. Floyd's method makes  $n \log n + O(n)$  comparisons on average.*

**Proof.** Given  $n$  keys, there are  $n!$  ( $\approx n^n e^{-n} \sqrt{2\pi n}$  by Stirling's formula) permutations. Hence we can choose a permutation  $p$  of  $n$  keys such that

$$C(p|n) \geq n \log n - 2n, \quad (6.16)$$

justified by Theorem 2.2.1, page 109. In fact, most permutations satisfy Equation 6.16.

**Claim 6.6.1** Let  $h$  be the heap constructed by the **Heapify** step with input  $p$  that satisfies Equation 6.16. Then

$$C(h|n) \geq n \log n - 6n. \quad (6.17)$$

**Proof.** Assume the contrary,  $C(h|n) < n \log n - 6n$ . Then we show how to describe  $p$ , using  $h$  and  $n$ , in fewer than  $n \log n - 2n$  bits as follows. We will encode the **Heapify** process that constructs  $h$  from  $p$ . At each loop, when we push  $k = A[i]$  down the subtree, we record the path that key  $k$  traveled: 0 indicates a left branch, 1 means a right branch, 2 means halt. In total, this requires  $n \log 3 \sum_i i / 2^{i+1} \leq 2n \log 3$  bits. Given the final heap  $h$  and the above description of updating paths, we can reverse the procedure of **Heapify** and reconstruct  $p$ . Hence,  $C(p|n) < C(h|n) + 2n \log 3 + O(1) < n \log n - 2n$ , which is a contradiction. (The term  $6n$  above can be reduced by a more careful encoding and calculation.)  $\square$

We give a description of  $h$  using the history of the  $n - 1$  heap rearrangements during the Sort step. We only need to record, for  $i := n - 1, \dots, 2$ , at the  $(n - i + 1)$ st round of the Sort step, the final position where  $A[i]$  is inserted into the heap. Both algorithms insert  $A[i]$  into the same slot using the same number of data moves, but a different number of comparisons.

We encode such a final position by describing the path from the root to the position. A path can be represented by a sequence  $s$  of 0's and 1's,

with 0 indicating a left branch and 1 indicating a right branch. Each path  $i$  is encoded in self-delimiting form by giving the value  $\delta_i = \log n - l(s_i)$  encoded in self-delimiting binary form, followed by the literal binary sequence  $s_i$  encoding the actual path. This description requires at most

$$l(s_i) + 2 \log \delta_i \quad (6.18)$$

bits. Concatenate the descriptions of all these paths into sequence  $H$ .

**Claim 6.6.2** We can effectively reconstruct heap  $h$  from  $H$  and  $n$ .

**Proof.** Assume  $H$  is known and the fact that  $h$  is a heap on  $n$  different keys. We simulate the Sort step in reverse. Initially,  $A[1..n]$  contains a sorted list with the least element in  $A[1]$ .

**for**  $i := 2, \dots, n - 1$  **do:** {now  $A[1..i - 1]$  contains the partially constructed heap and  $A[i..n]$  contains the remaining sorted list with the least element in  $A[i]$ } Put the key of  $A[i]$  into  $A[1]$ , while shifting every key on the  $(n - i)$ th path in  $H$  one position down starting from the root at  $A[1]$ . The last key on this path has nowhere to go and is put in the empty slot in  $A[i]$ .

**termination** {Array  $A[1..n]$  contains heap  $h$ .}

□

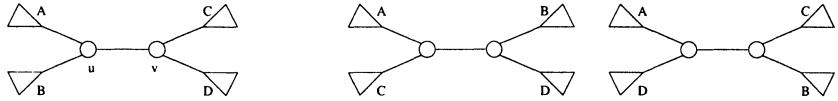
It follows from Claim 6.6.2 that  $C(h|n) \leq l(H) + O(1)$ . Therefore, by Equation 6.17, we have  $l(H) \geq n \log n - 6n$ . By the description in Equation 6.18, we have

$$\sum_{i=1}^n (l(s_i) + 2 \log \delta_i) = \sum_{i=1}^n ((\log n) - \delta_i + 2 \log \delta_i) \leq n \log n - 6n.$$

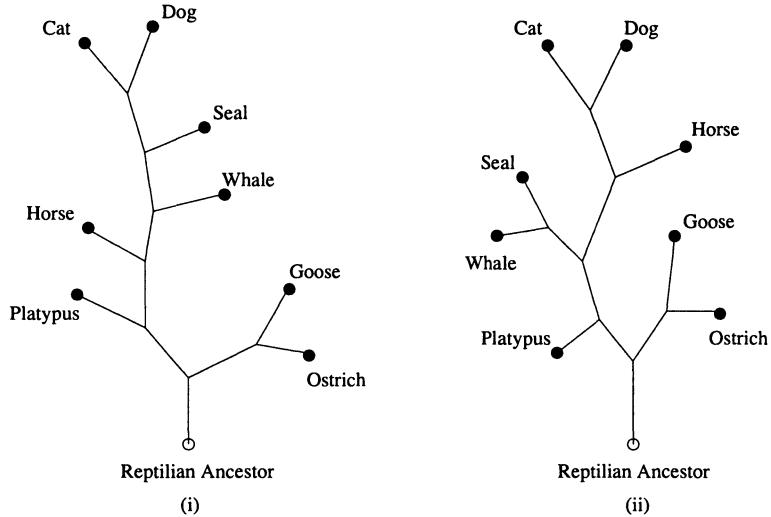
It follows that  $\sum_{i=1}^n (\delta_i - 2 \log \delta_i) \leq 6n$ . This is only possible if  $\sum_{i=1}^n \delta_i = O(n)$ . Therefore, the average path length is at least  $\log n - c$ , for some fixed constant  $c$ . In each round of the Sort step the path length equals the number of data moves. The combined total path length is at least  $n \log n - nc$ .

It follows that starting with heap  $h$ , Heapsort performs at least  $n \log n - O(n)$  data moves. Trivially, the number of data moves is at most  $n \log n$ . Together this shows that Williams's method makes  $2n \log n - O(n)$  key comparisons, and Floyd's method makes  $n \log n + O(n)$  key comparisons.

Since most permutations are Kolmogorov random, these bounds for one random permutation  $p$  also hold for all permutations *on average*. But we can make a stronger statement. We have taken  $C(p|n)$  at least  $\epsilon n$  below the possible maximum, for some constant  $\epsilon > 0$ . Hence, a fraction of at least  $1 - 2^{-\epsilon n}$  of all permutations on  $n$  keys will satisfy the above bounds. □



**FIGURE 6.2.** The two possible nni's on  $(u, v)$ : swap  $B \leftrightarrow C$  or  $B \leftrightarrow D$



**FIGURE 6.3.** The nni distance between (i) and (ii) is 2

## Exercises

**6.6.1.** [25] Improve the  $n \log n - 6n$  bound in Equation 6.17, page 414, by reducing  $6n$  via a better encoding and more precise calculation.

**6.6.2.** [40] In computational biology, evolutionary trees are represented by unrooted unordered binary trees with uniquely labeled leaves and unlabeled internal nodes. Measuring the distance between such trees is useful in biology. A *nearest neighbor interchange (nni)* operation swaps two subtrees that are separated by an internal edge  $(u, v)$ , as shown in Figure 6.2.

- Show that in Figure 6.3 it takes 2 nni moves to convert (i) to (ii).
- Show that  $n \log n + O(n)$  nni moves are sufficient to transform a tree of  $n$  leaves to any other tree with the same set of leaves.
- Prove an  $\Omega(n \log n)$  lower bound for Item (b), using the incompressibility method.

*Comments.* Item (b) is from [K. Culik II and D. Wood, *Inform. Process. Lett.*, 15(1982), 39–42; M. Li, J. Tromp, and L. Zhang, *J. Theoret. Biology*, 182(1996), 463–467]. The latter paper contains principal references related to the nni metric. Item (c) is by D. Sleator, R. Tarjan, and

W. Thurston [*SIAM J. Discr. Math.*, 5(1992), 428–450], who proved the  $\Omega(n \log n)$  lower bound for a more general graph transformation system.

## 6.7 Longest Common Subsequence

Certain problems concerning subsequences and supersequences of a given set of sequences arise naturally in quite practical situations. For example, in molecular biology, a longest common subsequence (of some DNA sequences) is commonly used as a measure of similarity in the comparison of biological sequences. Other applications of longest common subsequences include data compression and syntactic pattern recognition.

**Definition 6.7.1** If  $s = s_1 \cdots s_m$  and  $t = t_1 \cdots t_n$  are two sequences, then  $s$  is a *subsequence* of  $t$ , and equivalently,  $t$  is a *supersequence* of  $s$ , if for some sequence of indices  $i_1 < \cdots < i_m$ , we have  $s_j = t_{i_j}$  for all  $j$  ( $1 \leq j \leq m$ ). Given a finite set of sequences  $S$ , a *shortest common supersequence* (SCS) of  $S$  is a shortest sequence  $s$  such that each sequence in  $S$  is a subsequence of  $s$ . A *longest common subsequence* (LCS) of  $S$  is a longest sequence  $s$  such that each sequence in  $S$  is a supersequence of  $s$ .

It is well known that the SCS and LCS problems are NP-hard. In the worst case, the SCS and LCS problems cannot even be efficiently approximated unless  $P = NP$ . For example, the following is known for the LCS problem. If there is a polynomial-time algorithm that on some input sequences always finds a common subsequence of length  $c > 0$  times the length of the longest common subsequence, then  $P = NP$ . This holds also for the problem as stated in the theorem below. However, many simple heuristic algorithms for SCS and LCS turn out to work well in practice. An incompressibility argument shows that indeed these algorithms perform well *on the average*.

**Definition 6.7.2** Consider LCS problems on an alphabet  $\Sigma = \{a_1, \dots, a_k\}$ . Let  $\text{lcs}(S)$  denote the length of an LCS of a set  $S \subseteq \Sigma^*$  of sequences.

### Algorithm Long-Run

**Step 1** Determine the greatest  $m$  such that  $a^m$  is a common subsequence of all input sequences, for some  $a \in \Sigma$ .

**Step 2** Output  $a^m$  as a common subsequence.

**Theorem 6.7.1** Assume the notation above. Let  $S \subseteq \Sigma^*$  be a set of  $n$  sequences each of length  $n$ , and let  $\epsilon > 0$  be a constant. Algorithm Long-Run outputs a common subsequence of  $S$  of length  $\text{lcs}(S) - O(\text{lcs}(S)^{1/2+\epsilon})$  for a fraction of least  $1 - 1/n^2$  of all inputs, and hence on average.

**Proof.** Assume the notation above. Fix a string  $x$  of length  $n^2$  over  $\Sigma$  with

$$C(x) \geq (n^2 - 2 \log n) \log k. \quad (6.19)$$

Divide  $x$  into  $n$  equal length segments  $x_1, \dots, x_n$ . Choose the set  $S$  in the theorem as  $S = \{x_1, \dots, x_n\}$ .

The following claim is a corollary of the proof of Theorem 2.6.1, page 162, counting each letter as a block of size  $\log k$ , assuming  $k$  is a power of 2.

**Claim 6.7.1** Let  $a \in \Sigma, x_i \in S$ , and let  $\epsilon > 0$  be a constant. Denote the number of occurrences of  $a$  in  $x_i$  by  $m$ . If  $|m - n/k| > n^{1/2+\epsilon}$ , then there is a constant  $\delta > 0$  such that

$$C(x_i|k) \leq (n - \delta n^{2\epsilon}) \log k.$$

A direct proof of this claim is also easy. There are only  $D = \binom{n}{m}(k-1)^{n-m}$  strings of length  $n$  with  $m$  occurrences of  $a$ . Therefore, one can specify  $x_i$  by  $n, k, m$  and its index  $j$ , with  $l(j) = \log D$ , in this ensemble. An elementary estimate by Stirling's formula yields, for some  $\delta > 0$ ,

$$\log \binom{n}{m} (k-1)^{n-m} \leq (n - \delta n^{2\epsilon}) \log k.$$

**Claim 6.7.2**

$$\text{lcs}(S) < \frac{n}{k} + n^{1/2+\epsilon}.$$

**Proof.** Let  $s$  be an LCS of  $S$ . Then  $l(s) = \text{lcs}(S) \leq n$  by definition. Assume, by way of contradiction, that  $l(s)$  is greater than claimed in the lemma. We give a short description of  $x$ , for some fixed  $\delta > 0$ , by saving  $n^\delta \log k$  bits on the description of each  $x_i$  through the use of  $s$ .

Let  $s = s_1 s_2 \dots s_p$ , with  $p = l(s)$ . Fix an  $x_i$ . We will do another encoding of  $x_i$ . We align  $s$  with the corresponding letters in  $x_i$  as far to the left as possible, and rewrite

$$x_i = \alpha_1 s_1 \alpha_2 s_2 \dots \alpha_p s_p z.$$

Here  $\alpha_1$  is the longest prefix of  $x_i$  containing no  $s_1$ ;  $\alpha_2$  is the longest substring of  $x_i$  starting from  $s_1$  containing no  $s_2$ ; and so on. The string  $z$  is the remaining part of  $x_i$  after  $s_p$ . This way,  $\alpha_j$  does not contain an occurrence of letter  $s_j$ , for  $j = 1, \dots, p$ . That is, each  $\alpha_j$  contains only letters in  $\Sigma - s_j$ .

Then  $x_i$  can be considerably compressed with help of  $s$ . Divide  $x_i = yz$  such that the prefix  $y$  is

$$y = \alpha_1 s_1 \alpha_2 s_2 \dots \alpha_p s_p.$$

From  $s$  we can reconstruct which  $k - 1$  letters from  $\Sigma$  appear in  $\alpha_i$ , for each  $i$ . We map  $y$  to an equally long string  $y'$  as follows: For  $i = 1, \dots, p$ , change  $s_i$  to  $a_k$ , in  $y$ . Moreover, change each occurrence of  $a_k$  in  $\alpha_j$  in the letter  $s_j$ . We can map  $y'$  back to  $y$ , using  $s$ , by reversing this process (because the original  $\alpha_j$  did not contain an occurrence of  $s_j$ ).

The letter  $a_k$  occurs at least  $(n/k) + n^{1/2+\epsilon}$  times in  $y'$ , since  $l(s)$  is at least this long. Then, by Claim 6.7.1, for some constant  $\delta > 0$ , we have

$$C(y'|k) \leq (l(y') - \delta n^{2\epsilon}) \log k.$$

From  $y', s, z, k$  we can reconstruct  $x_i$ . (We have defined  $x_i = yz$ .) Giving also the lengths of  $y', s, z$  in self-delimiting format in  $O(\log n)$  bits, we can describe  $x_i$ , given  $k$  and  $s$ , by the number of bits in the right side of the equation below (using  $l(y') + l(z) \leq n$ ):

$$C(x_i|k, s) \leq (n - \delta n^{2\epsilon}) \log k + O(\log n). \quad (6.20)$$

We repeat this for every  $x_i$ . In total, we save more than  $\Omega(n^{1+2\epsilon} \log k)$  bits to encode  $x$ . Thus,

$$C(x|k) / \log k \leq n^2 - \Omega(n^{1+2\epsilon}) + l(s) + O(n \log n) < n^2 - 2 \log n.$$

This contradicts the incompressibility of  $x$  asserted in Equation 6.19.  $\square$

It follows from Claim 6.7.1 and Equation 6.19, by repeating the argument following Equation 6.20 in the proof of Claim 6.7.2, that for some  $\epsilon > 0$ , each  $a \in \Sigma$  occurs in each  $x_1, \dots, x_n$  at least  $(n/k) - O(n^{1/2+\epsilon})$  times. This means that  $a^m$  with

$$l(m) = \frac{n}{k} - O(n^{1/2+\epsilon}) \quad (6.21)$$

is a common subsequence of  $x_1, \dots, x_n$ . By Claim 6.7.2,  $\text{lcs}(S) - l(m) = O(n^{1/2+\epsilon})$ .

Altogether, we have shown that the statement in the theorem is true for this particular input  $x_1, \dots, x_n$ . The fraction of strings of length  $n^2$  satisfying the theorem is at least  $1 - 1/n^2$ , since that many strings satisfy Equation 6.19. The theorem follows by taking the average.  $\square$

## Exercises

---

**6.7.1.** [39] Prove that the expected length of the longest common subsequence of two random binary sequences of length  $n$  is bounded above by  $0.867n$ .

Comments. Hint: use the same encoding scheme as in Section 6.7 and count the number of encodings. The number 0.867 is roughly the root of the following equation  $x - 2x \log x - 2(1-x) \log(1-x) = 2$ . Source: T. Jiang, e-mail, Feb. 24, 1993. This bound was first proved by V. Chvátal and D. Sankoff in [*J. Appl. Probab.*, 12(1975), 306–315]. The precise bound on the expected length is still open.

**6.7.2.** [39] Consider the SCS problem defined in Section 6.7. Prove by incompressibility the following: Let  $S \subseteq \Sigma^*$  be a set of  $n$  sequences of length  $n$ , and let  $\delta = \sqrt{2}/2 \approx 0.707$ . Let  $\text{scs}(S)$  be the length of a SCS of  $S$ . The algorithm Majority-Merge below produces a common supersequence of length  $\text{scs}(S) + O(\text{scs}(S)^\delta)$  on the average.

**Algorithm Majority-Merge** {Input:  $n$  sequences, each of length  $n$ .}

**Step 1** Set supersequence  $s := \epsilon$ ; { $\epsilon$  is the null string.}

**Step 2** {Let the letters  $a$  form a majority among the leftmost letters of the remaining sequences.} Set  $s := sa$  and delete the front  $a$  from these sequences. Repeat this step until no sequences are left.

**Step 3** Output  $s$ .

Comments. Source: T. Jiang and M. Li, *SIAM J. Comput.*, 24:5(1995), 1122–1139. Part of the proof was from [D. Foulser, M. Li, and Q. Yang, *Artificial Intelligence*, 57(1992), 143–181].

## 6.8 Formal Language Theory

---

Part of formal language theory consists in establishing a hierarchy of language families. The main division is the Chomsky hierarchy, with regular languages, context-free languages, context-sensitive languages, and recursively enumerable languages.

A “pumping” lemma (for regular languages) shows that some languages are not regular, but often does not decide which languages are regular and which languages are not. There are many different pumping lemmas, each of them appropriate for limited use. Therefore, some effort has been made to present “pumping” lemmas that are exhaustive, in the sense that they *characterize* the regular languages [J. Jaffe, *SIGACT News*, 10:2(1978), 48–49; D. Stanat and S. Weiss, *SIGACT News*, 14:1(1982), 36–37; A. Ehrenfeucht, R. Parikh, and G. Rozenberg, *SIAM J. Comput.*, 10(1981), 536–541]. These pumping lemmas are complicated and hard to use, while the last reference uses Ramsey theory. Using incompressibility

we find an elegant characterization of the regular languages that makes our intuition about the “finite state”-ness of these languages rigorous and that is easy to apply.

**Definition 6.8.1** Let  $\Sigma$  be a finite nonempty alphabet, and let  $Q$  be a (possibly infinite) nonempty set of states. A *transition function* is a function  $\delta : \Sigma \times Q \rightarrow Q$ . We extend  $\delta$  to  $\delta'$  on  $\Sigma^*$  by  $\delta'(\epsilon, q) = q$  and

$$\delta'(a_1 \dots a_n, q) = \delta(a_n, \delta'(a_1 \dots a_{n-1}, q)).$$

Clearly, if  $\delta'$  is not one-to-one, then the automaton “forgets” because some  $x$  and  $y$  from  $\Sigma^*$  drive  $\delta'$  into the same memory state. An *automaton*  $A$  is a quintuple  $(\Sigma, Q, \delta, q_0, q_f)$  where everything is as above, and  $q_0, q_f \in Q$  are distinguished *initial state* and *final state*, respectively. We call  $A$  a *finite automaton* (FA) if  $Q$  is finite.

An alternative way of looking at it is as follows: We denote “indistinguishability” of a pair of histories  $x, y \in \Sigma^*$  by  $x \sim y$ , defined as  $\delta'(x, q_0) = \delta'(y, q_0)$ . “Indistinguishability” of strings is reflexive, symmetric, transitive, and right-invariant ( $\delta'(xz, q_0) = \delta'(yz, q_0)$  for all  $z$ ). Thus, “indistinguishability” is a right-invariant equivalence relation on  $\Sigma^*$ . It is a simple matter to ascertain this formally.

**Definition 6.8.2** The language accepted by automaton  $A$  as above is the set  $L = \{x : \delta'(x, q_0) = q_f\}$ . A *regular language* is a language accepted by a finite automaton.

It is a straightforward exercise to verify from the definitions the following fact (which will be used later):

**Theorem 6.8.1 (Myhill, Nerode)** *The following statements are equivalent.*

- (i)  $L \subseteq \Sigma^*$  is accepted by some finite automaton.
- (ii)  $L$  is the union of equivalence classes of a right-invariant equivalence relation of finite index on  $\Sigma^*$ .
- (iii) For all  $x, y \in \Sigma^*$  define right-invariant equivalence  $x \sim y$  by the following: for all  $z \in \Sigma^*$  we have  $xz \in L$  iff  $yz \in L$ . Then the number of  $\sim$ -equivalence classes is finite.

Subsequently, closure of finite automaton languages under complement, union, and intersection follow by simple construction of the appropriate  $\delta$  functions from given ones. Details can be found in any textbook on the subject.

**Example 6.8.1** Consider the language  $\{0^k 1^k : k \geq 1\}$ . If it were regular, then the state  $q$  of the accepting finite automaton  $A$ , subsequent to processing  $0^k$ , together with  $A$ , is a description of  $k$ . Namely, by running  $A$ , initialized in state  $q$ , on input consisting of only 1's, the first time  $A$  enters an accepting state is after precisely  $k$  consecutive 1's. The size of the description of  $A$  and  $q$  is bounded by a constant, say  $c$ , that is independent of  $k$ . Altogether, it follows that  $C(k) \leq c + O(1)$ . But choosing  $k$  with  $C(k) \geq \log k$  we obtain a contradiction for all large enough  $k$ . We generalize this observation in the lemma below.  $\diamond$

**Lemma 6.8.1 (KC-Regularity)** *Let  $L \subseteq \Sigma^*$  be regular,  $L_x = \{y : xy \in L\}$ . There is a constant  $c$  such that for each  $x$ , if  $y$  is the  $n$ th string in  $L_x$ , then  $C(y) \leq C(n) + c$ .*

**Proof.** Let  $L$  be a regular language. A string  $y$  such that  $xy \in L$  for some  $x$  can be described by

- this discussion, and a description of the automaton that accepts  $L$ ;
- the state of the automaton after processing  $x$ , and the number  $n$ .

The first item requires  $O(1)$  bits. Thus  $C(y) \leq C(n) + O(1)$ .  $\square$

**Example 6.8.2** Prove that  $\{1^p : p \text{ is not prime}\}$  is not regular. Consider the string  $xy = 1^p$  with  $p$  the  $(k+1)$ th prime. Set  $x = 1^{p'}$ , with  $p'$  the  $k$ th prime. Then  $y = 1^{p-p'}$ , and  $y$  is the lexicographical first element in  $L_x$ . Hence, by Lemma 6.8.1,  $C(p-p') = O(1)$ . But the difference between two consecutive primes grows unbounded. Since there are only  $O(1)$  descriptions of length  $O(1)$ , we have a contradiction.  $\diamond$

**Example 6.8.3** Prove that  $L = \{xx^Rw : x, w \in \{0, 1\}^* - \{\epsilon\}\}$  is not regular (if  $x = x_1 \dots x_m$ , then  $x^R = x_m \dots x_1$ ). Set  $x = (01)^m$ , where  $C(m) \geq \log m$ . Then the lexicographically first word in  $L_x$  is  $y = (10)^m 0$ . Hence,  $C(y) = \Omega(\log m)$ , contradicting the KC-Regularity Lemma.  $\diamond$

**Example 6.8.4** Prove that  $L = \{0^i 1^j : \gcd(i, j) = 1\}$  is not regular. Set  $x = 0^{(p-1)!} 1$ , where  $p > 3$  is a prime,  $l(p) = n$  and  $C(p) \geq \log n - \log \log n$ . Then the lexicographically first word in  $L_x$  is  $1^{p-1}$ , contradicting the KC-regularity lemma.  $\diamond$

**Example 6.8.5** Prove that  $\{p : p \text{ is the standard binary representation of a prime}\}$  is not regular. Suppose the contrary, and  $p_i$  denotes the  $i$ th prime,  $i \geq 1$ . Consider the least binary  $p_m = uv$  ( $= u2^{l(v)} + v$ ), with  $u = \prod_{i < k} p_i$  and  $v$  not in  $\{0\}^*\{1\}$ . Such a prime  $p_m$  exists since each interval  $[n, n +$

$n^{11/20}$ ] of the natural numbers contains a prime [D. Heath-Brown and H. Iwaniec, *Invent. Math.*, 55(1979), 49–69].

Consider  $p_m$  now as an integer,  $p_m = 2^{l(v)} \prod_{i < k} p_i + v$ . Since integer  $v > 1$ , and  $v$  is not divided by any prime less than  $p_k$  (because  $p_m$  is prime), the binary length  $l(v) \geq l(p_k)$ . Because  $p_k$  goes to infinity with  $k$ , the value  $C(v) \geq C(l(v))$  also goes to infinity with  $k$ . But since  $v$  is the lexicographical first suffix, with integer  $v > 1$ , such that  $uv \in L$ , we have  $C(v) = O(1)$  by the KC-Regularity Lemma, which is a contradiction.  $\diamond$

Characterizations (such as the Myhill-Nerode Theorem 6.8.1) of regular languages seem to be practically useful only to show regularity. The need for pumping lemmas stems from the fact that characterizations tend to be very hard to use to show nonregularity. In contrast, the compressibility characterization below is useful for both purposes.

**Definition 6.8.3** Enumerate  $\Sigma^* = \{y_1, y_2, \dots\}$  with  $y_i$  the  $i$ th element in the total order. For  $L \subseteq \Sigma^*$  and  $x \in \Sigma^*$ , let  $\chi = \chi_1 \chi_2 \dots$  be the *characteristic sequence* of  $L_x = \{y : xy \in L\}$ , defined by  $\chi_i = 1$  if  $xy_i \in L$ , and  $\chi_i = 0$  otherwise. We denote  $\chi_1 \dots \chi_n$  by  $\chi_{1:n}$ .

**Theorem 6.8.2 (Regular KC-Characterization)** *There is a constant  $c_L$  depending only on  $L \subseteq \Sigma^*$  such that the following statements are equivalent:*

- (i)  $L$  is regular;
- (ii) for all  $x \in \Sigma^*$ , for all  $n$ ,  $C(\chi_{1:n}|n) \leq c_L$ ;
- (iii) for all  $x \in \Sigma^*$ , for all  $n$ ,  $C(\chi_{1:n}) \leq C(n) + c_L$ ;
- (iv) for all  $x \in \Sigma^*$ , for all  $n$ ,  $C(\chi_{1:n}) \leq \log n + c_L$ .

**Proof.** (i)  $\rightarrow$  (ii). By similar proof as for the KC-Regularity Lemma.

(ii)  $\rightarrow$  (iii)  $\rightarrow$  (iv). Trivial.

(iv)  $\rightarrow$  (i): By (iv) and Claim 6.8.1 below, there are only finitely many distinct  $\chi$ 's associated with the  $x$ 's in  $\Sigma^*$ . Define the right-invariant equivalence relation  $\sim$  by  $x \sim x'$  if  $\chi = \chi'$ . This relation induces a partition of  $\Sigma^*$  into equivalence classes  $[x] = \{y : y \sim x\}$ . Since there is a 1-1 correspondence between the  $[x]$ 's and the  $\chi$ 's, and there are only finitely many distinct  $\chi$ 's, there are also only finitely many  $[x]$ 's. This implies that  $L$  is regular by the Myhill-Nerode theorem: define a finite automaton using one state for each equivalent class, and define the transition function accordingly. The proof of the theorem is finished, apart from proving Claim 6.8.1.

**Claim 6.8.1** For each constant  $c$  there are only finitely many sequences  $\omega \in \{0, 1\}^\infty$  such that for all  $n$ , we have  $C(\omega_{1:n}) \leq \log n + c$ .

This claim is a weaker version of Exercise 2.3.4, page 124, which claims that there are only finitely many such  $\omega$ 's and all of them are recursive reals. D.W. Loveland in [*Inform. Contr.*, 15(1969), 510–526] credits the following result to A.R. Meyer: For each constant  $c$  there are only finitely many  $\omega \in \{0, 1\}^\infty$  with  $C(\omega_{1:n}|n) \leq c$  for all  $n$  and each such  $\omega$  is a recursive real. G.J. Chaitin [*Theoret. Comput. Sci.*, 2(1976), 45–48] improves the condition first to  $C(\omega_{1:n}) \leq C(n) + c$ , and then further to  $C(\omega_{1:n}) \leq \log n + c$ . We provide an alternative and simpler proof, which is sufficient for our purpose, avoiding establishing that the  $\omega$ 's are recursive reals.

**Proof.** Let  $c$  be a positive constant, and let

$$\begin{aligned} A_n &= \{x \in \{0, 1\}^n : C(x) \leq \log n + c\}, \\ A &= \{\omega \in \{0, 1\}^\infty : \forall_{n \in \mathcal{N}} [C(\omega_{1:n}) \leq \log n + c]\}. \end{aligned} \quad (6.22)$$

If the cardinality  $d(A_n)$  of  $A_n$  dips below a fixed constant  $c'$  for infinitely many  $n$ , then  $c'$  is an upper bound on  $d(A)$ . This is because it is an upper bound on the cardinality of the set of prefixes of length  $n$  of the elements in  $A$  for *all*  $n$ .

Fix any  $l \in \mathcal{N}$ . Choose a binary string  $y$  of length  $2l + c + 1$  satisfying

$$C(y) \geq 2l + c + 1. \quad (6.23)$$

Choose  $i$  maximum such that for division of  $y$  in  $y = mn$  with  $l(m) = i$  we have

$$m \leq d(A_n). \quad (6.24)$$

(This holds at least for  $i = 0 = m$ .) Define similarly a division  $y = sr$  with  $l(s) = i + 1$ . By maximality of  $i$ , we have  $s > d(A_r)$ . From the easily proven  $s \leq 2m + 1$ , it then follows that

$$d(A_r) \leq 2m. \quad (6.25)$$

We prove  $l(r) \geq l$ . Since by Equations 6.24 and 6.22 we have

$$m \leq d(A_n) \leq 2^c n,$$

it follows that  $l(m) \leq l(n) + c$ . Therefore,

$$2l + c + 1 = l(y) = l(n) + l(m) \leq 2l(n) + c,$$

which implies that  $l(n) > l$ . Consequently,  $l(r) = l(n) - 1 \geq l$ .

We prove  $d(A_r) = O(1)$ . By dovetailing the computations of the reference universal machine  $U$  (Theorem 2.1.1, page 97) for all programs  $p$  with  $l(p) \leq \log n + c$ , we can enumerate all elements of  $A_n$ . We can reconstruct  $y$  from the  $m$ th element, say  $y_0$ , of this enumeration. Namely,

from  $y_0$  we reconstruct  $n$  since  $l(y_0) = n$ , and we obtain  $m$  by enumerating  $A_n$  until  $y_0$  is generated. By concatenation we obtain  $y = mn$ . Therefore,

$$C(y) \leq C(y_0) + O(1) \leq \log n + c + O(1). \quad (6.26)$$

From Equation 6.23 we have

$$C(y) \geq \log n + \log m. \quad (6.27)$$

Combining Equations 6.26 and 6.27, it follows that  $\log m \leq c + O(1)$ . Therefore, by Equation 6.25,

$$d(A_r) \leq 2^{c+O(1)}.$$

Here,  $c$  is a fixed constant independent of  $n$  and  $m$ . Since  $l(r) \geq l$  and we can choose  $l$  arbitrarily,  $d(A_r) \leq c_0$  for a fixed constant  $c_0$  and infinitely many  $r$ , which implies  $d(A) \leq c_0$ , and hence the claim.  $\square \quad \square$

The KC-Regularity Lemma may be viewed as a corollary of the KC-Characterization Theorem. If  $L$  is regular, then trivially  $L_x$  is regular. It follows immediately that there are only finitely many associated  $\chi$ 's, and each can be specified in at most  $c$  bits,  $c$  a constant depending only on  $L$ . If  $y$  is, say, the  $m$ th string in  $L_x$ , then we can specify  $y$  as the string corresponding to the  $m$ th “1” in  $\chi$ , using only  $C(m) + O(1)$  bits to specify  $y$  (absorbing  $c$  in the  $O(1)$  term). Hence,  $C(y) \leq C(m) + O(1)$ .

## Exercises

**6.8.1.** [10] The KC-Regularity Lemma can be generalized in several ways. Prove the following version. Let  $L$  be regular and  $L_x = \{y : xy \in L\}$ . Let  $\phi$  be a partial recursive function depending only on  $L$  that enumerates strings in  $\Sigma^*$ . For each  $x$ , if  $y$  is the  $n$ th string in the complement of  $L_x$  enumerated by  $\phi$ , then  $C(y) \leq C(n) + c$ , with  $c$  a constant depending only on  $L$ . Use this generalization to give an alternative proof of Example 6.8.4.

*Comments.* Source: M. Li and P. Vitányi, *SIAM J. Comput.*, 24:2(1995), 398–410.

**6.8.2.** [10] Prove that  $\{0^n 1^m : m > n\}$  is not regular.

**6.8.3.** [18] Prove that  $L = \{x \# y : x \text{ appears (possibly nonconsecutively) in } y\}$  is not regular.

**6.8.4.** [20] Prove that  $L = \{x \# y : \text{at least half of } x \text{ is a substring in } y\}$  is not regular.

**6.8.5.** [20] Prove that  $L = \{x \# y \# z : xy = z\}$  is not regular.

**6.8.6.** [37] A DCFL language is a language that is accepted by a deterministic pushdown automaton.

(a) Show that  $\{xx^R : x \in \Sigma^*\}$  and  $\{xx : x \in \Sigma^*\}$  are not DCFL languages, using an incompressibility argument.

(b) Similar to Lemma 6.8.1, the following is a criterion separating DCFL from CFL. Prove it. Let  $L \subseteq \Sigma^*$  be a DCFL and  $c$  be a constant. Let  $x$  and  $y$  be fixed finite words over  $\Sigma$  and  $\omega$  be a recursive sequence over  $\Sigma$ . Let  $u$  be a suffix of  $yx\dots yx$ ,  $v$  a prefix of  $\omega$ , and  $w \in \Sigma^*$  such that

1.  $v$  can be described in  $c$  bits given  $L_u$  in lexicographical order;
2.  $w$  can be described in  $c$  bits given  $L_{uv}$  in lexicographical order; and
3.  $C(v) \geq 2 \log \log l(u)$ .

Then there is a constant  $c'$  depending only on  $L, c, x, y, \omega$  such that  $C(w) \leq c'$ .

(c) Use (b) to prove (a).

*Comments.* Source: M. Li and P. Vitányi, *SIAM J. Comput.*, 24:2(1995), 398–410. In this paper, an incompressibility criterion more general than Item (b) is given for separating DCFL from CFL. See also [S. Yu, *Inform. Process. Lett.*, 31(1989), 47–51] and [M.A. Harrison, *Introduction to formal language theory*, Addison-Wesley, 1978] for basics of formal language theory and traditional approaches to this problem such as iteration lemmas.

**6.8.7.** [27] We have characterized the regular languages using Kolmogorov complexity. It is immediately obvious how to characterize recursive languages in terms of Kolmogorov complexity. If  $L \subseteq \Sigma^*$  and  $\Sigma^* = \{v_1, v_2, \dots\}$  is an effective enumeration, then we define the characteristic sequence  $\chi = \chi_1 \chi_2 \dots$  of  $L$  by  $\chi_i = 1$  if  $v_i \in L$  and  $\chi_i = 0$  otherwise. A language  $L$  is recursive if  $\chi$  is a recursive sequence.

(a) A set  $L \subseteq \Sigma^*$  is recursive iff there exists a constant  $c_L$  (depending only on  $L$ ) such that for all  $n$ , we have  $C(\chi_{1:n}|n) < c_L$ .

(b) If  $L$  is recursively enumerable, then there is a constant  $c_L$  such that for all  $n$ , we have  $C(\chi_{1:n}|n) \leq \log n + c_L$ .

(c) There exists a recursively enumerable set  $L$  such that  $C(\chi_{1:n}) > \log n$ , for all  $n$ .

*Comments.* Items (b) and (c) are Barzdins's Lemma, Theorem 2.7.2, restated. It quantitatively characterizes all recursively enumerable languages in terms of Kolmogorov complexity. Hint for Item (c): Exercise 2.3.4. With  $L$  as in Item (c),  $\Sigma^* - L$  also satisfies Item (b), so Item

(b) cannot be extended to a Kolmogorov complexity characterization of recursively enumerable sets.

**6.8.8.** [23] Assume the terminology in Exercise 6.8.7. Consider  $\chi$  defined in the proof for Item (ii) of Barzdins's Lemma, Theorem 2.7.2. Essentially,  $\chi_i = 1$  if the  $i$ th program started on the  $i$ th input string halts and outputs 0, and  $\chi_i = 0$  otherwise. Let  $A$  be the language with  $\chi$  as its characteristic sequence.

(a) Show that  $A$  is a recursively enumerable set and its characteristic sequence satisfies  $C(\chi_{1:n}) \geq \log n$ , for all  $n$ .

(b) Let  $\chi$  be as in Item (a). Define a sequence  $h$  by

$$h = \chi_1 0^2 \chi_2 0^{2^2} \dots \chi_i 0^{2^i} \chi_{i+1} \dots$$

Show:  $C(h_{1:n}) = O(C(n)) + \Theta(\log \log n)$ . Therefore, if  $h$  is the characteristic sequence of a set  $B$ , then  $B$  is not recursive, but more "sparsely" nonrecursive, as is  $A$ .

*Comments* Item (a) follows from the proof of Barzdins's Lemma, Theorem 2.7.2. Source: J.M. Barzdins, *Soviet. Math. Dokl.*, 9(1968), 1251–1254; D.W. Loveland, *Proc. 1st ACM Symp. Theory Comput.*, 1969, pp. 61–66.

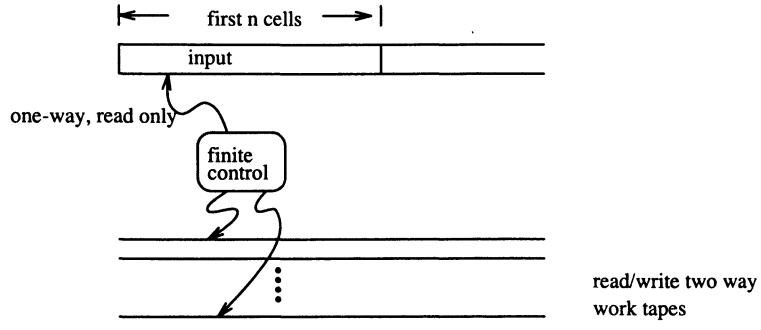
**6.8.9.** [19] The probability that the universal prefix machine  $U$  halts on self-delimiting binary input  $p$ , randomly supplied by tosses of a perfect coin, is  $\Omega$ ,  $0 < \Omega < 1$ . Let  $v_1, v_2, \dots$  be an effective enumeration without repetitions of  $\Sigma^*$ . Define  $L \subseteq \Sigma^*$  such that  $v_i \in L$  iff  $\Omega_i = 1$ . Section 3.6.2 implies that  $K(\Omega_{1:n}) \geq n$  for all but finitely many  $n$ . Show that  $L$  and its complement are not recursively enumerable.

*Comments.* It can be proved that  $L \in \Delta_2^0 - (\Sigma_1^0 \cup \Pi_1^0)$ , in the arithmetic hierarchy. See Section 3.6.2, page 217, and Exercise 1.7.21, page 46.

## 6.9 Online CFL Recognition

The elegant incompressibility proof below demonstrates a lower bound on the time for language recognition by a multitape Turing machine, as shown in Figure 6.4. A multitape Turing machine recognizes a language *online* if before reading each successive input symbol, it decides whether the partial input string scanned so far belongs to the language.

A context-free language is *linear* if it is generated by a linear context-free grammar in which no production rule contains more than one nonterminal symbol on the right-hand side. The known upper bound on the time required for online recognition of a linear context-free language by a multitape Turing machine is  $O(n^2)$ , even if only one work tape is available. We prove a corresponding  $\Omega(n^2/\log n)$  lower bound. Let  $x_i^R$  denote



**FIGURE 6.4.** Multitape Turing machine

$x_i$  written in reverse, and let  $y, x_1, \dots, x_k \in \{0, 1\}^*$ . Define a language  $L$  as

$$L = \{y \# x_1 @ x_2 @ \dots @ x_k : \text{for some } i, x_i^R \text{ is a substring of } y\}.$$

The language  $L$  is linear context-free since it is generated by the following linear grammar, with starting symbol  $S$ :  $S \rightarrow S_1 | S @ | S0 | S1; S_1 \rightarrow 0S_1 | 1S_1 | S_2; S_2 \rightarrow 0S_20 | 1S_21 | S_3 @; S_3 \rightarrow S_30 | S_31 | S_3 @ | S_4 \#; S_4 \rightarrow 0S_4 | 1S_4 | \epsilon$ .

**Theorem 6.9.1** *A multitape Turing machine that online recognizes  $L$  requires time  $\Omega(n^2 / \log n)$ .*

**Proof.** Assume that a multitape Turing machine  $T$  accepts  $L$  online in  $o(n^2 / \log n)$  steps. Choose  $y$  such that  $C(y) \geq l(y) = n$ . Using  $y$ , we will construct a hard input of length  $O(n)$  for  $T$ . The idea of the proof is to construct an input

$$y \# x_1 @ \dots @ x_k @$$

such that no  $x_l$  is a reverse of a substring of  $y$  and yet each  $x_l$  is hard enough to make  $T$  use  $\epsilon n$  steps, for some  $\epsilon > 0$  not depending on  $n$ . If  $k = \Omega(n / \log n)$ , then  $T$  will be forced to take  $\Omega(n^2 / \log n)$  steps. Our task is to prove the existence of such  $x_l$ 's. We need two lemmas:

**Lemma 6.9.1** *Let  $n = l(x)$ , and let  $p$  be a program described in the proof below. Assume that  $C(x|n, p) \geq n$ . Then no substring of length longer than  $2 \log n$  occurs, possibly overlapping, in  $x$  more than once.*

**Proof.** Let  $x = uvw$ , with  $v$  of length greater than  $2 \log n$  occurring exactly twice in  $uv$ . Let this discussion be formulated in terms of a program  $p$  that reconstructs  $x$  from the description below using the

value of  $n$  (given for free). To describe  $x$ , given  $p$  and  $n$ , we only need to concatenate the following information:

- the locations of the start bits of the two occurrences of  $v$  in  $uv$  using  $\log n(n - 1)$  bits;
- the literal word  $uw$ , using exactly  $l(uw)$  bits.

Altogether, this description requires  $n - 2 \log n + \log n(n - 1)$  bits. Since  $C(x|n, p)$  is the shortest such description, we have  $C(x|n, p) < n$ . This contradicts the assumption in the lemma.  $\square$

**Lemma 6.9.2** *If a string has no repetition of length  $m$ , then it is uniquely determined by the set of its substrings of length  $m+1$ ; that is, it is the unique string with no repetition of length  $m$  and with precisely this set of substrings of length  $m+1$ .*

**Proof.** Let  $S$  be the set of substrings of  $x$  of length  $m+1$ . Let  $a, b \in \{0, 1\}$ , and  $u, v, w \in \{0, 1\}^*$ . The prefix of  $x$  of length  $m+1$  corresponds uniquely to the  $ua \in S$  such that for no  $b$  is  $bu$  in  $S$ . For any prefix  $vw$  of  $x$  with  $l(w) = m$ , there is a unique  $b$  such that  $wb \in S$ . Hence, the unique prefix of length  $l(vw) + 1$  is  $vwb$ . The lemma follows by induction.  $\square$

We continue to prove the theorem. By Lemmas 6.9.1 and 6.9.2 we let  $m = 3 \log n$  so that  $y$  is uniquely determined by its set of substrings of length  $m$ . For  $i = 1, \dots, k$ , assume inductively that  $x_1, \dots, x_{i-1}$ , each of length  $m$ , have been chosen so that the input prefix  $y \# x_1 @ \dots @ x_{i-1} @$  does not yet belong to  $L$ , and  $T$  spends at least  $\epsilon n$  steps on each  $x_j @$  block for  $j < i$ .

We claim that for each  $i$  ( $1 \leq i \leq k$ ), there is an  $x_i$  of length  $m$  that is not a reverse substring of  $y$  such that appending  $x_i @$  to the input requires at least  $t = \epsilon n$  additional steps by  $T$ , where  $\epsilon > 0$  does not depend on  $n$ . Setting  $k = n/m$ , this proves the theorem.

Assume, by way of contradiction, that this is not the case. We devise a short description of all the length  $m$  substrings of  $y$ , and hence, by Lemma 6.9.2, of  $y$ . Simulate  $T$  with input  $y \# x_1 @ \dots @ x_{i-1} @$ , and record the following information at the time  $t_0$  when  $T$  reads the last  $@$  sign:

- this discussion;
- the work tape contents within distance  $t = \epsilon n$  of the tape heads;
- the specification of  $T$ , length  $n$ , current state of  $T$ , and locations of  $T$ 's heads.

With this information, one can easily search for all  $x_i$ 's such that  $l(x_i) = m$  and  $x_i^R$  is a substring of  $y$  as follows: Simulate  $T$  from time  $t_0$ , using the above information and with input suffix  $x_i@$ , for  $t$  steps. By assumption, if  $T$  accepts or uses more than  $t$  steps, then  $x_i$  is a reverse substring of  $y$ . If  $\epsilon$  is sufficiently small and  $n$  is large, then all the above information adds up to fewer than  $n$  bits, a contradiction.  $\square$

## Exercises

**6.9.1.** [28] A *k-head deterministic finite automaton*, denoted as  $k$ -DFA, is similar to a deterministic finite automaton except that it has  $k$ , rather than one, one-way read-only input heads. In each step, depending on the current state and the  $k$  symbols read by the  $k$  heads, the machine changes its state and moves some of its heads one step to the right. It stops when all heads reach the end of input, and at this time it accepts the input if it is in a final state. Use the incompressibility method to show that for each  $k \geq 1$  there is a language  $L$  that is accepted by a  $(k+1)$ -DFA but is not accepted by any  $k$ -DFA.

*Comments.* Hint: use the language

$$L_b = \{w_1 \# \cdots \# w_b @ w_b \# \cdots \# w_1 : w_i \in \{0, 1\}^*\}$$

with  $b = \binom{k}{2} + 1$ . Intuitively, when  $w_i$ 's are all random, for each pair of  $w_i$ 's we must have two heads “matching” them concurrently. But a  $k$ -DFA can match only  $\binom{k}{2}$  pairs. This result was first conjectured in 1965 by A. Rosenberg [*IBM J. Res. Develop.*, 10(1966), 388–394]. The case  $k = 2$  was settled by H. Sudborough [*Inform. Contr.*, 30(1976), 1–20] and O.H. Ibarra and C.E. Kim [*Acta Informatica*, 4(1975), 193–200]. The case  $k > 2$  took a decade to settle [A. Yao and R. Rivest, *J. ACM*, 25(1978), 337–340; C.G. Nelson, *Technical Report*, 14-76(1976), Harvard University]. These proofs use more complicated counting arguments. The proof by Kolmogorov complexity is folklore; it can be found in Section 6.4.3 of the first edition of this book.

**6.9.2.** [40/O45] Refer to Exercise 6.9.1 for the definition of  $k$ -DFA's, prove the following. Let  $L = \{x \# y : x \text{ is a substring of } y\}$ .

- (a) No 2-DFA can do string-matching, that is, no 2-DFA accepts  $L$ .
- (b) No 3-DFA accepts  $L$ .
- (c) No  $k$ -DFA accepts  $L$ , for any integer  $k$ .
- (d) [Open] No  $k$ -DFA with “sensing” heads accepts  $L$ , for any  $k$ , where the term *sensing* means the heads can detect each other when they meet.

*Comments.* The results in this exercise were motivated by a conjecture of Z. Galil and J. Seiferas [*J. Comput. System Sci.*, 26:3(1983), 280–294]

that no  $k$ -DFA can do string-matching, for any  $k$ . Galil and Seiferas proved that 6-head two-way DFA can do string-matching in linear time. Item (a) was first proved in [M. Li and Y. Yesha, *Inform. Process. Lett.*, 22(1986), 231–235]; item (b) in [M. Geréb-Graus and M. Li, *J. Comput. System Sci.*, 48(1994), 1–8]. Both of these papers provided useful tools for the final solution, item (c), by T. Jiang and M. Li [*Proc. 25th ACM Symp. Theory Comput.*, 1993, pp. 62–70].

**6.9.3.** [38] A  $k$ -head PDA ( $k$ -PDA) is similar to a pushdown automaton except that it has  $k$  input heads. Prove that  $k + 1$  heads are better than  $k$  heads for PDAs. That is, prove that there is a language that is accepted by a  $(k + 1)$ -PDA but not by any  $k$ -PDA.

*Comments.* Conjectured by M.A. Harrison and O.H. Ibarra [*Inform. Contr.*, 13(1968), 433–470] in analogy to Exercise 6.9.1. Partial solutions were obtained by S. Miyano [*Acta Informatica*, 17(1982), 63–67; *J. Comput. System Sci.*, 27(1983), 116–124]; and M. Chrobak [*Theoret. Comput. Sci.*, 48(1986), 153–181]. The complete solution, using incompressibility, is in [M. Chrobak and M. Li, *J. Comput. System Sci.*, 37:2(1988), 144–155].

**6.9.4.** [35] (a) A  $k$ -pass DFA is just like a usual DFA except that the input head reads the input  $k$  times, from the first symbol to the last symbol, moving right only during each pass. Use incompressibility to show that a  $k$ -pass DFA is exponentially more succinct than a  $(k - 1)$ -pass DFA. In other words, for each  $k$ , there is a language  $L_k$  such that  $L_k$  can be accepted by a  $k$ -pass DFA with  $O(kn)$  states, but the smallest  $(k - 1)$ -pass DFA accepting  $L_k$  requires  $\Omega(2^n)$  states.

(b) A sweeping two-way DFA is again just like a usual finite automaton except that its input head may move in two directions with the restriction that it can only reverse directions at the two ends of the input. If a sweeping two-way DFA makes  $k - 1$  reversals during its computation, we call it a  $k$ -sweep two-way DFA. Show that there is a language  $R_k$  that can be accepted by a  $2k$ -sweep two-way DFA with  $p(k, n)$  states for some polynomial  $p$ , but the smallest  $(2k - 1)$ -sweep two-way DFA accepting  $R_k$  requires an exponential number of states in terms of  $k$  and  $n$ .

*Comments.* W.J. Sakoda and M. Sipser studied sweeping automata in [*Proc. 10th ACM Symp. Theory Comput.*, 1978, pp. 275–286]. They called the  $k$ -pass DFA by the name of “ $k$ -component series FA.” Source: T. Jiang, e-mail, 1992.

## 6.10 Turing Machine Time Complexity

---

It is remarkable that lower bounds on the running time of various computation models can be proved with incompressibility arguments. This method has been quite successfully applied in solving open problems of long standing and to obtain elegant simplifications of existing proofs. In this section, we give one such example proving a lower bound on the time required to simulate a multitape Turing machine by a 1-(work)tape Turing machine (with a one-way input tape).

### Definition 6.10.1

A *k-tape Turing machine with one-way input*, as shown in Figure 6.4, has  $k$  work tapes and a one-way read-only input tape that contains the input. Initially, the input is written on the leftmost tape segment, one symbol per tape square, and the input head scans the leftmost input symbol. The end of the input is delimited by a distinguished endmarker.

Observe that this model with  $k = 1$  is far more powerful than the single-tape Turing machine model of Figure 6.1, page 381, where the single tape serves both as input tape and work tape. For instance, a Turing machine with one work tape apart from the input tape can recognize  $L = \{w\#w^R : w \in \{0, 1\}^*\}$  in real time, in contrast to the  $\Omega(n^2)$  lower bound required in Section 6.1.1 of Section 6.1. The additional marker  $\#$  allows the positive result and does not change the lower bound result.

In the literature, an *online* model is also used. In this model, the input tape is one-way, and on every prefix of the input the Turing machine writes the output, accepting or rejecting the prefix, on a write-only one-way output tape. Proving lower bounds is easier with the online model than with the one-way input model we use in this section. We use the latter model and thus prove stronger results.

A basic question in Turing machine complexity is whether additional work tapes add power. It is known that one tape can online simulate  $n$  steps of  $k$  tapes in  $O(n^2)$  steps. It has been a two-decade-long open question as to whether the known simulation is tight. Kolmogorov complexity has helped to settle this question by an  $\Omega(n^2)$  lower bound; see Exercise 6.10.2.

The tight  $\Omega(n^2)$  lower bound requires a lengthy proof. We provide a weaker form of the result whose simpler proof yet captures the central ideas of the full version.

### Theorem 6.10.1

*It requires  $\Omega(n^{3/2}/\log n)$  time to deterministically simulate a linear time 2-tape Turing machine with one-way input by a 1-tape Turing machine with one-way input.*

**Proof.** We first prove a useful lemma. For a 1-tape Turing machine  $T$ , call  $T$ 's input tape head  $h_1$  and work tape head  $h_2$ . Let  $s$  be a segment of  $T$ 's input and  $R$  be a tape segment on its work tape. We say that  $T$

*maps s into R* if  $h_2$  never leaves tape segment  $R$  while  $h_1$  is reading  $s$ . We say  $T$  maps  $s$  *onto*  $R$  if  $h_2$  traverses the *entire* tape segment  $R$  while  $h_1$  reads  $s$ . The *c.s.* at position  $p$  of the work tape is a sequence of pairs of the form

(state of  $T$ , position of  $h_1$ ),

which records the status of  $T$  when  $h_2$  enters  $p$  each time.

The following lemma states that a tape segment bordered by short *c.s.*'s cannot receive a lot of information without losing some. We assume the following situation: Let the input string start with  $x\#$  where  $x = x_1x_2\dots x_k$  with  $l(x_i) = l(x)/k$  for all  $i$ . Let  $R$  be a segment of  $T$ 's storage tape such that  $T$  maps all blocks in  $S = \{x_{i_1}, \dots, x_{i_l}\}$  into tape segment  $R$ , where  $S \subseteq \{x_i : 1 \leq i \leq k\}$ .

- Lemma 6.10.1 (Jamming Lemma)** *The contents of the storage tape of  $T$  at the time when  $h_1$  moves to the  $\#$  marker can be reconstructed by using only the sequence of blocks  $\bar{S} = \{x_i : 1 \leq i \leq k\} - S$ , the final contents of  $R$ , the two final *c.s.*'s on the left and right boundaries of  $R$ , a description of  $T$ , and a description of this discussion.*

Roughly speaking, if the number of missing bits  $\sum_{j=1}^l l(x_{i_j})$  exceeds the number of added description bits (those for  $R$  and the two crossing sequences around  $R$ ), then the Jamming Lemma implies that either  $x = x_1\dots x_k$  is not incompressible or some information about  $x$  has been lost.

**Proof.** Let the two positions at the left boundary and the right boundary of  $R$  be  $l_R$  and  $r_R$ , respectively. Subdivide the input tape into  $c$ -sized slots with  $c = l(x)/k$ . Put the blocks  $x_j$  of  $\bar{S}$  in their correct positions on the input tape in the following way: In the simulation  $h_1$  reads the input from left to right without backing up. We have a list  $\bar{S}$  of  $c$ -sized blocks available. We put the consecutive blocks of  $\bar{S}$  on the  $c$ -sized slots on the input tape such that the slots, which are traversed by  $h_1$  with  $h_2$  all the time positioned on  $R$ , are left empty. This can be easily determined from the left and right crossing sequences of  $R$ .

Simulate  $T$  with  $h_2$  staying to the left of  $R$  using the *c.s.* at  $l_R$  to construct the work tape contents to the left of  $R$ . Also simulate  $T$  with  $h_2$  staying to the right of  $R$  using the *c.s.* at  $r_R$  to construct the work tape contents to the right of  $R$ . Such a simulation is standard.

We now have obtained the contents of  $T$ 's work tape at the end of processing  $x\#$ , apart from the contents of  $R$ . The final contents of  $R$  is given and is put in position. Together, we now have  $T$ 's work tape contents at the time when  $h_1$  reaches  $\#$ .

Notice that although there are many unknown  $x_i$ 's (in  $S$ ), they are never read since  $h_1$  skips over them because  $h_2$  never goes into  $R$ .  $\square$

To prove the theorem we use the witness language  $L$  defined by

$$L = \{x_1 @ x_2 @ \dots @ x_k \# y_1 @ \dots @ y_l \# 0^i 1^j : x_i = y_j\}. \quad (6.28)$$

Clearly,  $L$  can be accepted in linear time by a 2-tape machine. Assume, by way of contradiction, that a deterministic 1-tape machine  $T$  accepts  $L$  in  $T(n) < c^{-5}n^{3/2}/\log n$  time, for some fixed new constant  $c$  and  $n$  large enough. We derive a contradiction by showing that then some incompressible string must have too short a description.

Assume, without loss of generality, that  $T$  writes only 0's and 1's in its work squares and that  $l(T) = O(1)$  is the number of states of  $T$ . Fix the new constant  $c$  and take the word length  $n$  as large as needed to derive the desired contradictions below and such that the formulas in the sequel are meaningful.

First, choose an incompressible string  $x \in \{0, 1\}^*$  of length  $l(x) = n$  and  $C(x) \geq n$ . Let  $x$  consist of the concatenation of  $k = \sqrt{n}$  substrings,  $x_1, x_2, \dots, x_k$ , each substring  $\sqrt{n}$  bits long. Let

$$x_1 @ x_2 @ \dots @ x_k \#$$

be the initial input segment on  $T$ 's input tape. Let time  $t_{\#}$  be the step at which  $h_1$  reads  $\#$ . If more than  $k/2$  of the  $x_i$ 's are mapped onto a contiguous tape segment of size at least  $n/c^3$ , then  $T$  requires  $\Omega(n^{3/2}/\log n)$  time, which is a contradiction. Therefore, there is a set  $S$  consisting of  $k/2$  blocks  $x_i$ , such that for every  $x_i \in S$  there is a tape segment of  $\leq n/c^3$  contiguous tape squares into which  $x_i$  is mapped. In the remainder of the proof we restrict attention to the  $x_i$ 's in this set  $S$ . Order the elements of  $S$  according to the order of the left boundaries of the tape segments into which they are mapped. Let  $x_m$  be the median.

The idea of the remainder of the proof is as follows: Intuitively, the only thing  $T$  can do before input head  $h_1$  crosses  $\#$  is somehow “copy” the  $x_i$ 's onto its work tape, and afterwards “copy” the  $y_j$ 's onto the work tape. There must be a pair of these  $x_i$  and  $y_j$  that are separated by  $\Omega(n)$  distance, since all these blocks together must occupy  $\Omega(n)$  space. At this time head  $h_1$  still has to read the  $0^i 1^j$  part of the tape. Hence, we can force  $T$  to check whether  $x_i = y_j$ , which means it has to spend about  $\Omega(n^{3/2}/\log n)$  time. To convert this intuition into a rigorous proof we distinguish two cases:

In the first case we assume that many  $x_i$ 's in  $S$  are mapped (jammed) into a small tape segment  $R$ . That is, when  $h_1$  (the input tape head) is reading them,  $h_2$  (the work tape head) is always in this small tape

segment  $R$ . We show that then, contrary to assumption,  $x$  can be compressed (by the Jamming Lemma). Intuitively, some information must have been lost.

In the second case, we assume there is no such “jammed” tape segment and that the records of the  $x_i \in S$  are “spread evenly” over the work tape. In that case, we will arrange the  $y_j$ ’s so that there exists a pair  $(x_i, y_j)$  such that  $x_i = y_j$  and  $x_i$  and  $y_j$  are mapped into tape segments that are far apart, at distance  $\Omega(n)$ . Then we complete  $T$ ’s input with final index  $0^i 1^j$  so as to force  $T$  to match  $x_i$  against  $y_j$ . As in Section 6.1.1, page 380, either  $T$  spends too much time or we can compress  $x$  again, yielding a second contradiction and proving, for large enough  $n$ ,

$$T(n) \geq c^{-5} n^{3/2} / \log n.$$

**CASE 1 (JAMMED)** Assume there are  $k/c$  blocks  $x_i \in S$  and a fixed tape segment  $R$  of length  $n/c^2$  on the work tape such that  $T$  maps all of these  $x_i$ ’s into  $R$ . Let  $S'$  be the set of such blocks.

We will construct a short program that prints  $x$ . Consider the two tape segments of length  $l(R)$  to the left and to the right of  $R$  on the work tape. Call them  $R_l$  and  $R_r$ , respectively. Choose positions  $p_l$  in  $R_l$  and  $p_r$  in  $R_r$  with the shortest c.s.’s in their respective tape segments. Both c.s.’s must be shorter than  $\sqrt{n}/(c^2 \log n)$ . Namely, if the shortest c.s. in either tape segment is at least  $\sqrt{n}/(c^2 \log n)$  long, then  $T$  uses at least

$$\sqrt{n}/(c^2 \log n) \cdot n/c^2$$

steps, and there is nothing to prove. Let tape segment  $R_l'$  ( $R_r'$ ) be the portion of  $R_l$  ( $R_r$ ) right (left) of  $p_l$  ( $p_r$ ).

Now, using the description of

- this discussion (including the text of the program below) and simulator  $T$  in  $O(1)$  bits;
- the values of  $n$ ,  $k$ ,  $c$ , and the positions of  $p_l$ ,  $p_r$  in  $O(\log n)$  bits;
- the literal concatenated list  $\{x_1, \dots, x_k\} - S'$ , using  $n - n/c$  bits;
- the state of  $T$  and the position of  $h_2$  at time  $t_{\#}$  in  $O(\log n)$  bits;
- the two c.s.’s at positions  $p_r$  and  $p_l$  at time  $t_{\#}$  in at most  $2\sqrt{n}(l(T) + O(\log n))$  bits; and
- the contents at time  $t_{\#}$  of tape segment  $R_l'RR_r'$  in at most  $3n/c^2 + O(\log n)$  bits;

we can construct a program to check whether a candidate string  $y$  equals  $x$  by running  $T$  as follows:

Check if  $l(y) = l(x)$ . By the Jamming Lemma (using the above information related to  $T$ 's processing of the initial input segment  $x_1 @ \dots @ x_k \#$ ) reconstruct the contents of  $T$ 's work tape at time  $t_\#$ , the time  $h_1$  gets to the first  $\#$  sign. Divide  $y$  into  $k$  equal pieces and form  $y_1 @ \dots @ y_k$ . Run  $T$ , initialized in the appropriate state, head positions, and work tape contents (at time  $t_\#$ ), as the starting configuration, on each input suffix of the form

$$y_1 @ \dots @ y_k \# 0^i 1^i.$$

By definition of  $L$ , the machine  $T$  can only accept for all  $i$  if and only if  $y = x$ .

This description of  $x$  requires not more than

$$n - \frac{n}{c} + \frac{3n}{c^2} + O(\sqrt{n} \log n) + O(\log n) \leq \gamma n$$

bits, for some constant  $0 < \gamma < 1$  and large enough  $c$  and  $n$ . However, this contradicts the incompressibility of  $x$  ( $C(x) \geq n$ ).

**CASE 2 (NOT JAMMED)** Assume that for each fixed tape segment  $R$ , with  $l(R) = n/c^2$ , there are at most  $k/c$  blocks  $x_i \in S$  mapped into  $R$ .

Fix a tape segment of length  $n/c^2$  into which the median  $x_m$  is mapped. Call this segment  $R_m$ . Then, at most  $k/c$  strings  $x_i$  in set  $S$  are mapped into  $R_m$ . Therefore, for large enough  $c$  (and  $c > 3$ ), at least  $k/6$  of the  $x_i$ 's in  $S$  are mapped into the tape right of  $R_m$ . Let the set of those  $x_i$ 's be  $S_r = \{x_{i_1}, \dots, x_{i_{k/6}}\} \subset S$ . Similarly, let  $S_l = \{x_{j_1}, \dots, x_{j_{k/6}}\} \subset S$  consist of  $k/6$  strings  $x_i$  that are mapped into the tape left of  $R_m$ . Without loss of generality, assume  $i_1 < i_2 < \dots < i_{k/6}$ , and  $j_1 < j_2 < \dots < j_{k/6}$ .

Set  $y_1 = x_{i_1}$ ,  $y_2 = x_{j_1}$ ,  $y_3 = x_{i_2}$ ,  $y_4 = x_{j_2}$ , and so forth. In general, for all integers  $s$ ,  $1 \leq s \leq k/6$ ,

$$y_{2s} = x_{j_s} \quad \text{and} \quad y_{2s-1} = x_{i_s}. \tag{6.29}$$

Using this relationship, we now define an input prefix for  $T$  to be

$$x_1 @ \dots @ x_k \# y_1 @ \dots @ y_{k/3} \#. \tag{6.30}$$

There exists a pair  $y_{2i-1} @ y_{2i}$  that is mapped into a segment of size less than  $n/(4c^2)$ . Otherwise,  $T$  uses at least

$$(k/6)(n/(4c^2)) = n^{3/2}/(24c^2)$$

steps, and there is nothing to prove.

Now this pair  $y_{2i-1}, y_{2i}$  is mapped into a segment with distance at least  $n/c^3$  either to  $x_{is}$  or to  $x_{js}$ . Without loss of generality, let  $y_{2s-1}, y_{2s}$  be mapped  $n/c^3$  away from  $x_{is}$ . So  $y_{2s-1}$  and  $x_{is}$  are separated by a region  $R$  of size  $n/c^3$ . Attach suffix  $0^{is}1^{2s-1}$  to the initial input segment of Equation 6.30 to complete the input to  $T$  to

$$x_1 @ \dots @ x_k \# y_1 @ \dots @ y_{k/3} \# 1^{is} 1^{2s-1}. \quad (6.31)$$

So at the time when  $T$  reads the second  $\#$  sign,  $x_{is}$  is mapped into the tape left of  $R$ , and  $y_{2s-1}$ , which is equal to  $x_{is}$ , is mapped into the tape right of  $R$ .

Determine position  $p$  in  $R$  that has the shortest *c.s.* of  $T$ 's computation on the input of Equation 6.31. If this *c.s.* is longer than  $\sqrt{n}/(c^2 \log n)$ , then  $T$  uses at least

$$(n/c^3) \cdot \sqrt{n}/(c^2 \log n)$$

steps, and there is nothing to prove. Therefore, assume the shortest *c.s.* has length at most  $\sqrt{n}/(c^2 \log n)$ . Then again we can construct a short program  $P$  to accept *only*  $x$  by a “cut-and-paste” argument, and show that it yields too short a description of  $x$ . Using the description of

- this discussion (including the text of the program  $P$  below) and simulator  $T$  in  $O(1)$  bits;
- the values of  $n, k, c$ , and the position  $p$  in  $O(\log n)$  bits;
- $n - \sqrt{n}$  bits for  $S - \{x_{is}\}$ ;
- $O(\log n)$  bits for the index  $i_s$  of  $x_{is}$  to place it correctly on the input tape; and
- $\leq \sqrt{n}/c$  bits to describe the *c.s.* of length  $\sqrt{n}/(c^2 \log n)$  at  $p$  (assuming  $c \gg l(T)$ );

we can construct a program to reconstruct  $x$  as follows:

Construct the input of Equation 6.31 on  $T$ 's input tape with the two blocks  $x_{is}$  and  $y_{2s-1}$  filled with blanks. Now we search for  $x_{is}$  as follows. For each candidate  $z$  with  $l(z) = \sqrt{n}$  put  $z$  in the  $y_{2s-1}$ 's position and do the following simulation:

Using the *c.s.* at point  $p$  we run  $T$  such that  $h_2$  always stays at the right of  $p$  ( $y_{2s-1}$ 's side). Whenever  $h_2$  encounters  $p$ , we check whether the current status matches the corresponding ID in the *c.s.* If it does, then we use the next ID of the *c.s.* to continue. If in the course of this

simulation process  $T$  rejects or there is a mismatch (that is, when  $h_2$  gets to  $p$ , machine  $T$  is not in the same state or  $h_1$ 's position is not as indicated in the *c.s.*), then  $z \neq x_{i_s}$ . If the crossing sequence at  $p$  of  $T$ 's computation for candidate  $z$  matches the prescribed *c.s.*, then we know that  $T$  would accept the input of Equation 6.31 with  $y_{2s-1}$  replaced by  $z$ . Therefore,  $z = x_{i_s}$ .

The description of  $x$  requires not more than

$$n - \sqrt{n} + \sqrt{n}/c + O(\log n) \leq n - \gamma\sqrt{n}$$

bits for some positive  $\gamma > 0$  and large enough  $c$  and  $n$ . This contradicts the incompressibility of  $x$  ( $C(x) \geq n$ ) again.

Case 1 and Case 2 complete the proof that  $T(n) \geq c^{-5}n^{3/2}/\log n$ .  $\square$

## Exercises

---

**6.10.1.** [33] Consider the 1-tape Turing machine as in Section 6.1.1, page 380. Let the input be  $n/\log n$  integers each of size  $O(\log n)$ , separated by # signs. The Element Distinctness problem is to decide whether all these integers are distinct. Prove that the Element Distinctness problem requires  $\Omega(n^2/\log n)$  time on such 1-tape Turing machine.

*Comments.* A similar bound also holds for 1-tape nondeterministic Turing machines. Source: [A. López-Ortiz, *Inform. Process. Lett.*, 51:6(1994), 311–314].

**6.10.2.** [42] Extend the proof of Theorem 6.10.1 to prove the following: Simulating a linear-time 2-tape deterministic Turing machine by a 1-tape deterministic Turing machine requires  $\Omega(n^2)$  time. (Both machines of the one-way input model.)

*Comments.* Hint: set the block size for  $x_i$  to be a large constant, and modify the language to one that requires comparison of  $\Omega(n)$  pairs of  $x_i$ 's and  $y_j$ 's. The lower bound is optimal since it meets the  $O(n^2)$  upper bound of [J. Hartmanis and R. Stearns, *Trans. Amer. Math. Soc.*, 117(1969), 285–306]. Source: W. Maass, *Trans. Amer. Math. Soc.*, 292(1985), 675–693; M. Li and P. Vitányi, *Inform. Comput.*, 78(1988), 56–85.

**6.10.3.** [38] A  $k$ -pushdown store machine is similar to a  $k$ -tape Turing machine with one-way input except that the  $k$  work tapes are replaced by  $k$  pushdown stores. Prove: simulating a linear-time 2-pushdown store deterministic machine with one-way input by a 1-tape nondeterministic Turing machine with one-way input requires  $\Omega(n^{3/2}/\sqrt{\log n})$  time.

*Comments.* Source: M. Li and P. Vitányi, *Inform. Comput.*, 78(1988), 56–85. This bound is optimal since it is known that simulating a linear-time 2-pushdown store deterministic machine with one-way input by

a 1-tape nondeterministic Turing machine with one-way input can be done in  $O(n^{3/2}\sqrt{\log n})$  time [M. Li, *J. Comput. System Sci.*, 7:1(1988), 101–116].

**6.10.4.** [44] Show that simulating a linear time 2-tape deterministic Turing machine with one-way input by a 1-tape nondeterministic Turing machine with one-way input requires  $\Omega(n^2/((\log n)^2 \log \log n))$  time.

*Comments.* Hint: let  $S$  be a sequence of numbers from  $\{0, \dots, k-1\}$ , where  $k = 2^l$  for some  $l$ . Assume that each number  $b \in \{0, \dots, k-1\}$  is somewhere in  $S$  adjacent to the number  $2b \pmod{k}$  and  $2b+1 \pmod{k}$ . Then for every partition of  $\{0, \dots, k-1\}$  into two sets  $G$  and  $R$  such that  $d(G), d(R) > k/4$  there are at least  $k/(c \log k)$  (for some fixed  $c$ ) elements of  $G$  that occur somewhere in  $S$  adjacent to a number from  $R$ . Subsequently prove the lower bound using the language  $L$  defined by: All members in  $L$  are strings in  $\{0, 1\}^*$  and are constructed as follows: Let  $u = u_1 \dots u_k$ , where  $u_i$ 's are of equal length. Form  $uu = u_1 \dots u_{2k}$  with  $u_{k+i} = u_i$ . Then inserting  $u_i$  between  $u_{2i-1}$  and  $u_{2i}$  for  $1 \leq i \leq k$  results in a member in  $L$ . Source: W. Maass, *Trans. Amer. Math. Soc.*, 292(1985), 675–693. The language  $L$  defined in this hint will not allow us to obtain an  $\Omega(n^2)$  lower bound. Let  $G = \langle Z_n, E_{ab} \rangle$ , where  $Z_n = \{0, 1, \dots, n-1\}$ ,  $E_{ab} = \{(i, j) : j \equiv (ai + b) \pmod{n} \text{ for } i \in Z_n\}$ , and  $a$  and  $b$  are fixed positive integers. Then  $G$  has a separator, a set of nodes whose removal separates  $G$  into two disconnected, roughly equal-sized components of size  $O(n/\sqrt{\log_a n})$ . Using such a separator,  $L$  can be accepted in subquadratic time by a 1-tape online deterministic machine [M. Li, *J. Comput. System Sci.*, 7:1(1988), 101–116].

**6.10.5.** [46] Prove: simulating a linear time 2-tape deterministic Turing machine with one-way input by a 1-tape nondeterministic Turing machine with one-way input requires  $\Omega(n^2/\log^{(k)} n)$  time, for any  $k$ , where  $\log^{(k)} = \log \log \dots \log$  is the  $k$ -fold iterated logarithm. This improves the result in Exercise 6.10.4.

*Comments.* Source: Z. Galil, R. Kannan, and E. Szemerédi [*J. Comput. System Sci.*, 38(1989), 134–149; *Combinatorica*, 9(1989), 9–19].

**6.10.6.** [O47] Does simulating a linear time 2-tape deterministic Turing machine with one-way input by a 1-tape nondeterministic Turing machine with one-way input require  $\Omega(n^2)$  time?

**6.10.7.** [46] A  $k$ -queue machine is similar to a  $k$ -tape Turing machine with one-way input except with the  $k$  work tapes replaced by  $k$  work queues. A queue is a first-in last-out (FIFO) device. Prove (with one-way input understood):

- (a) simulating a linear time 1-queue machine by a 1-tape Turing machine requires  $\Omega(n^2)$  time.

- (b) simulating a linear time 1-queue machine by a 1-tape nondeterministic Turing machine requires  $\Omega(n^{4/3}/\log^{2/3} n)$  time.
- (c) simulating a linear time 1-pushdown store machine (which accepts precisely CFLs) by a 1-queue machine, deterministically or nondeterministically, requires  $\Omega(n^{4/3}/\log n)$  time.

*Comments.* Items (a) and (b) are from [M. Li and P.M.B. Vitányi, *Inform. Comput.*, 78(1988), 56–85]; Item (c) is from [M. Li, L. Longpré, and P.M.B. Vitányi, *SIAM J. Comput.*, 21:4(1992), 697–712]. The bound in Item (a) is tight. The bound in Item (b) is not tight; the best upper bound is  $O(n^{3/2}\sqrt{\log n})$ , in [M. Li, *J. Comput. System Sci.*, 7:1(1988), 101–116]. The bound in Item (c) is not tight; the upper bound is known to be  $O(n^2)$  (also to simulate a 2-pushdown store machine).

**6.10.8.** [43] Use the terminology of Exercise 6.10.7, with one-way input understood.

- (a) Show that simulating a linear time deterministic 2-queue machine by a deterministic 1-queue machine takes  $\Omega(n^2)$  time.
- (b) Show that simulating a linear time deterministic 2-queue machine by a nondeterministic 1-queue machine takes  $\Omega(n^2/(\log^2 n \log \log n))$  time.
- (c) Show that simulating a linear time deterministic 2-tape Turing machine by nondeterministic 1-queue machine takes  $\Omega(n^2/\log^2 n \log \log n)$  time.

*Comments.* The upper bounds in all cases are  $O(n^2)$  time. Source: M. Li, L. Longpré, and P.M.B. Vitányi, *SIAM J. Comput.*, 21:4(1992), 697–712. For additional results on simulating  $(k+1)$ -queue machines and 2-tape or multitape machines by  $k$ -queue machines see [M. Hühne, *Theoret. Comput. Sci.*, 113:1(1993), 75–91].

**6.10.9.** [38] Consider the stronger offline deterministic Turing machine model with a two-way read-only input tape. Given an  $l \times l$  matrix  $A$ , with  $l = \sqrt{n/\log n}$  and element size  $O(\log n)$ , arranged in row-major order on the two-way (1-dimensional) input tape,

- (a) Show that one can transpose  $A$  (that is, write  $A^T$  on a work tape in row-major form) in  $O(n \log n)$  time on such a Turing machine with two work tapes.
- (b) Show that it requires  $\Omega(n^{3/2}/\sqrt{\log n})$  time on such a Turing machine with one work tape to transpose  $A$ .
- (c) From Items (a) and (b), obtain a lower bound on simulating two work tapes by one work tape for the above machines.

*Comments.* Source: M. Dietzfelbinger, W. Maass and G. Schnitger, *Theoret. Comput. Sci.*, 82:1(1991), 113–129.

**6.10.10.** [37] We analyze the speed of copying strings for Turing machines with a two-way input tape and one or more work tapes.

- (a) Show that such a Turing machine with one work tape can copy a string of length  $s$ , initially positioned on the work tape, to a work tape segment that is  $d$  tape cells removed from the original position in  $O(d + sd/\log \min(n, d))$  steps. Here  $n$  denotes the length of the input.
- (b) Show (by the incompressibility method) that the upper bound in Item (a) is optimal. For  $d = \Omega(\log n)$ , such a Turing machine with one work tape requires  $\Omega(sd/\log \min(n, d))$  steps to copy a string of length  $s$  across  $d$  tape cells.
- (c) Use Item (a) to show that such Turing machines can simulate  $f(n)$ -time bounded multitape Turing machines in  $O(f(n)^2/\log n)$  steps. This is faster by a multiplicative factor  $\log n$  than the straightforward simulations.

*Comments.* Source: M. Dietzfelbinger, *Inform. Process. Lett.*, 33(1989/1990), 83–90.

**6.10.11.** [38] Show that it takes  $\Theta(n^{5/4})$  time to transpose a Boolean matrix on a Turing machine with a two-way read-only input tape, a work tape, and a one-way write-only output tape. That is, the input is a  $\sqrt{n} \times \sqrt{n}$  matrix  $A$  that is initially given on the input tape in row-major order. The Turing machine must output  $A$  in column-major order on the output tape.

*Comments.* Hint: for the upper bound, partition columns of  $A$  into  $n^{1/4}$  groups each with  $n^{1/4}$  consecutive columns. Process each group separately (you need to put each group into a smaller region first). The lower bound proof is harder. Fix a random matrix  $A$ . Let the simulation time be  $T$ . Split  $T$  into  $O(n^{3/4})$  printing intervals. Within each interval  $O(n^{1/4})$  entries of  $A$  are printed and half such intervals last fewer than  $n^{1/2}$  steps. Also, split the work tape into (disjoint) intervals of size  $O(n^{1/2})$  so that one quarter of the printing intervals do not overlap with two work tape intervals. Say, an input bit is mapped to a work tape interval if while the input head is reading that bit, the work tape head is in this interval. A work tape interval is *underinformed* if many bits in the printing interval it corresponds to are not mapped into this work tape interval before they are printed. Show that if there are many underinformed work tape intervals,  $A$  is compressible. Then show that if there are not many underinformed intervals, there must be many “overburdened intervals,” that is, more bits than the length of such intervals are mapped in to each interval. This also implies the compressibility of  $A$ .

Source: M. Dietzfelbinger and W. Maass, *Theoret. Comput. Sci.*, 108(1993), 271–290.

**6.10.12.** [O44] Obtain a tight bound for simulating two work tapes by one work tape for Turing machines with a two-way input tape.

*Comments.* W. Maass, G. Schnitger, E. Szemerédi, and G. Turan, [*Computational Complexity*, 3(1993), pp. 392–401] proved (not using Kolmogorov complexity) the following: Let  $L = \{A\#B : A = B^t \text{ and } a_{ij} \neq 0 \text{ only when } i, j \equiv 0 \pmod{\log m} \text{ where } m = 2^k, \text{ for some } k, \text{ is the size of matrices}\}$ . Accepting  $L$  requires  $\Omega(n \log n)$  time on a Turing machine with a two-way input tape and one work tape. Since  $L$  can be accepted in  $O(n)$  time by a similar machine with two work tapes, this result implies that two tapes are better than one for deterministic Turing machines with a two-way input tape. An upper bound to this question is given in Exercise 6.10.10 Item (c).

**6.10.13.** [40] Consider an online deterministic Turing machine with a one-way input tape, some work tapes/pushdown stores and a one-way output tape. The result of computation is written on the output tape. “Online simulation” means that after reading a new input symbol the simulating machine must write down precisely the output of the simulated machine for the processed initial input segment before it goes on to read the next input symbol. Show the following:

- (a) It requires  $\Omega(n(\log n)^{1/(k+1)})$  time to online simulate  $k+1$  pushdown stores by  $k$  tapes.
- (b) On-line simulating one tape plus  $k-1$  pushdown stores by  $k$  pushdown stores requires  $\Omega(n(\log n)^{1/(k+1)})$  time.
- (c) Each of the above lower bounds holds also for a probabilistic simulation where the probabilistic simulator flips a random coin to decide the next move. (No error is allowed. The simulation time is the average taken over all coin-tossing sequences.)

*Comments.* Item (a) is from [W.J. Paul, *Inform. Contr.*, 53(1982), 1–8]. Item (b) is due to P. Dúriš, Z. Galil, W.J. Paul, and R. Reischuk [*Inform. Contr.*, 60(1984), 1–11]. Item (c) is from [R. Paturi, J. Simon, R. Newman-Wolfe, and J. Seiferas, *Inform. Comput.*, 88(1990), 88–104]. The latter paper also includes proofs for Items (a) and (b).

**6.10.14.** [40] Consider the machine model in Exercise 6.10.13, except that the work tapes are two-dimensional. Such a machine works in real time if at each step it reads a new input symbol and is online. (Then it processes and decides each initial  $m$ -length segment in precisely  $m$  steps.) Show that for such machines, two work tapes with one head each cannot real-time simulate one work tape with two independent heads.

*Comments.* Source: W.J. Paul, *Theoret. Comput. Sci.*, 28(1984), 1–12.

**6.10.15.** [48] Like in Exercise 6.10.14, consider the Turing machine model of Exercise 6.10.13 but this time with 1-dimensional tapes. Show

that a Turing machine with two single-head one-dimensional tapes cannot recognize the set  $\{x2x' : x \in \{0,1\}^*\text{ and }x'\text{ is a prefix of }x\}$ , in real time, although it can do so with three tapes, two two-dimensional tapes, or one two-head tape, or in linear time with just one tape.

*Comments.* This is considerably more difficult than the problem in Exercise 6.10.14. In particular, this settles the longstanding conjecture that a two-head Turing machine can recognize more languages in real time if its heads are on the *same* one-dimensional tape than if they are on *separate* one-dimensional tapes. Source: partial results in [W.J. Paul, *Ibid.*; P.M.B. Vitányi, *J. Comput. System Sci.*, 29(1984), 303–311]. This thirty year open question was finally settled by T. Jiang, J.I. Seiferas, and P.M.B. Vitányi in [*J. ACM*, to appear].

**6.10.16.** [38] A *tree work tape* is a complete, infinite, rooted binary tree used as storage medium (instead of a linear tape). A work tape head starts at the root and can in each step move to the direct ancestor of the currently scanned node (if it is not the root) or to either one of the direct descendants. A *multihead tree machine* is a Turing machine with a one-way linear input tape, one-way linear output tape, and several tree work tapes each with  $k \geq 1$  heads. We assume that the finite control knows whether two work tape heads are on the same node or not. A *d-dimensional work tape* consists of nodes corresponding to  $d$ -tuples of integers, and a work tape head can in each step move from its current node to a node with each coordinate  $\pm 1$  of the current coordinates. Each work tape head starts at the origin, which is the  $d$ -tuple with all zeros. A *multihead d-dimensional machine* is like the multihead tree machine but with  $d$ -dimensional work tapes.

- (a) Show that simulating a multihead tree machine online by a multihead  $d$ -dimensional machine requires time  $\Omega(n^{1+1/d}/\log n)$  in the worst case. Hint: prove this for a tree machine with one tree tape with a single head that runs in real time.
- (b) Show the same lower bound as in Item (a) where the multihead  $d$ -dimensional machine is made more powerful by allowing the work tape heads also to move from their current node to the current node of any other work tape head in a single step.

*Comments.* Source: M.C. Loui, *SIAM J. Comput.*, 12(1983), 463–472. The lower bound in Item (a) is optimal since it can be shown that every multihead tree machine of time complexity  $t(n)$  can be simulated online by a multihead  $d$ -dimensional machine in time  $O(t(n)^{1+1/d}/\log t(n))$ . It is known that every log-cost RAM (Exercise 6.10.17) can be simulated online in real time by a tree machine with one multihead tree tape [W.J. Paul and R. Reischuk, *J. Comput. System Sci.*, 22(1981), 312–327]. Hence, we can simulate RAMs online by  $d$ -dimensional machines

in time that is bounded above and below by the same bounds as the simulation of tree machines. See also [M.C. Loui, *J. Comput. System Sci.*, 28(1984), 359–378].

**6.10.17.** [37] A log-cost random access machine (log-cost RAM) has the following components: an infinite number of registers each capable of holding an integer and a finite sequence of labeled instructions including “output,” “branch,” “load/store,” “add/subtract between two registers.” The time cost for execution of each instruction is the sum of the lengths of the integers involved.

(a) Every tree machine with several tree tapes, each with one head, of time complexity  $t$  can be simulated online by a log-cost RAM of time complexity  $O(t \log t / \log \log t)$ . Show that this is optimal.

(b) Show that online simulating a linear time log-cost RAM by a  $d$ -dimensional Turing machine requires  $\Omega(n^{1+1/d} / \log n (\log \log n)^{1+1/d})$ .

*Comments.* Source: D.R. Luginbuhl, Ph.D. Thesis, 1990, University of Illinois at Urbana-Champaign; M.C. Loui and D.R. Luginbuhl, *SIAM J. Comput.* 21:5(1992), 959–971; M.C. Loui and D.R. Luginbuhl, *Math. Systems Theory*, 25:4(1992), 293–308.

**6.10.18.** [38] Consider the machine models in Exercise 6.10.13 Item (c). All machines below have one multidimensional tape with one head.

(a) Show that an  $l$ -dimensional machine running in time  $T$  can be simulated by a probabilistic  $k$ -dimensional machine running in time  $O(T^r (\log T)^{1/k})$ , where  $r = 1 + 1/k - 1/l$ .

(b) Show that a probabilistic  $k$ -dimensional machine requires time  $\Omega(T^r)$  to simulate an  $l$ -dimensional machine running in time  $T$ , where  $r = 1 + 1/k - 1/l$ .

*Comments.* Source: [N. Pippenger, *Proc. 14th ACM Symp. Theory Comput.*, 1982, pp. 17–26]. Pippenger used Shannon’s information measure to prove Item (b).

**6.10.19.** [30] Prove that if the number of states is fixed, then a 1-tape nondeterministic Turing machine with no separate input tape (with only one read/write two-way tape) can accept more sets within time bound  $a_2 n^a$  than within  $a_1 n^a$ , for  $0 < a_1 < a_2$  and  $1 < a < 2$ .

*Comments.* Source: K. Kobayashi, *Theoret. Comput. Sci.*, 40(1985), 175–193.

## 6.11 Parallel Computation

---

We prove a tight lower bound on the number of parallel steps needed to sum  $n$  integers on a widely used parallel computing model, using an incompressibility argument.

A *parallel random access machine* (PRAM), also called a “concurrent-read and concurrent-write priority PRAM,” consists of a finite number of processors (computers, each with an infinite local memory and infinite computing power), indexed as  $P(1), P(2), P(3), \dots$ , and an infinite number of shared memory cells  $c(i)$ ,  $i = 1, 2, \dots$ , each capable of holding any integer. Initially, the input is contained in the first  $n$  memory cells. The number of processors is polynomial in  $n$ . Each step of the computation consists of all processors in parallel executing three phases as follows. Each processor

1. reads from a shared memory cell,
2. performs any deterministic computation, and
3. may attempt writing into some shared memory cell.

At each step each processor is in some *state*. The actions and the next state of each processor at each step depend on the current state and the value read. In case of a *write conflict*, that is, more than one processor tries to write to the same memory cell, we choose to let the processor with the minimum index succeed in writing. By the end of computing, the shared memory contains the output. This model is often used to analyze parallel algorithms.

Addition or multiplication of  $n$  numbers in parallel on such PRAM’s is fundamental. We present a tight  $\Omega(\log n)$  lower bound on the number of parallel steps required to add or multiply  $n$  numbers, using Kolmogorov complexity. At first thought, the  $\log n$  lower bound may look trivial. However, the proof is a little subtler than it appears, as Exercise 6.11.1 demonstrates.

**Theorem 6.11.1** *Adding (or multiplying)  $n$  integers, each of at least  $n^\epsilon$  bits for a fixed  $\epsilon > 0$ , requires  $\Omega(\log n)$  parallel steps on a PRAM.*

**Proof.** Suppose, by way of contradiction, that a PRAM  $M$  with  $P = n^d$  processors adds  $n$  integers in  $T < (\epsilon/2) \log n$  parallel steps for infinitely many  $n$ . Consider large enough  $n$  such that  $n^{\epsilon/2} d(\log n)(\log \log n) < n^\epsilon/2$  and  $\log n$  is greater than the  $O(1)$  constant (not depending on  $n$ ) to be introduced later in the proof. The programs (possibly infinite) of  $M$  can be encoded into an oracle  $A$ . The oracle  $A$ , when queried with  $(i, l)$ , returns the initial section of length  $l$  of the program for  $P(i)$ . Fix a binary string  $x$  of length  $n^{1+\epsilon}$  such that

$$C(x|A, T, P, n) \geq n^{1+\epsilon}.$$

Partition  $x$  equally into  $n$  parts  $x_1, x_2, \dots, x_n$ . Consider the computation of  $M$  on input  $(x_1, \dots, x_n)$ , with  $x_i$  initially in  $c(i)$  ( $1 \leq i \leq n$ ). We define inductively (with respect to this input) a processor to be *alive* at step  $t$  in this computation if

- it writes the output; or
- it succeeds in writing something at some step  $t' \geq t$  that is read at some step  $t'' \geq t'$  by a processor that is alive at step  $t''$ .

Informally, a processor is alive at step  $t$  if it has the potential of “influencing” the outcome of the computation. An input is *useful* if it is read at some step  $t$  by a processor alive at step  $t$ . By induction on the number of steps we find that for a  $T$  step computation, the number of useful inputs and the number of processors that are alive at any time are both at most  $2^T$ .

It is not difficult to see that given all the useful inputs and the set  $\text{ALIVE} = \{(P(i), t_i) : P(i) \text{ was alive until step } t_i > 0\}$ , using oracle  $A$  we can simulate  $M$  to uniquely reconstruct the output  $\sum_{i=1}^n x_i$ . Since  $T < (\epsilon/2) \log n$  we have  $2^T < n^{\epsilon/2} < n$ . Since there are  $n$  inputs  $x_i$ , there must be an input  $x_{i_0}$  that is not useful. We now compress the description of  $x$  to below  $C(x)$ .

We only need  $2^T(\log P)(\log T) \leq n^{\epsilon/2}d(\log n)(\log \log n) < n^\epsilon/2$  bits to represent  $\text{ALIVE}$ . Let  $x'$  consist of a self-delimiting description of index  $i_0$  concatenated with the literal ordered concatenation of all  $x_i$ 's with  $i \neq i_0$ . This  $x'$  requires at most  $n^{1+\epsilon} - n^\epsilon + 2 \log n$  bits. Therefore, the total number of bits involved in self-delimiting descriptions of  $\text{ALIVE}$  and the above literal string  $x'$  is not more than

$$J = n^{1+\epsilon} - n^\epsilon + n^\epsilon/2 + O(\log n) + O(1).$$

The following program reconstructs  $x$ . We can recover  $\sum_{i=1}^n x_i$  by simulating  $M$  using the oracle  $A$ , and then reconstruct  $x_{i_0}$  from  $\sum_{i=1}^n x_i$  and the set  $\{x_i : x_i \neq x_{i_0}\}$  recovered from  $x'$ . Subsequently, concatenate the  $x_i$ 's in appropriate order according to  $x'$  to obtain  $x$ . But then

$$C(x|A, T, P, n) \leq J + O(1) < n^{1+\epsilon},$$

which contradicts the incompressibility of  $x$ .  $\square$

## Exercises

---

**6.11.1.** [22] Consider the following “proof” for Theorem 6.11.1 without using Kolmogorov complexity: assume that a PRAM  $M$  adds  $n$  numbers in  $o(\log n)$  time. Take any input  $x_1, \dots, x_n$ . Then there is an input  $x_k$  that is not useful (as defined in the proof of Theorem 6.11.1). Then if we change  $x_k$  to  $x_k + 1$ , the output should still be the same since  $x_k$  is not useful, a contradiction. What is wrong with this proof?

**6.11.2.** [26] A function  $f(x_1, \dots, x_n)$  is called *invertible* if for each  $i$ , argument  $x_i$  can be computed from  $\{x_1, \dots, x_n\} - \{x_i\}$  and  $f(x_1, \dots, x_n)$ . Use the PRAM model with  $q$  processors defined in this section. Show that it requires  $\Omega(\min\{\log(b(n)/\log q), \log n\})$  time to compute any invertible function  $f(x_1, \dots, x_n)$ , where  $l(x_i) \leq b(n)$ , for all  $i$ , and  $\log n = o(b(n))$ .

*Comments.* Source: M. Li and Y. Yesha, *J. ACM*, 36:3(1989), 671–680.

**6.11.3.** [36] [Computing minimum index] Modify the PRAM model as follows. We now have  $n$  processors  $P(1), \dots, P(n)$  and only one shared memory cell,  $c(1)$ . Each processor knows one input bit. If several processors attempt to write into  $c(1)$  at the same time, then they must all write the same data, otherwise each write fails. This PRAM version requires  $\Omega(\log n)$  time to find the smallest index  $i$  such that  $P(i)$  has input bit 1. Can you give two proofs, one using incompressibility arguments and the other not?

*Comments.* The original proof without using Kolmogorov complexity is due to F. Fich, P. Ragde, and A. Wigderson [*SIAM J. Comput.*, 17:3(1988), 606–627].

**6.11.4.** [30] Many VLSI lower bound proofs are based on crossing sequence type arguments similar to that of Section 6.1.1 on page 380. The VLSI model has the following components.

- A  $n$ -input, 1-output Boolean function  $f(x_1, x_2, \dots, x_n)$  that is to be computed;
- A synchronous circuit  $C$ , which computes  $f$  that contains AND, OR, NOT gates. There are also  $n$  fixed input gates (which are called *where-oblivious*) that are possibly not on the boundary of the layout of  $C$ . The time an input arrives may depend on the data value.
- A VLSI (for convenience: rectangle) layout  $V$  that realizes  $C$ , where wires are of unit width and gates (nodes) occupy unit squares.

Suppose we want to find  $C$  that computes a given  $f$  in time  $T$  and a VLSI layout of  $C$  with area  $A$ , minimizing the  $AT^2$  cost measure. The method used to prove  $AT^2 = \Omega(n^2)$  lower bounds is roughly as

follows: Draw a line to divide the layout into two parts, with about half of the inputs on each part. If the line cuts through at least  $\omega$  wires and nodes, then  $A > \Omega(\omega^2)$ . Further, since for each time unit only one bit of information can flow through a wire,  $T > I/\omega$ , where  $I$  is the *amount* of information that has to be passed between the two parts. Then for each specific problem one only needs to show that  $I = \Omega(n)$  for any division.

Use incompressibility and crossing sequence arguments to prove the  $AT^2 = \Omega(n^2)$  lower bound for the following predicates:

- (a) (Pattern-matching) Given an input containing a binary text string,  $t$ , of  $(1 - \alpha)n$  bits and a pattern,  $p$ , of  $\alpha n$  bits, with  $\alpha < \frac{1}{2}$ , determine whether the pattern occurs in the text.
- (b) (Selection/Equality testing) Given  $2n$  input bits, divide them into halves of  $n$  bits each. The first half is used for *selection*: it has  $n/2$  zeros and  $n/2$  ones. The value of the predicate is 1 if the  $n/2$  bit number obtained by selecting those bits in the second half of the input at positions corresponding to the zero bits in the selection mask is equal to the  $n/2$  bit number obtained in the same way from the one bit positions in the selection mask.
- (c) (Deterministic context-free language) The value of the predicate is 1 if the inputs belong to the DCFL  $xcx^R$ , with  $x$  a word from  $\{0, 1, *\}^*$ , subsequent to deletion of any  $*$ 's (don't cares) in the input.
- (d) (Factor verification 1) Given binary numbers  $x, y$ , and  $z$  of  $\alpha n$ ,  $\beta n$ , and  $(1 - \alpha - \beta)n$  bits each ( $0 < \alpha, \beta < 1$ ), respectively, determine whether or not  $xy = z$ .
- (e) (Factor verification 2) Given binary numbers  $x$  and  $y$  of  $\alpha n$  and  $(1 - \alpha)n$  bits, respectively for  $0 < \alpha < 1$ , determine whether or not  $x = y^2$ .

*Comments.* The VLSI model and the results in this exercise are from [R.J. Lipton and R. Sedgewick, *Proc. 13th ACM Symp. Theory Comput.*, 1981, 300–307]. Their proofs do not use incompressibility. The original VLSI model is from [C.D. Thompson, *Proc. 11th ACM Symp. Theory Comput.*, 1979, pp. 81–88].

## 6.12 Switching Lemma

---

A key lemma in the study of circuit complexity is the so-called Håstad's Switching Lemma. It is used to separate depth- $k$  and depth- $(k + 1)$  circuit classes, and to construct oracles relative to which the polynomial hierarchy is infinite and properly contained in PSPACE. The traditional proof of this lemma uses sophisticated probabilistic arguments. We give a simple elementary proof using the incompressibility method.

According to Definition 5.4.3 on page 344, a  $k$ -DNF formula is a disjunction of conjunctions with each conjunct (or term) containing at most  $k$  literals. A  $k$ -CNF is a conjunction of disjunctions with each disjunct (or clause) containing at most  $k$  literals.

**Definition 6.12.1** A *restriction*  $\rho$  is a function from a set of variables to  $\{0, 1, \star\}$ . Given a Boolean function  $f$ ,  $f|_{\rho}$  is the restriction of  $f$  in the natural way:  $x_i$  is free if  $\rho(x_i) = \star$  and  $x_i$  takes on the value  $\rho(x_i)$  otherwise. The *domain* of a restriction  $\rho$ ,  $\text{dom}(\rho)$ , is the set of variables mapped to 0 or 1 by  $\rho$ .

We can also naturally view a restriction  $\rho$  as a *term* of  $f$  if  $f|_{\rho} = 1$ . A *minterm* is a restriction such that no proper subset of the variables set by the restriction forms a term. Let  $R_l$  be the set of restrictions on  $n$  variables that leave  $l$  variables free. Obviously,  $d(R_l) = \binom{n}{l} 2^{n-l}$ .

**Lemma 6.12.1 (Switching Lemma)** *Let  $f$  be a  $t$ -CNF on  $n$  variables,  $\rho$  a random restriction  $\rho \in R_l$  and  $\alpha = 12tl/n \leq 1$ . Then the probability that  $f|_{\rho}$  is an  $s$ -DNF is at least  $1 - \alpha^s$ .*

**Proof.** Fix a  $t$ -CNF  $f$  on  $n$  variables, and integers  $s$  and  $l < n$ . Note,  $f|_{\rho}$  is  $l$ -DNF; we can assume  $s \leq l$ . In this proof, we will use conditional complexity  $C(\cdot|\mathbf{x})$ , where  $\mathbf{x}$  denotes the list of fixed values of  $f, t, l, n, s$  and several (fixed) programs needed later.

**Claim 6.12.1** For any  $\rho \in R_l$  such that  $f|_{\rho}$  is not  $s$ -DNF,  $\rho$  can be effectively described by some  $\rho' \in R_{l-s}$ , a string  $\sigma \in \{0, 1, \star\}^{st}$  such that  $\sigma$  has  $s$  non- $\star$  positions, and  $\mathbf{x}$ . That is,  $C(\rho|\rho', \sigma, \mathbf{x}) = O(1)$ .

Before proving Claim 6.12.1, we show that it implies the switching lemma. Fix a random restriction  $\rho \in R_l$  with

$$C(\rho|\mathbf{x}) \geq \log(d(R_l)\alpha^s), \quad (6.32)$$

where  $\alpha = 12tl/n \leq 1$ . If we show that  $f|_{\rho}$  is an  $s$ -DNF, then since there are at least  $d(R_l)(1 - \alpha^s)$   $\rho$ 's in  $R_l$  satisfying Equation 6.32 by Theorem 2.2.1, this will imply the lemma.

Assume that  $f|_{\rho}$  is not an  $s$ -DNF and  $\rho' \in R_{l-s}$  as in Claim 6.12.1. Obviously  $C(\rho'|\mathbf{x}) \leq \log d(R_{l-s})$ . Since  $l(\sigma) = st$  and  $\sigma$  has  $s$  non- $\star$

positions, we have

$$C(\sigma|\mathbf{x}) \leq \log \binom{st}{s} + s \leq s \log et + s = s \log 2et,$$

by standard estimation (Stirling's approximation), where  $e = 2.718\dots$ . By Claim 6.12.1, we have

$$C(\rho|\mathbf{x}) \leq C(\rho'|\mathbf{x}) + C(\sigma|\mathbf{x}) \leq \log d(R_{l-s}) + s \log 2et. \quad (6.33)$$

By Equations 6.32 and 6.33, we have

$$d(R_l)\alpha^s \leq d(R_{l-s})2^{s \log 2et}.$$

Substituting  $\binom{n}{l}2^{n-l}$  for  $d(R_l)$  and  $\binom{n}{l-s}2^{n-l+s}$  for  $d(R_{l-s})$ , and using the fact  $\binom{n}{l} / \binom{n}{l-s} \geq ((n-l+s)/l)^s$ , we obtain

$$\frac{12tl}{n} \leq \frac{4etl}{n-l+s}.$$

But the above formula cannot hold simultaneously with  $12tl/n \leq 1$ , a contradiction. This proves the switching lemma.

**Proof.** (For Claim 6.12.1) Given a  $t$ -CNF  $f$  on  $n$  variables and a restriction  $\rho \in R_l$  such that  $f|_\rho$  is not  $s$ -DNF. Let

$$f = \bigwedge_{j=1}^k D_j, \quad (6.34)$$

where each  $D_j$  is a disjunct of size at most  $t$ . We also can write  $f|_\rho$  as a DNF:  $f|_\rho = \vee_j C_j$ , where each  $C_j$  (the corresponding restriction) is a minterm of  $f|_\rho$ . Since  $f|_\rho$  is not an  $s$ -DNF, there must be a minterm  $\pi$  that contains at least  $s+1$  variables. We will extend  $\rho$  to  $\rho'$  using  $s$  of the variables of  $\pi$ .

First we split  $\pi$  into "subrestrictions." Assume that  $\pi_1, \dots, \pi_{i-1}$  have already been defined and  $\text{dom}(\pi) \setminus \text{dom}(\pi_1 \dots \pi_{i-1}) \neq \emptyset$ . Choose the first disjunct  $D_j$  in Equation 6.34 that is not already 1 under restriction  $\rho\pi_1 \dots \pi_{i-1}$ . Let  $S$  be the set of variables that appear both in  $D_j$  and in  $\text{dom}(\pi) \setminus \text{dom}(\pi_1 \dots \pi_{i-1})$ . Define  $\pi_i$  as

$$\pi_i(x) = \begin{cases} \pi(x) & \text{if } x \in S, \\ * & \text{otherwise.} \end{cases}$$

Because  $\pi$  is a minterm, it must force each disjunct to 1, and no subrestriction of  $\pi$  (namely  $\pi_1 \dots \pi_{i-1}$ ) will. Thus the above process is always possible. Let  $k$  be the least integer such that  $\pi_1 \dots \pi_k$  sets at least  $s$  variables. "Trim"  $\pi_k$  so that  $\pi_1 \dots \pi_k$  sets exactly  $s$  variables.

Change  $\pi_i$  to  $\tilde{\pi}_i$ : for each variable  $x \in \text{dom}(\pi_i)$ , if it appears in the corresponding  $D_j$  as  $x$  then  $\tilde{\pi}_i(x) = 0$ ; if it appears in  $D_j$  as  $\bar{x}$  then  $\tilde{\pi}_i(x) = 1$ . Thus  $\pi_i \neq \tilde{\pi}_i$  since  $\rho\pi_1 \dots \pi_{i-1}\pi_i$  forces  $D_j$  to be 1 but  $\rho\pi_1 \dots \pi_{i-1}\tilde{\pi}_i$  does not. If  $x$  is the  $m$ th variable in  $D_j$ , the  $m$ th digit of  $\sigma^{(i)}$  is

$$\sigma_m^{(i)} = \begin{cases} \pi_i(x) & \text{if } x \in \text{dom}(\pi_i) (= \text{dom}(\tilde{\pi}_i)), \\ * & \text{otherwise.} \end{cases}$$

Since  $D_j$  is of size at most  $t$ ,  $l(\sigma^{(i)}) = l(D_j) \leq t$ . Let

$$\rho' = \rho\tilde{\pi}_1 \dots \tilde{\pi}_k, \quad \sigma = \sigma^{(1)} \dots \sigma^{(k)} \star^{st-k}.$$

Pad  $\sigma$  with  $*$ 's so that  $l(\sigma) = st$ . Since  $\pi_1 \dots \pi_k$  sets exactly  $s$  variables,  $\sigma$  has  $s$  non- $*$  positions.

Now we show how to recover  $\pi_1, \dots, \pi_k$ , hence  $\rho$ , from  $\sigma$  and  $\rho'$  (given  $\mathbf{x}$ ). Assume we have already recovered  $\pi_1, \dots, \pi_{i-1}$ , from which we can infer  $\rho\pi_1 \dots \pi_{i-1}\tilde{\pi}_i \dots \tilde{\pi}_k$ , using  $\rho'$ . Recall that  $\pi_i$  was defined by choosing the first clause  $D_j$  not already forced to 1 by  $\rho\pi_1 \dots \pi_{i-1}$ . Since  $\tilde{\pi}_i$  does not force  $D_j$  to be 1 and  $\tilde{\pi}_{i+1} \dots \tilde{\pi}_k$  are defined on variables not contained in  $D_j$ , we simply identify  $D_j$  from  $f|_{\rho\pi_1 \dots \pi_{i-1}\tilde{\pi}_i \dots \tilde{\pi}_k}$ , as the first non-1 clause. Given  $D_j$ , recover  $\pi_i$  using  $\sigma^{(i)}$ . With  $\rho\pi_1 \dots \pi_k$  and  $\pi_1 \dots \pi_k$ , we can recover  $\rho$ . This proves Claim 6.12.1, and the theorem.  $\square \quad \square$

Let us summarize the central ideas in the above proof. When we choose a random  $\rho$ , the number of bits needed to specify each extra variable is roughly  $O(\log n)$ . However, the fact that  $f|_\rho$  is not an  $s$ -DNF implies that it has a long minterm, and this allows us to construct a  $\sigma$ , together with  $\rho'$ , specifying  $s$  extra variables at the expense of roughly  $\log 2et$  bits per variable. So a large term is a kind of regularity a random restriction  $\rho$  does not produce.

**Example 6.12.1** Lemma 6.12.1 is a powerful lemma in circuit complexity. Let's define a depth- $k$  (unbounded fan-in Boolean) circuit as follows: The input to the circuit is  $I = \{x_1, \dots, x_n, \bar{x}_1, \dots, \bar{x}_n\}$ . The circuit has  $k$  alternating levels of AND and OR gates, each with unbounded fan-in. The  $k$ th (top) level contains just one gate, which gives the output of the circuit. Each gate in the  $i$ th level gets an arbitrary number of inputs from the outputs of the  $i - 1$ st level, assuming  $I$  is at the zeroth level. A *parity* function  $f(x_1 \dots x_n) = 1$  if and only if an odd number of  $x_i$ 's are 1's. It is easy to show that a polynomial-size depth-2 circuit cannot compute parity. Assume that this is the case for  $k - 1$ . For a depth- $k$  circuit, if we apply a random restriction to it, then by Lemma 6.12.1, with high probability, we can switch the bottom two levels, say, from AND-OR to OR-AND. Then the 2nd level OR can merge with the 3rd level OR, hence reducing

the circuit depth to  $k - 1$ . Note that a restriction of a parity function remains a parity function. Making this kind of induction precise, one can prove the following: there is a constant  $c > 0$  such that a depth- $k$  circuit with at most  $2^{c^{k/k-1}n^{1/k-1}}$  gates cannot compute parity.  $\diamond$

## 6.13 History and References

---

Apparently, W.J. Paul [*Proc. Int. Conf. Fund. Comput. Theory*, L. Budach (Ed.), 1979, pp. 325–334] is the pioneer of using the incompressibility method, and he proved several lower bounds with it. R.V. Freivalds ['On the running time of deterministic and nondeterministic Turing machines,' *Latv. Mat. Ezhegodnik*, 23(1979) 158–165 (in Russian)] proved a lower bound on the time of Turing machine computations for a certain problem, implicitly using a Kolmogorov complexity argument in a veiled form of "optimal enumerations," justified by the Invariance Theorem 2.1.1. He did not use the incompressibility method.

The most influential paper is probably the paper by W.J. Paul, J. Seiferas, and J. Simon, [*J. Comput. System Sci.*, 23:2(1981), 108–126]. This was partly because the paper by W.J. Paul, which contains Section 6.1.1, page 380, was not widely circulated.

The aim of the paper by Paul, Seiferas, and Simon was "to promote the approach" of applying Kolmogorov complexity to obtain lower bounds. In that paper, using incompressibility arguments, the authors greatly simplified the proof of a difficult theorem proved by S.O. Aanderaa [pp. 75–96 in: *Complexity of Computation*, R. Karp (Ed.), Amer. Math. Soc., 1974], which states: *real time* simulation of  $k$  tapes by  $k - 1$  tapes is impossible for deterministic Turing machines. Earlier, M.O. Rabin [*Israel J. Math.*, 1(1963), 203–211] proved the particular case  $k = 2$  of this result. In 1982 W.J. Paul [*Inform. Contr.*, 53(1982), 1–8] further improved this to a nonlinear lower bound by incompressibility arguments. In the same year, S. Reisch and G. Schnitger [*Proc. 23rd IEEE Found. Comput. Sci.*, 1982, pp. 45–52] published a paper giving three applications of incompressibility in areas other than Turing machine computational complexity. (The authors later lost contact with each other and they have never written up a journal version of this paper.) Subsequently, incompressibility arguments started to be applied to an ever-increasing variety of problems.

Lemma 6.1.1 in Section 6.1.1 was first proved by F.C. Hennie using a counting argument in [*Inform. Contr.*, 8:6(1965), 553–578]. The proof we give here is due to W.J. Paul. Section 6.1.2 is based on [R. Beigel, W. Gasarch, M. Li, and L. Zhang, 'Addition in  $\log_2 n + O(1)$  steps on average: a simple analysis,' *Manuscript*, 1996]. The original probabilistic analysis is in [A.W. Burks, H.H. Goldstine, and J. von Neumann,

‘Preliminary discussion of the logical design of an electronic computing instrument. Institute for Advanced Studies,’ Report (1946). Reprinted in *John von Neumann Collected Works*, vol 5 (1961)]. Improved probabilistic analysis can be found in [B.E. Briley, *IEEE Trans. Computers*, C-22:5(1973)] and [G. Schay, *Amer. Math. Monthly*, 102:8(1995), 725–730]. Background material on adder design can be found in [K. Hwang, *Computer arithmetic: principles, architecture, and design*, Wiley, New York, 1979]. Lemma 6.1.3 in Section 6.1.3 is due to J. Seiferas and Y. Yesha [Personal communication, 1986]. The idea of proving a lower time bound for palindrome recognition by a probabilistic Turing machine, as mentioned in the comment at the end of Section 6.1 and in Exercise 6.10.13 Item (c), is due to R. Paturi, J. Simon, R. Newman-Wolfe, and J. Seiferas, [*Inform. Comput.*, 88(1990), 88–104].

The discussion in Section 6.2 on the quantitative relation between high-probability properties of finite objects and individual randomness of finite objects is taken from [H. Buhrman, M. Li, and P.M.B. Vitányi, ‘Kolmogorov random graphs,’ *Manuscript*, CWI, Amsterdam, November 1995]. With respect to infinite binary sequences the distinction between laws of probability (that hold with probability one) and individual random sequences is discussed in [M. van Lambalgen, *Random Sequences*, PhD Thesis, Universiteit van Amsterdam, Amsterdam, 1987] and [V.V. V'yugin, ‘Ergodic theorem for individual random sequences,’ *Manuscript*, Inst. Inform. Transm. Probl., Russ. Acad. Sci., Moscow, Russia, 1996]. Section 6.3 on combinatorics follows [M. Li and P.M.B. Vitányi, *J. Comb. Theory, Ser. A*, 66:2(1994), 226–236]. The problems on tournaments in Section 6.3 are from [P. Erdős and J. Spencer, *Probabilistic methods in combinatorics*, Academic Press, 1974]. See also [N. Alon, J. Spencer, and P. Erdős, *The probabilistic method*, Wiley, 1992] for the probabilistic method. Section 6.3.3 was suggested by W. Gasarch; the lower bound on the Ramsey numbers was originally proved in [P. Erdős, *Bull. Amer. Math. Soc.*, 53(1947), 292–294]; see [P. Erdős and J. Spencer, *Ibid.*]. The lower bound for the coin-weighing problem in Theorem 6.3.4 has been established, using probabilistic or information theoretic methods, by P. Erdős and A. Rényi [*Publ. Hungar. Acad. Sci.*, 8(1963), 241–254], L. Moser [*Combinatorial structures and their applications*, Gordon and Breach, 1970, pp. 283–384], and N. Pippenger [*J. Combinat. Theory, Ser. A*, 23(1977), 105–115]. The latter paper contains proofs, by entropy methods, of Theorem 6.3.4 on page 393 and Exercise 6.3.3 on page 395. Recently, entropy methods have also been used quite successfully in proving lower bounds on parallel sorting [J. Kahn and J. Kim, *Proc. 24th ACM Symp. Theory Comput.*, 1992, pp. 178–187], perfect hashing [I. Newman, P. Ragde, and A. Wigderson, *5th IEEE Conf. Structure in Complexity Theory*, 1990, pp. 78–87], lower

bounds on parallel computation [R. Boppana, *Proc. 21st ACM Symp. Theory Comput.*, 1989, pp. 320–326].

Section 6.4 is based on [H. Buhrman, M. Li, P. Vitányi, ‘Kolmogorov random graphs,’ Manuscript, CWI, Amsterdam, November 1995]. For random graphs in a probabilistic sense see for example [B. Bollobás, *Random graphs*, Academic Press, 1985]. The statistics of subgraphs of high complexity graphs, Theorem 6.4.1, has a corresponding counterpart in quasi-random graphs, and a similar expression is satisfied almost surely by random graphs [N. Alon, J. Spencer, and P. Erdős, *The probabilistic method*, Wiley, 1992] pp. 125–140, see especially Property  $P_1(s)$  on page 126. The latter property may be weaker in terms of quantification of “almost surely” and the  $o(\cdot)$  and  $O(\cdot)$  estimates involved than the result we present here.

The results in Section 6.5 are from [H. Buhrman, J.H. Hoepman, and P.M.B. Vitányi, *Proc. 15th ACM Symp. Principles Distribute Comput.*, 1996, pp. 134–142]. Related research (see exercises) appears in [E. Kranakis and D. Krizanc, ‘Boolean Routing on Cayley Networks,’ *Proc. 3rd Int'l Colloq. Structure Inform. Communication Complexity*, Siena, Italy, 1996, to appear; E. Kranakis and D. Krizanc, *Proc. 13th Symp. Theoret. Aspects Comput. Sci.*, 1996, pp. 529–540; E. Kranakis, D. Krizanc, and F. Luccio, *Proc. 13th Symp. Math. Found. Comput. Sci.*, 1995, pp. 392–401].

Heapsort was originally discovered by J.W.J. Williams [*Comm. ACM*, 7(1964), 347–348]. R.W. Floyd [*Comm. ACM*, 7(1964), 701] subsequently improved the algorithm. Researchers had previously tried to analyze the precise average-case complexity of Heapsort with no success. For example, the analysis typically works only for the first step; after one step the heap changes and certain properties such as that all heaps are equally likely no longer hold. Section 6.6 is based on an explanation by I. Munro on a summer evening in 1992. The solution to the average-case complexity of Heapsort was first obtained by R. Schaffer and R. Sedgewick [*J. Algorithms*, 15(1993), 76–100]. The proof in the form given in Section 6.6 is due to I. Munro. Claim 6.6.1, that the “Heapify” procedure produces a random heap from a random input, was observed by T. Jiang, at a 1993 Dagstuhl seminar, and I. Munro.

Section 6.7 is from [T. Jiang and M. Li, *SIAM J. Comput.*, 24:5(1995), 1122–1139]. Section 6.8 follows [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 24:2(1995), 398–410]. For the history and an introduction of formal language theory, see [M.A. Harrison, *Introduction to formal language theory*, Addison-Wesley, 1978; J.E. Hopcroft and J.D. Ullman, *Introduction to automata theory, languages, and computation*, Addison-Wesley, 1979].

The proof in Section 6.9 of the lower bound on the time required for linear context-free language recognition is due to J. Seiferas [*Inform. Contr.*, 69(1986), 255–260] who simplified the original proof of H. Gallaire [*Inform. Contr.*, 15(1969), 288–295]. The latter paper improved a  $\log n$  multiplicative factor weaker result by F.C. Hennie. Gallaire’s proof uses a complicated counting argument and de Bruijn sequences. T. Kasami [*Inform. Contr.*, 10(1967), 209–214] proved that linear context-free languages can be online recognized in  $O(n^2)$  by a one work tape Turing machine.

It is fair to say that the solutions to the conjectures on  $k$ -PDA in Exercise 6.9.3, string-matching in Exercise 6.9.2, Theorem 6.10.1, and many exercises in Section 6.10 would not have been possible, at least not proved in such a short period of time, without the use of incompressibility arguments. This seems all the more true about the results in Section 6.10 concerning whether or not an extra tape adds computational power in various Turing machine models and especially for the “two heads are better than two tapes” result in Exercise 6.10.15. These results were open for decades before they were solved using the incompressibility method and some other Kolmogorov complexity-related techniques, and no other proofs are known.

It is well known that if a  $k$ -tape Turing machine runs in  $O(T(n))$  time, then it can be simulated by a 1-tape Turing machine in  $O(T^2(n))$  time [J. Hartmanis and R. Stearns, *Trans. Amer. Math. Soc.*, 117(1969), 285–306] and by a 2-tape Turing machine in  $O(T(n) \log T(n))$  time [F.C. Hennie and R. Stearns, *J. ACM*, 4(1966), 533–546]. For years, only several weak lower bounds were known with complicated proofs, such as M.O. Rabin’s paper from 1963 and S.O. Aanderaa’s paper of 1974 above. These papers consider the restricted online model with an extra output tape. For the more general model used in Theorem 6.10.1, P. Dúriš, Z. Galil, W.J. Paul, and R. Reischuk [*Inform. Contr.*, 60(1984), 1–11] proved that it requires  $\Omega(n \log n)$  time to simulate two tapes by one.

Research advanced quickly only after the incompressibility argument was invented. W.J. Paul [*Inform. Contr.*, 53(1982), 1–8] proved Exercise 6.10.13 Item (a) on page 442, improving Aanderaa’s result. Around 1983/1984, independently and in chronological order, Wolfgang Maass at UC Berkeley, one of us [ML] at Cornell, and the other one [PV] at CWI Amsterdam, obtained an  $\Omega(n^2)$  lower bound on the time to simulate two tapes by one tape (deterministically), and thereby closed the gap between 1 tape versus  $k \geq 2$  tapes (Exercise 6.10.2 on page 438). All three relied on Kolmogorov complexity, and actually proved more in various ways. (One of us [PV], at first not realizing how to use incompressibility, reported in [P.M.B. Vitányi, *Theoret. Comput. Sci.*, 34(1984), 157–168] an  $\Omega(n^{3/2})$  lower bound on the time to simulate a single push-down store online by one *oblivious* tape unit. However, after being en-

lightened by J. Seiferas about how to use incompressibility with respect to another result, he realized how to apply it to the 1-tape versus 2-tape problem without the oblivious restriction [P.M.B. Vitányi, *Inform. Process. Lett.*, 21(1985), 87–91 and 147–152], and the optimal results cited below.) W. Maass also obtained a nearly optimal (almost square) lower bound for nondeterministic simulation (Exercise 6.10.4, page 439) Maass's lower bound on nondeterministic simulation was later improved by Z. Galil, R. Kannan, and E. Szemerédi [*Proc. 18th ACM Symp. Theory Comput.*, 1986, pp. 39–49] to  $\Omega(n^2 / \log^{(k)} n)$  by constructing a language whose computation graph does not have small separators (Exercise 6.10.5, page 439). See the exercises for many more lower bounds proved in this direction. This section is based on [W. Maass, *Trans. Amer. Math. Soc.*, 292(1985), 675–693; M. Li and P.M.B. Vitányi, *Inform. Comput.*, 78(1988), 56–85]. The latter also contains results on tapes versus stacks and queues. Many lower bounds using Kolmogorov complexity for various models of computation, such as machines with extra two-way input tapes, machines with queues, random access machines, machines with many heads on a tape, machines with tree tapes, machines with  $k$ -dimensional tapes, and probabilistic machines, have since been proved. We have tried to cover these results in the exercises, where also the references are given.

Theorem 6.11.1 in Section 6.11 was independently proved by P. Beame [*Inform. Comput.*, 76(1988), 13–28] without using Kolmogorov complexity and by M. Li and Y. Yesha [*J. ACM*, 36:3(1989), 671–680], which contains the current proof. Slightly weaker versions of Theorem 6.11.1 were proved by F. Meyer auf der Heide and A. Wigderson [*SIAM J. Comput.*, 16(1987), 100–107] using a Ramsey theorem, by A. Israeli and S. Moran [private communication, 1985] and by I. Parberry [Ph.D. Thesis, 1984, Warwick University]. In the latter three proofs, one needs to assume that the integers have arbitrarily (or exponentially) many bits.

The proof of Lemma 6.12.1 in Section 6.12 is based on a paper by L. Fortnow and S. Laplante [*Inform. Comput.*, 123(1995), 121–126], which in turn was based on a proof by A. Razborov [pp. 344–386 in *Feasible Mathematics II*, P. Clote, J. Remmel (Eds), 1995]. This lemma was originally proved by J. Håstad [pp. 143–170 in *Randomness and Computation*, S. Micali (Ed), JAI Press, 1989] for the purpose of simplifying and improving Yao's lower bound on unbounded circuits [A. Yao, *Proc. 26th IEEE Symp. Found. Comput. Sci.*, 1985, pp. 1–10]. Note that in Lemma 6.12.1, in order to simplify the proof, we have  $\alpha = 12tl/n$  instead of Håstad's  $\alpha = 5tl/n$  or Fortnow and Laplante's  $\alpha = 5.44tl/n$ .

There are applications of Kolmogorov complexity we do not cover. This is because some applications require lengthy discussions of computational models and preliminary facts; and some others are indirect applications.

A.M. Ben-Amram and Z. Galil [*J. ACM*, 39:3(1992)] use Kolmogorov complexity to formalize the concept of incompressibility for general data types and prove a general lower bound for incompressible data types. J. Shallit and his coauthors in a series of papers study a variation of descriptive complexity, “automaticity,” where the description device is restricted to finite automata, see [J. Shallit and Y. Breitbart, *Proc. 11th Symp. Theoret. Aspects Comput. Sci.*, 1994, pp. 619–630, part of this paper to appear in *J. Comput. System Sci.*]. U. Vazirani and V. Vazirani [*Theoret. Comput. Sci.*, 24(1983), 291–300] studied probabilistic polynomial time reductions. It is possible to do their reduction by Kolmogorov complexity. Kolmogorov complexity has also been studied in relation to the tradeoff of table size and number of probes in hashing by H.G. Mairson [*Proc. 24th IEEE Found. Comput. Sci.*, 1983, pp. 40–47]. See also [K. Mehlhorn, *Proc. 23rd IEEE Found. Comput. Sci.*, 1982, 170–175]. D. Hammer and A.Kh. Shen’ [A strange application of Kolmogorov complexity, *Math. Systems Theory*, to appear] use complexity to derive a geometric relation, and a geometric relation to derive a property of complexity. Namely, from  $2C(a, b, c) \leq C(a, b) + C(b, c) + C(c, a) + O(\log n)$  one can derive  $\|V\|^2 \leq \|S_{xy}\| \cdot \|S_{yz}\| \cdot \|S_{zx}\|$ . Here  $V$  is a set in three-dimensional space,  $S_{xy}, S_{yz}, S_{zx}$  are its two-dimensional projections, and  $\|\cdot\|$  denotes volume. Moreover, from the well-known Cauchy-Schwartz Inequality one can derive  $2K(a, b, c) \leq K(a, b) + K(b, c) + K(c, a) + O(1)$ . The incompressibility method has been applied to logical definability by M. Zimand [*Inform. Process. Lett.*, 57(1996), 59–64] and to finite-model theory and database query languages by J. Tyszkiewicz in [‘The Kolmogorov expression complexity of logics,’ in *Proc. Comput. Sci. Logic Conf.*, 1995] and [*Proc. 8th Intern. Conf. Database Theory, Lecture Notes in Computer Science*, Vol 893, G. Gottlob and M.Y. Vardi (Eds), Springer-Verlag, 1995, pp. 97–110]. M. Zimand [*Ibid.*] studies a ‘high-low Kolmogorov complexity law’ equivalent to a 0-1 law in logic. See also [R. Book, *SIAM J. Comput.*, 23(1994), 1275–1282]. K.W. Regan [*Proc. 10th IEEE Conf. Structure in Complexity Theory*, 1995, pp. 50–64] uses Kolmogorov complexity to prove superlinear lower bounds for some problems in a type of hierarchical memory model that charges higher cost for nonlocal communication.

# Resource-Bounded Complexity

Recursion theory has a resource-bounded version in computational complexity theory. Similarly, Kolmogorov complexity has resource-bounded Kolmogorov complexity (also known as generalized Kolmogorov complexity). Several authors suggested early on the possibility of restricting the power of the device used to (de)compress strings. Says Kolmogorov:

“The concept discussed . . . does not allow for the ‘difficulty’ of preparing a program  $p$  for passing from an object  $x$  to an object  $y$ . [. . . some] object permitting a very simple program, i.e., with very small complexity  $K(x)$  can be restored by short programs only as the result of computations of a thoroughly unreal nature. [. . . this concerns] the relationship between the necessary complexity of a program and its permissible difficulty  $t$ . The complexity  $K(x)$  that was obtained [before] is, in this case, the minimum of  $K^t(x)$  on the removal of the constraints on  $t$ .”

Additional restrictions on resource bounds yield a rich mathematical theory and abundant applications. Many statements that used to be trivial or easily provable become nontrivial or very difficult questions with resource bounds. Even the definition of Kolmogorov complexity itself becomes nonstandard.

Allowing unlimited computational resources such as time and space, it does not matter whether we define the complexity of  $x$  as the length of the shortest program that *prints*  $x$ , or that *accepts only*  $x$ ; both definitions turn out to be equivalent. However, it is not known whether the two definitions are still equivalent under polynomial-time restriction of the computational processes involved. We will specifically study two parameters of the “permissible difficulty” of Kolmogorov, the *time* and *space* complexity.

## 7.1 Mathematical Theory

---

Consider Turing machines with a separate read-only input tape, a fixed number of work tapes on which the computation takes place, and a write-only output tape. The input is a finite binary string. Initially, it is placed on the input tape, delimited by distinguished endmarkers. A machine with  $k$  work tapes is called a  *$k$ -tape Turing machine*. Machines with  $k \geq 0$  are called *multitape Turing machines*.

Let  $\phi_1, \phi_2, \dots$  be an effective enumeration of partial recursive functions. Let  $T_\phi$  be a multitape Turing machine that computes  $\phi$ . Let  $n = l(x)$ . If  $T_\phi(y) = x$  in  $t(n)$  steps (time) and  $s(n)$  tape cells (space), then we also write  $\phi^{t,s}(y) = x$ . We identify the natural numbers  $\mathcal{N}$  and the set of finite binary sequences  $\{0, 1\}^*$  as in Equation 1.1, page 12.

### Definition 7.1.1

Let  $x, y, p \in \mathcal{N}$ . Any recursive function  $\phi$  together with  $p, y$ , such that  $\phi(p, y) = x$ , is a *description* of  $x$ . The *resource-bounded Kolmogorov complexity*  $C_\phi^{t,s}$  of  $x$ , *conditional* to  $\phi$  and  $y$ , is defined by

$$C_\phi^{t,s}(x|y) = \min\{l(p) : \phi^{t,s}(p, y) = x\},$$

and  $C_\phi^{t,s}(x|y) = \infty$  if there is no such  $p$ . Define the unconditional resource-bounded Kolmogorov complexity of  $x$  as  $C_\phi^{t,s}(x) = C_\phi^{t,s}(x|\epsilon)$ .

Obviously, for total recursive  $t(n), s(n)$ , the function  $C_\phi^{t,s}$  is total recursive as well. We prove an *Invariance Theorem* for the resource-bounded Kolmogorov complexity. As usual, this invariance condition is the basis of the theories and applications to follow. The price we pay for adding resource bounds is that the invariance property becomes considerably weaker.

### Theorem 7.1.1

*There exists a universal partial recursive function  $\phi_0$ , such that for any other partial recursive function  $\phi$ , there is a constant  $c$  such that*

$$C_{\phi_0}^{ct \log t, cs}(x|y) \leq C_\phi^{t,s}(x|y) + c,$$

*for all  $x$  and  $y$ . The constant  $c$  depends on  $\phi$  but not on  $x$  and  $y$ .*

**Proof.** Consider a standard enumeration of multitape Turing machines  $T_1, T_2, \dots$ . Let  $\phi_i$  denote the partial recursive function computed by  $T_i$ . Let  $n = l(y)$ . Let  $\langle y, p \rangle = 1^{l(y)}0yp$  be the standard linear time and space invertible pairing bijection. Let  $\phi_0$  be the partial recursive function computed by a universal Turing machine  $U$  with two worktapes, such that  $U$  simulates all  $T_i$ 's according to the well-known simulation of F.C. Hennie and R. Stearns [J. ACM, 13(1966), 533–546]. It follows from that simulation that for each  $i$  there is a constant  $c$  such that  $U(\langle y, \langle i, p \rangle \rangle) = T_i(y, p)$ . Suppose that  $T_i$ , started on input  $\langle y, p \rangle$ , computes  $x$  in  $t(n)$  time

and  $s(n)$  space. Then the computation of  $U$ , started on input  $\langle y, \langle i, p \rangle \rangle$ , takes at most time  $ct(n) \log t(n)$  and space  $cs(n)$ . That is,

$$\phi_0^{ct \log t, cs}(\langle y, \langle i, p \rangle \rangle) = \phi_i^{t, s}(y, p).$$

Choosing  $c$  large enough to also exceed the length of the encoding of  $T_i$  in the input for  $U$ , say  $c \geq 2i + 1$ , proves the theorem.  $\square$

The Invariance Theorem allows us to drop the subscript  $\phi$  in “ $C_\phi^{t,s}$ ” and write “ $C^{t,s}$ ,” provided the statement we make is not sensitive to the additive constant term in the complexity of each string, the multiplicative logarithmic factor in the time complexity, and the multiplicative constant factor in the space complexity.

**Definition 7.1.2** Let  $U$  compute  $\phi_0$  of Theorem 7.1.1. The  $t$ -time bounded and  $s$ -space bounded version of  $C(x|y)$  is defined by

$$C^{t,s}(x|y) = \min\{l(p) : U(\langle y, p \rangle) = x \text{ in } t(n) \text{ steps}\}.$$

Resource-bounded Kolmogorov complexity behaves differently from the unbounded case. The complexity of string  $x$  can be defined as the length of either the shortest program that *generates*  $x$ , or the shortest program that *accepts only*  $x$ . These definitions are clearly equivalent when there are no resource bounds. It is unknown whether they are the same, for example, for the polynomial-time-bounded versions. There may be a short program that accepts precisely  $x$  in time  $t(n)$ , but one may have trouble finding an equally short program that prints  $x$  in time  $t(n)$ . In Definition 7.1.1 we have defined the resource-bounded complexity in terms of the shortest program that generates  $x$ . We also require the other type.

**Definition 7.1.3** A *predicate* is a 0–1 valued function.

Let  $\psi_1, \psi_2, \dots$  be an effective enumeration of partial recursive predicates. Let  $T_\psi$  be the multitape Turing machine which computes  $\psi$ . The machine  $T_\psi(x)$  outputs 0 or 1. If  $T_\psi$  uses at most  $t(n)$  time and  $s(n)$  space on any input of length  $n$ , then we write  $\psi^{t,s}(x)$ .

**Definition 7.1.4** Let  $x, y, p \in \mathcal{N}$ . Let  $\psi$  be a partial recursive predicate. The *CD* complexity of  $x$ , *conditional* to  $\psi$  and  $y$ , is defined as

$$CD_\psi^{t,s}(x|y) = \min\{l(p) : \forall v, \psi^{t,s}(v, p, y) = 1 \text{ iff } v = x\},$$

and  $CD_\psi^{t,s}(x|y) = \infty$  if there is no such  $p$ . When  $y = \epsilon$ ,  $CD^{t,s}(x) = CD^{t,s}(x|\epsilon)$  is the unconditional *CD*-complexity of  $x$ .

Statement and proof of the invariance theorem for  $CD$ -complexity are omitted. They are completely analogous to Theorem 7.1.1. In cases where this is justified by the Invariance Theorem, we drop the index  $\psi$  in  $CD_\psi^{t,s}$ .

Therefore, we distinguish between  $C^{t,s}(x)$  as the length of the shortest program generating  $x$  of length  $n$  in  $t(n)$  time and  $s(n)$  space, and  $CD^{t,s}(x)$  as the length of the shortest program accepting precisely  $x$  in  $t(n)$  time and  $s(n)$  space.

**Definition 7.1.5** Let  $S \subseteq \mathcal{N}$ . We define  $C^{t,s}(x|S)$  or  $CD^{t,s}(x|S)$  similar to the above, except that the underlying Turing machine is now equipped with a distinguished *oracle tape*. In the course of its computation the Turing machine can write questions of the form “is  $z \in S$ ?” on the oracle tape. After the question is written, which takes at least  $l(z)$  steps, the oracle gives the correct answer “yes” or “no” in one step.

In the notation  $C^{t,s}$  and  $CD^{t,s}$ , we always use the first parameter  $t$  for time and second parameter  $s$  for space. Considering only time complexity we write  $C^t$  and  $CD^t$ ; considering only space complexity, we write  $C^s$  and  $CD^s$ . When neither  $t$  nor  $s$  is required, that is, when  $t = \infty$  and  $s = \infty$ , then  $C^{t,s}$  and  $CD^{t,s}$  both converge to the original plain Kolmogorov complexity  $C$ .

**Definition 7.1.6** Assume the notation above. We define the following descriptional complexity classes.

$$\begin{aligned} C_\phi[f(n), t(n), s(n)|y] &\stackrel{\text{def}}{=} \{x : C_\phi^{t,s}(x|y) \leq f(n), n = l(x)\}, \\ CD_\psi[f(n), t(n), s(n)|y] &\stackrel{\text{def}}{=} \{x : CD_\psi^{t,s}(x|y) \leq f(n), n = l(x)\}, \end{aligned}$$

where the subscripts  $\phi$  and  $\psi$  are dropped in case they are superfluous. If  $y = \epsilon$ , then we also drop the conditional “ $|\epsilon$ .”

**Theorem 7.1.2** Let  $p, q$  be polynomial time bounds, and let  $A$  be an NP-complete set.

- (i) For all  $p$ , there exist  $q$  such that  $CD^q(x) \leq C^p(x) + O(1)$ .
- (ii) For all  $p$ , there exist  $q$  such that  $C^q(x|A) \leq CD^p(x) + O(1)$ .

**Proof.** Item (i) is immediate. To prove Item (ii), use the NP oracle repeatedly to determine the successive bits of  $x = x_1 \dots x_n$  as follows: Let  $T$  be a Turing machine that accepts only  $x$  and runs in time  $p(n)$ . Define a Turing machine  $T^A$  that generates  $x$  using an appropriate oracle  $A$ . Assume  $T^A$  has already determined  $x_1 \dots x_i$ . At the next step,  $T^A$  asks the oracle  $A$  whether  $T$  accepts a string with prefix  $x_1 \dots x_i 0$ . If the

answer is “yes,” then  $T^A$  sets  $x_{i+1} := 0$ , otherwise it sets  $x_{i+1} := 1$ . The oracle can answer these questions, since the set

$$\{\langle T, y, 1^t, 1^n \rangle : T \text{ accepts } yz \text{ in time } t \text{ for some } z \text{ with } l(yz) = n\}$$

is in NP. The oracle machine generates all of  $x$  this way, taking time at most polynomial in  $p(n)$ .  $\square$

The fact that it is not known how to improve Theorem 7.1.2 may suggest that  $C^{t,s}$  and  $CD^{t,s}$  are different measures, at least for polynomial time bounds. This difference disappears, however, when we have enough time available. Namely, in exponential time the machine can search through all strings of length  $n$ .

The resource-bounded prefix complexities  $K^{t,s}$  and  $KD^{t,s}$  are defined similarly to  $C^{t,s}$ , and  $CD^{t,s}$ . The definitions and the proofs of the similar invariance theorems are omitted.

### 7.1.1 Computable Majorants

Here, consider only time-bounded Kolmogorov complexity and write  $C^t$  or  $CD^t$ . Similar results to the ones obtained here can be derived for space-bounded complexities. Co-enumerable functions (below) are the functions that can be approximated from above by recursive functions, Section 4.1. We repeat the definition in Example 4.1.1.

**Definition 7.1.7** A *majorant* of Kolmogorov complexity is a co-enumerable function  $\phi$  such that  $C(x) < \phi(x) + c$  for all  $x$ , where  $c$  is a constant depending on  $\phi$  but not on  $x$ .

**Example 7.1.1** Let  $t$  be a total recursive function. Then the function  $C^t$  is a total recursive majorant of  $C$ . We first prove  $C^t$  is a majorant. For each  $x$ , we can run the universal machine on all programs of size at most  $l(x) + c$  for  $t(l(x))$  steps each. Then  $C^t(x)$  is the length of the shortest program that halts with output  $x$ . Clearly,  $C^t(x) \geq C(x) - O(1)$ . This also shows that  $C^t$  is total recursive.  $\diamond$

If  $C^t$  is a total recursive majorant, then  $C(x)$  and  $C^t(x)$  may be exponentially different. That is, if we impose any recursive time limit, the length of the shortest description of some strings can sharply increase. We prove this fact in terms of characteristic sequences. Let  $A$  be a set, and  $\chi = \chi_1\chi_2\dots$  be its characteristic sequence defined by  $\chi_i = 1$  if  $x_i \in A$  and  $\chi_i = 0$  otherwise. By Barzdins’s Lemma, Theorem 2.7.2 on page 171, if  $A$  is recursively enumerable, then  $C(\chi_{1:n}|n) \leq \log n + O(1)$ . The following *blow-up* theorem is one of the very first results in “time-limited” Kolmogorov complexity.

**Theorem 7.1.3** *There is a recursively enumerable set  $A$  with characteristic sequence  $\chi$  such that for all total recursive  $t$  and all  $n$  we have  $C^t(\chi_{1:n}|n) \geq c_t n$ , where  $0 < c_t < 1$  is a constant independent of  $n$  (but dependent on  $t$ ).*

**Proof.** Let  $\mathbf{T} = T_1, T_2, \dots$  be a standard enumeration of Turing machines. Let  $\mathbf{M} = M_1, M_2, \dots$  be another enumeration of Turing machines such that  $T_k = M_i$  with  $i = 2^{k-1} + j2^k$ , for all  $k \geq 1$  and  $j \geq 0$ . Stated differently,  $k$  is the largest integer such that  $2^{k-1}$  divides  $i$ . That is, in the following sequence, if  $k$  is the index in position  $i$ , then  $M_i$  is  $T_k$ :

$$1213121412131215121312141213121612\dots$$

Thus,  $T_1$  occurs every second machine in the list,  $T_2$  appears every four machines in the list, and  $T_k$  occurs every  $2^k$  machines in the list.

Given any total recursive time bound  $t$ , there is some  $T_k$  that computes the value  $t(n)$  from input  $n$ . We know that  $T_k$  occurs with intervals of  $2^k$  consecutive machines in the  $\mathbf{M}$  list, starting with the  $(3 \cdot 2^{k-1})$ th machine.

The following process enumerates an infinite sequence  $\chi = \chi_1\chi_2\dots$  by diagonalization. Later we show that there is a constant  $c_t > 0$  such that  $C^t(\chi_{1:n}|n) \geq c_t n$ .

**Step 1** Set  $i = 1$  and  $\chi_1 := 0$ .

**For** all  $i > 1$  dovetail the following computations: Set  $n := 2^{i-1}$ .

{This step will enumerate the segment  $\chi_{n+1:2n}$ .}

Set  $n' := 2^{2^k} n$ , where  $k$  is the index of  $M_i$  on the  $\mathbf{T}$  list. This means that  $M_i = T_k$  computes the partial recursive function, say,  $t$ . Simulate the computation of  $M_i(n')$ . If  $M_i(n')$  terminates, then its output is time bound  $t(n')$ .

Simulate all binary programs of size up to  $n - 1$  for  $t(n')$  steps on the reference universal Turing machine. Choose  $\chi_{n+1:2n}$  such that it is different from the segment from position  $n + 1$  to  $2n$  of the output of any of these simulations. Since there are only  $2^n - 1$  programs of size at most  $n - 1$ , and there are  $2^n$  different candidates for  $\chi_{n+1:2n}$ , this is always possible. If  $M_i(n')$  does not terminate, then define  $\chi_{n+1:2n} := 0^n$ .

The above procedure does not give an effective construction for  $\chi$ . If  $M_i(n')$  does not terminate, then we have no value of  $t(n')$  to use; hence, we cannot decide that  $\chi_{n+1:2n} = 0^n$  for  $n = 2^{i-1}$ .

**Claim 7.1.1** Let  $A \subseteq \mathcal{N}$  be defined by  $x \in A$  iff  $\chi_x = 1$ . Then  $A$  is recursively enumerable.

**Proof.** For all  $x$  simultaneously dovetail the following computation: First, determine  $i$  such that  $n + 1 < x \leq 2n$  for  $n = 2^{i-1}$ . Enumerate  $\chi_{n+1:2n}$  by simulating  $M_i$ . If  $M_i(n')$  never terminates, then  $\chi_{n+1:2n} = 0^n$ . That is, both  $x \notin A$ , and our simulation of  $M_i(n')$  doesn't halt. If  $M_i(n')$  terminates, then we effectively construct  $\chi_{n+1:2n}$ . We enumerate  $x$  as an element of  $A$  iff  $\chi_x = 1$ .  $\square$

**Claim 7.1.2** Let  $\chi$  be as above. It is incompressible as stated in the theorem, where without loss of generality we take  $t$  to be a monotonic increasing total recursive function.

**Proof.** Let  $t$  be any monotonic increasing total recursive function. Then  $t$  is computed by some machine  $T_k$ . By enumeration of  $\mathbf{M}$ , the machine  $T_k$  occurs for the first time as  $M_{k_0}$  with  $k_0 = 2^{k-1}$ . Subsequently,  $T_k$  returns in the list in fixed intervals of  $2^k$  machines. Let  $i$  run over the infinite sequence of values

$$k_0, k_0 + 2^k, k_0 + 2 \times 2^k, k_0 + 3 \times 2^k, \dots$$

Consider the special sequence of values of  $n$  defined by  $n = 2^{i-1}$ . Denote this sequence by  $S$ . Since  $t$  is total, the machine  $T_k$  halts for all inputs. Enumerating  $\chi$ , for each  $n \in S$  we used a copy of  $T_k$  for a terminating computation of segment  $\chi_{n+1:2n}$ . This segment was determined differently from all corresponding segments of outputs of programs of length at most  $n - 1$  and halting within time  $t(n')$  on the universal Turing machine.

For each large enough  $n \in S$ , there is an  $m \in S$  such that  $n/2^{2^k} \ll m \leq n/2$  and no program of length at most  $m - 1$  does output an initial segment  $\chi_{1:2m}$  in  $t(m') = t(2^{2^k}m) \gg t(n)$  steps. (This is why we need to use  $t(n')$  rather than  $t(n)$  in the definition of  $\chi$ .) Therefore,

$$C^{t(m')}(\chi_{1:2m}) \geq m. \quad (7.1)$$

By way of contradiction, assume that

$$C^{t(n)}(\chi_{1:n}|n) \leq n/2^{2^k+1}.$$

Then, given  $n$ , we can compute  $\chi_{1:n}$  in time  $t(n) \ll t(m')$  from a program of length at most  $m/2$ . Subsequently, given an additional description of  $m$  in  $O(\log m)$  bits, we can output the prefix  $\chi_{1:2m}$  in linear time. Since we have altogether used fewer than  $t(m')$  steps, and a program of length less than  $m$ , we contradict Equation 7.1.  $\square$

Set the constant  $c_t := 1/2^{2^k+1}$  with  $k$  an index of a Turing machine computing  $t$ . The two claims prove the theorem.  $\square$

This lower bound is approximately optimal. Namely, let  $A$  be a recursively enumerable set. For the characteristic sequence  $\chi$  of  $A$  and each constant  $c > 0$ , there exists a total recursive function  $t$  such that for infinitely many  $n$ ,  $C^t(\chi_{1:n}|n) \leq cn$ . This result, as well as Theorem 7.1.3, and Barzdins's Lemma, Theorem 2.7.2 on page 171, are due to J.M. Barzdins [*Soviet Math. Dokl.*, 9(1968), 1251-1254]. See also Exercise 7.1.5 on page 472.

The theorem can be generalized to computable majorants. There is a recursively enumerable set  $A$ , with characteristic sequence  $\chi$ , such that for each computable majorant  $\phi$  and each  $n$ , we have  $\phi(\chi_{1:n}|n) \geq c_\phi n$ , with  $c_\phi$  a constant independent of  $n$  (but dependent on  $\phi$ ) (Exercise 7.1.4, page 472).

If a set  $A$  is recursive, then its characteristic sequence  $\chi$  is also recursive, that is, there is a program that outputs  $\chi_{1:n}$  on input  $n$ . By Exercise 2.3.4 on page 124, the following statements are equivalent:

- $\chi$  is an infinite recursive sequence.
- There exists a constant  $c$  such that for all  $n$ , we have  $C(\chi_{1:n}|n) \leq c$ .
- There exists a constant  $c$  such that for all  $n$ , we have  $C(\chi_{1:n}) \leq \log n + c$ .

Adding a recursive time bound, we obtain a blow-up theorem again. It is stated for  $C^t$ , but the same proof also works for  $K^t$ .

**Theorem 7.1.4** *Let  $f$  and  $t$  be unbounded total recursive functions. There is a recursive sequence  $\chi$  such that  $C^t(\chi_{1:n}|n) \geq n - f(n)$  infinitely often.*

**Proof.** Without loss of generality, let  $f(n) \leq n$  and  $\lim_{n \rightarrow \infty} f(n) = \infty$ . The sequence  $\chi$  will be constructed by diagonalization. Define a function  $g$  inductively by  $g(1) = 1$ , and  $g(n+1) = \min\{m : f(m) > g(n)\}$  for  $n > 1$ . Function  $g$  is total recursive, nondecreasing, and unbounded.

**Step 1** Set  $\chi_1 := 0$ .

**For**  $n := 2, 3, \dots$  **do** Compute  $\chi_{g(n-1)+1:g(n)}$  as follows: Simulate all programs of size less than  $g(n) - g(n-1)$ , each for  $t(g(n))$  steps. Extend  $\chi_{1:g(n-1)}$  to  $\chi_{1:g(n)}$  so that  $\chi_{1:g(n)}$  is not the initial segment of an output of any of the above simulations. {Since there are at most  $2^{g(n)-g(n-1)} - 1$  programs of length less than  $g(n) - g(n-1)$ , extending  $\chi$  in this way is always possible.}

Thus, by construction, for almost all  $n$  we have  $C^t(\chi_{1:g(n)}|g(n)) \geq g(n) - g(n-1)$ . By definition of  $g(n)$ , we have  $f(g(n)) \geq g(n-1)$ . Then for infinitely many  $n$ , we have  $C^t(\chi_{1:n}|n) \geq n - f(n)$ .  $\square$

In Theorem 2.5.5, page 146, it was stated that almost all infinite sequences have maximal Kolmogorov complexity. More precisely, with probability 1 in the sense of the uniform measure, for each infinite sequence  $\chi$  there is a constant  $c$  such that

$$C(\chi_{1:n}|n) \geq n - c$$

infinitely often. The above theorem shows that any recursive time bound blows up the Kolmogorov complexity of some recursive sequences to nearly maximal. But, similar to the complexity oscillations in pure Kolmogorov complexity for random sequences (Theorems 2.5.1 and 2.5.4, page 136 and page 145, respectively), recursive sequences cannot achieve maximal Kolmogorov complexity, even when they are restricted by a linear recursive time bound. In this sense, Theorem 7.1.4 is optimal.

**Theorem 7.1.5** *Let  $\chi$  be a recursive sequence, and let  $t(n) = \Omega(n)$  be an unbounded total recursive function. There is an unbounded total function  $f$  such that for all  $n$ , we have*

$$C^t(\chi_{1:n}|n) \leq n - f(n).$$

**Proof.** Since  $\chi$  is recursive, there is a Turing machine  $T$  that outputs  $\chi_n$  on input  $n$ . We can trivially design a “short” program  $q$  for the universal Turing machine  $U$  such that  $U(\langle p, n \rangle) = \chi_{1:n}$ . The computation of  $U$  on input  $\langle p, n \rangle$  simulates  $T$  on input 0 until it halts, then on input 1 until it halts, and so on, for a *total* of  $n$  steps. Let  $m$  be the last input for which  $T$  halted in this simulation. Then  $U$  has already obtained  $\chi_{1:m}$ . Therefore,  $q$  only needs to additionally “store”  $\chi_{m+1:n}$  of length  $n-m+1$  to be able to reconstruct  $\chi_{1:n}$  totally. This means that  $l(q) = n - m + c$  for some constant  $c$  (the index of  $T$ ). Moreover, the computation of  $\chi_{1:n}$  by  $U$  runs in linear time. Since  $m \rightarrow \infty$  as  $n \rightarrow \infty$ , the theorem follows.  $\square$

An infinite sequence  $\chi$  may have maximal Kolmogorov complexity in the sense that for some  $c$ , for all  $n$ ,

$$K(\chi_{1:n}|n) \geq n - c,$$

where  $K$  is the prefix complexity; see Theorem 3.9.1 on page 232. It is unknown to the authors whether this holds with  $\chi$  is recursive and  $K$  replaced by  $K^t$  with  $t$  a linear recursive time bound, or whether Theorem 7.1.5 holds with prefix complexity  $K$  instead of plain complexity  $C$ , or something in between.

**Example 7.1.2** We considered sequences that are incompressible with respect to time-limited Kolmogorov complexity. Sequences that are incompressible with respect to pure Kolmogorov complexity are shown to be random in Theorem 2.5.5, page 146, and Theorem 3.6.1, page 212. In Section 2.5, sequential (Martin-Löf) tests were used to define randomness of a sequence.

An infinite sequence is random iff it can withstand all sequential tests (Definition 2.5.2). A sequential test is defined by a recursive function (Definition 2.5.1).

Random numbers cannot be generated by arithmetic means. In many applications, such as cryptography, it is sufficient to generate numbers that appear to be random if we do not inspect them too closely. Obviously, the notion of sequential tests can be scaled down by permitting only time-bounded (or space-bounded) tests. In particular, one is often interested in sequences that are just random enough to pass sequential tests that run in polynomial time.

Let  $\{0,1\}^\infty$  be the set of infinite binary sequences with the uniform (Lebesgue) measure. This measure assigns probability  $2^{-l(x)}$  to the set of infinite sequences starting with  $x$ .

**Definition 7.1.8** A sequential test  $\delta$  is called a *sequential ptime- (pspace-) test* if  $\delta$  is polynomial time (space) computable. An infinite sequence  $\omega$  is *ptime- (pspace-) pseudorandom* if for every sequential ptime- (pspace-) test  $\delta$ , and every polynomial  $p$ , it satisfies  $\delta(\omega_{1:p(n)}) < n$  for all but finitely many  $n$ .

The condition on  $\omega$  that  $\delta(\omega_{1:p(n)}) < n$  guarantees that no evidence of nonrandomness of  $\omega$  at significance level  $2^{-n}$  can be found by  $\delta$  unless an initial segment of length greater than  $p(n)$  is examined. For the terminology of testing, consult Sections 2.4 and 2.5. Contrary to the case of unbounded sequential tests, it appears to be the case that there is no universal ptime sequential test. This is because a universal sequential ptime-test must simulate all other sequential ptime-tests, and it is unknown how to make such a universal test run in polynomial time. By diagonalization, on the other hand, one can show that there exist pspace-pseudorandom sequences that are double exponential space computable. We give a criterion in terms of space-bounded complexity for a sequence to be pspace-random.

**Lemma 7.1.1** *Let  $p$  be a polynomial. Let  $C^s$  denote the  $s$ -space bounded generating complexity. If  $C^q(\omega_{1:n}) > n - p(\log n)$  for all polynomials  $q$  and all but finitely many  $n$ , then  $\omega$  is pspace-pseudorandom.*

All pspace-pseudorandom sequences are also ptime-pseudorandom sequences, simply because each ptime-sequential test cannot use more than polynomial space. See also (Exercise 7.1.10 on page 474; Ker-I Ko, *Theoret. Comput. Sci.*, 48(1986), 9-33.)  $\diamond$

### 7.1.2 Resource-Bounded Hierarchies

Given a fixed reference universal machine, the three parameters in

$$C[f(n), t(n), s(n)] = \{x : C^{t,s}(x) \leq f(n), n = l(x)\}$$

characterize classes of strings of different resource-bounded Kolmogorov complexity. A shortest program  $p$  of length  $C^{t,s}(x)$ , from which we can compute  $x$  with  $l(x) = n$  in  $t(n)$  time and  $s(n)$  space simultaneously, is called a *seed*. Intuitively, with a longer seed, more time, or more space, we should be able to obtain some string that is not obtainable with these parameters. Conversely, with a shorter seed, less time, or less space, we may not be able to obtain strings that are obtainable with the current parameters.

As usual, everything needs to be related to one fixed reference universal Turing machine. For any reasonable  $f, s, n, x$  with  $l(x) = n$ , there is always a universal Turing machine  $U$  such that  $x \in C_U[f(n), n, s(n)]$ . This is similar to the case of pure complexity.

The linear space and time speedup theorems, which were useful tools in proving Turing machine space and time hierarchies, do not apply here because compressing time or space results in the increase of program length.

**Theorem 7.1.6** *Let  $f$  and  $g$  be unbounded total recursive functions such that  $f(n) + g(n) \leq n$ , and for each  $k$  we can compute the least  $n$  such that  $f(n) = k$  in space  $s(n) \geq \log n$  and time  $t(n) - n$ . Then for large  $n$ ,*

$$C[f(n) + g(n), t(n), s(n)] - C[f(n), \infty, \infty] \neq \emptyset.$$

**Proof.** It suffices to prove that for the fixed reference universal Turing machine  $U$ , there is a constant  $c$  depending only on  $U$  such that

$$C[f(n) + c, t(n), s(n)] - C[f(n), \infty, \infty] \neq \emptyset.$$

Now let  $n$  and constant  $c$  be large enough so that the following formulas are true. Fix a string  $x$  of length  $l(x) = f(n) + c/2$  such that  $C(x) \geq l(x)$ . We will show that

$$x0^{n-l(x)} \in C[f(n) + c, t(n), s(n)] - C[f(n), \infty, \infty].$$

Note that  $l(x0^{n-l(x)}) = n$ . Let  $p = qx$  be a program that computes  $x0^{n-l(x)}$  as follows:

- Compute the first  $n$  such that  $f(n) + c \geq l(p)$  in  $t(n) - n$  time and  $s(n)$  space.
- Print  $x$  followed by  $0^{n-l(x)}$ .

The length of the program is  $l(p) \leq f(n) + c$ . Hence,  $C^{t,s}(x0^{n-l(x)}) \leq f(n) + c$ . Therefore,  $x0^{n-l(x)} \in C[f(n) + c, t(n), s(n)]$ .

We show that  $x0^{n-l(x)} \notin C[f(n), \infty, \infty]$ . Suppose the contrary. But then we can reconstruct  $x$  as follows: Generate  $x0^{n-l(x)}$  using  $f(n)$  bits. In order to retrieve  $x$ , we need to remove the  $n - l(x)$  length suffix. We know that  $l(x) := f(n) + c/2$ . The length of a self-delimiting program to compute  $c$  is at most  $2 \log c$ . We know  $f(n)$  because we remembered this value as the length of the initial program. We can compute  $n$  from  $f(n)$  in some constant  $d$  bits by the assumption in the theorem. Delete the suffix  $0^{n-l(x)}$  to retrieve  $x$ . Altogether, this description of  $x$  takes  $f(n) + 2 \log c + d$  bits. For  $c$  large enough we obtain  $C(x) < f(n) + c/2 = l(x)$ , a contradiction.  $\square$

**Example 7.1.3** Let  $f(n) = \sqrt{n}$ , which is computable in  $O(n)$  time, and  $g(n) = \log \log n$ , which is also computable in  $O(n)$  time. By the above theorem, for the reference universal Turing machine  $U$ , for  $n$  large enough,

$$C_U[\sqrt{n} + \log \log n, cn, \infty] - C_U[\sqrt{n}, \infty, \infty] \neq \emptyset.$$

Since  $C_U[\sqrt{n}, cn, \infty]$  is by definition contained in both these classes, its containment in  $C_U[\sqrt{n} + \log \log n, cn, \infty]$  is proper.  $\diamond$

**Theorem 7.1.7** Let  $s'(n) \geq 2n + s(n) + c$  and let  $s(n)$  be nondecreasing and computable in space  $s'(n)$ . Let  $f(n) \leq n$  be an unbounded nondecreasing function computable in space  $s'(n) - \log n$ . For  $n$  large enough, we have

$$C[f(n), \infty, s'(n)] - C[n - 1, \infty, s(n)] \neq \emptyset.$$

**Proof.** The idea is to find the first string not in  $C[n - 1, \infty, s(n)]$ . The following program  $p$  for the reference universal Turing machine  $U$  prints a string  $x \in C[f(n), \infty, s'(n)] - C[n - 1, \infty, s(n)]$ . Program  $p$  first computes the first  $n$  such that  $f(n) > l(p)$ , using  $s'(n)$  space, including  $\log n$  space to count to  $n$ . Then  $p$  marks off  $s(n)$  tape cells. Subsequently,  $p$  simulates all  $s(n)$  space-bounded computations of  $U$ , with inputs of sizes up to  $n - 1$ . If  $x$  is the lexicographically first string of length  $n$  not generated by any of these  $s(n)$  space-bounded computations, then  $p$  prints  $x$ . We can find this  $x$  because we can simply repeat all  $s(n)$  space-bounded computations for each next candidate in the lexicographical enumeration of all strings of length  $n$ . Since there are  $2^n$  candidates, and only  $2^n - 1$  inputs of sizes up to  $n - 1$ , such an  $x$  exists. With input  $p$ , reference machine  $U$  uses  $s(n') \geq 2n + s(n) + c$  space. Here  $c$  is a constant,  $n$  tape cells are needed to generate all strings of length  $n$ , another  $n$  tape cells are needed to generate all inputs for the simulated computation, and  $s(n)$  space is needed to carry out the simulation. Initially,  $s(n')$  space is needed to compute  $f(n)$ .

Program  $p$  has size less than  $f(n)$ . It uses  $s'(n)$  space and generates a string  $x$  not in  $C[n - 1, \infty, s(n)]$ .  $\square$

**Example 7.1.4** Theorems 7.1.6 and 7.1.7 are still true with  $C$  replaced by  $CD$  everywhere. For example, for large enough  $n$ , we have  $CD[\log n, \infty, n^2] \subset CD[\log n, \infty, n^2 \log n]$  and  $CD[\log n, \infty, n^2] \subset CD[2 \log n, \infty, n^2]$ .  $\diamond$

A straightforward generalization to time-bounded Kolmogorov complexity leaves an exponential gap between two time classes. We leave this to the Exercises section. The nontrivial time counterpart of the above theorem is still an open problem. The above proof method fails, since time is not reusable. The time- or space-constructibility assumptions in Theorems 7.1.6, 7.1.7 are necessary. Indeed, as in computational complexity theory, we can prove the gap theorems.

**Theorem 7.1.8** *Given any total recursive function  $g(n) \geq n$  there exists a total recursive function  $s(n)$  such that for any  $f(n) < n$ , for large enough  $n$ ,*

$$C[f(n), \infty, s(n)] = C[f(n), \infty, g(s(n))].$$

**Proof.** Consider the reference universal Turing machine  $U$ . Given  $n$ , we define  $s(n)$  to be the first  $i$  such that no program of size up to  $n - 1$  (starting with empty input) halts using between  $i + 1$  and  $g(i)$  space on  $U$ . This  $s(n)$  can always be recursively found because there are only a finite number of programs of size up to  $n - 1$  and we can simulate them recursively in dovetailing style.

Let  $x \in C[f(n), \infty, g(s(n))]$  and let  $p$  be a program of length  $f(n)$  printing  $x$  using at most  $g(s(n))$  space. Because  $f(n) < n$ , and from our construction, we know that  $p$  does not use an amount of space between  $s(n) + 1$  and  $g(s(n))$ . Therefore,  $p$  uses at most  $s(n)$  space, and hence  $x \in C[f(n), \infty, s(n)]$ .  $\square$

Let  $Q$  be a collection of subsets in  $\{0, 1\}^*$ . A set is  $Q$ -immune if it is infinite and does not have infinite subsets belonging to  $Q$ . It is easy to see that for any unbounded total recursive  $f(n) < n$ , the set  $\{0, 1\}^* - C[f(n), \infty, \infty]$  is recursively enumerable-immune. For suppose  $\{0, 1\}^* - C[f(n), \infty, \infty]$  contains an infinite recursively enumerable subset accepted by a Turing machine  $T$ . Then we can construct the following machine  $T'$ : The machine  $T'$  enumerates the subset by simulating  $T$ . While enumerating, it eventually will find a string  $x$  such that  $f(l(x)) > l(T')$ . Since  $x$  is enumerated by a program  $T'$  that is shorter than  $f(n)$  with  $n = l(x)$ , we have  $x \in C[f(n), \infty, \infty]$ . But the subset enumerated by  $T'$  was in the complement of  $C[f(n), \infty, \infty]$ , which gives the required contradiction. We extend this idea in the following example.

**Example 7.1.5** Let  $\lim_{n \rightarrow \infty} s(n)/s'(n) \rightarrow 0$ . Let  $U$  be the reference universal Turing machine. Let  $s'(n) \geq n$  be a nondecreasing function. Let  $f(n)$  be an unbounded nondecreasing function computable in space  $s(n)$  by  $U$ . We prove that  $\{0, 1\}^* - C[f(n), \infty, s'(n)]$  is  $\text{DSPACE}[s(n)]$ -immune.

Assume, by way of contradiction, that  $\{0, 1\}^* - C[f(n), \infty, s'(n)]$  contains an infinite subset  $A \in \text{DSPACE}[s(n)]$ . Let  $T_A$  be a Turing machine that accepts  $A$  using  $O(s(n))$  space. We exhibit a program  $p$  for the reference universal Turing machine  $U$  that prints a long string in  $A$ .

**Step 1** Find the first  $i$  such that  $f(i) > l(p)$ .

**Step 2** Find the first  $x \in A$  {by simulating  $T_A$ } such that  $l(x) \geq i$ . Print  $x$ . {This step uses  $O(s(n))$  space, where  $n = l(x)$ .}

For  $n$  large enough, by choice of  $p$ , the reference universal Turing machine prints  $x$  in space  $s'(n)$  with the program  $p$ , where  $l(p)$  is at most  $f(n)$ . Therefore,  $x \in C[f(n), \infty, s'(n)]$ , a contradiction.  $\diamond$

The sets  $C[c \log n, t, s]$  will be especially useful in certain applications. Such sets belong to the wider class of sparse sets. A set is sparse if it has only a polynomial number of elements for each length.

## Exercises

**7.1.1.** [20] Define  $K^{t,s}$  and  $KD^{t,s}$ , and prove the invariance theorems for them.

**7.1.2.** [20] Prove an invariance theorem for  $CD^{t,s}$ .

**7.1.3.** [25] Prove the results of Theorem 7.1.2 for  $KD^s$ , the space-bounded prefix complexity.

**7.1.4.** [39] Prove that there is a recursively enumerable set  $A$  with characteristic sequence  $\chi$  such that for any total recursive majorant  $\phi$  of  $C$  and any  $n$ , we have  $\phi(\chi_{1:n}|n) \geq c_\phi n$ , where  $c_\phi$  is a constant independent of  $n$  (but dependent on  $\phi$ ).

*Comments.* Use the proof of Theorem 7.1.3 and also Example 4.1.1. Source: A.K. Zvonkin and L.A. Levin, *Russ. Math. Surveys*, 25:6(1970), 83–124], attributed to J.M. Barzdins and N.V. Petri. See also Example 4.1.1.

**7.1.5.** [35] Show that there is a recursively enumerable set  $A$  with characteristic sequence  $\chi$  such that for all total recursive functions  $t$ , and for all  $0 < c < 1$ , there exist infinitely many  $n$  such that  $C^t(\chi_{1:n}) > cn$ .

*Comments.* Compare this exercise with Theorem 7.1.3. R.P. Daley [*Inform. Contr.*, 23(1973), 301–312] proved a more general form of this result using the uniform complexity of Exercise 2.3.3, page 124.

**7.1.6.** [40] Uniform complexity was defined in Exercise 2.3.3, page 124. Here we define the time-bounded uniform complexity. For an infinite string  $\omega$  and the reference universal Turing machine  $U$ ,

$$C^t(\omega_{1:n}; n) = \min\{l(p) : \forall i \leq n [U^t(p, i) = \omega_{1:i}]\}.$$

Define

$$CU[f, t, \infty] = \{\omega : \forall^\infty n [C^t(\omega_{1:n}; n) \leq f(n)]\}.$$

Let  $\omega$  be a Mises-Wald-Church random sequence, but with total recursive admissible place-selection instead of the standard definition with partial recursive admissible place-selection rules. Show that for any unbounded, nondecreasing, total recursive time bound  $t$  we have  $\omega \notin CU[n/2, t, \infty]$ . This holds a fortiori for the more restricted set of Mises-Wald-Church random sequences using partial recursive place-selection rules.

*Comments.* This exercise shows that there is a Mises-Wald-Church random sequence, as defined on page 53, with limiting frequency  $\frac{1}{2}$ , that has very low uniform Kolmogorov complexity, but high time-bounded uniform Kolmogorov complexity. In Exercise 2.5.13, Item (a), page 153, we have proved that there is a Mises-Wald-Church random sequence  $\omega$ , with partial recursive place-selection rules as defined on page 53 and limiting frequency  $\frac{1}{2}$ , such that  $C(\omega_{1:n}; n) \leq f(n) \log n + O(1)$  for every unbounded, nondecreasing, total recursive function  $f$ . In Exercise 2.5.13, Item (c), page 153, we have proved that there is a Mises-Wald-Church random sequence  $\omega$ , with total recursive place-selection rules and limiting frequency  $\frac{1}{2}$ , such that  $C(\omega_{1:n}; n) \leq f(n) + O(1)$  for every unbounded, nondecreasing, total recursive function  $f$ . Now we have shown that if there is a recursive time bound, then all Mises-Wald-Church random sequences (even with respect to only total recursive place-selection rules) indeed look quite random.

One should be able to prove similar theorems for  $C$  and  $K$  versions of this exercise and Exercise 7.1.7. Source: R.P. Daley, *Inform. Contr.*, 23(1973), 301–312. The fact that some Mises-Wald-Church random sequences have low Kolmogorov complexity is from [R.P. Daley, *Math. Systems Theory*, 9(1975), 83–94].

**7.1.7.** [40] Use the uniform complexity of Exercise 7.1.6. We consider a time-information tradeoff theorem for resource-bounded uniform complexity. Let  $f_i = \lceil n^{1/i} \rceil$ , for  $i = 1, 2, \dots$ . Construct a recursive infinite sequence  $\omega$  and a set of total recursive, nondecreasing, unbounded, functions  $\{t_i\}$ , where the  $t_i$ 's are recursively enumerated as  $t_1, t_2, \dots$ , such that the following holds:

- (a) For all  $i > 1$ , we have  $\omega \notin CU[f_i, t_i, \infty]$ , where  $CU[\cdot]$  is defined in Exercise 7.1.6.

- (b) For all  $i, \omega \in CU[f_i, t_{i+1}, \infty]$ .

*Comments.* Source: R.P. Daley, *J. ACM*, 20:4(1973), 687–695. Daley proved a more general statement with a general characterization of  $\{f_i\}$ .

**7.1.8.** [29] Call  $\chi$  polynomial time computable if for some polynomial  $p$  and Turing machine  $T$ , the machine  $T$  on input  $n$  outputs  $\chi_{1:n}$  in time  $p(n)$ . Prove that  $\chi$  is polynomial time computable if and only if for some polynomial  $p$  and constant  $c$  for all  $n$  we have  $C^p(\chi_{1:n}; n) \leq c$ , where  $C(\cdot; \cdot)$  is the uniform complexity of Exercise 7.1.6.

*Comments.* This is a polynomial time version of the results about recursive infinite sequences in Exercise 2.3.4. Source: Ker-I Ko, *Theoret. Comput. Sci.*, 48(1986), 9–33.

**7.1.9.** [30] Let  $f(n)$  be computable in polynomial time and satisfy the condition  $\sum 2^{-f(n)} = \infty$ . Then for every infinite sequence  $\omega$ , there is a polynomial  $p$  such that for infinitely many  $n$ , we have  $C^p(\omega_{1:n}|n) \leq n - f(n)$ .

*Comments.* This is a polynomial time version of Theorem 2.5.1, page 136. Source: Ker-I Ko, *Theoret. Comput. Sci.*, 48(1986), 9–33.

**7.1.10.** [38] Use the terminology of Example 7.1.2. Prove that there exists a pspace-pseudorandom infinite sequence  $\omega$  such that  $\omega_{1:n}$  is computable in  $2^{2^n}$  space.

*Comments.* Hint: prove this in two steps. First, prove that there is an infinite sequence  $\omega$ , computable in double exponential space, such that for all polynomial  $p$  the  $p$  space-bounded *uniform* Kolmogorov complexity satisfies, for infinitely many  $n$ ,  $C^{\infty,p}[\omega_{1:n}; n] \geq n - 3 \log n$ . Second, prove that any such  $\omega$  satisfying the above condition is pspace-pseudorandom. Third, conclude that  $\omega$  is also ptime-pseudorandom. Source: Ker-I Ko, *Theoret. Comput. Sci.*, 48(1986), 9–33.

**7.1.11.** • [27/O48] In Section 3.9.1, we proved various versions of a symmetry of information theorem, with no resource bounds. For example, up to an  $O(\log l(xy))$  additive term,  $C(x, y) = C(x) + C(y|x)$ . Here  $C(x, y) = C(\langle x, y \rangle)$ , where  $\langle \cdot, \cdot \rangle$  is the standard pairing function.

(a) The symmetry of information holds for space-bounded complexity. Let  $s(n) \geq n$  be nondecreasing and let  $f(n)$  be a nondecreasing function computable in  $s(n)$  space. Show that for all  $x, y$  such that  $l(y) = f(l(x))$ , to within an additive term of  $O(\log l(xy))$ ,

$$C^{\infty, O(s(n))}(x, y | l(xy)) = C^{\infty, O(s(n))}(x | l(x)) + C^{\infty, O(s(n))}(y | x).$$

(b) Use the assumptions in Item (a), except that now  $f(n)$  is computable in exponential time. Let  $E = \bigcup_c \text{DTIME}(2^{cn})$ . Show that to within an

additive term of  $O(\log l(xy))$ ,

$$C^{E,\infty}(x,y|l(xy)) = C^{E,\infty}(x|l(x)) + C^{E,\infty}(y|x).$$

(c) (Open) Use the assumptions of Item (a) except that now  $f(n)$  is computable in nondecreasing polynomial time. Prove, or disprove,

$$C^{Poly,\infty}(x,y|l(xy)) = C^{Poly,\infty}(x|l(x)) + C^{Poly,\infty}(y|x).$$

*Comments.* Here and in Exercise 7.1.12 we consider the important question of *resource-bounded symmetry of information*. The main problem in making the proof of Item (b) work for Item (c) is a construction involving the enumeration of a set: one needs to find in polynomial time the  $i$ th element in an exponential size  $P$  set. In general, is there a similar symmetry of information theorem for subexponential time-bounded Kolmogorov complexity? Source: L. Longpré and S. Mocas, *Inform. Process. Lett.*, 46:2(1993), 95–100. They proved that if a certain version of (c) holds, then one-way functions do not exist.

**7.1.12.** [36] We investigate the open problem of Exercise 7.1.11, Item (c), further. We call this problem “polynomial time symmetry of information.” Let  $p(n)$  be a polynomial. We say a subset of  $\{0,1\}^*$  contains “almost all strings” (of  $\{0,1\}^*$ ) if for each  $n$  it contains a fraction of at least  $1 - 1/p(n)$  of  $\{0,1\}^{\leq n}$ . We call a string  $w$  a *polynomial time description* of  $x$  if the universal reference machine  $U$  computes  $x$  from  $w$  in polynomial time.

- (a) If the polynomial time symmetry of information holds, then there is a polynomial time algorithm that computes the shortest  $p$ -time description of a string for “almost all” strings.
- (b) If the polynomial time symmetry of information holds, then every polynomial time computable function is probabilistic polynomial time invertible for “almost all” strings in its domain.
- (c) If  $P = NP$  (which means that every polynomial time computable function is polynomial time invertible), then the polynomial time symmetry of information holds.

*Comments.* The prime question is whether or not symmetry of information holds in a polynomial-time-bounded environment. Intuitively, this problem is related the complexity of inverting a polynomial-time-computable function. The above evidence supports this intuition. Source: L. Longpré and O. Watanabe, *Inform. Computation*, 121:1(1995), 14–22.

**7.1.13.** [27/O35] Let  $T(n), f(n)$  be functions both computable in time  $T(n)$ . Prove  $C[f(n), T(n), \infty] \subset C[f(n), c2^{f(n)}T(n), \infty]$ , for some constant  $c$ . Open Problem: can we establish a tighter version of this time hierarchy without an exponential time gap in the hierarchy?

*Comments.* Source: L. Longpré, Ph.D. Thesis, Cornell University, 1986.

**7.1.14.** [25] Show that for some  $c$ ,  $C[n, \infty, c] \cap C[\log n, \infty, c \log n] \neq \emptyset$  for large  $n$ . Can you formulate other tradeoff (between complexity, time, and space) results?

**7.1.15.** [25] Kolmogorov complexity arguments may be used to replace diagonalization in computational complexity. Prove the following using Kolmogorov complexity:

- (a) If  $\lim_{n \rightarrow \infty} s(n)/s'(n) \rightarrow 0$ , and  $s'(n) \geq \log n$  computable in space  $s'(n)$ , then  $\text{DSPACE}[s'(n)] - \text{DSPACE}[s(n)] \neq \emptyset$ .
- (b) If  $\lim_{n \rightarrow \infty} s(n)/s'(n) \rightarrow 0$ , with  $s'(n)$  computable in space  $s'(n)$  and  $s'(n) \geq 3n$ , then there is a language  $L \in \text{DSPACE}[s'(n)]$  such that  $L$  is  $\text{DSPACE}[s(n)]$ -immune.
- (c) Let  $r(n)$  be any total recursive function. There exists a recursive language  $L$  such that for any Turing machine  $T_i$  accepting  $L$  in space  $S_i(n)$ , there exists a Turing machine  $T_j$  accepting  $L$  in space  $S_j(n)$  such that  $r(S_j(n)) \leq S_i(n)$ , for almost all  $n$ .
- (d) Exhibit a Turing machine that accepts an infinite set containing no infinite regular sets.

*Comments.* Source: suggested by B.K. Natarajan. In Item (a) we consider the *DSPACE hierarchy*. The original space hierarchy theorem was studied by R. Stearns, J. Hartmanis, and P. Lewis II [6th IEEE Symp. Switching Circuit Theory and Logical Design, 1965, pp. 179–190]. Item (c) is the *Blum speedup theorem* from [M. Blum, J. ACM, 14:2(1967), 322–336].

## 7.2 Language Compression

If  $A$  is a recursive set and  $x$  is lexicographically the  $i$ th element in  $A$ , then obviously  $C(x) \leq \log i + c_A$  for some constant  $c_A$  depending only on  $A$ . If further, the membership of  $A$  can be determined in polynomial time  $p(n)$ , and  $A^{=n} = \{x \in A : l(x) = n\}$ , then we are inclined to conjecture,

There exists a constant  $c_A$  such that for all  $x \in A^{=n}$ ,  
we have  $C^p(x|n) \leq \log d(A^{=n}) + c_A$ .

But since a Turing machine cannot search through all  $2^n$  strings in polynomial time  $p(n)$ , the argument that worked in the recursive case does not apply in the polynomial time setting. This conjecture has interesting consequences in language compression. Results in this section all stem from this question. Although whether or not the conjecture is true is

still an important open problem in time-bounded Kolmogorov complexity, we will try to provide some partial answers to it. In this section, we will deal only with time-bounded Kolmogorov complexity  $C^t$  and  $CD^t$ .

- Definition 7.2.1**
- (i) Let  $\Sigma = \{0, 1\}$ . A function  $f: \Sigma^* \rightarrow \Sigma^*$  is a *compression* of language  $L \subseteq \Sigma^*$  if  $f$  is one-to-one on  $L$  and for all except finitely many  $x \in L$ , we have  $l(f(x)) < l(x)$ .
  - (ii) For a function  $f$ , the inverse of  $f$  is  $f^{-1}: f(L) \rightarrow L$ , such that for every  $x \in L$ , we have  $f^{-1}(f(x)) = x$ . A language  $L$  is *compressible* in time  $T(n)$  if there is a compression function  $f$  for  $L$  that can be computed in time  $T(n)$ , and its inverse  $f^{-1}$  can also be computed in time  $T(n)$ .
  - (iii) A compression function  $f$  *optimally compresses* a language  $L$  if for every  $x \in L$  of length  $n$ ,  $l(f(x)) \leq \lceil \log \sum_{i=0}^n d(L^{=i}) \rceil$ .
  - (iv) A natural and optimal form of compression is *ranking*. The ranking function  $r_L: L \rightarrow \mathcal{N}$  maps  $x \in L$  to its index in a lexicographical ordering of  $L$ .

Obviously, language compression is closely related to the Kolmogorov complexity of the elements in the language. Efficient language compression is closely related to the time-bounded Kolmogorov complexity of the elements of the language. Using a ranking function on a recursive set, we achieve the optimal Kolmogorov complexity for most elements in the set. In this sense, ranking optimally compresses a recursive set. When there is no resource bound, this compression is trivial. Our purpose is to study the polynomial time setting of the problem. This is far from trivial.

### 7.2.1 Decision Compression

Given a polynomial time computable set  $A$  and  $x \in A^{=n}$ , can we compress  $x$  by representing  $x$  in  $\log d(A^{=n}) + O(1)$  bits? In order to apply the non-resource-bounded Kolmogorov complexity in previous chapters, we have often implicitly relied on Theorem 2.1.3 on page 103. The theorem states that given a recursively enumerable set  $A$ , for every  $x \in A^{=n}$  we have  $C(x|n) \leq \log d(A^{=n}) + O(1)$ . The short program for  $x$  is simply its index in the enumeration of elements in  $A^{=n}$ . Similarly, we also have,  $CD(x|n) \leq \log d(A^{=n}) + O(1)$ .

Even for  $A$  in P, a program that generates  $x$  this way must search through all  $2^n$  strings over  $\Sigma^n$  and test whether they belong to  $A^{=n}$ . Such a process takes at least exponential time. We are interested in the compression achievable in polynomial time. Remarkably, for  $CD$  near optimal, compression is achievable in polynomial time.

- Theorem 7.2.1** *Let  $A$  be any set. There is a constant  $c$  and a polynomial  $p$  such that for all strings  $x \in A^{=n}$ ,*

$$CD^p(x|A^{=n}) \leq 2 \log d(A^{=n}) + 2 \log n + c.$$

**Proof.** Let  $A$  be a set and  $d = d(A^{=n})$ .

**Lemma 7.2.1** *Let  $S = \{x_1, \dots, x_d\} \subseteq \{0, \dots, 2^n - 1\}$ . For all  $x_i \in S$  there exists a prime  $p_i \leq dn$  such that for all  $j \neq i$ ,  $x_i \not\equiv x_j \pmod{p_i}$ .*

**Proof.** Let  $N = 2^n$ . Consider prime numbers between  $c$  and  $2c$ . For each pair  $x_i, x_j \in S$  there are at most  $\log_c N = \log N / \log c$  different prime numbers  $p$  such that  $c \leq p \leq 2c$  and  $x_i \equiv x_j \pmod{p}$ . Since there are only  $d-1$  pairs of strings in  $S$  containing  $x_i$  it follows that there exists a prime number  $p_i$  among  $(d-1) \log N / \log c + 1$  prime numbers such that for all  $j \neq i$ ,  $x_i \not\equiv x_j \pmod{p_i}$ . The Prime Number Theorem, Exercise 1.5.8 on page 17, implies that there are at least  $c / \log c$  prime numbers between  $c$  and  $2c$ ; hence taking  $c > (d-1) \log N$  yields that  $p_i \leq 2d \log N = 2dn$ .  $\square$

Let  $A$  play the role of  $S$  in Lemma 7.2.1. For  $x \in A$ , apply Lemma 7.2.1 to get  $p_x$ . The *CD* program for  $x$  works as follows:

**Step 1** input  $y$ ;

**Step 2** If  $y \notin A^{=n}$  then reject  $y$

else if  $y \equiv x \pmod{p_x}$  then accept  $y$  else reject  $y$

The size of the above program is  $l(p_x) + l(x \pmod{p_x}) + O(1)$ . This is  $2 \log d(A^{=n}) + 2 \log n + O(1)$ . Notice that by padding, we can make the descriptions of  $p_x$  and  $x \pmod{p_x}$  equally long, precisely  $\log d + \log n$  bits. Concatenating the two descriptions we can parse them in the middle, saving an  $O(\log \log n)$  bit overhead for a delimiter to separate them. It is clear that the program runs in polynomial time, and only accepts  $x$ .  $\square$

**Corollary 7.2.1** *Let  $A$  be a set in  $\mathbf{P}$ . There is some polynomial  $p$  such that for every string  $x \in A^{=n}$  it holds that*

$$CD^p(x|n) \leq 2 \log d(A^{=n}) + 2 \log n + O(1).$$

Can we improve Theorem 7.2.1 and Corollary 7.2.1? The next theorem states that the compression can be made almost optimal at the penalty of adding a lot of information to the conditional. That is, we are given a magical string  $s$  depending only on  $A^{=n}$  for free.

**Theorem 7.2.2** *There is a polynomial  $p(n)$  such that for every set  $A$  and every large  $n$ , if  $x \in A^{=n}$ , then*

$$CD^p(x|A^{=n}, s) \leq \log d(A^{=n}) + \log \log d(A^{=n}) + O(1),$$

where  $A^{=n}$  is given as an oracle, and string  $s$  depends only on  $A^{=n}$  and has length about  $n \log d(A^{=n})$ .

In this theorem, we do not need to assume that  $A$  is accepted in polynomial time. If, for example, for some  $c$  and for each  $n$ , the set  $A^{=n}$  is accepted by a circuit of size  $n^c$ , then the oracle  $A^{=n}$  in the condition can be replaced by a finite string, of polynomial length, describing the circuit. In the proof, if we need to decide whether or not an element  $x$  is in  $A^{=n}$ , then we just simulate the circuit in polynomial time. The same applies to Theorem 7.2.4.

**Proof.** In order to prove the above theorem, we need a coding lemma. Let  $k = d(A^{=n})$ , and  $m = 1 + \lceil \log k \rceil$ . Let  $h : \Sigma^n \rightarrow \Sigma^m$  be a linear transformation given by a randomly chosen  $m \times n$  binary matrix  $R = \{r_{ij}\}$ . That is, for each  $x \in \Sigma^n$ , we have  $Rx$  is a string  $y \in \Sigma^m$  where  $y_i \equiv (\sum_j r_{ij} \cdot x_j) \bmod 2$ .

Let  $H$  be a collection of such functions. Let  $B, C \subseteq \Sigma^n$  and  $x \in \Sigma^n$ . The mapping  $h$  separates  $x$  within  $B$  if for every  $y \in B$ , different from  $x$ , it maps  $h(y) \neq h(x)$ . The mapping  $h$  separates  $C$  within  $B$  if it separates each  $x \in C$  within  $B$ . The set of mappings  $H$  separates  $C$  within  $B$  if for each  $x \in C$  some  $h \in H$  separates  $x$  within  $B$ . In order to give each element in  $B$  a (logarithmic) short code, we randomly hash elements of  $B$  into short codes. If collisions can be avoided, then elements of  $B$  can be described by short programs. We now state and prove the Coding Lemma.

**Lemma 7.2.2 (Coding Lemma)** *Let  $B \subseteq \Sigma^n$ , where  $d(B) = k$ . Let  $m = 1 + \lceil \log k \rceil$ . There is a collection  $H$  of  $m$  random linear transformations  $\Sigma^n \rightarrow \Sigma^m$  such that  $H$  separates  $B$  within  $B$ .*

**Proof.** Fix a random string  $z$  of length  $nm^2$  such that  $C(z|B, P, m, n) \geq l(z)$ , where  $P$  is the program for describing  $z$  (independent of the actual  $z$ ) in the following discussion. Cut  $z$  into  $m$  equal pieces. Use the  $nm$  bits from each piece to form an  $n \times m$  binary matrix in the obvious way. Thus, we have constructed a set  $H$  of  $m$  random matrices. We claim that  $H$  separates  $B$  within  $B$ .

Assume the contrary. That is, for some  $x \in B$ , no  $h \in H$  separates  $x$  within  $B$ . Then there exist  $y_1, \dots, y_m \in B$  such that  $h_i(x) = h_i(y_i)$ . Hence,  $h_i(x - y_i) = 0$ . (Here  $x - y_i$  means mod 2 element-wise subtraction.)

Since  $x - y_i \neq 0$ , the first column of  $h_i$  corresponding to a “1” in vector  $x - y_i$  can be expressed by the rest of columns using  $x - y_i$ . Now we can describe  $z$  by a fixed program  $P$  that uses the following data:

- index of  $x$  in  $B$ , using  $\lceil \log k \rceil$  bits;
- indices of  $y_1, \dots, y_m$ , in at most  $m\lceil \log k \rceil$  bits;

- matrices  $h_1, \dots, h_m$  each minus its redundant column, in  $m^2n - m^2$  bits.

From this description of the nonredundant columns of the  $h_i$ 's, and  $x$  and the  $y_i$ 's, the program  $P$  given  $B$  (as an oracle) will reconstruct  $h_i$ . The total length of the description is only

$$m^2n - m^2 + \lceil \log k \rceil + m(\lceil \log k \rceil) \leq m^2n - 1.$$

Hence,  $C(z|B, P, m, n) < l(z)$ , a contradiction.  $\square$

We continue the proof of the theorem. Let  $H$  be a collection of functions as given by the Coding Lemma 7.2.2. Let  $s$  be an encoding of  $H$ . For each  $z \in A^{=n}$  there is some  $h_i \in H$  that separates  $z$  within  $A^{=n}$ . Given oracle  $A^{=n}$  and  $s$ , a short program with input  $ih_i(z)$  can accept exactly  $z$  as follows:

Check in lexicographical order for each candidate  $x$  whether  $x \in A^{=n}$ . If  $x \in A^{=n}$ , then decode  $ih_i(z)$  into  $i$  and  $h_i(z)$ , and find  $h_i$  using  $i$  and  $s$ . Compute  $h_i(x)$ . Accept  $x$  iff  $h_i(x) = h_i(z)$ , using the fact that  $h_i$  separates  $z$  within  $A^{=n}$ .

Therefore,  $z$  can be described by  $i$  followed by  $h_i(z)$ , where  $i$  uses precisely  $\lceil \log m \rceil$  bits by padding and  $h_i(z)$  requires  $m$  bits. All of this can be done in some polynomial time  $p(n)$ . Hence,

$$\begin{aligned} CD^p(z|A^{=n}, s) &\leq m + \log m + O(1) \\ &= \log d(A^{=n}) + \log \log d(A^{=n}) + O(1). \end{aligned}$$

$\square$

**Example 7.2.1** In this example, we provide another nice application of the Coding Lemma 7.2.2. A PTM (*probabilistic Turing machine*) is a Turing machine that can flip a fair coin to determine its next move. A PTM is also called a random algorithm. The class BPP is the set of all decision problems solvable by a polynomial time PTM such that the “yes–no” answer always has probability at least  $\frac{1}{2} + \delta$  of being correct for some fixed  $\delta > 0$ . One can make the error probability less than  $1/2^n$  by running such an algorithm a polynomial number of times and taking the majority answer.

**Theorem 7.2.3**  $BPP \subseteq \Sigma_2^p \cap \Pi_2^p$ .

**Proof.** Let  $B \in BPP$  be accepted by a PTM  $T$  using  $m = n^k$  coin flips (random bits), with error probability  $\epsilon \leq 1/2^n$  on inputs of length  $n$ . Let  $E_x \subset \Sigma^m$  be the collection of coin-flip sequences of length  $m$  on which

$T$  rejects  $x$ . If  $x \in B$ , then we must have  $d(E_x) \leq 2^{m-n}$  since otherwise we would have  $\epsilon > 1/2^n$ . Setting  $l = 1 + m - n$ , the Coding Lemma 7.2.2 states that there is a collection  $H$  of  $l$  linear transformations from  $\Sigma^m$  to  $\Sigma^l$  separating  $E_x$  within  $E_x$ . If  $x \notin B$ , then  $d(E_x) > 2^{m-1}$ , since otherwise  $\epsilon > 1/2^n$  again. But then by the pigeon-hole principle, no such collection  $H$  as above exists. Hence,  $x \in B$  if and only if such an  $H$  exists. The latter can be expressed as

$$(\exists H)(\forall e \in E_x)(\exists h \in H)(\forall e' \in E_x)[e \neq e' \Rightarrow h(e) \neq h(e')].$$

The second existential quantifier has polynomial range, and hence can be absorbed into the polynomial time computation. So,  $\text{BPP} \subseteq \Sigma_2^p$ . Since  $\text{BPP}$  is closed under complement,  $\text{BPP} \subseteq \Pi_2^p$  and  $\text{BPP} \subseteq \Sigma_2^p \cap \Pi_2^p$ .  $\square$

Denote the set of languages accepted in random polynomial time by  $R$ . Let  $R_2 = R^{\text{NP}}$  be the collection of languages accepted in random polynomial time with an oracle for an NP-complete set. The above proof also yields  $\text{BPP} \subseteq R_2 \cap \text{CoR}_2$ .

$\diamond$

### 7.2.2 Description Compression

In this section, we discuss the compression issues with respect to  $C^t$ . Let  $A$  be total recursive and let  $d(A^{=n}) \leq p(n)$  for some polynomial  $p$  for all  $n$ . Then by Theorem 2.1.3, page 103, there exists a constant  $c$  such that for all  $x \in A^{=n}$

$$C(x|n) \leq c \log n + c. \quad (7.2)$$

The next theorem states that this is true up to a  $\log \log d(A^{=n})$  additive factor in a polynomial time setting with the help of a  $\Sigma_2^p$  oracle.

**Theorem 7.2.4** *There is a polynomial  $p(n)$  such that for each set  $A$  and each large  $n$ , if  $x \in A^{=n}$ , then*

$$C^p(x|A^{=n}, s, \text{NP}^A) \leq \log d(A^{=n}) + \log \log d(A^{=n}) + O(1), \quad (7.3)$$

$$C^p(x|A^{=n}, \Sigma_2^{p,A}) \leq \log d(A^{=n}) + \log \log d(A^{=n}) + O(1), \quad (7.4)$$

where  $A^{=n}$ , an  $\text{NP}^A$ -complete set, and a  $\Sigma_2^{p,A}$ -complete set are given as oracles, and  $s$  is a string of length about  $n \log d(A^{=n})$ .

**Proof.** Equation 7.3 follows from Theorems 7.1.2 and 7.2.2. To prove Equation 7.4, we use the  $\Sigma_2^{p,A}$  oracle  $B$  to find  $s$ , the encoding of  $H$  in the proof of Theorem 7.2.2. By that proof there exists an  $H$  such that for all  $x \in A^{=n}$ , some  $h_i$  in  $H$  separates  $x$  within  $A^{=n}$  ( $H$  separates  $A^{=n}$  within  $A^{=n}$ ). The idea is to reconstruct  $s$  by asking questions to  $B$ . Since  $H$  is not unique, we construct  $s$  bit by bit. Assume we have

constructed the prefix  $s_1 \dots s_i$  of  $s$ . We extend this with one more bit  $s_{i+1}$ , where  $s_{i+1} = 1$  if oracle  $B$  answers “yes” to the question “is there an  $H$  with prefix  $s_1 \dots s_i s_{i+1}$  such that for all  $x \in A^{=n}$ , some  $h_i \in H$  separates  $x$  within  $A^{=n}$ ?”.  $\square$

We can get rid of the oracle at the cost of obtaining much weaker inequalities than Equation 7.2 or Theorem 7.2.4. For a given language  $L$ , define the density function of  $L$  to be  $\mu_L(n) = d(L^{=n})/2^n$ .

**Theorem 7.2.5** *If  $L$  is acceptable in polynomial time,  $k > 3$ , and  $\mu_L \leq n^{-k}$ , then  $L$  can be compressed in probabilistic polynomial time. The compression function maps strings of length  $n$  to strings of length at most  $n - (k-3)\log n + c$ .*

**Proof.** Complicated and omitted. See [A. Goldberg and M. Sipser, *SIAM J. Comput.*, 20(1991), 524–536].  $\square$

Theorem 7.2.5 is far from optimal in two ways:

- If a language  $L$  is very sparse, say  $\mu_L \leq 2^{-n/2}$ , then one expects to compress  $n/2$  bits instead of only  $O(\log n)$  bits given by the theorem; can this be improved?
- The current compression algorithm is probabilistic; can this be made deterministic?

In computational complexity, oracles sometimes help us to understand the possibility of proving a new theorem. The following result shows that if  $S$ , the language to be compressed, is not restricted to be polynomial time acceptable, and the membership query of  $S$  is given by an oracle, then compression by only a logarithmic term is optimal. The point here is that polynomial time decidability of language membership is a stronger requirement than membership decidability by just oracle questions. The technique used in the following proof will be used again in the proof for Theorem 7.3.3.

**Definition 7.2.2** A set  $A$  is called a *sparse set* if for some constant  $c$ , for all  $n$ , the cardinality  $d(A^{=n}) \leq n^c + c$ .

**Theorem 7.2.6** *There is a sparse language  $L$  such that if  $L$  is compressed by a probabilistic polynomial time machine with an oracle for  $L$ , then the compression function maps strings of length  $n$  to strings of length  $n - O(\log n)$ .*

**Proof.** Let  $L$  be a language that contains exactly one string  $x$  for each length  $n$ , which is a tower of 2's:

$$n = 2^{2^{\dots^2}}.$$

The language  $L$  contains no strings other than those just described. Each string  $x \in L$  is maximally complex:

$$C(x) \geq l(x). \quad (7.5)$$

Let  $T$  be a probabilistic machine with an oracle for  $L$  that runs in time  $n^k$ . By way of contradiction, let  $f$  be the compression function such that for large  $n$  and  $x \in L$  of length  $n$ , and constant  $c_1$  to be chosen later,  $l(f(x)) \leq n - c_1 - (k + c_1) \log n$ . It is enough to show that  $T$  cannot restore  $x$  from  $f(x)$  with probability at least  $\frac{1}{2}$ .

Assume that  $T$  on input  $f(x)$  outputs  $x$  with probability  $\frac{1}{2}$ . Hence, on half of the  $2^{n^k}$  coin-toss sequences of length  $n^k$ ,  $T$  should output the correct  $x$ . Let  $R$  be the set of such coin-toss sequences that lead to correct  $x$ . Choose  $r \in R$  such that

$$C(r|x) \geq l(r) - 1. \quad (7.6)$$

Since  $C(r) \leq l(r) + O(1)$  (Equations 7.5,7.6) we have by the Symmetry of Information Theorem 2.8.2 on page 182 that for some constant  $c_2$

$$C(x|r) \geq n - c_2 \log n. \quad (7.7)$$

In order to give a short description of  $x$ , relative to  $r$ , we will describe in  $(k + 2) \log n + c_1$  bits all information needed to answer the oracle queries of  $T$  on input  $f(x)$ .

$L$  is so sparse that on input of length of  $n$  (which is a tower of 2's)  $T$  has no time to write down a string in  $L$  of length greater than  $n$  because any such string has length at least  $2^n$ . That is, for any query “ $x' \in L?$ ,” if  $l(x') > n$ , then the answer is “no.” In addition, we can write down all the strings in  $L$  of lengths less than  $n$  in less than  $2 \log n$  bits. This enables us to answer queries of lengths less than  $n$ .

Since  $T$ , with respect to the coin-toss sequence  $r$ , outputs  $x$ , we can redescribe  $x$  by simulating  $T$  on input  $f(x)$ , using the random sequence  $r$  and the following information:

- This discussion, in at most  $c_1$  bits;
- $2 \log n$  bits to encode strings in  $L$  of length less than  $n$ ;
- If  $T$  queries “ $x \in L?$ ,” then we use  $k \log n$  bits to indicate the first step when  $T$  makes such query. Before this step, if  $T$  queries about a string of length  $n$ , then the answer is “no.”

Therefore,

$$C(x|r) \leq c_1 + 2 \log n + k \log n + l(f(x)).$$

Choosing  $c_1 > c_2 + 2$ , the right side of the inequality is less than  $n - c_2 \log n$ , contradicting Inequality 7.7.  $\square$

There is a language  $L$ , of density  $\mu_L < 2^{-n/2}$ , that cannot be compressed by any deterministic polynomial time machine that uses an oracle for  $L$ . Further, the following results can be shown by diagonalization: (a) There is an exponential time language that cannot be compressed in deterministic polynomial time. (b) There is a double exponential time language that cannot be compressed in probabilistic polynomial time. See A. Goldberg and M. Sipser, *SIAM J. Comput.*, 20(1991), 524–536.

### 7.2.3 Ranking

Ranking is a special and optimal case of compression. Recall that the ranking function  $r_L$  maps the strings in  $L$  to their indices in the lexicographical ordering of  $L$ . If  $r_L : L \rightarrow \mathcal{N}$  is polynomial time computable, then so is  $r_L^{-1} : \mathcal{N} \rightarrow L$  (Exercise 7.2.6, page 486). We are only interested in polynomial time computable ranking functions. In the previous sections we met several hard-to-compress languages. But, there are also quite natural language classes that are easy to compress.

A one-way log-space Turing machine is a deterministic Turing machine with a separate one-way input tape on which the input is delimited between distinguished endmarkers, and a fixed number of separate work tapes. During its computation on an input of length  $n$ , the machine uses only  $O(\log n)$  space on its work tapes.

**Theorem 7.2.7** *If a language  $L$  is accepted by a one-way log-space Turing machine, then  $r_L$  can be computed in polynomial time.*

**Proof.** Let  $T$  be a one-way log-space Turing machine that accepts  $L$ . We want to compute  $r_L(x)$  for each string  $x$  of length  $n$ . Write  $y \leq x$  if  $y$  precedes  $x$  in the length-increasing lexicographical order. Let  $L_x = \{y \in L : y \leq x\}$ . Trivially,  $r_L(x) = d(L_x)$ . By storing  $x$  in the internal memory of  $T$ , we obtain a machine  $T_x$  accepting  $L_x$ . That is,  $T_x$  on each input simulates  $T$  and compares the input with  $x$  at the same time. The machine  $T_x$  accepts iff  $T$  accepts the input and the input is less than or equal to  $x$ .

$T_x$  has size polynomial in  $n = l(x)$  and uses  $\log n$  space. We can construct a directed computation graph  $G = (V, E)$  for  $T_x$  as follows: Let an ID of  $T_x$  be a triple,

(state, input head position, work tape content).

The set of nodes  $V$  contains all possible IDs of  $T_x$ . The set of edges  $E$  contains directed edges that represent moves of  $T_x$ . The sets  $V$  and

$E$  are polynomially bounded in  $n$  because  $T_x$  has a polynomial number of states,  $n$  input positions, and  $\log n$  work tape cells. Since  $T_x$  is deterministic and runs in bounded time,  $G$  is loop-free.

In order to compute  $r_L$ , we only need to calculate  $d(L_x)$  by counting the number of accepting paths in  $G$  from the initial ID to the final ID. This latter task can be easily done by dynamic programming since  $G$  does not contain loops.  $\square$

## Exercises

**7.2.1.** [20] (a) Show that a set  $S$  is sparse iff for all  $x \in S$ ,  $CD^p(x|S) \leq O(\log l(x))$  for some polynomial  $p$ .

(b) Show that set  $S \in P$  is sparse iff for all  $x \in S$ ,  $CD^p(x) \leq O(\log l(x))$  for some polynomial  $p$ .

*Comments.* Use Theorem 7.2.1. Source: H.M. Buhrman and L. Fortnow, *Proc. 14th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comput. Sci., Springer-Verlag, Heidelberg, 1997. The authors also define a non-deterministic version of  $CD$  and prove several results.

**7.2.2.** [30] (a) Show that for  $0 \leq x < y < 2^n$  there are at most  $n$  primes  $p$  such that  $x \equiv y \pmod{p}$ .

(b) Prove that there is some polynomial  $p$  such that for all formulas  $\phi(x_1, \dots, x_n) \in SAT$  and all  $r$  such that  $l(r) = p(l(\phi))$  and  $C(r) \geq l(r)$ , there is some satisfying assignment  $a$  of  $(x_1, \dots, x_n)$  such that  $CD^p(a|\phi, r) \leq O(\log n)$ .

*Comments.* Hint: see the proof of Theorem 7.2.1. Source: H.M. Buhrman and L. Fortnow, *Ibid.*, where a connection is given between this exercise and [L. Valiant and V. Vazirani, *Proc. 17th ACM Symp. Theory Comput.*, 1985, pp. 458–463].

**7.2.3.** [O43] Improve the upper bounds in Theorems 7.2.1 and 7.2.2. Can we achieve the optimal bound  $\log d(A^{=n}) + O(1)$ ?

*Comments.* H.M. Buhrman and L. Fortnow, *Ibid.*, observe that proving Corollary 7.2.1 to be optimal implies  $P = P^{\#P}$  (see Exercise 7.2.7 for the definition of  $\#P$ ).

**7.2.4.** [30] There is an infinite set  $A$  such that for every polynomial  $p$ ,  $CD^p(x|A) \geq l(x)/5$  for almost all  $x \in A$ .

*Comments.* Source: L. Fortnow and M. Kummer, *Theoret. Comput. Sci. A*, 161(1996), 123–140.

**7.2.5.** [28] Let  $f$  be a function on the natural numbers. Let  $\Sigma = \{0, 1\}$ . A set  $A$  belongs to the class  $P/f$  if there exist another set  $B \in P$  and a function  $h: \mathcal{N} \rightarrow \Sigma^*$  such that for all  $n$ , we have  $l(h(n)) \leq f(n)$ ; and for

all  $x$ , we have  $x \in A$  iff  $\langle x, h(l(x)) \rangle \in B$ . Define  $P/\text{poly} = \bigcup_{c>0} P/n^c$ . Prove that  $BPP \subseteq P/\text{poly}$ , using Kolmogorov complexity.

*Comments.* Let  $T$  be a probabilistic machine such that the error probability is  $1/2^{n^2}$  and each path is of length  $n^k$ . A Kolmogorov random string of length  $n^k$  will always give a correct path. Source: W. Gasarch, e-mail, July 16, 1991. The class  $P/f$  was defined by R.M. Karp and R.J. Lipton in [*L'Enseignement Mathématique*, 28(1982), 191–209].

**7.2.6.** [28] Let  $r_L$  be the ranking function of  $L$ . Show that if  $r_L$  is polynomial time computable, then so is  $r_L^{-1}$ .

*Comments.* Source: A. Goldberg and M. Sipser, *SIAM J. Comput.*, 20(1991), 524–536.

**7.2.7.** [30] A problem is in  $\#P$  if there is a nondeterministic Turing machine such that for each input, the number of distinct accepting paths of the Turing machine is precisely the number of solutions for the problem for this input.  $\#P$ -complete problems are defined (analogous to NP-complete problems) as follows: the problem is in  $\#P$ , and each  $\#P$  problem can be reduced to it by a polynomial time deterministic Turing machine computation. As an example, counting the number of satisfying truth assignments for SAT is  $\#P$ -complete. Let  $\mathbf{C}$  be a class of languages. Say  $\mathbf{C}$  is P-rankable if for all  $L \in \mathbf{C}$ , the ranking function  $r_L$  is polynomial time computable. Prove:

- (a)  $P$  is P-rankable iff  $NP$  is P-rankable.
- (b)  $P$  is P-rankable iff  $P = P^{\#P}$ .
- (c)  $\text{PSPACE}$  is P-rankable iff  $P = \text{PSPACE}$ .
- (d)  $P/\text{poly}$  is not P-rankable, where the class  $P/\text{poly}$  is defined in Exercise 7.2.5.
- (e) Languages accepted by two-way deterministic pushdown automata are P-rankable iff  $P = \#P$ .
- (f) Languages accepted by one-way multihead DFA's are P-rankable iff  $P = \#P$ .
- (g) Languages accepted by one-way log-space bounded nondeterministic Turing machines are P-rankable iff  $P = \#P$ .

*Comments.* Source: Items (a)–(d) are from [L. Hemachandra and S. Rudich, *J. Comput. System Sci.*, 41:2(1990), 251–271]. Items (e)–(g) are from [D.T. Huynh, *Math. Systems Theory*, 23(1990), 1–19]. The  $\#P$  class was introduced by L. Valiant [*Theoret. Comput. Sci.*, 8(1979), 189–201].

**7.2.8.** [40] Is it possible that  $\#P$  problems have solutions of low time-bounded Kolmogorov complexity relative to the input? Prove that if

there is a polynomial time Turing machine that on input Boolean formula  $f$ , prints out a polynomial-sized list of numbers among which one number is the number of solutions of  $f$  (though we may not know which one), then  $P = P^{\#P}$ . This gives strong evidence that in general the number of solutions of a Boolean formula has high time-bounded Kolmogorov complexity relative to the formula.

*Comments.* Source: J.-Y. Cai and L. Hemachandra, *Inform. Process. Lett.*, 38(1991), 215–219. Hint: use A. Shamir’s polynomial interpolation technique in [*Proc. 31st IEEE Found. Comput. Sci.*, 1990, pp. 11–15].

**7.2.9.** [35] Let  $\chi_L = \chi_1\chi_2\dots$  be the characteristic sequence for language  $L \subseteq \{0,1\}^*$  such that  $\chi_i = 1$  iff the lexicographically  $i$ th word  $w_i$  is in  $L$ . As before,  $L^{< n}$  is the set of strings in  $L$  with length less than  $n$ . Define the time-space-bounded Kolmogorov complexity of  $L$  as  $C^{t,s}(\chi_{L^{< n}})$ . Let  $t(n) = 2^{n^{O(1)}}$  and  $s(n, \epsilon) = c^{n^\epsilon}$ . Prove by diagonalization:

(a) There is language  $L \in \text{DTIME}[2^{2^{O(n)}}]$  such that  $t(n)$  time-bounded Kolmogorov complexity of  $L$  is “exponential” almost everywhere (note,  $\chi_{L^{< n}}$  has exponential length). That is, for all but finitely many  $n$ ,

$$C^{t(n),\infty}(\chi_{L^{< n}}|n) > s(n, \epsilon), \quad \text{for some } c > 1, \epsilon > 0.$$

(b) Use Item (a) to show: if  $L$  is  $\text{DTIME}[2^{2^{O(n)}}]$ -hard under polynomial time Turing reduction, then the  $t(n)$  time-bounded Kolmogorov complexity of  $L$  is “exponential” almost everywhere. That is, for all but finitely many  $n$ ,

$$C^{t(n),\infty}(\chi_{L^{< n}}|n) > s(n, \epsilon), \quad \text{for some } c > 1, \epsilon > 0.$$

(c) There is a language  $L \in \text{SPACE}[2^{O(n)}]$  such that for all but finitely many  $n$ , we have  $C^{\infty, 2^n}(\chi_{L^{< n}}|n) > 2^{n-2}$ .

(d) Use Item (c) to show: if  $L$  is  $\text{SPACE}[2^{O(n)}]$ -hard under polynomial time Turing reduction, then there exists a constant  $\epsilon > 0$  such that for all but finitely many  $n$ , we have  $C^{\infty, s(n, \epsilon)}(\chi_{L^{< n}}|n) > s(n, \epsilon)$ .

(e) There is a language  $L \in \text{SPACE}[2^{O(n)}]$  such that for large enough  $n$ ,  $C^{\infty, 2^n}(\chi_{L^{< n}}) > 2^n - n$ .

(f) Consider the  $P/\text{poly}$  class defined in Exercise 7.2.5. It is not known whether  $\text{DTIME}[t(n)]$  is contained in  $P/\text{poly}$ . A nonsparse language  $L$  is said to be  $P/\text{poly}$  *immune* if it has only sparse subsets in  $P/\text{poly}$ . A language  $L$  is said to be  $P/\text{poly}$  *bi-immune* if  $L$  and its complement are both  $P/\text{poly}$  immune. Use Item (e) to show that there exists a language  $L \in \text{SPACE}[2^{O(n)}]$  that is  $P/\text{poly}$  bi-immune. (Use the language  $L$  in (e).)

*Comments.* Source: D.T. Huynh, in [Proc. 1st Conf. Structure Complexity Theory, 1986, pp. 184–195; Theoret. Comput. Sci., 96(1992), 305–324; Inform. Comput., 90(1991), 67–85; Inform. Process. Lett., 37(1991), 165–169].

## 7.3 Computational Complexity

### 7.3.1 Constructing Oracles

Time- and space-bounded Kolmogorov complexities are natural tools in the study of time- and space-bounded computations and the study of the structures of the corresponding complexity classes. As in Chapter 6, they provide powerful techniques and help to make intuitive arguments rigorous.

With a simple input, a polynomial time bounded Turing machine cannot compute too complicated strings. Therefore, it cannot ask complicated questions to an oracle. Kolmogorov complexity enables us to formalize this intuition.

**Definition 7.3.1** A function  $f : \{0, 1\}^* \rightarrow \{0, 1\}^*$  is *honest* if there exists a  $k$  such that for every  $x \in \{0, 1\}^*$ ,

$$l(f(x)) \leq l(x)^k + k \quad \text{and} \quad l(x) \leq l(f(x))^k + k. \quad (7.8)$$

That is,  $f$  neither shrinks nor stretches  $x$  more than polynomially in its length. The following result is used to prove Theorems 7.3.1 and 7.3.2.

**Lemma 7.3.1 (Honesty Lemma)** *Let  $f$  be an honest function computable in polynomial time. For all  $t \geq 1$ , and all but finitely many  $n$ ,*

$$f(C[\log \log n, n^t, \infty]) \subseteq C[\log n, n^{\log \log n}, \infty].$$

**Proof.** Let  $n = l(x)$ . By assumption, there is a  $k$  such that Equation 7.8 holds. Let machine  $T$  compute  $f$  and run in time polynomial in  $n$ . For each  $x$  in  $C[\log \log n, n^t, \infty]$  there exists a  $y$ , with  $l(y) \leq \log \log n$ , such that some  $T'$  computes  $x$  from  $y$  and runs in time polynomial in  $n$ . So  $f(x)$  is computable from  $y$ ,  $T$ , and  $T'$ . This description has fewer than  $\log n$  bits for large  $n$ . Moreover, the computation takes polynomial time, which is less than  $n^{\log \log n}$  for large  $n$ . So, for all but finitely many  $n$ ,

$$f(C[\log \log n, n^t, \infty]) \subseteq C[\log n, n^{\log \log n}, \infty].$$

□

We present two examples applying this lemma. We first construct the so-called Baker-Gill-Solovay oracle.

**Theorem 7.3.1** *There is a recursive oracle  $A$  such that  $P^A \neq NP^A$ .*

**Proof.** By diagonalization. Define  $f$  inductively by  $f(1) = 2$  and  $f(k) = 2^{f(k-1)}$  for  $k > 1$ . Choose  $B \subseteq \{1^{f(k)} : k \geq 1\}$  with  $B \in \text{DTIME}[n^{\log n}] - P$ . Use  $B$  to construct  $A$  as follows: For every  $k$  such that  $1^{f(k)} \in B$ , put the first string of length  $f(k)$  from

$$C[\log n, n^{\log n}, \infty] - C[\log n, n^{\log \log n}, \infty] \quad (7.9)$$

in  $A$ . The set in Equation 7.9 is nonempty by Exercise 7.1.13. Clearly,  $A$  is recursive and  $B \in NP^A$ .

Suppose we present a polynomial time bounded Turing machine with an input  $1^{f(k)}$ . This input has length  $n = f(k)$ . The machine can only compute strings not in  $A$  or of length less than  $\log n$  in  $A$  by Lemma 7.3.1. Therefore, a  $P^A$ -oracle machine can in polynomial time only query about strings not in  $A$  or of length less than  $\log n$ . The former strings are not in  $A$  anyway. The membership of the latter strings in  $A$  can be determined in polynomial time. Therefore, the oracle is useless. Consequently, if  $B$  is in  $P^A$ , then also  $B \in P$ . But this contradicts our assumption that  $B \notin P$ .  $\square$

In recursion theory, all recursively enumerable-complete sets are recursively isomorphic. We define a similar concept in complexity theory.

**Definition 7.3.2** Two sets are said to be *P-isomorphic* if there is a polynomial time computable bijection between the two sets.

The Berman-Hartmanis conjecture states that all NP-complete sets are P-isomorphic. At this time of writing, all known natural NP-complete sets are P-isomorphic. The following theorem attempts to tackle the conjecture. The SAT decision problem was defined in Definition 1.7.8 and Definition 5.4.3.

**Theorem 7.3.2** *If there exists a set  $L \in P$  such that  $L \subset SAT$  and*

$$C[\log n, n^{\log n}, \infty] \cap SAT \subseteq L,$$

*then  $SAT - L$  is an NP-complete set that is not P-isomorphic to SAT.*

**Proof.** To see that  $SAT - L$  is  $\leq_m^P$ -complete for NP we reduce SAT to  $SAT - L$  by reducing all elements in  $L$  to a fixed element in  $SAT - L$  and all other strings to themselves. This is possible since  $L \in P$ .

At the same time, SAT and  $SAT - L$  cannot be P-isomorphic since an isomorphism  $h$  is an honest function in P and therefore cannot map the

simple strings in  $L$  onto more complex strings in  $\text{SAT} - L$ . That is, by the Honesty Lemma 7.3.1, for all but finitely many  $n$ ,

$$h(\text{SAT} \cap C[\log \log n, n^3, \infty]) \subseteq C[\log n, n^{\log n}, \infty].$$

Since  $\text{SAT} \cap C[\log \log n, n^3, \infty] \neq \emptyset$ , it follows that

$$h(\text{SAT} \cap C[\log \log n, n^3, \infty]) \not\subseteq \text{SAT} - L.$$

That is,  $\text{SAT}$  is not  $\text{P}$ -isomorphic to  $\text{SAT} - L$ .  $\square$

It turns out that resource-bounded Kolmogorov complexity provides rigor to intuition for oracle constructions. We give another slightly more involved example.

**Definition 7.3.3** A set  $A$  is *exponentially low* if  $\text{E}^A = \text{E}$ , where  $\text{E} = \bigcup_{c \in \mathcal{N}} \text{DTIME}[2^{cn}]$ . Obviously, for any  $A \in \text{P}$ , the set  $A$  is exponentially low.

**Theorem 7.3.3** *There is an exponentially low sparse set  $A$  that is not in  $\text{P}$ .*

**Proof.** Let  $B = C[n/2, 2^{3n}, \infty]$ . The set  $\overline{B}$  is the complement of  $B$ . Let  $A = \{x : x \text{ is a lexicographically least element of } \overline{B} \text{ of length } 2^{2^{m-2}} \text{ (tower of } m \text{ stacked } 2\text{'s), } m > 0\}$ . Here  $m$  is not meant to be constant. Trivially,  $A \in \text{E}$ . Furthermore,  $A$  is not in  $\text{P}$ . To prove this, assume that  $T$  accepts  $A$  in polynomial time. Then we can simulate  $T$  to find an  $x$  of length  $n$  such that  $n > 2(l(T) + \log n)$  and  $x \in A$ , all in less than  $2^{3n}$  time using fewer than  $l(T) + \log n$  bits. Hence,  $x \in B$ , a contradiction.

We also need to show that  $\text{E}^A = \text{E}$ . We do this by simulating an  $\text{E}^A$  oracle machine  $T^A$  by an  $\text{E}$  machine. Let the machine  $T^A$  run in  $2^{cn}$  time. An idea from the proof of Theorem 7.2.6 is useful:

**Claim 7.3.1** For constant  $c' > 3c + 3$ , machine  $T^A$  cannot ask any query “ $y \in A?$ ” with  $l(y) \geq c'n$ , for  $y$  is really in  $A$ .

**Proof.** Assume, by way of contradiction, that  $T^A$  can do so. Suppose further that  $y$  is the first string in  $A$  of length greater than  $c'n$  queried by  $T^A$  on input  $x$  of length  $n$ . We will show that  $y \in B$ , and hence  $y \notin A \subseteq \overline{B}$ , a contradiction. In order to show  $y \in B$ , we use the fact that  $A$  is so sparse that we can actually describe the list of all strings in  $A$  of lengths less than  $l(y)$  in fewer than  $2 \log(c'n)$  bits.

Assume that  $T^A$  with input  $x$  of length  $n$  queries  $y$  at step  $t < 2^{cn}$ . We can now reconstruct  $y$  by simulating  $T^A$  on input  $x$ , and stop  $T^A$  at time  $t$ . If  $T^A$  queries  $A$  before step  $t$ , by assumption on the computation time to produce a string in  $A$ , either the query is shorter than  $l(y)$ , or the

queried string is not in  $A$ . For short queries, we can supply the answer from an exhaustive description of  $A^{< c'n}$  of fewer than  $2 \log(c'n)$  bits. At time  $t$ , we can recover  $y$  from the query tape. This way we can describe  $y$  by the following items:

- $O(1)$  bits to describe this discussion and  $T^A$ ;
- $n$  bits to describe  $x$ ;
- $2 \log(c'n)$  bits to describe  $A^{< c'n}$ ;
- $cn$  bits to describe  $t$ .

In total, we use fewer than  $c'n/2 \leq l(y)/2$  bits. The time required for the simulation is at most  $2^{cn} < 2^{c'n} \leq 2^{l(y)}$ . Hence,

$$y \in B = C[l(y)/2, 2^{3l(y)}, \infty].$$

This contradicts the assumption on  $A$ . Hence,  $T^A$  queries no string in  $A$  of length greater or equal to  $c'n$ .  $\square$

Now we can simulate  $T^A$  without using an oracle. Whenever we meet an oracle query about a string that is longer than  $c'n$ , we answer “no” for  $A$ . Whenever we meet an oracle query about a string  $y$  such that  $l(y) \leq c'n$ , we simply perform an exhaustive search to decide whether  $y \in A$  using exponential time. This way, we simulate an  $E^A$  machine by an  $E$  machine without using oracle  $A$ . Hence,  $E^A = E$ .  $\square$

### 7.3.2 P-Printability

We want to characterize the sets  $C[k \log n, n^k, \infty]$ , for constant  $k$ . These sets are said to have small time-bounded Kolmogorov complexity.

**Definition 7.3.4** A set  $L$  is *polynomial-time printable* (P-printable) if for some integer  $k$  all the elements of  $L$  up to size  $n$  can be printed by a Turing machine in time  $n^k + k$ . Clearly every P-printable set is a sparse set in P. Let  $y \leq x$  denote that  $y$  precedes  $x$  in the length-increasing lexicographical order. The ranking function, as defined in Definition 7.2.1, for a language  $L$ , denoted by  $r_L$ , is defined as  $r_L(x) = d(\{y \in L : y \leq x\})$ . A *tally set* is a set over a one-letter alphabet.

**Theorem 7.3.4** Let  $L \subseteq \{0, 1\}^*$ . The following statements are equivalent:

- (i)  $L$  is P-printable;
- (ii)  $L$  is sparse and has a polynomial time computable ranking function;
- (iii)  $L$  is P-isomorphic to some tally set in P; and
- (iv)  $L \subseteq C[k \log n, n^k, \infty]$  for some constant  $k$  and  $L \in P$ .

**Proof.** (i)  $\rightarrow$  (ii) is immediate.

(ii)  $\rightarrow$  (iii). Let  $L$  have a polynomial time computable ranking function  $r_1$  and  $d(L^{\leq n}) < p(n)$  for some (nondecreasing) polynomial  $p(n)$ . Then  $r_2(x) = 1x - r_1(x)$  is a ranking function for the complement of  $L$ , where  $1x$  is treated as a binary number. The set

$$T = \{0^{np(n)+i} : r_1(1^{n-1}) < i \leq r_1(1^n)\}$$

is a tally set in  $P$ . Let  $r_3$  be a ranking function for  $\{0\}^* - T$ . By Exercise 7.2.6, all such ranking functions have inverses that are computable in time polynomial in the length of their output. It is now easy to see that the function that maps  $x$  of length  $n$  to  $0^{np(n)+r_1(x)}$  if  $x \in L$ , and to  $r_3^{-1}(r_2(x))$  if  $x \notin L$ , is a  $P$ -isomorphism that maps  $L$  one-to-one onto  $T$ .

(iii)  $\rightarrow$  (iv). By assumption, there is a  $P$ -isomorphism  $f$ , both  $f$  and  $f^{-1}$  computable in time  $n^c$  for some  $c$ , that maps  $L$  one-to-one onto a tally set  $T \subseteq \{0\}^*$ , where  $T \in P$ . Trivially,  $L \in P$ .

Since  $f$  is computable in time  $n^c$ , we have  $l(f(x)) \leq n^c$  for  $n = l(x)$ . Hence, the binary representation of  $f(x)$  has length at most  $c \log n$ . Then,  $x$  can be represented by  $f(x)$  in binary using only  $c \log n$  bits. To compute  $x$  from  $f(x)$ , we simply compute  $f^{-1}(0^{f(x)})$ , which takes polynomial time by assumption. Thus,  $L \subseteq C[k \log n, n^k, \infty]$  for some constant  $k$ .

(iv)  $\rightarrow$  (i). Assume that  $L \in P$  and that for some  $k$

$$L \subseteq C_U[k \log n, n^k, \infty].$$

We show how to print elements of  $L$  up to size  $n$  in polynomial time. Given  $n$ , simulate the reference universal machine  $U$ , using each string up to length  $k \log n$  as input, for at most  $n^k$  steps. For each output  $x$ , print  $x$  if the computation halted in  $l^k(x)$  steps,  $l(x) = n$ , and  $x \in L$ . Clearly, this can be done in time polynomial in  $n$ .  $\square$

**Corollary 7.3.1** Let  $L \subseteq \{0, 1\}^*$ . Then,  $L \subseteq C[k \log n, n^k, \infty]$  for some constant  $k$  if and only if  $A$  is  $P$ -isomorphic to a tally set.

## Exercises

**7.3.1.** [28/O43] Let SAT be the set of satisfiable Boolean formulas. By Definition 7.2.2, a set  $A$  is sparse if there is a constant  $c$  such that for all  $n$  we have  $d(A^{\leq n}) \leq n^c + c$ . In Section 1.7.4 we defined that a set  $B$  is *polynomial time Turing reducible* to set  $C$ , denoted by  $B \leq_T^P C$ , if there is a polynomial time oracle Turing machine that accepts  $B$ , using oracle  $C$ . One might suspect that formulas with low Kolmogorov complexity are easy to solve. Prove the following:

- (a) If  $A \in \text{NP}$  and  $A$  is sparse, then  $A \leq_T^P C[\log n, n^2, \infty] \cap \text{SAT}$ . Thus,  $\text{SAT} \cap C[\log n, n^2, \infty]$  is a complete set for all other sparse sets in NP under polynomial Turing reduction.
- (b)  $\text{SAT} \cap C[\log n, n^2, \infty] \in \text{P}$  iff there are no sparse sets in  $\text{NP} - \text{P}$ .
- (c) (Open)  $\text{SAT} \cap C[\log n, n^2, \infty] \in \text{P}?$

*Comments.* Source: this exercise and the next seven exercises are taken from J. Hartmanis, *Proc. 24th IEEE Found. Comput. Sci.*, 1983, pp. 439–445.

**7.3.2.** [29] Let  $g(n) \leq n$  be any unbounded, monotonically increasing function and let  $G(n)$  be such that for every  $k$ ,

$$\lim_{n \rightarrow \infty} n^k / G(n) = 0.$$

Show that  $C[g(n), G(n), \infty] \cap \text{SAT} \subseteq A_0 \subset \text{SAT}$  and  $A_0 \in \text{P}$  implies that SAT is not P-isomorphic to  $\text{SAT} - A_0$ . Also show that  $\text{SAT} - A_0$  is NP-complete.

**7.3.3.** [26] If  $C[\log n, n^{\log n}, \infty] \cap \text{SAT} \subseteq A_0 \subseteq \text{SAT}$  and  $A_0 \in \text{P}$ , then  $\text{E} = \text{NE}$ .

**7.3.4.** [40] Show the following:

- (a) There is a sparse set in  $\text{NP} - \text{P}$  iff  $\text{NE} \neq \text{E}$ .
- (b) Define  $\Delta_2^{\text{E}}$  analogous to  $\Delta_2^{\text{P}}$ . If  $\text{NE} = \Delta_2^{\text{E}}$ , and every sparse set in NP is polynomial time many-to-one reducible (see page 38 in Section 1.7.4) to  $\text{SAT} \cap C[\log n, n^2, \infty]$ , then  $\text{NE} = \text{E}$ . Therefore, all sparse sets in NP are in P.

*Comments.* For Item (b), use the Berman-Mahaney tree labeling technique [S. Mahaney, *J. Comput. System Sci.*, 25(1982), 130–143].

**7.3.5.** [30] Show that  $\text{P} = \text{NP}$  iff for all oracles  $A \subseteq C[\log n, \infty, n^2]$ , we have  $\text{P}^A = \text{NP}^A$ .

**7.3.6.** [29] Use Kolmogorov complexity to show that there exists a recursive oracle  $A$  such that  $\text{NP}^A$  has  $\text{P}^A$ -immune sets.

**7.3.7.** [33] Show that the set  $\{0,1\}^* - C[\log n, \infty, 2^n]$  is  $\text{SPACE}[2^{cn}]$ -immune for any  $c < 1$ .

*Comments.* Compare this exercise with Theorem 2.7.1.

**7.3.8.** [34] We construct a sparse random oracle set  $A$  as follows: For every  $n$ ,  $n = 1, 2, \dots$ , toss a fair coin. If the result is “tails,” then we do not include any string of length  $n$  in  $A$ ; if the result is “heads,” then we toss the coin  $n$  times and place the resulting binary string (the  $i$ th bit is 1 iff the  $i$ th toss is heads) in  $A$ . Prove that  $\Pr(\text{NP}^A \neq \text{P}^A) = 1$ .

**7.3.9.** [32] The method using generalized Kolmogorov complexity to construct oracles in Section 7.3.1 can be used to obtain many more oracles. Define

$$\text{EXPTIME} = \bigcup_{c \in \mathcal{N}} \text{DTIME}[2^{n^c}].$$

Notice that EXPTIME is different from class E. Let NSPACE stand for the nondeterministic version of PSPACE. Let us consider oracle Turing machines with an unbounded oracle tape. Prove the following:

- (a) There is an oracle  $A$  such that  $\text{NSPACE}^A \not\subset \text{EXPTIME}^A$ .
- (b) There is an oracle  $B$  such that  $\text{PSPACE}^B \subset \text{EXPTIME}^B$ , where the containment is proper.
- (c) There is an oracle  $C$  such that  $\text{EXPTIME}^C \not\subset \text{NSPACE}^C$ .
- (d) There is an oracle  $D$  with  $\text{PSPACE}^D \subset \text{NSPACE}^D = \text{EXPTIME}^D$ , where the containment is proper.

*Comments.* Item (d) implies that Savitch's Theorem [W.J. Savitch, *J. Comput. System Sci.* 4:2(1972), 177–192] does not relativize with unbounded oracle tape. Source: R. Gavaldà, L. Torenvliet, O. Watanabe, and J.L. Balcázar, *Proc. 15th Math. Found. Comput. Sci. Conf.*, 1991, pp. 269–276. For a comprehensive study in this direction, see [R. Gavaldà, Ph.D. Thesis, Universitat Politècnica de Catalunya, 1992].

**7.3.10.** [39] Show that if  $A$  is a set whose characteristic sequence is a random infinite binary sequence in the sense of Martin-Löf (see Section 2.5), then  $\text{P}^A \neq \text{NP}^A$ .

*Comments.* This result is presented in a more general setting using results of Section 2.5 by R.V. Book, J.H. Lutz, and K.W. Wagner in [*Math. Systems Theory*, 27(1994), 201–209]. Such a set  $A$  can be used to establish other known probability one oracle separations such as  $\text{PH}^A \neq \text{PSPACE}^A$ , where PH is the polynomial hierarchy. See also [R.V. Book and O. Watanabe, *Inform. Comput.*, 125(1996), 70–76].

**7.3.11.** [33] Show that for all  $t \geq 2$ , the set  $C[t \log n, n^t, \infty]$  is P-isomorphic to  $\{0\}^*$ .

*Comments.* See [E. Allender, O. Watanabe, *Inform. Comput.*, 86(1990), 160–178]. In this paper, the authors also use sets  $C[t \log n, n^t, \infty]$  to study the equivalent classes of tally sets under various types of reductions. See the papers [S. Tang and R.V. Book, *Proc. 15th Int. Colloq. Automata, Languages, Prog., Lect. Notes Comp. Sci.*, vol. 317, Springer-Verlag, 1988, pp. 519–599; R. Gavaldà and O. Watanabe, *SIAM J. Comput.*, 22(6) (1993), 1257–1275; H. Buhrman, E. Hemaspaandra, and L.

Longpré, *Proc. 8th IEEE Conference on Structure in Complexity Theory*, 1993, pp. 208–214].

**7.3.12.** [32] We say that  $A$  is *truth-table* reducible to  $B$  if there are some functions  $g_1, \dots, g_m$  and a Boolean function  $f$  where  $y_i$  is true iff  $g_i(a) \in B$ , and  $f(y_1, \dots, y_m)$  is true iff  $a \in A$ . We consider only polynomial time truth-table reductions where the  $f$  and  $g_i$ 's are computable in polynomial time. It is clear that if  $A$  many-to-one reduces to  $B$ , then  $A$  also truth-table reduces to  $B$ ; and if  $A$  truth-table reduces to  $B$ , then  $A$  also Turing reduces to  $B$ . Use time-space-bounded Kolmogorov complexity to construct a set  $D$  that is complete for  $E$  under polynomial Turing reduction but not complete for  $E$  under polynomial truth-table reduction.

*Comments.* Source: O. Watanabe, *Theoret. Comput. Sci.*, 54(1987), 249–265. Improved in [B. Fu, *SIAM J. Comput.*, 24:5(1995), 1082–1090].

## 7.4 Instance Complexity

In computational complexity, it is traditional to study the intractability of a decision problem by treating a set collectively. With Kolmogorov complexity, it is possible to also study the complexity of individuals in a set, and this is called *instance complexity*.

Consider a partial program  $p$  that gives answers 1 (accept), 0 (reject), or  $\perp$  (don't know). For set  $A$  we will write  $A(x)$  as the characteristic function of  $A$ , with  $A(x) = 1$  iff  $x \in A$ . We say that function  $p$  is *consistent* with a set  $A$  if  $p(x) = A(x)$  when  $p(x) \neq \perp$ . The computing time of machine  $T$  on input  $y$ , using program  $p$ , is denoted as  $\text{time}_T(p, y)$ . In the definition below,  $A$  is an arbitrary set, possibly nonrecursively enumerable, and  $t$  is an arbitrary function, possibly nonrecursive.

**Definition 7.4.1** Let  $T$  be a Turing machine. Given a set  $A$  and time bound  $t$ , define the ( $t$ -bounded) *instance complexity* of  $x$  with respect to  $T$  and  $A$  as

$$\text{ic}_T^t(x : A) = \min\{l(p) : T(p, x) \neq \perp, \text{ and } \forall y \text{ } T(p, y) \neq \perp \text{ implies } \text{time}_T(p, y) \leq t(l(y)), T(p, y) = A(y)\},$$

and is  $\infty$  if no such  $p$  exists.

By this definition, also  $x \notin A$  can have instance complexity with respect to  $A$ . Intuitively, the instance complexity of a string  $x$ , with respect to a set  $A$ , measures the length of the shortest program that correctly decides whether  $x$  is in  $A$ . This program doesn't need to decide about other strings as long as it does not contradict  $A$  when it makes a decision. The goal is to identify the “hard” instances that make a language “hard.” We state an invariance theorem for instance complexity.

**Theorem 7.4.1 (Invariance for Instance Complexity)** *There exists a universal Turing machine  $U$  such that for each Turing machine  $T$  there is a constant  $c$  such that for all  $A$  and  $t$  and  $x$ ,*

$$\text{ic}_U^{t'}(x : A) \leq \text{ic}_T^t(x : A) + c,$$

where  $t'(n) = ct(n) \log t(n) + c$ .

The proof for this invariance theorem is similar to that for time-bounded Kolmogorov complexity, and is left to the reader. Using the above invariance theorem, we fix a reference universal machine and drop the index  $U$  in  $\text{ic}_U^t(x : A)$ .

**Example 7.4.1** The relationship between time-bounded Kolmogorov complexity and instance complexity is a fundamental question. Obviously,  $CD^t$  is a special case of  $\text{ic}^t$  complexity because

$$CD^t(x) = \text{ic}^t(x : \{x\}).$$

If  $c$  is a large enough constant and  $t'(n) = ct(n) \log t(n) + c$ , it is not difficult to prove the following for any set  $A$  (see Exercise 7.4.1):

$$\begin{aligned}\text{ic}^{t'}(x : A) &\leq C^t(x) + c, \text{ and} \\ \text{ic}^{t'}(x : A) &\leq CD^t(x) + c.\end{aligned}$$

It is clear that for some set  $A$ , such as  $\{0, 1\}^*$ , the equalities do not hold for the above two inequalities. This leads us to the following *instance complexity conjecture*: Let  $t(n)$  be computable in time  $t(n)$ ,  $t'$  as above, and  $A$  recursive. If  $A \notin \text{DTIME}[t(n)]$ , then there is a constant  $c$  and infinitely many  $x$  such that

$$\text{ic}^t(x : A) \geq C^{t'}(x) - c.$$

Exercises 7.4.5, 7.4.7, and 7.4.8, contain solutions to various special cases of this conjecture.  $\diamond$

Instance complexity provides pleasant and simple characterizations for many fundamental complexity-theoretic properties.

**Lemma 7.4.1** *A set  $A$  is in  $\text{P}$  if and only if there exist a polynomial  $t$  and a constant  $c$  such that for all  $x$ , we have  $\text{ic}^t(x : A) \leq c$ .*

**Proof. (ONLY IF)** If  $A$  is in  $\text{P}$ , then for every  $x$  the polynomial time program that decides  $A$  is trivially a consistent program for  $A$  and decides whether or not  $x \in A$ . This gives constant instance complexity for every string  $x$ .

(IF) By assumption, for some  $c$ , for all  $x$ , we have  $\text{ic}^t(x : A) \leq c$ . Let  $B$  be the set of all programs, consistent with  $A$ , of size at most  $c$ , with running time bounded by  $t(n)$ . Decide whether or not  $x \in A$ , by simulating all the programs in  $B$  and accept  $x$  if and only if some program in  $B$  accepts  $x$ .  $\square$

**Definition 7.4.2** Let  $A$  be a recursive set. An infinite set  $C$  (not necessarily a subset of  $A$ ) is called a *polynomial complexity core* for  $A$  if for every total program  $p$  that decides  $A$  and polynomial  $t$ , we have  $\text{time}_p(x) > t(l(x))$  for all but finitely many  $x$  in  $C$ .

**Lemma 7.4.2** *A set  $C$  is a polynomial complexity core for  $A$  if and only if for every polynomial  $t$  and constant  $c$  we have  $\text{ic}^t(x : A) > c$  for all but finitely many  $x$  in  $C$ .*

**Proof.** (IF) Assume that there are infinitely many  $x$  in  $C$  such that  $\text{ic}^t(x : A) \leq c$  for some polynomial  $t$  and constant  $c$ . Let  $B$  be the set of programs, running in time  $t$ , that are consistent with  $A$  and of length less than  $c$ . Then there are infinitely many elements  $x$  in  $C$  that we can accept in polynomial time by simulating all programs in  $B$  simultaneously in dovetailing style. Thus,  $C$  cannot be a polynomial complexity core for  $A$ .

(ONLY IF) Assume that  $C$  is not a polynomial complexity core for  $A$ . Then there is a total program  $p$  for  $A$  and a polynomial  $t$  such that for infinitely many  $x$  in  $C$  we have  $\text{time}_p(x) \leq t(l(x))$ . The program  $p$  is consistent with  $A$  and it accepts infinitely many  $x$  in time  $t(l(x))$ . We modify  $p$  by adding a step counter of  $t(n)$  to it. When the computation exceeds  $t(n)$  steps, we halt  $p$ . The modified  $p$  is consistent with  $A$ , runs in polynomial time, and decides correctly infinitely many  $x$ . Thus, for infinitely many  $x$ , we have  $\text{ic}^t(x : A) \leq c$ .  $\square$

Among the strings of various instance complexities, one class is particularly interesting. This is the class of strings of logarithmic instance complexity computable in polynomial time.

**Definition 7.4.3**  $\text{IC}[\log, \text{poly}]$  is the class of sets  $A$  for which there exists a polynomial  $t$  and a constant  $c$  such that for all  $x$ ,  $\text{ic}^t(x : A) \leq c \log l(x) + c$ .

The SAT decision problem, Definition 1.7.8, page 39, is the following: Given a Boolean formula  $\phi(x_1, \dots, x_k)$  in conjunctive normal form of binary description length  $n$ , we want to decide whether or not there is a truth assignment to the variables  $x_1, \dots, x_k$  that makes the formula true. We always assume that  $k$  and  $n$  are polynomially related. The set SAT is the set of  $\phi$ 's for which there is a truth assignment that makes  $\phi$  true.

**Theorem 7.4.2** *If  $\text{SAT} \in \text{IC}[\log, \text{poly}]$ , then  $\text{NP} = \text{P}$ .*

**Proof.** Assume that  $\text{SAT} \in \text{IC}[\log, \text{poly}]$ . Then for some polynomial  $t$  and constant  $c$ , we have

$$\text{ic}^t(\phi : \text{SAT}) \leq c \log l(\phi) + c, \quad (7.10)$$

for all  $\phi$ . We treat each input as a disjunctive normal form  $\phi$ . If the input is of wrong format, then it can be rejected right away. Thus, for each  $\phi$ , there is a program  $p$  of length  $c \log l(\phi) + c$  such that  $p$  decides whether  $\phi \in \text{SAT}$  correctly and  $p$  is consistent with SAT for all other inputs. We show that if Equation 7.10 is true, then SAT can be decided in polynomial time.

Given input  $\phi(x_1, \dots, x_k)$ , let  $l(\phi(x_1, \dots, x_k)) = n$ . First, set  $A := \{p : l(p) \leq c \log n + c\}$ . Clearly,  $d(A)$  is of polynomial size. By Equation 7.10, for each input of length less than or equal to  $n$ , there is a program in  $A$  running in polynomial time  $t(n)$  that decides correctly whether this particular input is in SAT, and on other inputs it makes decisions consistent with SAT. In particular, there is a program  $p_0$  running in polynomial time  $t(n)$  that decides correctly whether or not  $\phi(x_1, \dots, x_k) \in \text{SAT}$ . Moreover,  $p_0$  is consistent with SAT on all other inputs. The program  $p_0$  may just output  $\perp$  signs for the latter two inputs.

In the procedure below we will, in a polynomial number of steps, either decide  $\phi(x_1, \dots, x_k) \in \text{SAT}$  correctly or delete at least one program  $p$ , with  $p \neq p_0$  and  $p$  inconsistent with SAT, from  $A$ . Repeating this procedure  $d(A) - 1$  times, we will either decide  $\phi(x_1, \dots, x_k)$  correctly or find  $p_0$  as the last remaining element by elimination. Since  $d(A)$  is polynomial, the entire process takes polynomially many steps. The procedure is as follows:

**Step 1** For all  $p \in A$ , simulate  $p$  for  $t(n)$  steps on input  $\phi(x_1, \dots, x_k)$ .

- If no program  $p$  rejects  $\phi(x_1, \dots, x_k)$  {implying that  $p_0$  has accepted  $\phi(x_1, \dots, x_k)$ ; some programs may output  $\perp$ } then we also accept  $\phi(x_1, \dots, x_k)$  and **exit** the procedure.
- If no  $p$  accepts  $\phi(x_1, \dots, x_k)$  {implying that  $p_0$  must have rejected  $\phi(x_1, \dots, x_k)$ } then reject  $\phi(x_1, \dots, x_k)$  and **exit** the procedure.

**Step 2** If some program in  $A$  accepts  $\phi(x_1, \dots, x_k)$  and some other program in  $A$  rejects  $\phi(x_1, \dots, x_k)$  then **for**  $i := 1, \dots, k$  **do**

Suppose the binary values  $b_1, \dots, b_{i-1}$  are determined in the previous loops and  $p \in A$  accepts  $\phi = \phi(b_1, \dots, b_{i-1}, x_i, \dots, x_k)$ .

Simulate all programs in  $A$  with inputs

$$\begin{aligned}\phi_0 &= \phi(b_1, \dots, b_{i-1}, 0, x_{i+1}, \dots, x_k), \\ \phi_1 &= \phi(b_1, \dots, b_{i-1}, 1, x_{i+1}, \dots, x_k).\end{aligned}$$

{Formulas such as  $\phi$ ,  $\phi_0$ , and  $\phi_1$  are all different instances.}

If neither of the above inputs is accepted by some program in  $A$  {this means  $\phi(b_1, \dots, b_{i-1}, x_i, \dots, x_k) \notin \text{SAT}$ , because there are programs  $q_0, q_1 \in A$  that are consistent with SAT and decide correctly whether  $\phi_0, \phi_1 \in \text{SAT}$ ; consequently  $p$  is not consistent with SAT and  $p \neq p_0$ } then delete  $p$  from  $A$  and **exit** the procedure.

If one of the inputs, say with  $x_i = 0$ , is accepted by some program in  $A$  then set  $b_i := 0$ .

The process either exits in Step 1 deciding  $\phi$ , or it exits in Step 2 with a program  $p$  inconsistent with SAT deleted from  $A$ , or it ends with a truth assignment  $(b_1, \dots, b_k)$ . In the latter case, we can check (in polynomial time) whether the truth assignment  $(b_1, \dots, b_k)$  really satisfies  $\phi(x_1, \dots, x_k)$ . If it does, then we accept  $\phi(x_1, \dots, x_k)$ . Otherwise, we delete  $p$  from  $A$  since it is inconsistent with SAT. (This way, we never delete programs consistent with SAT.) Therefore, in Step 2, we either finish with a satisfying assignment or by deleting an inconsistent program from  $A$ .

Repeating the procedure for at most  $d(A) - 1$  times we are guaranteed to decide  $\phi$ , or to find a satisfying assignment, or we have only  $p_0$  left in  $A$  at the last step.  $\square$

## Exercises

**7.4.1.** [25] Let  $t(n)$  be computable in time  $t(n)$ ,  $t'(n) = ct(n) \log t(n) + c$ , for some constant  $c$ . For any set  $A$ , string  $x$ , and some constant  $c$ , prove:

- (a)  $\text{ic}^{t'}(x : A) \leq C^t(x) + c$ , and
- (b)  $\text{ic}^{t'}(x : A) \leq CD^t(x) + c$ .

*Comments.* Item (a) and the next three exercises are from [P. Orponen, K. Ko, U. Schöning, and O. Watanabe, *J. ACM*, 41(1994), 96–121]. Item (b) was suggested by L. Fortnow.

**7.4.2.** [30] The proof of Theorem 7.4.2 depends on the so-called “self-reducibility” of the SAT problem. A set is *self-reducible* if the membership question for any element can be reduced in polynomial time to the membership question for a number of *shorter* elements. For example,

SAT is self-reducible since an arbitrary Boolean formula  $\phi(x_1, x_2, \dots, x_n)$  is satisfiable if and only if at least one of the two shorter Boolean formulas  $\phi(0, x_2, \dots, x_n)$  or  $\phi(1, x_2, \dots, x_n)$  is satisfiable. Prove the following generalization of Theorem 7.4.2. If a set is self-reducible and it is in  $\text{IC}[\log, \text{poly}]$ , then it is also in P.

*Comments.* See the survey by D. Joseph and P. Young in [*Complexity Theory Retrospective*, A. Selman (Ed.), Springer-Verlag, 1990, pp. 82–107] for more information on self-reducibility.

**7.4.3.** [32] The class P/poly is defined in Exercise 7.2.5, page 485. Use that exercise to define the class  $P/\log = \bigcup_{c>0} P/c\log n$  analogously. Show that  $P/\log$  is properly contained in  $\text{IC}[\log, \text{poly}]$ , which is properly contained in P/poly.

**7.4.4.** [31] There are sets with hard instance complexity everywhere. In particular, prove that there is a set  $A$  computable in  $2^{O(n)}$  time such that for some constant  $c$  and for all  $x$ ,

$$\text{ic}^{\text{exp}}(x : A) \geq C^{\text{exp}'}(x) - 2 \log C^{\text{exp}'}(x) - c,$$

where  $\text{exp}(n) = 2^n$  and  $\text{exp}'(n) = O(n2^{2n})$ .

**7.4.5.** [38/O43] In light of Exercise 7.4.4 and Lemma 7.4.2, let's say that a set  $A$  has *p-hard instances* if for every polynomial  $t$  there exists a polynomial  $t'$  and a constant  $c$  such that for infinitely many  $x$   $\text{ic}^t(x : A) \geq C^{t'} - c$ .

(a) (Open) Prove or disprove: every recursive set  $A \notin \text{P}$  has p-hard instances. This may be regarded as a polynomial version of the instance complexity conjecture in Example 7.4.1.

(b) Every recursive tally set  $A \notin \text{P}$  has p-hard instances.

(c) Prove the following claim: Let a recursive set  $A$  be NP-hard with respect to the polynomial-time Turing reduction in which the length of any query is not shorter than a fixed polynomial of the length of the input. Then  $A$  has p-hard instances unless  $A \in \text{P}$ . In particular, this holds for all natural NP-hard problems.

(d) (Open) We can also state a CD version of Item (a). Prove or disprove: for every recursive set  $A \notin \text{P}$  and every polynomial  $t$  there is a polynomial  $t'$  and a constant  $c$  such that for infinitely many  $x$ ,  $\text{ic}^t(x : A) \geq CD^{t'} - c$ .

*Comments.* Source: L. Fortnow and M. Kummer, *Theoret. Comput. Sci. A*, 161(1996), 123–140. They have also shown that the instance complexity conjecture, Item (a), and Item (d) all fail relative to some oracles. Item (a) also holds relative to some oracle.

**7.4.6.** [35] Let  $t$  be a recursive time bound. There is a recursive set  $A$  such that  $f(x) = \text{ic}^t(x : A)$  is not recursive.

*Comments.* This result, due to L. Fortnow and M. Kummer [*Ibid.*], was originally conjectured by P. Orponen, K. Ko, U. Schöning, and O. Watanabe, in [*J. ACM*, 41(1994), 96–121].

**7.4.7.** [33] In Definition 7.4.1, when we allow  $t$  to be arbitrary finite time, we will remove  $t$  from  $\text{ic}^t$  and simply write  $\text{ic}$ .

(a) Let  $R = \{x : C(x) \geq l(x)\}$ . We know that  $R$  is infinite and that it contains at least one string of each length. Show that there is a constant  $c$  such that for every  $x$  in  $R$ ,  $\text{ic}(x : R) \geq C(x) - c$ . That is,  $R$  contains only “hard instances.”

(b) Strings with high Kolmogorov complexity are individually hard to recognize by bounded computations. We say that a set  $C$  is dense if there is an  $\epsilon$  such that  $d(C \cap \Sigma^n) \geq 2^{\epsilon n}$ . Let set  $A$  be complete for the class  $\bigcup_{c \geq 0} \text{DTIME}[2^{cn}]$  with respect to polynomial time reductions. Show that there is a dense subset  $C \subseteq A$  such that for every nondecreasing polynomial  $t(n) \geq n \log n$ , for each  $x \in C$ ,  $\text{ic}^t(x : A) \geq C^t(x) - c$ , for some constant  $c$ .

(c) For every recursively enumerable set  $A$  with an infinite complement, there exists a constant  $c$  such that for infinitely many  $x \notin A$ ,  $\text{ic}(x : A) \geq C(x) - c$ .

(d) Show that there is a recursively enumerable set  $A$  with  $\text{ic}(x : A) \geq l(x)$  for infinitely many  $x \in A$ .

*Comments.* Hint for Item (a). Let  $L(p)$  be the set of strings accepted by  $p$  in the  $\text{ic}$  sense. Observe that there is a constant  $d$  such that for any total  $p$  for which  $L(p) \subseteq R$ , and for any  $x \in L(p)$ , we have  $l(x) \leq l(p) + d$ . (This is similar to Corollary 2.7.2, page 169.) Now the result follows easily using the definition of  $\text{ic}(\cdot)$ . Item (b) may be regarded as another special case of the instance complexity conjecture in Example 7.4.1. Items (a)–(c) are from [H.M. Buhrman and P. Orponen, *J. Comput. System Sci.*, 53:2(1996), 261–266]. Item (d) and an extensive study of related topics can be found in [M. Kummer, *SIAM J. Comput.*, 25:6(1996), 1123–1143].

**7.4.8.** [39] Use the definition of  $\text{ic}$  in Exercise 7.4.7.

(a) Show that if  $\text{ic}(x : A) \leq \log C(x) - 1$  for all but finitely many  $x$ , then  $A$  is recursive.

(b) There is a nonrecursive recursively enumerable set  $A$  and a constant  $c$  such that  $\text{ic}(x : A) \leq \log C(x) + c$  for all but finitely many  $x$ .

*Comments.* Item (b) resolves an open question in the first edition of this book: it refutes the unbounded version of the *instance complexity*

*conjecture* in Example 7.4.1. Originally proposed by P. Orponen, K. Ko, U. Schöning, and O. Watanabe [*J. ACM*, 41(1994), 96–121]. The solution is due to M. Kummer [*Ibid.*], where Item (a) is attributed to J. Tromp.

## 7.5 $Kt$ Complexity and Universal Optimal Search

---

It is meaningful to consider the *age* of strings. Loosely stated, the age of a string corresponds to the time we need to generate that string starting with a program of constant size. In this sense, the age of a random string  $x$  should be at least  $2^{l(x)}$ , considering the number of steps needed to generate  $x$  by enumerating all strings in length-increasing lexicographical order. This is also the expected time for a constant-size probabilistic program to generate  $x$  by fair coin flips. For random  $x$ ,  $\text{age}(x) = \Omega(2^{l(x)})$ . The case of nonrandom strings is more interesting. Generate programs length-increasing lexicographically and simulate them dovetailing style. This leads to a definition of the *age* as

$$\text{age}(x) = \min_p \{2^{l(p)} t : U(p) = x \text{ in } t \text{ steps}\}.$$

This way, *age* is dominated by the total time needed for a string to appear out of nothing, enumerated by a constant-size program. Taking the logarithm, we arrive at Levin's  $Kt$  complexity:  $Kt(x) = \log \text{age}(x)$ . Let us now define  $Kt$  complexity formally. We use the monotone machines of Definition 4.5.2, page 276, or the prefix machines of Example 3.1.1, page 193. The two models are completely equivalent for use in Definition 7.5.1. Such a machine is a Turing machine with a one-way read-only input tape, a one-way write-only output tape and a work tape. Initially, the input tape contains a one-way infinite binary sequence.

### Definition 7.5.1

Let the universal monotonic machine  $U$  scan an initial input segment  $p$  before it prints  $x$  (without necessarily halting after it does so). Let  $t(p, x)$  be the number of steps taken until  $x$  is printed. Then  $Kt$  is defined by

$$Kt_U(x) = \min_p \{l(p) + \log t(p, x)\}.$$

Since an invariance theorem such as Theorem 2.1.1, page 97, can be proved in the standard way, we will drop the subscript  $U$  and write  $Kt$ . For  $x \in \{0, 1\}^*$ , if  $x \in K[m - \log t, t, \infty]$ , with  $m$  minimal, then  $x$  has  $Kt$  complexity  $Kt(x) = m$ . Namely,  $x$  can be computed (and hence enumerated) by a self-delimiting program of length  $m - \log t$  in  $t$  steps.

### 7.5.1 Universal Optimal Search

There is a *universal search* method that will solve all problems of a certain class of inverting problems in time that is optimal but for a multiplicative constant. Despite its simplicity, the idea of universal optimal search is a powerful one. Let  $T_1, T_2, \dots$  be a standard enumeration of

prefix machines, and let  $\phi_1, \phi_2, \dots$  be the corresponding enumeration of partial recursive functions. If  $\phi$  is a partial recursive function and  $\phi(y) = x$ , then  $y$  is a  $\phi$ -witness for  $x$ .

**Definition 7.5.2** An algorithm  $A$  *inverts problem*  $\phi$  if given some  $x$  in the range of  $\phi$ , algorithm  $A$  computes a  $\phi$ -witness  $y$  for  $x$  and checks that  $\phi(y) = x$ . Algorithm  $A$  diverges outside the range of  $\phi$ .

**Example 7.5.1** Many computational problems consist in finding feasible algorithms to invert functions. Given a composite natural number, we are required to find a factorization. Once a splitting (partial factorization) is found, we check whether it is correct. A solution to this inversion problem does not solve the corresponding decision problem of whether or not a natural number is prime.

Given a satisfiable Boolean formula, we want to find a truth assignment that satisfies it, Definition 1.7.8. Once an assignment is found, we check whether it makes the formula true. A solution to this inversion problem does not imply a solution to the corresponding decision problem SAT of whether or not a given Boolean formula is satisfiable.

To decide whether or not we can node color a given graph on  $n$  nodes with  $k$  colors, such that no edge connects two nodes of the same color, is an NP-complete problem. We know we can color the graph with  $n$  colors. Hence, there is a least number of colors  $\chi$  ( $1 \leq \chi \leq n$ ) to color the graph. This  $\chi$  is the chromatic number of the graph. Suppose we have an algorithm that finds a witness  $k$ -coloring if there is one in time  $t(n)$ . We can run this algorithm for  $1 \leq k \leq n$  dovetail style. The dovtailed algorithm finds  $k_0$  together with a  $k_0$ -coloring in  $nt(n)$  steps. This solves the  $k$ -coloring problem for all  $k$ . The drawback is that we need to know  $t(n)$  to make the actual decision. A polynomial time solution to the inversion problem of finding a graph  $k$ -coloring, together with its running time  $t$ , also gives a polynomial time solution for the corresponding NP-complete decision problem.  $\diamond$

To solve inverting problems naively usually requires us to search through exponentially many candidate witnesses. But we can also take the universal optimal inversion algorithm.

**Lemma 7.5.1** *Let  $\phi$  be an inverting problem. If there is an algorithm  $A$  that inverts  $\phi$  in time  $t(n)$ , then the SIMPLE algorithm below inverts  $\phi$  in time  $ct(n)$ , where  $c$  is a constant depending only on  $A$ .*

**Proof.** We describe Algorithm SIMPLE. Run all machines  $T_i$  one step at a time according to the following scheme:  $T_1$  every second step,  $T_2$  every second step in the remaining unused steps,  $T_3$  every second step in the

remaining unused steps, and so on. That is, according to the sequence of indices,

$$1213121412131215121312141213121612\dots$$

If  $T_k$  inverts  $\phi$  on  $x$  in  $t$  steps, then this procedure will do the same in  $2^k t + 2^{k-1}$  steps. Choosing  $c = 2^{k+1}$  proves the lemma.  $\square$

A similar universal optimal search procedure  $A$  was developed by L.A. Levin using  $Kt$ . It has the advantage that the multiplicative constant  $c = 2^{k+1}$  in Lemma 7.5.1 is reduced to  $2^{K(T_k)+1}$ .

**Theorem 7.5.1** *Let  $\phi$  be an inverting problem. If there is an algorithm  $A$  that inverts  $\phi$  in time  $t(n)$ , then the SEARCH algorithm below inverts  $\phi$  in time  $ct(n)$ , where  $c$  is a constant depending only on  $A$ .*

**Proof.** We need a conditional and modified  $Kt$ . Define  $Kt'(w|x, \phi) = \min\{l(p) + \log t(p, w) : \text{given } x, \text{ program } p \text{ prints } w \text{ and tests whether or not } \phi(w) = x \text{ in time } t(p, x)\}$ . Given  $x$  and  $\phi$ , algorithm SEARCH will generate all strings  $w$  in order of increasing  $Kt'(w|x, \phi)$ , and try whether  $\phi(w) = x$  until it has found a witness that inverts  $\phi$ .

The algorithm SEARCH is as follows: The universal prefix machine  $U$  lexicographically runs all self-delimiting programs  $p$  (of length less than  $i$ ) for  $2^i 2^{-l(p)}$  steps in phase  $i$ ,  $i = 1, 2, \dots$ , until it has inverted  $\phi$  on  $x$ .

**Claim 7.5.1** All strings  $w$  of  $Kt'$  complexity less than or equal to  $k$  are generated and tested in  $2^{k+1}$  steps.

**Proof.** If  $Kt'(w|x, \phi) \leq i$ , then  $l(p) + \log t(p, x) \leq i$ . That is,  $t(p, x) \leq 2^{i-l(p)}$ , which is precisely the allotted time for this program in phase  $i$ . Since  $U$  is a prefix machine, we have by the Kraft Inequality, Theorem 1.11.1, that  $\sum 2^{-l(p)} \leq 1$ , with the sum taken over all  $p$  for which the computation of  $U$  with input  $p$  terminates. Consequently,

$$\sum_{1 \leq i \leq k} \sum_{0 < i - l(p)} 2^{i-l(p)} \leq \sum_{U(p) < \infty} 2^{-l(p)} \sum_{1 \leq i \leq k} 2^i \leq 2^{k+1}.$$

$\square$

Let  $m = \min\{Kt'(w|x, \phi) : w \text{ is a } \phi\text{-witness for } x\}$ . Suppose there exists a prefix machine  $T$  that inverts  $\phi$  on  $x$  in time  $t(n)$  with  $n = l(x)$ . By definition,  $m \leq Kt'(T|x, \phi)$ . The SEARCH algorithm inverts  $\phi$  on  $x$  in  $2^{m+1}$  steps by the claim. By definition,  $Kt'(T|x, \phi) \leq K(T) + \log t(n)$ . Therefore, SEARCH uses a number of steps of at most

$$2^{K(T)+1} t(n).$$

Setting  $c = 2^{K(T)+1}$  we prove the theorem.  $\square$

**Example 7.5.2** Let  $T_k$  be an inversion algorithm running in time  $t(n)$ . The SEARCH algorithm will use time  $2^{K(k)+O(1)}t(n)$ , which is at most  $O(k(\log k)^2)t(n)$ . The SIMPLE algorithm uses  $2^{k+1}t(n)$  time. If the inversion algorithm  $T_k$  for a given inversion problem is very simple, for example,  $K(T_k) = \log \log k$ , then SEARCH runs in time  $O(\log k)t(n)$ . This is much better than the time used by SIMPLE, which stays at  $2^{k+1}t(n)$ .  $\diamond$

**Example 7.5.3** The enumeration used in the above optimal search algorithm leads to another variant of  $Kt$ . Let a monotonic machine  $T$  enumerate a sequence of strings. The *height*  $h_T(x)$  of  $x$  with respect to  $T$  is the logarithm of the shortest time by which  $T$  outputs  $x$ , without necessarily halting after  $x$ . Since an invariance theorem such as Theorem 2.1.1 can be proved in the usual manner, we drop the subscript and write  $h(x)$ . It is not difficult to show that  $h(x) = Kt(x) + O(1)$ .

For an arbitrary NP problem, let  $t(x)$  be the time it takes to find a witness by the optimal algorithm for input  $x$ . Then it easy to see that

$$\log t(x) = \min Kt'(w|x) + O(1),$$

where the minimum is taken over all witnesses  $w$  for  $x$ .  $\diamond$

## Exercises

**7.5.1.** [27] Similar to Definition 7.5.1, one can define the  $Ct$  version of  $Kt$ . Below we set  $n = l(x)$ .

- (a) Show that  $Ct(x) \leq s(n)$  implies  $x \in C[s(n), 2^{s(n)}, \infty]$ , and the last formula implies  $Ct(x) \leq 2s(n)$ .
- (b) For a context-free language  $L$ , let  $C_L(n) = \min\{Ct(x) : x \in L^{=n}\}$ . Show that  $C_L(n) = O(\log n)$ .
- (c) Define  $C^L(n) = \max\{Ct(x) : x \in L^{=n}\}$ . Show that  $L$  is P-printable iff  $L$  is in P and  $C^L(n) = O(\log n)$ .
- (d) Show that  $C_L(n) = O(\log n)$  for all  $L$  in P iff  $C_L(n) = O(\log n)$  for all  $L$  in NP.
- (e) Every nondeterministic exponential time computable predicate  $L$  is computable in deterministic exponential time iff  $C_L(n) = O(\log n)$ .

*Comments.* Source: [E. Allender, In: *Kolmogorov complexity and computational complexity*, O. Watanabe (Ed.), Springer-Verlag, 1992, pp. 4–22]. This reference contains applications of  $C_L(n)$  and  $C^L(n)$  in random oracle constructions, pseudorandom generators, and circuit complexity.

## 7.6 Time-Limited Universal Distributions

The universal distribution  $\mathbf{m}$  was applied to the average-case analysis of algorithms in Section 4.4 and learning in Section 5.4.3. A drawback of  $\mathbf{m}$  is that it is not computable. It is not difficult to scale the entire theory down to a more feasible domain. For each time bound  $t(n)$  we construct a function  $t'(n)$  defined by  $t'(n) = nt(n)$  such that  $\mathbf{m}^{t'}$  is computable in time  $nt(n)2^n$  and multiplicatively dominates each probability distribution  $P$  with a distribution function  $P^*$  that is computable in time  $t(n)$ .

It is convenient to formulate this section in terms of distribution functions  $P^* : \{1, 2, \dots\} \rightarrow [0, 1]$ , where  $P^*(x)$  is the summed probability of all elements not exceeding  $x$ . Its *density*  $P(x) = P^*(x) - P^*(x-1)$  is the usual probability of  $x$ .

A function  $f$  is computable in time  $t$  if there exists a Turing machine  $T$  that on input  $x$  computes output  $f(x)$  in at most  $t(n)$  steps, where  $n = l(x)$ . Let  $U$  be the reference universal prefix machine  $U$ . The  $t$ -time bounded version of  $K(x)$  is defined similarly to the one for  $C(\cdot)$  in Definition 7.1.2 on page 461:

$$K^t(x) = \min\{l(p) : U(p) = x \text{ in } t(n) \text{ steps}\}.$$

For all  $t$  and  $x$  we have  $K^t(x) \leq n + O(1)$ . The limiting value of  $K^t(x)$ , as  $t(\cdot)$  grows unboundedly, is  $K(x)$ .

**Definition 7.6.1** The  $t$ -time-bounded version of  $\mathbf{m}$ , denoted by  $\mathbf{m}^t$ , is defined as follows:

$$\begin{aligned} \mathbf{m}^t(x) &= 2^{-K^t(x)}, \\ \mathbf{m}^{*t}(x) &= \sum_{y \leq x} \mathbf{m}^t(y). \end{aligned}$$

Note that for all  $t$  and  $x$ , we have  $\mathbf{m}^t(x) \leq 2^{-l(x)+O(1)}$ , and in the limit, as  $t(\cdot)$  grows unboundedly,  $\mathbf{m}^t(x)$  goes to  $\mathbf{m}(x)$ .

Consider the class  $\mathcal{P}^{*t}$  of  $t$ -time computable probability distributions  $P^*$ . Let  $n = l(x)$ . We want to show that  $\mathbf{m}^{nt(n)}(x)$  multiplicatively dominates all probability density functions  $P$  with the corresponding  $P^* \in \mathcal{P}^{*t}$ .

We do not know whether  $\mathbf{m}^{t(n)}(x)$  or  $\mathbf{m}^{*t(n)}(x)$  themselves are  $t$ -time computable. The best we can do is to compute the density function  $\mathbf{m}^{t(n)}(x)$  in  $t(n)2^{n+1}$  time, namely, computing  $K^{t(n)}(x)$  by simulating all programs of length up to  $n$ , and determining the length of the shortest program that halts with output  $x$  in  $t(n)$  steps. Similarly, we compute  $\mathbf{m}^{*t(n)}(x)$  in at most  $t(n)2^{2n+2}$  time by computing the sum

$$\mathbf{m}^{*t(n)}(x) = \sum_{y \leq x} \mathbf{m}^{t(l(y))}(y).$$

**Theorem 7.6.1** *The distribution  $\mathbf{m}^{nt(n)}$  is universal for the class  $\mathcal{P}^{\ast t}$  in the sense that it multiplicatively dominates each  $P$  in that class. That is, there exists a constant  $c_P$  such that for all  $x$ , we have  $c_P \mathbf{m}^{nt(n)}(x) \geq P(x)$  where  $c_P$  depends on  $P$  but not on  $x$  or  $t$ .*

**Proof.** This follows immediately from the following claim.

**Claim 7.6.1** *If the probability distribution  $P^*$  is computable in time  $t(n)$ , then there is a constant  $c_P$  such that for all  $x$  with  $l(x) = n$ ,*

$$K^{nt(n)}(x) \leq -\log P(x) + c_P.$$

**Proof.** Without a time bound, a proof similar to that of the optimality of a constructive version of the Shannon-Fano code would be sufficient, as in the proof of Theorem 4.3.1 on page 247. But we have to deal with the time bound here.

Divide the real interval  $[0, 1)$  into subintervals such that the code word  $p(x)$  for source word  $x$  “occupies”  $[P^*(x-1), P^*(x))$ . There is room enough since  $\sum_x P(x) \leq 1$ . The *binary interval* determined by the finite binary string  $r$  is the half open interval  $[0.r, 0.r + 2^{-l(r)})$  corresponding to the set of reals (cylinder)  $\Gamma_r$  consisting of all reals  $0.r\dots$ . If  $\Gamma_r$  is the greatest binary interval contained in  $I_x = [P^*(x-1), P^*(x))$ , then  $x$  is encoded as  $p(x) := r$ . It is easy to show that the greatest binary interval  $I_r$  in any interval  $I_x$  of  $[0, 1)$  has size at least one-fourth of  $I_x$ . Since length  $|I_x| = P(x)$ , it follows that  $l(p(x)) \leq -\log P(x) + 2$  bits.

We have to give polynomial time encoding and decoding algorithms. The encoding algorithm is trivial. Since  $P^*$  is computable in time  $t(n)$ , given source word  $x$ ,  $l(x) = n$ , the code word  $p(x)$  can be computed from  $P^*(x-1)$  and  $P^*(x)$  in  $O(t(n))$  time. In order to compute  $p^{-1}$ , the decoding function, given a code word  $p(x)$ , we proceed as follows:

**Step 1** Set  $k := 1$ .

**Step 2 (Doubling)** Repeat set  $k := 2k$  until  $\Gamma_{p(x)}$  falls in or to the right of interval  $[P^*(k-1), P^*(k))$ . Set  $l := k/2$  and  $u := k$ .

**Step 3 (Binary search)** Set  $m := (u + l)/2$ . If  $\Gamma_{p(x)}$  is in  $[P^*(m-1), P^*(m))$ , then return  $x := m$  else set  $u := m$  if  $\Gamma_{p(x)}$  is to the left of  $P^*(m-1)$  and set  $l := m$  if  $\Gamma_{p(x)}$  is to the right of  $P^*(m)$ .

This *decoding algorithm* is similar to a binary search, and it takes at most  $\sum_{i=0}^n O(t(i)) = O(nt(n))$  time to find  $x$ . By construction  $l(p(x)) \leq -\log P(x) + 2$ . This completes our encoding/decoding of  $x$  using distribution  $P^*$ .

We can reconstruct  $x$  in  $O(nt(n))$  steps from the following description:

- A description of this discussion (including the Decoding Algorithm) in  $O(1)$  bits;
- a self-delimiting program  $q$  to compute  $P^*$ ; and
- the self-delimiting code word  $p(x)$ .

This description takes  $l(p) + l(q) + O(1)$  bits. Since  $K^{nt(n)}(x)$  is the shortest program from which  $x$  can be reconstructed in  $O(nt(n))$  steps, setting  $c_P = l(q) + O(1)$ , we have  $K^{nt(n)}(x) \leq -\log P(x) + c_P$ .  $\square \quad \square$

Clearly, one can choose  $c_P := K^{nt(n)}(P) + O(1)$  since program  $q$  computing  $P^*$  may be reconstructed in time  $O(n(t(n)))$  from a shorter description  $q'$ .

**Example 7.6.1** In many algorithms, we consider only inputs in a set  $D \subseteq \mathcal{N}$ , say strings of fixed length  $n$ . We can precompute the time-limited version  $\mathbf{m}^t(\cdot|D)$  in the form of an interval representation of the table once and for all, and use it to sample according to  $\mathbf{m}_t$  by means of a sequence of fair coin flips analogous to Example 4.4.3 on page 271. Such a table of the time-limited  $\mathbf{m}^t(\cdot|D)$  needs to be precomputed only *once*, and being available, can be used repeatedly by *any* randomized algorithm using  $\mathbf{m}^t(\cdot|D)$ . An application of this is in “simple pac learning,” Section 5.4.3. How do we use  $\mathbf{m}^t(\cdot|D)$  with, say,  $D = \{0, 1\}^n$ ? Divide  $[0, 1)$  into  $2^n$  half-open disjoint intervals  $I_y$  with

$$|I_y| = \frac{\mathbf{m}^t(y)}{\mathbf{m}^{*t}(x+2^n) - \mathbf{m}^{*t}(x)},$$

such that  $\bigcup_y I_y = [0, 1)$ ,  $y = x+1, \dots, x+2^n$ .  $\diamond$

**Example 7.6.2** The entire theory of simple pac-learning can be reformulated in terms of  $t$ -time limited simple distributions,  $\mathbf{m}^{nt(n)}$ , and  $K^{nt(n)}$ . Fix a low  $t$ , like  $t(n) = O(n^2)$ , and precompute once and for all the  $\mathbf{m}^{n^3}$  table. Let  $\mathcal{C}$  be a concept class that is polynomially learnable under  $\mathbf{m}^{n^3}$ . For example, suppose that  $\mathcal{C}$  is the class of  $n^2$ -simple DNF (analogous to simple DNF in Exercise 5.4.3 on page 349, using  $K^{n^2}$ ). We can use this table, together with random coin flips as explained in Example 4.4.3 on page 271, to polynomially learn  $n^2$ -simple DNF under all  $n^2$ -simple distributions.  $\diamond$

## Exercises

---

**7.6.1.** [35] Call a probability distribution  $P : \mathcal{N} \rightarrow \mathcal{R}$  *malign* for a class of algorithms if each algorithm in the class runs in  $P$ -average time (space) equal to worst-case time (space). In Section 4.4 it was shown that  $\mathbf{m}$  is malign for the recursive algorithms. Denote  $P^*(x) = \sum_{y \leq x} P(y)$ . Function  $P^*$  is computable in time  $t$  if there exists a  $t(n)$  time bounded Turing machine that on input  $\langle x, 1^i \rangle$  writes the truncated binary expansion  $p$  of  $P(x)$  such that  $|P^*(x) - p| \leq 2^{-i}$ .

- (a) Show that for each total recursive function  $f$  there exists a probability distribution  $P$  that is malign for the class of  $f$ -time bounded algorithms, and  $P^*$  is computable in polynomial time.

Define a probability ensemble as the set of probability distributions  $\{P_n\}$  defined by  $P_n(x) = P(x | l(x) = n)$ . It is polynomial time computable if each  $P_n^*$  is polynomial time computable.

- (b) Show that there exists an exponential time computable probability ensemble which is malign for the class of polynomial time algorithms.  
(c) Show there does not exist a polynomial time computable probability ensemble that is malign for the class of polynomial time algorithms.

*Comments.* Source: P.B. Miltersen, *SIAM J. Comput.*, 22:1(1993), 147–156. This contains many more results on malignness. K. Kobayashi studies the differences and relations between malign measures and universal distributions dominating classes of distributions in [*IEICE Trans. Inform. Systems*, E76-D:6(1993), 634–640; Transformations that preserve malignness of universal distributions, *Theoret. Comput. Sci.*, to appear.]. A. Jakoby, R. Reischuk, and C. Schindelhauer study malign distributions for average case circuit complexity in [*Proc. 12th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comp. Sci., vol. 900, Springer-Verlag, Heidelberg, 1995, 628–639].

**7.6.2.** [18] (a) Let  $l(x) = n$ . Show that probabilities  $\mathbf{m}^t(x) = 2^{-K^t(x)}$  for the set of all  $x$ 's with  $K^t(x) = O(\log n)$ , can be computed in time polynomial in  $t(n)$ .

- (b) Use Item (a) to show that one can precompute  $\mathbf{m}_t$  for the high probability  $x$ 's ( $K^t(x) = O(\log n)$ ) in polynomial time for  $t(n)$  is polynomial.

*Comments.* Source: M. Li and P. Vitányi, *SIAM J. Comput.*, 20:5(1991), 915–935.

## 7.7 Logical Depth

One may consider a book on number theory difficult, or “deep.” The book will list a number of difficult theorems of number theory. However, it has very low Kolmogorov complexity since all theorems are derivable from the initial few definitions. Our estimate of the difficulty, or “depth,” of the book is based on the fact that it takes a long time to reproduce the book from part of the information in it. The existence of a “deep” book is itself evidence of some long evolution preceding it. Currently, the sequence of primes is being broadcast to outer space since it is deemed deep enough to prove to aliens that it arose as a result of a long evolution.

From the point of view of an investigator, a sequence is deep if it yields its secrets only slowly: one will be able to discover all significant regularities in it only if one analyzes it long enough.

**Example 7.7.1** A suggestive example is provided by DNA sequences. Such a sequence is quite regular and has some 90% redundancy, possibly due to evolutionary history. A DNA sequence over an alphabet of four letters  $\{A, C, G, T\}$  looks like nothing but a super-long ( $3 \times 10^9$  characters for humans) computer program. A particular three-letter combination literally signifies “begin” of the encoding of a protein. Following the “begin” command, every next block of three consecutive letters encodes one of the 20 amino acids. At the end another three-letter combination signifies the “end” of the program for this protein. Such a sequence is not Kolmogorov random, and it encodes the structure of a living being.

DNA is much less random than, say, a typical configuration of gas in a container. On the other hand, DNA is more random than a crystal. Both gases and crystals are structurally trivial; the former is in complete chaos and the latter is in total order. Intuitively, DNA contains more useful information than both. A “deep” object, such as DNA, is something really simple but “disguised” by complicated manipulations of nature or computation by computer. To quote Charles Bennett:

“A structure is deep, if it is superficially random but subtly redundant, in other words, if almost all its algorithmic probability [ $m$ ] is contributed by slow-running programs [...] A priori the most probable explanation of ‘organized information’ such as the sequence of bases in a naturally occurring DNA molecule is that it is the product of an extremely long biological process.”

◇

Logical depth is the necessary number of steps in the deductive or causal path connecting an object with its plausible origin. Formally, it is the time required by a universal computer to compute the object from its compressed original description.

It turns out that it is quite subtle to give a formal definition of “depth” that satisfies our intuitive notion of it. After some attempts at a definition, we will settle for Definition 7.7.1. As usual, we write  $x^*$  to denote

the shortest self-delimiting program (of the reference universal prefix machine  $U$ ) for  $x$ . If there is more than one of the same length, then  $x^*$  is the first such program in a fixed enumeration.

**Attempt 1** The number of steps required to compute  $x$  from  $x^*$  is not a stable quantity since there might be a program of just a few more bits using substantially less time to generate  $x$ . That this can happen is shown by the hierarchy theorems of Section 7.1.2. Therefore, a proper definition of depth probably should “compromise” between the program size and computation time.

**Attempt 2** Relax the strict requirement of minimum program to *almost minimum* programs. Define that a string  $x$  has depth  $d$  within error  $2^{-b}$  if  $x$  can be computed in  $d$  steps by a program  $p$  of no more than  $b$  bits in excess of  $x^*$ . That is,  $2^{-l(p)}/2^{-K(x)} \geq 2^{-b}$ .

This definition is stable but is unsatisfactory because of the way it treats multiple programs of the same length. If  $2^b$  distinct programs of length  $m+b$  all compute  $x$ , then together they account for the same algorithmic probability

$$\sum \{2^{-l(p)} : U(p) = x, l(p) = m+b\}, \quad (7.11)$$

as one program of length  $m$  printing  $x$  does. That is, they are as likely to produce  $x$  as output of the universal reference prefix machine when its input is provided by fair coin tosses. But with the proposed definition,  $2^b$  programs of length  $m+b$  make the emerging of  $x$  no more probable than one program of length  $m+b$ .

We shall explicitly take the algorithmic probability into account. The universal prior probability of a string  $x$  is

$$Q_U(x) = \sum_{U(p)=x} 2^{-l(p)},$$

where  $U$  is the reference universal prefix machine. This is the probability that  $U$  would print  $x$  if its input were provided by random tosses of a fair coin. By Theorem 4.3.3, page 253, we can define

$$-\log Q_U(x) + O(1) = -\log \mathbf{m}(x) = K(x). \quad (7.12)$$

Theorem 4.3.3 is stated with an additional  $O(1)$  term in the second equality. It shows that  $2^{-K(x)}$  is a universal discrete semimeasure. This means that we are free to choose the reference universal semimeasure  $\mathbf{m}$  in Theorem 4.3.1 exactly equal to  $2^{-K(x)}$ .

Thus, weighing all possible causes of emergence of  $x$  appropriately, we are led to the following definition:

**Definition 7.7.1** The *depth* of a string  $x$  at significance level  $\epsilon = 2^{-b}$  is

$$\text{depth}_\epsilon(x) = \min\{t : Q_U^t(x)/Q_U(x) \geq \epsilon\},$$

where  $Q_U^t(x) = \sum_{U^t(p)=x} 2^{-l(p)}$  and  $U^t(p) = x$  means that  $U$  computes  $x$  within  $t$  steps and halts. A string  $x$  is  $(d, b)$ -deep if  $d = \text{depth}_\epsilon(x)$  and  $\epsilon = 2^{-b}$ .

If  $x$  is  $(d, b)$ -deep, then  $x$  receives an approximately  $1/2^{b \pm \delta}$  fraction of its algorithmic probability (for some small  $\delta$ ) from programs running in  $d$  steps. Below we formalize this statement and make  $\delta$  precise. A binary string  $x$  is  $b$ -compressible if  $l(x^*) \leq l(x) - b$ . Otherwise,  $x$  is  $b$ -incompressible.

**Theorem 7.7.1** A string  $x$  is  $(d, b)$ -deep ( $b$  up to precision  $K(d) + O(1)$ ) if and only if  $d$  is the least time needed by a  $b$ -incompressible program to print  $x$ .

**Proof.** The precise form of the statement of the theorem is

$$\frac{1}{2^{b+K(d)+O(1)}} \leq \frac{Q_U^d(x)}{Q_U(x)} \leq \frac{1}{2^{b-O(1)}}, \quad (7.13)$$

where  $Q_U^d$  is the  $d$  time-bounded analogue of  $Q_U$ , in the sense that the right inequality implies the “if” part and the left inequality implies the “only if” part.

(IF) We prove the right inequality of Equation 7.13. Suppose that  $d$  is the least time required to compute  $x$  by any  $b$ -incompressible program. In other words, each program computing  $x$  in fewer than  $d$  steps is compressible by  $b$  bits. For each such program  $p$ , there is a program  $p'$  such that  $U(p') = p$  and  $l(p') \leq l(p) - b$ . Let  $q$  be a program that simulates  $U$  on  $p'$  to obtain  $p$  and subsequently simulates  $p$  to obtain  $x$ . The length of  $q$  is about  $l(p) - b + O(1)$ . The relationship between  $p$  and  $q$  is 1-1.

The notation  $U^d(p) = x$  means that  $U$  computes  $x$  in  $d$  steps from  $p$ . Let  $\alpha = Q_U(x) - \sum_{U(q)=x} 2^{-l(q)}$ , where the sum is taken over all  $q$ ’s as defined above. By Equation 7.12, we have  $\alpha \geq 0$ . Then,

$$\begin{aligned} Q_U^d(x)/Q_U(x) &= \sum_{U^d(p)=x} 2^{-l(p)} \Bigg/ \left( \alpha + \sum_{U(q)=x} 2^{-l(q)} \right) \\ &\leq \sum_{U^d(p)=x} 2^{-l(p)} \Bigg/ \sum_{U(q)=x} 2^{-l(q)} \\ &\leq \sum_{U^d(p)=x} 2^{-l(p)} \Bigg/ \sum_{U(q)=x} 2^{-(l(p)-b+O(1))} \\ &\leq 1/2^{b-O(1)}. \end{aligned}$$

(ONLY IF) Assume, by way of contradiction, that the left inequality of Equation 7.13 is false. Then given  $x^*$  and  $d^*$  (the shortest self-delimiting programs for  $x$  and  $d$ ), we can enumerate the set  $A$  of all programs computing  $x$  in time at most  $d$  by simulating all programs of length less than  $l(x)$  for  $d$  steps. Thus, set  $A$  can be computed by a program of size

$$l = K(x) + K(d) + O(1). \quad (7.14)$$

If  $X$  is a set, then define  $L(X) = \sum_{p \in X} 2^{-l(p)}$ . By definition of  $A$ , we have  $L(A) = \sum_{p \in A} 2^{-l(p)} = Q_U^d(x)$ . By the contradictory assumption,  $2^{b+K(d)+O(1)} < Q_U(x)/Q_U^d(x)$ . By Theorem 4.3.3,  $Q_U(x) = 2^{-K(x)+O(1)}$ . Together, it follows that

$$L(A) = Q_U^d(x) < 2^{-K(x)-b-K(d)-O(1)}. \quad (7.15)$$

**Claim 7.7.1** Let  $B$  be a prefix-free set (no element is a proper prefix of any other element),  $\sum_{x \in B} 2^{-l(x)} < 2^{-m}$ , and  $B$  is enumerated by program  $s$ . Then every string in  $B$  can be effectively compressed by at least  $m - l(s) - O(1)$  bits.

**Proof.** Increasing the probability of each  $x$  in  $B$  by a multiplicative factor  $2^m$ , we still have  $\sum_{x \in B} 2^{-l(x)} 2^m < 1$ . Therefore, the elements of  $B$  can be coded by the Shannon-Fano code as in Lemma 4.3.3, page 253, of Coding Theorem 4.3.3. The code word for each  $x$  in  $B$  has length at most  $l(x) - m + 2$ . In order to make this coding effective, we use  $s$  to enumerate all and only strings in  $B$ . This takes an additional  $l(s) + O(1)$  bits in the code for each  $x \in B$ . Hence, each string in  $B$  is effectively compressed by  $m - l(s) - O(1)$  bits.  $\square$

By Equations 7.14 and 7.15 and Claim 7.7.1, every program in  $A$  can be compressed by  $b + K(x) + K(d) + O(1) - l - O(1)$  bits. By proper choice of the  $O(1)$  constants, we conclude that every program in  $A$  can be compressed by more than  $b$  bits, contradicting the fact that  $x$  is  $(d, b)$ -deep.  $\square$

Definition 7.7.1 is not equivalent to the seemingly close Attempt 2. The reason is that in trying to work out the equivalence, the Coding Theorem 4.3.3 must be used, and it brings in an exponential factor in time. The two are close when the depth we are interested in is much larger than exponential.

**Example 7.7.2** If  $x$  is  $(d, b)$ -deep, then the universal machine must take at least  $d$  steps to compute  $x$  from  $x^*$ . Otherwise,  $x$  would be less than  $d$ -deep at all significance levels. Note that if a string is  $d$ -deep at significance level  $2^{-b}$ , then it is also  $d$ -deep at significance level  $2^{-(b+1)}$ . A  $b$ -incompressible program printing  $x$  is an explanation of how the phenomenon could occur

with probability at least  $2^{-b}$ . Thus, with larger  $b$ , the explanation is less compelling and less likely. If the object is very deep even with a larger  $b$ , then this means that even a less compelling explanation takes a long time.  $\diamond$

**Definition 7.7.2** At any significance level,  $x$  is  $d$ -shallow if its depth does not exceed  $d$ . Any string  $x$  must take at least  $n = l(x)$  steps to be printed. If  $x$  is  $n \pm O(1)$ -shallow (at all significance levels), then we will simply call  $x$  shallow.

**Example 7.7.3** A random string  $x$  of length  $n$  is always shallow, because it can be printed by its shortest program  $x^*$  of length  $n \pm O(1)$  in  $n$  steps. Thus, a random string is shallow at all significance levels.

String  $1^n$  is shallow since a constant-sized incompressible program prints it in time  $n$  given input  $0^n$ . The strings  $(01)^n$  and  $01^n$  are likewise shallow.  $\diamond$

**Example 7.7.4** This example demonstrates the distinction between depth and information. Consider two sequences: the halting probability  $\Omega$  defined in Section 3.6.2, and  $\chi$  defined in the proof of Theorem 7.1.3. Although both encode the halting information of Turing machines,  $\chi$  is deep and  $\Omega$  shallow. Since  $\Omega$  encodes the halting problem with maximal density, it is recursively indistinguishable from a random string and practically useless, since  $K(\Omega_{1:n}) \geq n - O(1)$  by Section 3.6.2. Thus,  $\Omega$  is shallow.

On the other hand, since  $\chi$  is a characteristic sequence of a recursively enumerable set constructed in the proof of Theorem 7.1.3, we have  $C(\chi_{1:n}|n) \leq \log n$  by Barzdins's Lemma (Theorem 2.7.2). If  $t$  is the minimum time required to compute  $\chi_{1:n}|n$  from an  $O(\log n)$  sized incompressible program, then  $\chi_{1:n}$  is  $t$ -deep at all significance levels.

As a matter of fact,  $\chi_{1:n}$  is very deep since  $t$ , the time required to compute an initial segment  $\chi_{1:n}$  from a  $\log n$  sized program, increases faster than any computable function. Namely, Theorem 7.1.3 implies that  $\chi_{1:n}$  can be computed from a program logarithmic in  $n$  if unlimited time is allowed, but can only be computed from a program that is linear in  $n$  if any recursive bound is imposed on the decoding time.  $\diamond$

**Example 7.7.5** Depth is stable. That is, deep strings cannot be quickly computed from shallow ones (Exercise 7.7.3). In the genetic sense organisms evolve relatively slowly. This may be called the Slow Growth Law. There is a mathematical version of such a law. Consider any string  $x$  and two significance parameters  $s_2 > s_1$ . A random program generated by coin tossing has probability less than  $2^{-(s_2-s_1)+O(1)}$  of transforming  $x$  into an excessively deep output, one whose  $s_2$ -significance depth exceeds the

$s_1$ -significance depth of  $x$  plus the run time of the transforming program plus  $O(1)$  (Exercise 7.7.5).  $\diamond$

**Example 7.7.6** It is not known whether all functions computable by algorithms that use space at most polynomial in the input length (PSPACE) can also be computed by deterministic algorithms that use time polynomial in the input length (P). Suppose the space that can be used by a computation is fixed at polynomial in the input length. If  $P = \text{PSPACE}$ , then every string derivable from an input has depth at most polynomial. If a string were the result of an exponentially long computation, then there would be a shorter computation of polynomial length to obtain it from the input. That is, the maximal depth of any string computable in PSPACE would be polynomial. If, on the other hand,  $P \neq \text{PSPACE}$ , then there are strings computable in PSPACE whose depth is larger than polynomial.

Consider, for example, the development of DNA. It is a fair assumption that no computations in the evolution of DNA will exceed PSPACE. Then the maximal logical depth achievable is possibly exponential, such as  $2^n$ , for a seed of length  $n$ . But if  $P = \text{PSPACE}$ , then there is always a shorter computation of length at most polynomial in  $n$ , and the evolved DNA is at most polynomially deep.

Similar statements can be made for micro states in evolving thermodynamic systems; for example, if the container, temperature variation, and pressure variation for a sample of gas molecules are fixed.  $\diamond$

Apparently, depth is different from  $Kt$ , since it may be the case that a random string is *greater* than a nonrandom one of equal length in  $Kt$  but shallower in logical depth.

## Exercises

**7.7.1.** [22] Strengthen the first inequality of Equation 7.13 to

$$2^{-b-\min\{K(b), K(d)\}-O(1)} \leq Q_U^d(x)/Q_U(x).$$

*Comments.* Hint: given  $x^*$  and  $b^*$ , we can enumerate all programs in order of increasing running time and stop when the accumulated algorithmic probability measure exceeds  $2^{-K(x)+b}$ . Source: Lemma 3 in [C.H. Bennett, pp. 227–257 in: *The universal Turing machine: a half-century survey*, R. Herken (Ed.), Oxford University Press, 1988]. This paper is also the source of all other exercises, except Exercise 7.7.4, in this section.

**7.7.2.** [28] Deep strings are not easy to identify, but can be constructed by diagonalization. Prove that the following program finds a ( $t$ -) deep string: “Find all  $x$  of length  $n$  such that  $Q_U^t(x) > 2^n$ ; print the first string which is not in this set.”

**7.7.3.** [30] Prove the following “stability property.” Deep strings cannot be quickly computed from shallow ones. More precisely, there is a polynomial  $p$  and a constant  $c$ , both depending only on the universal machine  $U$ , such that if  $q$  is a program to compute  $x$  in time  $t$  and if  $q$  is less than  $(d, b)$ -deep, then  $x$  is less than  $(d + p(t), b + c)$ -deep.

**7.7.4.** [25] Depth is machine independent: If  $U, U'$  are two optimal universal machines, prove that there exists a polynomial  $p$  and a constant  $c$ , both depending only on  $U$  and  $U'$ , such that  $(p(d), b + c)$ -depth on either machine is a sufficient condition for  $(d, b)$ -depth on the other.

**7.7.5.** [30] Show the *Slow Growth Law*: Given a string  $x$  and two significance parameters  $s_2 > s_1$ , a random program generated by coin tossing has probability less than  $2^{-(s_2-s_1)+O(1)}$  of transforming  $x$  into an excessively deep output, one whose  $s_2$ -significance depth exceeds the  $s_1$ -significance depth of  $x$  plus the run time of the transforming program plus  $O(1)$ .

## 7.8 History and References

---

The earliest discussion of resource-bounded Kolmogorov complexity is the quotation of A.N. Kolmogorov at the outset of this chapter. The earliest result in resource-bounded Kolmogorov complexity the authors know of is Theorem 7.1.3, due to J.M. Barzdins [*Soviet Math. Dokl.*, 9(1968), 1251–1254]. In the early 1970s, R.P. Daley wrote his Ph.D. thesis on Kolmogorov complexity with D.W. Loveland at Carnegie-Mellon University. Parts of it were published as [*J. ACM*, 20:4(1973), 687–695; *Inform. Contr.*, 23(1973), 301–312]. This work concerns resource bounds for uniform complexity  $C(x; n)$  as defined by Loveland (Exercise 2.3.3, page 124). These papers, and [L.A. Levin, *Problems Inform. Transmission*, 9(1973), 265–266; R.P. Daley, *Theoret. Comput. Sci.*, 4(1977), 301–309; L. Adleman, ‘Time, space, and randomness,’ LCS Report TM-131, 1979, MIT], are significant early documents of resource-bounded Kolmogorov complexity. The papers [J. Hartmanis, *Proc. 24th IEEE Found. Comput. Sci.*, 1983, pp. 439–445; Ker-I Ko, *Theoret. Comput. Sci.*, 48(1986), 9–33; and M. Sipser, *Proc. 15th ACM Symp. Theory Comput.*, 1983, pp. 330–335] started the recent research wave on resource-bounded Kolmogorov complexity. Related work on polynomial resource-bounded notions of randomness is by Y. Wang [*Randomness and Complexity*, Ph.D. Thesis, Heidelberg, 1996; *Proc. 11th IEEE Conf. Structure in Complexity Theory*, 1996, pp. 180–189].

The *CD* notation in Definition 7.1.4 as well as Theorem 7.1.2 were introduced by M. Sipser in [*Ibid.*]. Theorems 7.1.4 and 7.1.5 and Definition 7.1.8 are from [Ker-I Ko, *Ibid.*]. The theory of co-enumerable and recursive majorants of complexity is found in A.K. Zvonkin and L.A.

Levin, *Russ. Math. Surveys*, 25:6(1970), 83–124]; see also Example 4.1.1. and Exercise 7.1.4, page 472.

The notation  $C[f(n), t(n), s(n)]$  in Section 7.1.2 is from [J. Hartmanis, *Proc. 24th IEEE Found. Comput. Sci.*, 1983, pp. 439–445]. This paper had major influence on the research applying resource-bounded Kolmogorov complexity to computational complexity theory. Earlier, Daley [*J. ACM*, 20:4(1973), 687–695] defined similar notions for infinite sequences with uniform complexity:  $\mathcal{C}[f|t] = \{x : (\forall^\infty n) C^t(x_{1:n}; n) \leq f(n)\}$ , Exercises 7.1.6, and 7.1.7. The hierarchy theorems in Section 7.1.2 were first developed by J. Hartmanis [*Ibid.*]. L. Longpré's [Ph.D. Thesis, Cornell University, 1986] contains Theorem 7.1.8, as well as the details of hierarchy theorems for complexity, space, and time.

Section 7.2 on language compression is based on [M. Sipser, *Proc. 15th ACM Symp. Theory Comput.*, 1983, pp. 330–335; H.M. Buhrman and L. Fortnow, *Proc. 14th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comput. Sci., Springer-Verlag, Heidelberg, 1997; A. Goldberg and M. Sipser, *SIAM J. Comput.*, 20(1991), 524–536]. The elegant Theorem 7.2.1, added to the second edition, is due to H.M. Buhrman and L. Fortnow, *Ibid.* Theorems 7.2.2 and 7.2.4 are from [M. Sipser, *Ibid.*]. Theorems 7.2.5 and 7.2.6 are from [A. Goldberg and M. Sipser, *Ibid.*]. The original proofs of the Coding Lemma and Theorem 7.2.6 use probabilistic arguments. Example 7.2.1 is one of the early applications of time-bounded Kolmogorov complexity. Originally, Sipser proves  $BPP \subseteq \Sigma_4^p \cap \Pi_4^p$  using time-bounded Kolmogorov complexity. The current version,  $BPP \subseteq \Sigma_2^p \cap \Pi_2^p$ , is due to P. Gács. Theorem 7.2.7 is independently proved in [E. Allender [Ph.D. Thesis, Georgia Inst. Tech., 1985; A. Goldberg and M. Sipser, *Ibid.*]. Section 7.3.1 is based on [J. Hartmanis, *Proc. 24th IEEE Found. Comput. Sci.*, 1983, pp. 439–445], which also contains Lemma 7.3.1, and Theorems 7.3.1 and 7.3.2. Theorem 7.3.3 is due to R.V. Book, P. Orponen, D. Russo, and O. Watanabe [*SIAM J. Comput.*, 17(1988), 504–516]. Theorem 7.3.4 is from [E. Allender and R. Rubinstein, *SIAM J. Comput.*, 17(1988), 1193–1202], while the equivalence of Items (i) and (iv) is independently due to J.L. Balcázar and R.V. Book [*Acta Informatica*, 23(1986), 679–688] and J. Hartmanis and L. Hemachandra [*Inform. Process. Lett.*, 28(1988), 291–295].

Application of resource-bounded Kolmogorov complexity in computational complexity is widespread. We did not cover [J.H. Lutz, *J. Comput. System Sci.*, 44(1992), 220–258] generalizing a uniform (Lebesgue) measure over  $\{0, 1\}^*$  to define *resource-bounded measure*  $\mu$  over  $\text{ESPACE} = \bigcup_{c \in \mathcal{N}} \text{DSPACE}[2^{cn}]$ . Note that here  $\text{ESPACE}$  is a countable set and under the usual uniform measure it has measure zero. Lutz uses a special-purpose notion of measure and shows that complete problems form a

negligibly small subclass in ESPACE. For each real number  $a > 1$ , if

$$X = \{L : C^{\infty,s}(\chi_{L=n}) \geq 2^n - an, \text{ for all but finitely many } n\},$$

where the space bound  $s = 2^{O(n)}$ , then  $\mu(X|\text{ESPACE}) = 1$ . That is, almost every language computable in  $2^{O(n)}$  space has high  $2^{O(n)}$  space-bounded Kolmogorov complexity, almost everywhere. Every problem in a complexity class is reducible to every complete problem in that class. Intuitively, such complete problems must have highly organized structures and hence have *low* Kolmogorov complexity. This is the case for ESPACE. It is also shown that for every hard language (under many-to-one polynomial reduction)  $L$  for ESPACE there exists  $\epsilon > 0$  such that  $C^{-,s}(\chi_{A=n}) < 2^n - 2^{n^\epsilon}$  infinitely often, with the space bound  $s = 2^{2n}$ . In [Theoret. Comput. Sci., 81(1991), 127–135] Lutz also showed that if  $P$  is properly contained in BPP, then  $\mu(E|\text{ESPACE}) = 0$ , where  $E = \bigcup_{c \in \mathcal{N}} \text{DTIME}[2^{cn}]$ . See also [R.V. Book and J.H. Lutz, SIAM J. Comput., 22:2(1993), 395–402]; the survey by D.W. Juedes and J.H. Lutz in [*Kolmogorov complexity and its relation to computational complexity theory*, O. Watanabe (Ed.), Springer-Verlag, 1992, pp. 43–65]; and the survey by J.H. Lutz in [*Proc. 8th IEEE Conf. Structure in Complexity Theory*, 1993, pp. 158–175]. H. Buhrman and L. Longpré in [*Proc. 13th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comp. Sci., vol. 1046, Springer-Verlag, Heidelberg, 1996, 13–24] treat a resource bounded version of compressibility of infinite sequences that is shown to be equivalent to the resource bounded martingales in [J.H. Lutz [*J. Comput. System Sci.*, 44(1992), 220–258]. K.W. Regan and J. Wang [*SIGACT News*, 25(1994), 106–113] and A. Naik, K.W. Regan and D. Sivakumar [*Theoret. Comput. Sci.*, 148(1995), 325–349] show how to construct oracles in a quasilinear time model.

Section 7.4 on instance complexity is based on [P. Orponen, K. Ko, U. Schöning, and O. Watanabe, *J. ACM*, 41(1994), 96–121]. Several open questions in the first edition of this book were solved in [M. Kummer, SIAM J. Comput., 25:6(1996), 1123–1143; H.M. Buhrman and P. Orponen, *J. Comput. System Sci.*, 53:2(1996), 261–266; and L. Fortnow and M. Kummer, *Theoret. Comput. Sci. A*, 161(1996), 123–140]. See Exercises 7.4.5, 7.4.8, 7.4.6, and 7.4.7.

$Kt$  complexity of Section 7.5 is from [L.A. Levin, *Problems Inform. Transmission*, 9:3(1973), 265–266]. Levin’s universal optimal search algorithm in Section 7.5.1 is Theorem 2 (without proof) of that paper. (Theorem 1 of the paper (without proof) is the independent discovery of NP-complete problems. Earlier, S.A. Cook proved similar results in [*Proc. 3rd ACM Symp. Theory Comput.*, 1971, pp. 151–158]. Apparently, Levin did not want to publish the NP-completeness part of the paper. Kolmogorov urged him to publish and recommended the paper to *Problems Inform. Transmission*. Levin agreed on the condition that

he could include the universal search part.) Our treatment partly follows [Yu. Gurevitch, *EATCS Bull.*, 35(1988), 71–82; R.J. Solomonoff, Optimum sequential search, Memorandum, 1984]. In the simulation of the search procedure, it is standard to use so-called Kolmogorov-Uspensky machines, which allow a universal machine to simulate each Kolmogorov-Uspensky machine in linear time. In contrast, Turing machines have a logarithmic factor of slowdown because the universal Turing machine has only a fixed number of tapes. We ignored this problem. See [L.A. Levin, *Inform. Contr.*, 61(1984), 17–37] for further results on  $Kt$  complexity. Although the constant in the universal search procedure is forbiddingly large (exponential in the size of the shortest optimal program being looked for) there are some ideas to deal with this. Solomonoff [pp. 473–491 in: L.N. Kanal and J.F. Lemmer, Eds., *Uncertainty in Artificial Intelligence*, Elsevier, 1986] proposes to learn first small chunks (subgoals) using the universal distribution, which could be done “relatively fast.” Then use the small chunks as building blocks to learn a more composite goal from the new copy of the universal distribution over the building blocks as the basic elements. This way one progresses fast to composite goals. Some computer experiments with universal optimal search are reported in [J. Schmidhuber, *Proc. 12th Int'l Conf. Machine Learning*, Morgan Kaufmann, 1995, 488–496] for learning simple threshold units with high generalization capacity; and [J. Schmidhuber and M. Wiering, *Proc. 13th Int'l Conf. Machine Learning*, Morgan Kaufmann, 1996, To appear] for solving partially observable mazes. Related work on a notion of “potential” is based on L. Adleman [*Ibid.*] and is not treated in this edition. In the first edition of this book we reformulated the ideas in terms of  $Kt$  complexity. L. Hemachandra and G. Wechsung [*Theoret. Comput. Sci.* 83(1991), 313–322] continue this approach. Our treatment of  $Kt$  complexity, universal optimal search, and potential has greatly profited from discussions with P. Gács and R. Solomonoff.

Other work on resource-bounded complexity includes: G. Peterson, *Proc. 21st IEEE Found. Comput. Sci.*, 1980, pp. 86–95; E. Allender, *J. Comput. System Sci.*, 39(1989), 101–124; J.H. Lutz, *J. Comput. System Sci.*, 41(1990), 307–320; M. Hermo and E. Mayordomo, *Math. Systems Theory*, 27(1994), 247–356; J.L. Balcázar and U. Schöning, *Theoret. Comput. Sci.*, 99(1992), 279–290.

Section 7.6 on computable versions of the universal distribution is partially based on [M. Li and P.M.B. Vitányi, *SIAM J. Comput.*, 20:5(1991), 915–935]. The original nonresource-bounded versions of the universal distribution are treated in Chapter 4. The computable version of the universal distribution can be used as the dominating distribution in simple pac-learning, Section 5.4. Another application is exhibited in Section 4.4. Related work on the maximal gain of average-case complexity over worst-case complexity for (algorithm, distribution) pairs and associated fam-

ilies like polynomial time algorithms and polynomial time computable distributions is [P.B. Miltersen, *SIAM J. Comput.*, 22:1(1993), 147–156; K. Kobayashi, *IEICE Trans. Inform. Systems*, E76-D:6(1993), 634–640; K. Kobayashi, Transformations that preserve malignness of universal distributions, *Theoret. Comput. Sci.*, to appear; A. Jakoby, R. Reischuk, and C. Schindelhauer, *Proc. 12th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comp. Sci., vol. 900, Springer-Verlag, Heidelberg, 1995, 628–639]. A.K. Jagota and K.W. Regan ['Performance of MAX-CLIQUE approximation heuristics under description-length weighed distributions,' UB-CS-TR 92-24, SUNY at Buffalo, 1992] have used a distribution  $\mathbf{q}(x)$  that is essentially a time bounded version of  $\mathbf{m}(x)$ . See also the History and Reference Section of Chapter 4.

Logical depth as in Section 7.7 was first described by G.J. Chaitin in [*IBM J. Res. Develop.*, 21(1977), 350–359] and studied at greater length by C.H. Bennett in [*Emerging Syntheses in Science*, D. Pines (Ed.), Addison-Wesley, 1987, pp. 215–233; *The universal Turing machine: a half-century survey*, R. Herken (Ed.), Oxford Univ. Press, 1988, pp. 227–258]. The material contained in Section 7.7 is primarily based on the latter article of C.H. Bennett. There one can find Definition 7.7.1 and a variant (Lemma 3 in the paper) of Theorem 7.7.1. The proof given in the article is rather sketchy and informal. Bennett also defines the depth for infinite strings, which we have not discussed. Our treatment has greatly profited from discussions with P. Gács. For material on logical depth and fault-tolerant evolution, see [P. Gács, 'Self-correcting two-dimensional arrays,' in *Randomness in Computation*, vol 5, S. Mi-cali (Ed.), JAI Press, 1989, pp. 223–326; an asynchronous version is in W.G. Wang, Ph.D. Thesis, Boston University, 1990]. D.W. Juedes, J.I. Lathrop, and J.H. Lutz [Computational depth and reducibility, To appear in *Theoret. Comput. Sci.*] study relations between logical depth of infinite sequences and recursively time-bounded reducibility.

The edited volume [*Kolmogorov complexity and computational complexity*, O. Watanabe (Ed.), Springer-Verlag, 1992] contains five survey papers related to several topics in this chapter. E. Allender presents applications of time-bounded Kolmogorov complexity to questions in complexity theory, see also Exercise 7.5.1 on page 505. R.V. Book surveys sets with small information content, relating Sections 7.3.1 and 7.3.2 and current research in computational complexity theory. L. Longpré studies the statistical properties and tests for resource-bounded Kolmogorov random strings, supplementing Section 7.1.1. D.W. Juedes and J.H. Lutz treat the Kolmogorov complexity of characteristic sequences of languages. The survey by V.A. Uspensky analyzes the definitional and quantitative relations between the various versions of (nonresource bounded) Kolmogorov complexity.

# Physics, Information, and Computation

Various issues in information theory and theoretical physics can be fruitfully analyzed by Kolmogorov complexity. This is the case for physical aspects of information processing and for application of complexity to physics issues. Physicists have used complexity arguments in a variety of settings like information distance, thermodynamics, chaos, biology, and philosophy. We touch briefly upon several themes, but focus on three main issues.

First, we analyze the relation between Shannon’s entropy and the expected prefix complexity and expected plain Kolmogorov complexity. Shannon’s entropy measures the uncertainty in a statistical ensemble of messages, while Kolmogorov complexity measures the algorithmic information in an individual message. In Section 2.8 it was observed that the  $P$ -expected value of  $C(\cdot)$  is asymptotic to the entropy  $H(P)$ . Using the prefix complexity and the theory of universal distributions we establish the precise relationships.

Second, we look at energy dissipation in computing. Continued advance in the miniaturization and mobilization of computing devices will eventually require (near) dissipationless computing. Apparently, there are no physical laws that require reversible computations to dissipate energy. It is for this reason that notions of reversible computers, both physical and logical, are rapidly gaining prominence. We treat reversible computing first because it is used in some of the later topics in this chapter. We give both a variety of physical implementations that are theoretically possible, and a logical model—the reversible Turing machine. We then define the Kolmogorov complexity variant associated with this machine.

While Kolmogorov complexity is an absolute measure of information in an object, it is desirable to have a similar absolute notion for the information distance between two objects. Among other things, absolute information distance is related to theories of pattern recognition where one wants to express the notion of “picture similarity.” We give several alternative definitions for “information distance,” based on considerations about either reversible (dissipationless) computations or irreversible computations. It turns out that these definitions define the same notion. Next, we formulate certain weak requirements that any reasonable “picture distance” ought to satisfy. It then turns out that our earlier notion is the optimal “picture distance.” A related distance measures the amount of nonreversibility of a computation. Finally, we treat the minimum dissipated energy in a computation.

Third, we look at an application of Kolmogorov complexity in statistical thermodynamics. There, one explains the classical theory of thermodynamics by statistical and information-theoretic analysis of an underlying deterministic model. It turns out that a complexity analysis using the powerful methods developed in the first few chapters gives a basis of an algorithmic theory for entropy. Some applications include a proof of an “entropy nondecrease over time” property, and “entropy stability” property, “entropy increase” for certain systems, and an analysis of Maxwell’s demon. We end with some miscellaneous topics.

## 8.1 Algorithmic Complexity and Shannon’s Entropy

There is a close relation between information theory and coding, Section 1.11. In particular, the quantity

$$H(P) = - \sum_x P(X = x) \log P(X = x)$$

is the entropy of a random variable  $X$  with probability  $P(X = x)$  of outcome  $x$  (Definition 1.11.1 on page 67). For convenience we abbreviate the notation “ $P(X = x)$ ” to “ $P(x)$ ,” conforming to our customary usage.

Let  $E$  be a prefix-code assigning codeword  $E(x)$  to source word  $x$ . The Noiseless Coding Theorem 1.11.2 on page 75 asserts that the minimal average codeword length  $L(P) = \sum_x P(x)l(E(x))$  among all such prefix-codes  $E$  satisfies

$$H(P) \leq L(P) \leq H(P) + 1.$$

By definition  $K(x) = l(x^*)$ , where  $x^*$  is the shortest self-delimiting program for  $x$  with respect to the reference prefix machine. If there is more than one such shortest program then  $x^*$  is the first one in the standard enumeration. Consider the mapping  $E$  defined by  $E(x) = x^*$ . This is a prefix-code, and by its definition a very parsimonious one. Suppose the

source words  $x$  are distributed as a random variable  $X$  with probability  $P(x)$ . The expected code word length of source words with respect to probability distribution  $P$  is  $\sum_x P(x)K(x)$ .

What we would like to know is the following: While  $K(x)$  is fixed for each  $x$  and *independent* of the probability distribution  $P$ , is  $K$  still so universal that its  $P$ -expected code word length  $\sum_x P(x)K(x)$  achieves the minimal average code word length  $H(P) = -\sum_x P(x) \log P(x)$ ?

This universality requirement contrasts with the Shannon-Fano code of Example 1.11.2 on page 67, which does achieve the  $H(P)$  bound at the cost of setting the code word length equal to the negative logarithm of the specific source word probability.

Surprisingly, under some mild restrictions on the distributions  $P$ , the expectation of  $K(x)$  achieves  $H(P)$ . We can view the  $K(x)$ 's as the code word length set of a “universal” Shannon-Fano code based on the universal probability, Theorem 1.11.2 on page 75. The expectation of  $K(x)$  differs from  $H(P)$  by a constant depending on  $P$ . Namely,  $H(P) = \sum_x P(x)K(x) + c_P$ , where the constant  $c_P$  depends on the length of the program to compute the distribution  $P$ . In Exercise 8.1.6 this dependence is removed.

If the set of outcomes is infinite, then it is possible that  $H(P)$  is infinite. For example, with  $x \in \mathcal{N}$  and  $P(x) = 1/(x \log x)$  we have  $H(P) > \sum_x 1/x$  which diverges. If the expected  $K(x)$  is close to  $H(P)$  then it diverges as well. Two diverging quantities are compared by looking at their quotient or difference. The latter allows us to express really small distinctions.

**Theorem 8.1.1** *Let  $H(P) = -\sum_x P(x) \log P(x)$  be the entropy of a recursive probability distribution  $P$ . Then*

$$0 \leq \left( \sum_x P(x)K(x) - H(P) \right) \leq c_P,$$

*where  $c_P$  is a constant that depends only on  $P$ .*

**Proof.** Since  $K(x)$  is the code word length of a prefix-code for  $x$ , the first inequality of the Noiseless Coding Theorem 1.11.2 on page 75 states that

$$H(P) \leq \sum_x P(x)K(x).$$

By Theorem 4.3.1 on page 247, if  $\mathbf{m}$  is the universal enumerable distribution, then  $P(x) \leq 2^{K(P)}\mathbf{m}(x)$ . By Theorem 4.3.3 on page 253, we

have  $-\log \mathbf{m}(x) = K(x) + O(1)$ . Together this shows that  $-\log P(x) \geq K(x) - K(P) + O(1)$ . It follows that

$$\sum_x P(x)K(x) \leq H(P) + K(P) + O(1).$$

Define the constant  $c_p$  by

$$c_p := K(P) + O(1),$$

and the theorem is proven. Note that the constant implied in the  $O(1)$  term depends on the lengths of the programs occurring in the proof of Theorem 4.3.3. These depend only on the reference universal prefix machine.  $\square$

In other words, statistical entropy of a random variable  $X$  taking outcomes in  $\mathcal{N}$  is equal, within an additive constant, to the expected value of complexity of elements in  $\mathcal{N}$ , or individual entropy.

Above we required  $P(\cdot)$  to be recursive. Actually, we only require that  $P$  be an enumerable function, which is a weaker requirement than recursivity. However, together with the condition that  $P(\cdot)$  is a probability distribution,  $\sum P(x) = 1$ , this means that  $P(\cdot)$  is recursive by Example 4.3.2.

What does this mean for the plain  $C(\cdot)$  complexity? Since  $K(x) \leq C(x) + K(C(x))$  by Example 3.1.3 on page 194, also  $-\log P(x)$  and  $C(x)$  are close to each other with high probability. Substituting  $K(x)$  in Theorem 8.1.1 we find

$$\begin{aligned} -c_p &\leq H(P) - \sum_x P(x)C(x) \\ &\leq \sum_x P(x)K(C(x)), \end{aligned} \tag{8.1}$$

which is bounded only if  $\sum_x P(x)K(C(x))$  converges. Let  $P$  run through a sequence of distributions  $P_1, P_2, \dots$  such that  $H(P_i) \rightarrow \infty$  for  $i \rightarrow \infty$  and  $\lim_{i \rightarrow \infty} K(P_i)/H(P_i) = 0$ . Then, we find that  $H(P)$  is asymptotically equal to the expected complexity  $\sum_x P(x)C(x)$ , as was also established in a weaker form in Example 2.8.1 on page 181.

**Theorem 8.1.2** *Assume the notation above.  $\lim_{H(P) \rightarrow \infty} \sum_x P(x)C(x)/H(P) = 1$ .*

**Proof.** Using Equation 8.1 and the definition of  $\lim_{H(P) \rightarrow \infty}$  above, we only need to establish

$$\lim_{H(P) \rightarrow \infty} \frac{\sum_x P(x)K(C(x))}{H(P)} = 0.$$

This follows from the fact that by Theorem 8.1.1

$$\lim_{H(P) \rightarrow \infty} \frac{\sum_x P(x)K(x)}{H(P)} = 1,$$

while  $C(x) \leq \log x + O(1)$  and  $\sum_{l(x)=n} P(x|l(x)=n)K(x) = \Theta(\log n)$ .  $\square$

Since the expected complexities  $C(\cdot)$  and  $K(\cdot)$  are asymptotically equal to the entropy, the intended interpretation of  $C(x)$ , and especially  $K(x)$ , as a measure of the information content of an individual object  $x$  is supported by a tight quantitative relationship to Shannon's probabilistic notion.

**Example 8.1.1** We examine the expected complexity over a set  $A \subseteq \mathcal{N}$  of source words. Assume that  $A$  is recursively enumerable. Applying the proof of Theorem 8.1.1 using Corollary 4.3.2 on page 255 yields

$$0 \leq \left( \sum_{x \in A} P(x|x \in A)K(x|A) \right) - H(P(\cdot|A)) \leq K(P) + O(1). \quad (8.2)$$

If  $A$  is finite, then the expected complexity and the entropy of the conditional probability  $P(x|x \in A)$  over the set of outcomes in  $A$  are always finite. A natural case is to consider the set  $A$  of all  $x$  of length  $n$  with  $A = \{0, 1\}^n$ . Then since  $K(x|A) = K(x|n) + O(1)$ , we find that the expected prefix complexity of strings of length  $n$  for each recursive distribution  $P$  equals the entropy of the conditional probability plus a term  $K(n)$ ,

$$0 \leq \left( \sum_{l(x)=n} P(x|l(x)=n)K(x|n) \right) - H(P(\cdot|l(x)=n)) \leq c_P, \quad (8.3)$$

where the constant  $c_P$  depends on  $P$  but not on  $n$ , that is,  $c_P \leq K(P) + O(1)$ . As a consequence, also

$$\begin{aligned} H(P(\cdot|l(x)=n)) &\leq \sum_{l(x)=n} P(x|l(x)=n)K(x) \\ &\leq H(P(\cdot|l(x)=n)) + K(n) + c_P, \end{aligned} \quad (8.4)$$

where the left inequality follows a fortiori from the Noiseless Coding Theorem 1.11.2 on page 75 since the  $K(x)$ 's are prefix-code word lengths. (Since for finite sets the entropy is finite we can here compare the quantities involved rather than bound their difference as above.) Note that  $K(n)$  can be very small for regular  $n$  and is bounded from above by  $\log n + 2 \log \log n$  for all  $n$ .  $\diamond$

**Exercises**

**8.1.1.** [28] Show that  $\lim_{H(P) \rightarrow \infty} \sum_x P(x)f(x)/H(P) = 1$ , for  $f(x) = \log x + 2\log\log x$  and  $P$  runs over all nondecreasing probability distributions over  $\mathcal{N}$ .

**8.1.2.** [25] Let  $P$  be a probability distribution.

(a) Show that  $H(P(\cdot|l(x) = n)) - \sum_{l(x)=n} P(x|l(x) = n)C(x) \leq \log n + 2\log\log n + O(1)$ , where the  $O(1)$  term is independent of  $P$  and  $n$ .

(b) Show that  $H(P(\cdot|l(x) = n)) - \sum_{l(x)=n} P(x|l(x) = n)C(x|n) \leq \log n + O(1)$ , where the  $O(1)$  term is independent of  $P$  and  $n$ .

*Comments.* Hint Item (a): by  $K(x) \leq C(x) + K(C(x))$  and the left side of Equation 8.4 we have  $H(P(\cdot|l(x) = n)) - \sum_{l(x)=n} P(x|l(x) = n)C(x) \leq \sum_{l(x)=n} P(x|l(x) = n)K(C(x))$ . Hint Item (b): with  $l(x) = n$  we have  $K(x|n) \leq C(x|n) + \log n + O(1)$ . Namely,  $C(x|n)$  can be given in precisely  $\log n$  bits, given  $n$ . Hence there is a self-delimiting description of  $x$  given  $n$  of length  $C(x|n) + \log n + O(1)$ . Source of Item (b): B. List [Personal communication 1996].

**8.1.3.** [19] Let  $P$  be a (not necessarily recursive) probability distribution and  $A \subseteq \mathcal{N}$  be recursively enumerable, such that for each pair  $x, x' \in A$  we have  $P(x|x \in A) \geq P(x'|x \in A)$  iff  $C(x|A) \leq C(x'|A)$ . Show that  $\sum_{x \in A} P(x|x \in A)C(x|A) \leq H(P(\cdot|x \in A)) + O(1)$ , where the  $O(1)$  term is independent of  $P$ .

*Comments.* Hint: let  $x^1, x^2, \dots$  be an enumeration of  $A$  by nondecreasing  $C(\cdot|A)$  complexity. Then,  $C(x^i|A) \leq \log i + O(1)$  follows immediately by observing that one out of the  $i+1$  first recursively enumerated strings of  $A$  must come after  $x^i$  by the pigeon hole principle. This gives  $\sum_{x \in A} P(x|x \in A)C(x|A) \leq \sum_{1 \leq i \leq d(A)} P(x^i|x \in A) \log i + O(1) \leq H(P(\cdot|x \in A)) + O(1)$  with the  $O(1)$  term independent of  $P$  and  $A$ . (The last inequality follows from  $P(x^i|x \in A) \leq \sum_{1 \leq j \leq i} P(x^j|x \in A)/i$  yielding  $\log(iP(x^i|x \in A)) \leq \log 1 = 0$ .) Source: generalization of an idea of B. List [Personal communication 1996].

**8.1.4.** [21] Let  $P$  be a recursive distribution of  $m$  identically distributed random variables. Show that the analogue of Theorem 8.1.1 holds. In this case, the entropy is proportional to  $m$ , but  $c_P$  is  $O(\log m)$ .

*Comments.* Hint: use  $K(P) = O(\log m)$  and a similar proof to Theorem 8.1.1. Source: T.M. Cover, P. Gács, and R.M. Gray, *Ann. Probab.*, 17:3(1989), 840–865. This is isomorphic to Exercise 8.1.2.

**8.1.5.** [21] Let  $P$  be a recursive distribution of  $m$  differently distributed random variables. Show we can derive the analogue of Theorem 8.1.1. In this case, the entropy is proportional to  $m$ , but  $c_P$  is  $O(m)$ .

**8.1.6.** [32] Let  $P$  be a (not necessarily recursive) distribution and  $A \subseteq \mathcal{N}$ .

- (a) Show that  $0 \leq [\sum_{x \in A} P(x|A)K(x|P(\cdot|x \in A))] - H(P(\cdot|x \in A)) = O(1)$ , where the  $O(1)$  term is independent of  $A$  and  $P$  and depends only on the reference prefix machine. This is a stronger result than Equation 8.2 where  $c_P$  is allowed to depend on  $P$ .
- (b) Show that  $0 \leq [\sum_{x \in A} P(x|x \in A)K(x|P(\cdot|x \in A))] - H(P(\cdot|x \in A)) < 1$  for all  $A$  and  $P$  for some appropriate reference prefix machines. This achieves exactly the optimum expected code word length of the Noiseless Coding Theorem 1.11.2 on page 75.

*Comments.* Note that Item (a) reduces to Equation 8.3 for  $A = \{0, 1\}^n$  and  $P$  the uniform distribution. Hint: in Item (a) use a universal prefix machine with a table of  $(x, P(x|x \in A))$  pairs for all  $x \in A$  on its auxiliary conditional information tape. With an  $O(1)$  program to compute a Shannon-Fano code, this machine when given an input  $p$  determines whether  $p$  is the Shannon-Fano code word for some  $x$ . By Lemma 4.3.3 on page 253 such a code word has length  $-\log P(x|x \in A) + O(1)$ . If this is the case, then the machine outputs  $x$ , otherwise it halts without output. Therefore,  $K(x|P(\cdot|x \in A)) \leq -\log P(x|x \in A) + O(1)$ . This shows the upper bound. The lower bound follows as usual from the Noiseless Coding Theorem 1.11.2. Item (b) follows by appropriate modification of the reference machine. Source: the basic Equation 8.2 appears in [P. Gács, *Lecture Notes on Descriptive Complexity and Randomness*, Manuscript, Boston University, 1987; T.M. Cover, P. Gács, and R.M. Gray, *Ann. Probab.*, 17:3(1989), 840–865; W.H. Zurek *Phys. Rev. A*, 40(1989), 4731–4751] but is presumably older. A complex argument for Item (a) in a physics setting is [C. Caves, pp. 91–115 in: W.H. Zurek, Ed., *Complexity, Entropy and the Physics of Information*, Addison-Wesley, 1991]. Item (b) is spelled out in [R. Schack, ‘Algorithmic information and simplicity in statistical physics,’ *Phys. Rev. E*, submitted]. The latter two papers use the nonstandard complexity variant of Example 3.9.1, but this is not required to prove the result.

**8.1.7.** [39] Examine the lossless compression of individual sequences over alphabet  $A$  by a class  $\Sigma$  of recursive bijective binary prefix codes (called “sequential codes” in the sequel) with a recursive inverse, defined as follows: A parsing rule  $\pi$  is defined as a sequence of pairs  $\{(W_n, \Psi_n) : n \in \mathcal{N}\}$  where  $W_n$  is a finite prefix subset of  $A^*$ , and  $\Psi_n$  is a mapping from  $W_n$  to  $\mathcal{N}$ . Applying  $\pi$  to a source sequence  $\omega$ , one obtains a sequence of selection indexes  $i_1, i_2, \dots$  in  $\mathcal{N}$  and a parse sequence  $u_1, u_2, \dots$  over  $A$ , such that  $\omega = u_1 u_2 \dots$  and the following holds. Initial index  $i_1 = 1$  forces segment  $u_1$  to be selected from  $W_1$ ; for all  $t > 1$  selection index  $i_t = \Psi_{i_{t-1}}(u_{t-1})$  forces segment  $u_t$  to be selected from  $W_{i_t}$ . A *sequential code* is a sequence of triples  $\sigma = \{(W_n, \Psi_n, \Phi_n) : n \in \mathcal{N}\}$  such that  $\pi_\sigma = \{(W_n, \Psi_n) : n \in \mathcal{N}\}$  is a parsing rule and for each  $n \in \mathcal{N}$ ,  $\Phi_n$  is a one-to-one mapping of  $W_n$  onto a prefix subset of  $\{0, 1\}^n$ . A sequential code  $\sigma$  encodes an individual sequence  $\omega$  into a sequence  $\zeta \in \{0, 1\}^*$

as follows: First the  $\pi_\sigma$ -parsing of  $\omega$  is generated as defined above. Then  $\zeta = \phi_{i_1}(u_1)\phi_{i_2}(u_2) \dots$ . This class of sequential codes includes all finite state codes and the Lempel-Ziv code. Given an infinite sequence  $\omega$  over an  $\alpha$ -ary alphabet, we are interested in the behavior of the OPTA code  $\rho(\omega) = \inf_{\sigma \in \Sigma} \rho(\omega|\sigma)$ ,  $\rho(\omega|\sigma) = \limsup_{n \rightarrow \infty} \sigma(\omega_1 \dots \omega_n)/n$ .

- (a) Show that this class does not contain a universal element that compresses each sequence at least as well as any other element in the class (at least as well as the OPTA code). That is, is there a universal code  $\sigma \in \Sigma$  such that  $\rho(\omega|\sigma) = \rho(x)$  for all  $\omega$ ?
- (b) Show that the OPTA compression function is bounded from below by the Kolmogorov complexity compression function, and that for some individual sequences the OPTA compression is strictly less than the Kolmogorov compression. That is, what is the relation between  $\rho(x)$  and  $k(\omega) = \limsup_{n \rightarrow \infty} K(\omega_1 \dots \omega_n)/n$  (where  $K(\cdot)$  is the Kolmogorov complexity function).

*Comments.* Source: J.C. Kieffer and E. Yang, *IEEE Trans. Inform. Theory*, IT-42:1(1996), 29–39.

## 8.2 Reversible Computation

Computers can be regarded as engines that must dissipate energy in order to process information. Von Neumann reputedly thought that a computer operating at temperature  $T$  must dissipate at least  $kT \ln 2$  joule per elementary bit operation, where  $k \approx 1.38 \times 10^{-23}$  joule/°kelvin is Boltzmann’s constant and  $T$  is the absolute temperature in °kelvin. This is about  $2.8 \times 10^{-21}$  joule at room temperature.

Around 1960, R. Landauer more thoroughly analyzed this question and concluded that it is only “logically irreversible” operations that must dissipate energy. An operation is *logically reversible* if its inputs can always be deduced from the outputs. Erasure of information is not reversible. Erasing each bit costs  $kT \ln 2$  energy when the computer operates at temperature  $T$ .

### 8.2.1 Energy Dissipation

Briefly, Landauer’s line of reasoning runs as follows: Distinct logical states of a computer must be represented by distinct physical states of the computer hardware. Each bit has one degree of freedom (0 or 1), corresponding to one or more degrees of freedom in the physical hardware. The  $n$  bits collectively have  $n$  degrees of freedom. This corresponds to  $2^n$  possible logical states and hence to at least  $2^n$  physical states of the hardware.

Suppose  $n$  bits are irreversibly erased (reset to zeros). Before the erasure operation, these  $n$  bits could be in any of  $2^n$  possible logical states. After the erasure, they are compressed to just one unique state. According

to the second law of thermodynamics, this loss of degrees of freedom of the physical system must be compensated for. If the information is not simply transmitted from the system to the environment, but cannot be retrieved any more, then the compensation must happen by an increase of temperature in the system and its environment, that is, by heat dissipation.

Thus, the only computer operations that are necessarily thermodynamically costly are those that are logically irreversible, that is, operations that map several distinct logical states of the computer onto a common successor, thereby throwing away information about the computer's previous state. Any computation that discards information irreversibly costs energy. For example, a logic gate with more input lines than output lines inevitably loses information, and hence is irreversible and therefore dissipative.

The ultimate limits of miniaturization of computing devices, and therefore the speed of computation, are governed by an unavoidable heat increase through energy dissipation. Such limits have already been reached by current high-density electronic chips. The question of how to reduce the energy dissipation of computation determines future advances in computing power. Since battery technology improves by only twenty percent every ten years, low-power computing will similarly govern advances in mobile communication and computing.

Over the last fifty years, the energy dissipation per logic operation has been reduced by approximately a factor of ten every five years. Such an operation in 1988 dissipated about 1/10th picojoule (1 picojoule is  $10^{-12}$  joule) versus around  $10^9$  picojoule in 1945 [R.W. Keyes, *IBM J. Res. Devel.*, 32(1988), 24–28].

Extrapolations of current trends suggest that reduction of the energy dissipation per logic operation below  $kT$  (thermal noise, on the order of  $10^{-8}$  picojoule at room temperature) becomes a relevant issue by 2015. For example, a computer using  $kT$  per operation at room temperature at a frequency of 1 gigahertz with  $10^{18}$  logic gates packed in a cubic centimeter will dissipate about 3 megawatts. This is difficult to cool.

The drive for ever greater computing power through more densely packed logic circuits will eventually require that we find methods other than cooling. An alternative way is to develop reversible logic that computes (almost) without energy dissipation.

Logical reversibility does not imply dissipation freeness. A computer may compute in a logically reversible manner and yet dissipate heat. But the laws of physics allow for technologies to make logical reversible computers operate in a dissipationless manner. Logically reversible are computers built from Fredkin gates, or the reversible Turing machine, discussed later. Thought experiments exhibit a computer that is both

reversible and dissipationless. An example is the billiard ball computer discussed below.

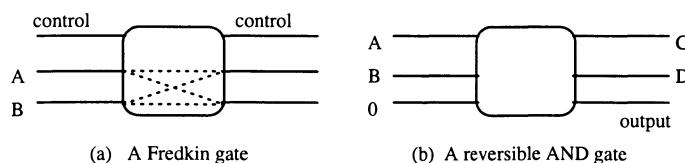
### 8.2.2 Reversible Logic Circuits

To build a computer we only need two types of logical gates: AND, NOT; all other gates can be built from AND and NOT. The NOT gate is already logically reversible. But the AND gate is not. The AND and OR gates have each two inputs but only one output. Hence, they must lose information.

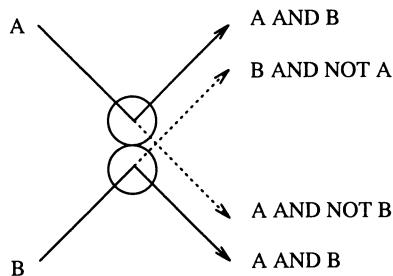
Can we design something like an AND gate that has the same number of input lines and output lines and is reversible? In principle, we only need to add garbage output lines. This is easy, and there are many such gates. However, we would like to find a universal reversible gate from which all Boolean gates can be synthesized. One such universal reversible gate is the so-called Fredkin gate as shown in Figure 8.1(a). If the control bit is 0, then the input values of  $A$  and  $B$  are transmitted unaltered to the outputs on the same level; and if the control bit is 1, then they are switched to the opposite output.

The Fredkin gate is universal in the sense that it can be used to construct all other Boolean gates in a reversible variant. For example, Figure 8.1(b) shows a reversible AND gate built from a Fredkin gate. Here, the top wire marked by  $A$  and  $C$  is the control wire of the Fredkin gate. The bottom wire marked “0” has a constant “0” input. The inputs to the AND gate are marked  $A$  and  $B$ , the single output wire is marked “output.” If  $A = 0$ , then the constant input 0 is transmitted unaltered to the output. If  $A = 1$ , then the input from  $B$  is transmitted unaltered to the output. Hence, the output is  $A$  AND  $B$ .

Fredkin’s gate and the AND gate in Figure 8.1 are reversible because we have added “garbage” output bits in order to ensure that the inputs can always be deduced from the outputs. Therefore, in principle they do not need to dissipate energy. Theoretically, building a dissipationless computer using Fredkin gates or reversible AND and NOT gates is possible. There is no *thermodynamic* reason that such a computer should cost any energy. A computation would be started by an energy jolt, say some amount of electricity, proceed without dissipating energy, and at



**FIGURE 8.1.** Reversible Boolean gates



**FIGURE 8.2.** Implementing reversible AND gate and NOT gate

the end of the computation simply return the same amount of electricity to the source.

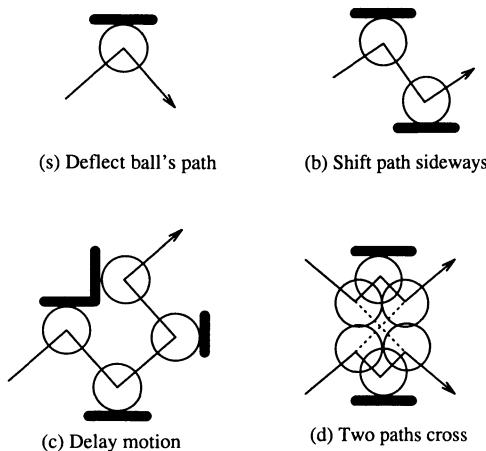
### 8.2.3 Reversible Ballistic Computer

Consider an idealized computer using elastic frictionless billiard balls. The presence of a ball represents a 1 and no ball represents a 0. The ballistic computer contains “mirrors” at some positions that perfectly elastically reflect the balls. All collisions between balls are perfectly elastic as well. Between the collisions, the balls travel in straight lines with constant speed, as in Newtonian mechanics. For example, we can think of the balls as idealized molecules. To start the computation, if an input bit is 1, we fire a ball, if an input bit is 0, we do not fire a ball. All input balls are fired simultaneously, and at the same speed.

Figure 8.2 implements an AND gate for inputs  $A$  and  $B$ . If we set  $B = 1$ , then we obtain a NOT gate for  $A$  (and setting  $A = 1$  yields a NOT gate for  $B$ ). We will also need the constructions in Figure 8.3 using mirrors to deflect a ball’s path, to shift a path, to delay the ball’s motion without changing its final direction, and to allow two lines of motion to cross.

It is possible to emulate any computation using the above gadgets. Suppose the setup lets all of the balls reach the output end simultaneously. After we observe the output, we can simply reflect back all the output balls, including the many “garbage balls,” to reverse the computation. The billiard balls will then come out of the ballistic computer exactly where we sent them in, with the same speed. The kinetic energy can be returned to the device that kicked the balls in. Then the computer is ready for a next round of dissipationless action. An example input-output+garbage scheme for a ballistic computer is shown in Figure 8.4.

A ballistic computer with molecules as balls (the molecular computer) is extremely unstable in principle. Any small initial error in the position or speed of the balls is amplified by a factor of 2 in each collision. Due to the quantum-mechanical uncertainty relation there are unavoidable errors of size at least about Planck’s constant. After a few dozen collisions, the whole computer deteriorates.



**FIGURE 8.3.** Controlling billiard ball movements

To prevent deterioration, one may use objects at the quantum level that are more stable. R.P. Feynman suggested the use of electron spin orientation in [*Int. J. Theoret. Phys.*, 21(1982), 467–488; *Optics News*, 11(1985), 11]. Other quantum-mechanics computer proposals are [P.A. Benioff, *Int. J. Theoret. Phys.*, 21(1982), 177–202; *Ann. New York Acad. Sci.*, 480(1986), 475–486; N. Margolus, in W.H. Zurek (Ed.), *Complexity, entropy and the physics of information*, Addison-Wesley, 1991, pp. 273–287].

C.H. Bennett [*Int. J. Theoret. Phys.*, 21:2(1982), 905–940] suggested a Brownian computer, where the computing particles are larger than molecules but small enough to be subject to Brownian motion. The idea is to use thermal noise to drive the computing balls. The computing balls can move on a fixed computation trajectory between start and finish of the computation. By Brownian motion they perform a random walk, back and forth along the trajectory, until they finally arrive at the finish. These computers do dissipate energy, but this can in principle be made arbitrarily small.

C.H. Bennett and R. Landauer [*Scientific American*, July 1985, 48–56] also described how to construct an idealized Fredkin gate and given such gates, another alternative way of constructing a billiard ball computer. All the balls are linked together and pushed forward by one mechanism. The balls move along pipes. The entire assembly is immersed in an ideal viscous fluid. The frictional forces that act on the balls will be proportional to their velocity; there will be no static friction.

If we move the balls very slowly, then the frictional force will be very weak. The energy that must be expended to work against friction is equal to the product of the frictional force and the distance the ball traveled. Thus, we can use as little energy as we wish simply by slowing down the computation. There is no minimum amount of energy that must be expended to perform any computation.

In [R.C. Merkle, *Nanotechnology*, 4(1993), 21–40], two methods to realize such reversible computations using electronic switching devices in conventional

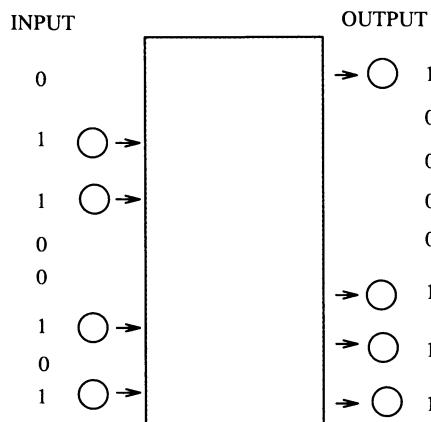


FIGURE 8.4. A billiard ball computer

technologies (like nMOS, CMOS, and Charge Coupled Devices) are proposed. In fact, such reversible almost dissipationless electronics have already been used in existing laboratory computers. (They happened to have been designed with other aims in mind.)

#### 8.2.4 Reversible Turing Machines

Reversibility of a Turing machine's transition function (Section 1.7) can be guaranteed by requiring disjointness of the ranges of the quadruples, just as determinism is guaranteed by requiring disjointness of their domains. To assure that the machine's global input-output relation is also one-to-one, it is necessary to impose a standard format on the initial and final instantaneous descriptions. In particular, one requires that all storage space used during the computation, other than that used for containing the input and output strings, be blank at the beginning and end of the computation.

Let  $T_1, T_2, \dots$  be the standard enumeration of Turing machines with input tape symbols 0 and 1 only (no blanks). Each such Turing machine has the property that the set of programs for which it halts is a prefix set (no element in the set is a proper prefix of another element in the set). This is the prefix machine defined in Chapter 3. These (in general irreversible) prefix machines compute all and only partial recursive prefix functions  $\phi_1, \phi_2, \dots$  (Definition 3.1.1 on page 192).

Consider the model of Turing machine in Section 1.7. The elementary operations were rules in quadruple format  $(p, s, a, q)$ . Quadruples are said to *overlap in domain* if they cause the machine in the same state and scanning the same symbol to perform different actions. A *deterministic Turing machine* is defined as a Turing machine with quadruples that pairwise do not overlap in domain.

Now consider the special format (deterministic) Turing machines using quadruples of two types: *read/write* quadruples and *move* quadruples. A *read/write* quadruple  $(p, a, b, q)$  causes the machine in state  $p$  scanning tape symbol  $a$  to write symbol  $b$  and enter state  $q$ . A *move* quadruple  $(p, 1, \sigma, q)$  causes the machine in state  $p$  to move its tape head by  $\sigma \in \{-1, 0, +1\}$  squares and enter state  $q$ , oblivious to the particular symbol in the currently scanned tape square. (Here “ $-1$ ” means “one square left,” “ $0$ ” means “no move” and “ $+1$ ” means “one square right.”) Quadruples are said to *overlap in range* if they cause the machine to enter the same state and either both write the same symbol or (at least) one of them moves the head. Said differently, quadruples that enter the same state overlap in range unless they write different symbols. A *reversible Turing machine* is a deterministic Turing machine with quadruples that pairwise do not overlap in range. A  $k$ -tape reversible Turing machine uses  $(2k+2)$  tuples that for each tape separately, select a *read/write* or *move* on that tape. Moreover, where the quadruples resulting from the restriction to each tape separately have pairwise nonoverlapping ranges.

To show that each partial recursive function can be computed by a reversible Turing machine one can proceed as follows. Take the standard irreversible Turing machine computing that function. We modify it by adding an auxiliary storage tape called the “history tape.” The quadruple rules are extended to 6-tuples to additionally manipulate the history tape. To be able to reversibly undo (retrace) the computation deterministically, the new 6-tuple rules have the effect that the machine keeps a record on the auxiliary history tape consisting of the sequence of quadruples executed on the original tape. Reversibly undoing a computation entails also erasing the record of its execution from the history tape.

This notion of reversible computation means that only one-to-one recursive functions can be computed. To reversibly simulate  $t$  steps of an irreversible computation from  $x$  to  $f(x)$  one reversibly computes from input  $x$  to output  $\langle x, f(x) \rangle$ . Say this takes  $t' = O(t)$  time. Since this reversible simulation at some time instant has to record the entire history of the irreversible computation, its space use increases linearly with the number of simulated steps  $t$ . That is, if the simulated irreversible computation uses  $s$  space, then for some constant  $c > 1$  the simulation uses  $t' \approx c + ct$  time and  $s' \approx c + c(s+t)$  space. After computing from  $x$  to  $f(x)$  the machine reversibly copies  $f(x)$ , reversibly undoes the computation from  $x$  to  $f(x)$  erasing its history tape in the process, and ends with one copy of  $x$  and one copy of  $f(x)$  in the format  $\langle x, f(x) \rangle$  and otherwise empty tapes.

We list some relevant properties. Let  $\psi_i$  be the partial recursive function computed by the  $i$ th such reversible prefix machine. Let  $\langle \cdot \rangle$  be a bijective recursive pairing mapping over the integers, that can be decoded from

left to right by a prefix machine. Among the more important properties of reversible prefix machines are the following:

**Universal reversible prefix machine** There exists a reversible prefix machine that is universal, say  $UR$ , computing  $\psi_0$ , such that for all  $k$  and  $x$ , we have that  $\psi_0(\langle k, x \rangle) = \langle k, \psi_k(x) \rangle$ .

**Irreversible to reversible** Two irreversible algorithms, one for computing  $y$  from  $x$  and the other for computing  $x$  from  $y$ , can be efficiently combined to obtain a reversible algorithm for computing  $y$  from  $x$ . More formally, for any two indices  $i$  and  $j$  one can effectively obtain an index  $k$  such that for any strings  $x$  and  $y$ , if  $\phi_i(x) = y$  and  $\phi_j(y) = x$ , then  $\psi_k(x) = y$ .

**Saving input copy** From any index  $i$  one may obtain an index  $k$  such that  $\psi_k$  has the same domain as  $\phi_i$  and for every  $x$ , we have  $\psi_k(x) = \langle x, \phi_i(x) \rangle$ . In other words, an arbitrary prefix machine can be simulated by a reversible one that saves a copy of the irreversible machine's input in order to assure a global one-to-one mapping.

**Efficiency** The above simulation can be performed rather efficiently. In particular, for any  $\epsilon > 0$  one can find a reversible simulating machine that runs in time  $O(t^{1+\epsilon})$  and space  $O(s \log t)$  compared to the time  $t$  and space  $s$  of the irreversible machine being simulated.

**One-to-one functions** From any index  $i$  one may effectively obtain an index  $k$  such that if  $\phi_i$  is one-to-one, then  $\psi_k = \phi_i$ . The reversible Turing machines  $\{\psi_k\}$ , therefore, provide an effective enumeration of all one-to-one partial recursive functions.

**Definition 8.2.1** Let  $\psi(p, x)$  be a partial recursive function satisfying the following: For every  $p$ , the function  $\psi(p, x)$  is one-to-one as a function of  $x$ ; for every  $x$ , the set of  $p$ 's such that  $\psi(p, x) < \infty$  is a prefix set; and for each  $y$  the set of  $p$ 's such that for some  $x$  we have  $\psi(p, x) = y$  is a prefix set. Then  $\psi$  is a *reversible prefix partial recursive function*.

Such a  $\psi$  may be thought of as the function computed by a reversible prefix machine that performs a one-to-one mapping on  $x \leftrightarrow y$  under the control of a program  $p$  that acts like a catalyst in that it remains on the program tape throughout the computation. Such a machine has a one-way read-only program tape initially containing program  $p$ , a read-only conditional data tape initially containing input  $x$ , and a one-way write-only output tape containing  $y$  on termination.

Any other work tapes used during the computation are supplied in blank condition at the beginning of the computation and must be left blank at

the end of the computation. The program tape's head begins and ends scanning the leftmost square of the program, which is self-delimiting both for forward computations from each input  $x$  as well as for backward computations from each output  $y$ .

A *universal reversible prefix machine UR* whose program size is minimal to within an additive constant can readily be shown to exist.

**Definition 8.2.2** The *reversible prefix complexity* is  $KR(y|x) := \min\{l(p) : UR(p, x) = y\}$ .

## Exercises

**8.2.1.** [20] Construct reversible OR, NOT, and XOR gates using Fredkin gates.

**8.2.2.** [31] Show the existence of a universal reversible Turing machine.

*Comments.* Source: C.H. Bennett, *IBM J. Res. Develop.*, 17(1973), 525-532. If the original computation takes  $t$  steps and uses  $s$  space then Bennett's simulation requires  $t' = \Theta(t)$  steps and  $s' = \Theta(st)$  space.

**8.2.3.** [34] Let an irreversible computation use  $t$  steps and  $s$  space.

(a) Show how to simulate it reversibly using  $t' = \Theta(t^{1+\epsilon}/s^\epsilon)$  steps and  $s' = \Theta(c(\epsilon)s(1 + \log t/s))$  space with  $c(\epsilon) = \epsilon 2^{1/\epsilon}$  for any  $\epsilon > 0$  using always the same simulation method with different parameters. Typically,  $\epsilon = \log 3$ .

(b) Show that Item (a) implies that each irreversible computation using  $s$  space can be simulated by a reversible computation using  $s^2$  space in  $t' = \Theta(t^{1+\epsilon})$  time.

*Comments.* Hint: in Item (a) do not save the entire history of the irreversible computation, but break up the simulated computation into segments of about  $s$  steps and save in a hierarchical manner *checkpoints* consisting of complete instantaneous descriptions of the simulated machine (entire tape contents, tape heads positions, state of the finite control). After a later checkpoint is reached and saved, the simulating machine reversibly undoes its intermediate computation reversibly erasing the intermediate history and reversibly canceling the previously saved checkpoint. Subsequently, the computation is resumed from the new checkpoint onwards. The reversible computation simulates  $k^n$  segments of length  $m$  of irreversible computation in  $(2k - 1)^n$  segments of length  $\Theta(m + s)$  of reversible computation using  $n(k - 1) + 1$  checkpoint registers using  $\Theta(m + S)$  space each, for each  $k, n, m$ . This way it is established that there are various tradeoffs possible in time-space in between  $t' = \Theta(t)$  and  $s' = \Theta(ts)$  at one extreme ( $k = 1, m = t, n = 1$ ) and the  $t'$  and  $s'$  in Item (a) using always the same simulation method but with different parameters  $k, n$  where  $\epsilon = \log_k(2k - 1)$  and  $m = \Theta(s)$ .

Typically, for  $k = 2$  we have  $\epsilon = \log 3$ . Since for  $t > 2^s$  the machine goes into a computational loop, we always have  $s \leq \log t$ . This shows Item (b). Source: for Items (a) and (b) see [C.H. Bennett, *SIAM J. Comput.*, 18(1989), 766–776; R.Y. Levine and A.T. Sherman, *SIAM J. Comput.*, 19(1990), 673–677]. In [M. Li and P.M.B. Vitányi, *Proc. Royal Soc. London, Ser. A*, 452(1996), 769–789], reversible simulations are investigated using reversible pebble games. It is shown that the simulations as in Item (a) are optimal in the sense of simulating the greatest number of steps using least space overhead  $s' - s$ . It is shown that the space overhead can be reduced at the cost of limited irreversible erasing. By Landauer’s principle this implies a space-energy tradeoff.

## 8.3 Information Distance

Kolmogorov complexity is a measure of absolute information content of individual objects. It is desirable to have a similar measure of absolute information distance between individual objects. Such a notion of universal informational “distance” between two strings is the minimal quantity of information sufficient to translate between  $x$  and  $y$ , generating either string effectively from the other. Such a notion should be asymptotically machine-independent and can serve as an absolute measure of the informational, or “picture,” distance between discrete objects  $x$  and  $y$ .

### 8.3.1 Max Distance

Intuitively, the minimal information distance between  $x$  and  $y$  is the length of the shortest program for a universal computer to transform  $x$  into  $y$  and  $y$  into  $x$ . This program then functions in a “catalytic” manner, being retained in the computer before, during, and after the computation. This measure will be shown to be, up to a logarithmic additive term, equal to the *maximum* of the conditional Kolmogorov complexities.

**Example 8.3.1** The conditional complexity  $K(y|x)$  itself is unsuitable as optimal information distance because it is asymmetric:  $K(\epsilon|x)$ , where  $\epsilon$  is the empty string, is small for all  $x$ , yet intuitively a long random string  $x$  is not close to the empty string. The asymmetry of the conditional complexity  $K(x|y)$  can be remedied by defining the algorithmic informational distance between  $x$  and  $y$  to be the sum of the relative complexities,  $K(y|x) + K(x|y)$ . The resulting metric will overestimate the information required to translate between  $x$  and  $y$  in case there is some redundancy between the information required to get from  $x$  to  $y$  and the information required to get from  $y$  to  $x$ .  $\diamond$

The example suggests investigating to what extent the information required to compute  $x$  from  $y$  can be made to overlap with that required to compute  $y$  from  $x$ . In some simple cases, complete overlap can be achieved, so that the same minimal program suffices to compute  $x$  from  $y$  as to compute  $y$  from  $x$ .

**Example 8.3.2** If  $x$  and  $y$  are independent random binary strings of the same length  $n$  (up to additive constants  $K(x|y) = K(y|x) = n$ ), then their bitwise exclusive-or  $x \oplus y$  serves as a minimal program for both computations. Similarly, if  $x = uv$  and  $y = vw$  where  $u$ ,  $v$ , and  $w$  are independent random strings of the same length, then  $u \oplus w$  is a minimal program to compute either string from the other.

Now suppose that more information is required for one of these computations than for the other, say,

$$K(y|x) > K(x|y).$$

Then the minimal programs cannot be made identical because they must be of different sizes. Nevertheless, in simple cases, the overlap can still be made complete, in the sense that the larger program (for  $y$  given  $x$ ) can be made to contain all the information in the smaller program, as well as some additional information. This is so when  $x$  and  $y$  are independent random strings of unequal length, for example  $u$  and  $vw$  above. Then  $u \oplus v$  serves as a minimal program for  $u$  from  $vw$ , and  $(u \oplus v)w$  serves as one for  $vw$  from  $u$ .

The programs for going between independent random  $x$  and  $y$  can, if one wishes, be made *completely independent*. For example use  $y$  to go from  $x$  to  $y$ , and  $x$  to go from  $y$  to  $x$ . This may be true, in general, at least to within logarithmic terms.

◇

Up to logarithmic error terms, the information required to translate between two strings can always be represented in the maximally overlapping way of Example 8.3.2. Namely, let

$$\begin{aligned} k_1 &= K(x|y), \\ k_2 &= K(y|x), \\ l &= k_2 - k_1. \end{aligned}$$

The next theorem shows that there is a string  $d$  of length  $k_1 + O(\log k_1)$  and a string  $q$  of length  $l + \log l$  such that  $d$  serves as the minimal program both from  $xq$  to  $y$  and from  $y$  to  $xq$ . This means that the information required to pass from  $x$  to  $y$  is always maximally correlated with the information required to get from  $y$  to  $x$ . It is therefore never the case

that a large amount of information is required to get from  $x$  to  $y$  and a large *but independent* amount of information is required to get from  $y$  to  $x$ . (It is very important here that the time of computation is completely ignored: this is why this result does not contradict the idea of one-way functions.)

The process of going from  $x$  to  $y$  may be broken into two stages. First, add the string  $q$ ; second, use the difference program  $d$  between  $qx$  and  $y$ . In the reverse direction, first use  $d$  to go from  $y$  to  $qx$ ; second, erase  $q$ . Thus, the computation from  $x$  to  $y$  needs both  $d$  and  $q$ , while the computation from  $y$  to  $x$  needs only  $d$  as program.

The foregoing is true of ordinary computations, but if one insists that the computation be performed reversibly, that is by a machine whose transition function is one-to-one, then the full program  $p = dq$  is needed to perform the computation in either direction. This is because reversible computers cannot get rid of unwanted information simply by erasing it as ordinary irreversible computers do. If they are to get rid of unwanted information at all, they must cancel it against equivalent information already present elsewhere in the computer.

Let  $T_1, T_2, \dots$  be the standard enumeration of prefix machines of Example 3.1.1, page 193. For a partial recursive function  $\phi$  computed by  $T$ , let

$$E_\phi(x, y) = \min\{l(p) : \phi(p, x) = y, \phi(p, y) = x\}.$$

**Lemma 8.3.1** *There is a universal prefix machine  $U$  computing  $\phi_0$  such that for each partial recursive prefix function  $\phi$  and all  $x, y$ ,*

$$E_{\phi_0}(x, y) \leq E_\phi(x, y) + c_\phi,$$

where  $c_\phi$  is a constant that depends on  $\phi$  but not on  $x$  and  $y$ .

**Proof.** This is a consequence of the existence of a universal prefix machine and is identical to the proof of Theorem 2.1.1, page 97.  $\square$

Then for every two universal prefix machines computing  $\phi_0$  and  $\psi_0$ , we have for all  $x, y$  that  $|E_{\phi_0}(x, y) - E_{\psi_0}(x, y)| \leq c$ , with  $c$  a constant depending on  $\phi_0$  and  $\psi_0$  but not on  $x$  and  $y$ . We fix a particular universal prefix machine  $U$  as reference machine and define

$$E_0(x, y) = \min\{l(p) : U(p, x) = y, U(p, y) = x\}.$$

We shall prove that up to an additive logarithmic term,  $E_0$  is equal to the “max distance” defined as follows.

**Definition 8.3.1** The *max distance* between  $x, y$  is  $E_1(x, y) = \max\{K(x|y), K(y|x)\}$ .

**Theorem 8.3.1** Let  $C(x|y) = k_1$  and  $C(y|x) = k_2$ , and  $l = k_2 - k_1 \geq 0$ . There is a string  $q$  of length  $l$  such that

$$E_0(qx, y) = K(k_1, k_2) + k_1 + O(1).$$

**Proof.** Given  $k_1, k_2$ , we can enumerate the set  $S = \{(x, y) : C(x|y) \leq k_1, C(y|x) \leq k_2\}$ . Without loss of generality, assume that  $S$  is enumerated without repetition, and with witnesses of length exactly  $k_1$  and  $k_2$ . Now consider a dynamic graph  $G = (V, E)$  where  $V$  is the set of binary strings, and  $E$  is a dynamically growing set of edges that starts out empty.

Whenever a pair  $(x, y)$  is enumerated, we add an edge  $e = \{xq, y\}$  to  $E$ . Here,  $q$  is chosen to be the  $(i2^{-k_1})$ th binary string of length  $l$ , where  $i$  is the number of times we have enumerated a pair with  $x$  as the first element. So the first  $2^{k_1}$  times we enumerate a pair  $(x, )$  we choose  $q = 0^l$ , for the next  $2^{k_1}$  times we choose  $q = 0^{l-1}1$ , etc. The condition  $C(y|x) \leq k_2$  implies that  $i < 2^{k_2}$  hence  $i2^{-k_1} < 2^l$ , so this choice is well-defined.

In addition, we “color” edge  $e$  with a binary string of length  $k_1 + 1$ . Call two edges *adjacent* if they have a common endpoint. If  $c$  is the minimum color not yet appearing on any edge adjacent to either  $xq$  or  $y$ , then  $e$  is colored  $c$ . Since the degree of any node is bounded by  $2^{k_1}$  (when acting as an  $xq$ ) plus  $2^{k_1}$  (when acting as a  $y$ ), a color is always available.

A *matching* is a set of nonadjacent edges. Note that the colors partition  $E$  into at most  $2^{k_1+1}$  matchings, since no edges of the same color are ever adjacent. Since the pair  $(x, y)$  in the statement of the theorem is necessarily enumerated, there is some  $q$  of length  $l$  and color  $c$  such that the edge  $\{xq, y\}$  is added to  $E$  with color  $c$ .

Knowing  $k_1, k_2, c$  and either of the nodes  $xq$  or  $y$ , one can dynamically reconstruct  $G$ , find the unique  $c$ -colored edge adjacent to this node, and output the neighbour. This shows that a self-delimiting program of size  $K(k_1, k_2) + k_1 + O(1)$  suffices to compute in either direction between  $xq$  and  $y$ .  $\square$

**Corollary 8.3.1** Up to a logarithmic additive term,  $E_0(x, y) = E_1(x, y)$ .

We may call this theorem the *Conversion Theorem* since it asserts the existence of a difference string  $d$  that converts both ways between  $x$  and  $y$  and at least one of these conversions is optimal. If  $k_1 = k_2$ , then the conversion is optimal in both directions.

### 8.3.2 Picture Distance

Every reasonable distance to measure similarity of pictures should be an effectively approximable, symmetric, positive function of  $x$  and  $y$

satisfying a reasonable normalization condition and obeying the triangle inequality. It turns out that  $E_1$  is minimal up to an additive constant among all such distances. Hence, it is a universal “picture distance” that accounts for any effective similarity between pictures.

Let us identify digitized black-and-white pictures with binary strings. There are many distances defined for binary strings, for example, the Hamming distance and the Euclidean distance. Such distances are sometimes appropriate. For instance, if we take a binary picture and change a few bits on that picture, then the changed and unchanged pictures have small Hamming or Euclidean distance, and they do look similar.

However, this is not always the case. The positive and negative prints of a photo have the largest possible Hamming and Euclidean distance, yet they look similar in our eyes. Also, if we shift a picture one bit to the right, again the Hamming distance may increase by a lot, but the two pictures remain similar.

Many approaches to pattern recognition try to define picture similarity. Section 8.3.1 has given evidence that  $E_1(x, y) = \max\{K(x|y), K(y|x)\}$  is a natural way to formalize a notion of the minimal algorithmic informational distance between  $x$  and  $y$ . Let us show that the distance  $E_1$  is, in a sense, minimal among all reasonable similarity measures.

A distance measure must be nonnegative for all  $x \neq y$ , zero for  $x = y$ , symmetric, and satisfy the triangle inequality. This is not sufficient since a distance measure like  $D(x, y) = 1$  for all  $x \neq y$  must be excluded. For each  $x$  and  $d$ , we want only finitely many elements  $y$  at a distance  $d$  from  $y$ . Exactly how fast we want the distances of the strings  $y$  from  $x$  to go to  $\infty$  is not important, it is only a matter of scaling. For convenience, we will require the following *normalization property*:

$$\sum_{y:y \neq x} 2^{-D(x,y)} < 1.$$

We consider only distances that are computable in some broad sense. This condition will not be seen as unduly restrictive. As a matter of fact, only co-enumerability (effective approximation from above) of  $D(x, y)$  will be required (Section 4.1). This is reasonable: as we have more and more time to process  $x$  and  $y$  we may discover newer and newer similarities among them, and thus may revise our upper bound on their distance. The co-enumerability means exactly that  $D(x, y)$  is the limit of a computable sequence of such upper bounds.

**Definition 8.3.2** An *admissible picture distance*,  $D(x, y)$ , is a total nonnegative function on the pairs  $x, y$  of binary strings that is 0 only if  $x = y$ , is symmetric, satisfies the triangle inequality, and is co-enumerable and normalized.

The following theorem shows that  $E_1$  is, in some sense, the optimal admissible picture distance. It is remarkable that this distance happens to also have a “physical” interpretation as the approximate length of the conversion program of Theorem 8.3.1 and, as shown in the next section, of the smallest program that transforms  $x$  into  $y$  on a reversible machine.

**Theorem 8.3.2** *For an appropriate constant  $c$ , let  $E(x, y) = E_1(x, y) + c$  if  $x \neq y$  and 0 otherwise. Then  $E(x, y)$  is an admissible distance function that is minimal in the sense that for every admissible distance function  $D(x, y)$ ,*

$$E(x, y) < D(x, y) + O(1).$$

**Proof.** The nonnegativity and symmetry properties are immediate from the definition. Trivially,  $K(x|z) \leq K(x|y) + K(y|z) + c$  for some nonnegative constant  $c$ . Therefore, by writing everything out, it follows that

$$E_1(x, z) \leq E_1(x, y) + E_1(y, z) + c.$$

Let this  $c$  be the one used in the statement of the theorem; then  $E(x, y)$  satisfies the triangle inequality without an additive constant. We prove the normalization property: For fixed  $x$ , the function

$$f_x(y) = 2^{-\max\{K(x|y), K(y|x)\}}$$

is enumerable (computably approximable from below, page 240). Trivially,  $f_x(y) \leq 2^{-K(y|x)}$ . Moreover,  $2^{-K(y|x)} \leq \sum 2^{-l(p)}$ , where the sum is taken over all programs  $p$  for which the reference prefix machine  $U$ , given  $x$ , computes  $y$ . The sum is the probability that  $U$ , given  $x$ , computes  $y$  from a program  $p$  generated bit-by-bit uniformly at random. Hence,  $\sum_y f_x(y) \leq 1$ . Therefore,  $E(x, y) = -\log f_x(y) + c$  is normalized, for  $c$  an appropriate constant.

We prove the minimality: Let  $x$  be fixed. If  $\sum_y 2^{-D(x,y)} < 1$ , then by Theorem 4.3.1 there is a constant  $c_1$ , independent of  $y$ , such that  $c_1 \mathbf{m}(y|x) \geq 2^{-D(x,y)}$ . Theorem 4.3.4 shows that  $-\log \mathbf{m}(y|x) = K(y|x) + O(1)$ . Repeating this argument with  $x$  and  $y$  interchanged, there is a constant  $c_2$  such that  $D(x, y) \geq \max\{K(x|y), K(y|x)\} - c_2$ .  $\square$

### 8.3.3 Reversible Distance

Another information distance is based on the idea that one should aim for dissipationless computations, and hence for reversible ones. Such an information distance is given by the length of the shortest reversible program that transforms  $x$  into  $y$  and  $y$  into  $x$  on a universal reversible computer. This distance turns out to be  $E_2(x, y) = KR(x|y) = KR(y|x)$ . It will be shown also that  $E_2 = E_1$ , up to a logarithmic additive term. It is remarkable that three so differently motivated definitions turn out to define one and the same notion.

In Theorem 8.3.1 it was demonstrated that for all strings  $x$  and  $y$ , there exists a conversion program  $p$ , of length at most logarithmically greater than  $E_1(x, y) = \max\{K(y|x), K(x|y)\}$ , such that  $U(p, x) = y$  and  $U(p, y) = x$ . We show that the length of this minimal conversion program is equal within a constant to the length of the minimal reversible program for transforming  $x$  into  $y$ .

**Theorem 8.3.3** *Up to an additive constant,*

$$KR(y|x) = \min\{l(p) : U(p, x) = y \text{ and } U(p, y) = x\}.$$

**Proof.** This proof is an example of the general technique for combining two irreversible programs, for  $y$  from  $x$  and for  $x$  from  $y$ , into a single reversible program for  $y$  from  $x$ . In this case the two irreversible programs are almost the same, since by Theorem 8.3.1 the minimal conversion program  $p$  is both a program for  $y$  given  $x$  and a program for  $x$  given  $y$ .

The computation proceeds by several stages, as shown in Figure 8.5. To illustrate motions of the head on the self-delimiting program tape, the program  $p$  is represented by the string “prog” in the table, with the head position indicated by a caret. Each of the stages can be accomplished without using any many-to-one operations.

In stage 1, the computation of  $y$  from  $x$ , which might otherwise involve irreversible steps, is rendered reversible by saving a history, on previously blank tape, of all the information that would have been thrown away.

In stage 2, making an extra copy of the output onto blank tape is an intrinsically reversible process, and therefore can be done without writing anything further in the history. Stage 3 exactly undoes the work of stage 1, which is possible because of the history generated in stage 1.

STAGE AND ACTION	PROGRAM	WORK TAPE	
0. Initial configuration	prog	$x$	
1. Compute $y$ , saving history	proĝ	$y$	( $y x$ )-history
2. Copy $y$ to blank region	proĝ	$y$	( $y x$ )-history
3. Undo comp. of $y$ from $x$	prog	$x$	$y$
4. Swap $x$ and $y$	prog	$y$	$x$
5. Compute $x$ , saving history	proĝ	$x$	( $x y$ )-history
6. Cancel extra $x$	proĝ	$x$	( $x y$ )-history
7. Undo comp. of $x$ from $y$	prog	$y$	

**FIGURE 8.5.** Combining irreversible computations of  $y$  from  $x$  and  $x$  from  $y$  to achieve a reversible computation of  $y$  from  $x$

Perhaps the most critical stage is stage 5, in which  $x$  is computed from  $y$  for the sole purpose of generating a history of that computation. Then, after the extra copy of  $x$  is reversibly disposed of in stage 6 by cancelation (the inverse of copying onto blank tape), stage 7 undoes stage 5, thereby disposing of the history and the remaining copy of  $x$ , while producing only the desired output  $y$ .

Not only are all its operations reversible, but the computations from  $x$  to  $y$  in stage 1 and from  $y$  to  $x$  in stage 5 take place in such a manner as to satisfy the requirements for a reversible prefix interpreter. Hence, the minimal irreversible conversion program  $p$ , with constant modification, can be used as a reversible program for  $UR$  to compute  $y$  from  $x$ .

Conversely, the minimal reversible program for  $y$  from  $x$ , with constant modification, serves as a program for  $y$  from  $x$  for the ordinary irreversible prefix machine  $U$ , because reversible prefix machines are a subset of ordinary prefix machines. This establishes the theorem.  $\square$

**Definition 8.3.3** The *reversible distance*  $E_2(x, y)$  between  $x$  and  $y$  is defined by

$$E_2(x, y) = KR(y|x) = \min\{l(p) : UR(p, x) = y\}.$$

As just proved, this is within an additive constant of the size of the minimal conversion program of Theorem 8.3.1. Although it may be logarithmically greater than the optimal distance  $E_1$ , it has the intuitive advantage of being the actual length of a concrete program for passing in either direction between  $x$  and  $y$ . The optimal distance  $E_1$ , on the other hand, is defined only as the greater of two one-way program sizes and we don't know whether it corresponds to the length of any two-way translation program.

$E_2$  may indeed be legitimately called a distance because it is symmetric and obeys the triangle inequality to within an additive constant (which can be removed by the additive renormalization technique described at the end of Section 8.3.2).

**Theorem 8.3.4**  $E_2(x, z) < E_2(x, y) + E_2(y, z) + O(1)$ .

**Proof.** We will show that given reversible  $UR$  programs  $p$  and  $q$  for computing  $(y|x)$  and  $(z|y)$ , respectively, a program of the form  $spq$ , where  $s$  is a constant supervisory routine, serves to compute  $z$  from  $x$  reversibly. Because the programs are self-delimiting, no punctuation is needed between them. If this were an ordinary irreversible  $U$  computation, the concatenated program  $spq$  could be executed in an entirely straightforward manner, first using  $p$  to go from  $x$  to  $y$ , then using  $q$  to go from  $y$  to  $z$ .

STAGE AND ACTION	PROGRAM	WORK TAPE
0. Initial configuration	$\hat{p}p\text{prog}q\text{prog}$	$x$
1. Compute $(y x)$ , transcribing pprog.	$\hat{p}p\text{prog}q\text{prog}$	$y$ pprog
2. Space forward to start of qprog.	$p\text{prog}\hat{q}p\text{rog}$	$y$ pprog
3. Compute $(z y)$ .	$p\text{prog}\hat{q}p\text{rog}$	$z$ pprog
4. Cancel extra pprog as head returns.	$\hat{p}p\text{rog}q\text{prog}$	$z$

**FIGURE 8.6.** Reversible execution of concatenated programs for  $(y|x)$  and  $(z|y)$  to transform  $x$  into  $z$

However, with reversible *UR* programs, after executing  $p$ , the head will be located at the beginning of the program tape, and so will not be ready to begin reading  $q$ . It is therefore necessary to remember the length of the first program segment  $p$  temporarily, to enable the program head to space forward to the beginning of  $q$ , but then cancel this information reversibly when it is no longer needed.

A scheme for doing this is shown in Figure 8.6, where the program tape's head position is indicated by a caret. To emphasize that the programs  $p$  and  $q$  are strings concatenated without any punctuation between them, they are represented respectively in the table by the expressions “pprog” and “qprog,” and their concatenation  $pq$  by “pprogqprog.”  $\square$

### 8.3.4 Sum Distance

Only the irreversible erasures of a computation need to dissipate energy. This raises the question of the minimal amount of irreversibility required in transforming string  $x$  into string  $y$ , that is, the number of bits we have to add to  $x$  at the beginning of a reversible computation from  $x$  to  $y$ , and the number of garbage bits left (apart from  $y$ ) at the end of the computation that must be irreversibly erased to obtain a “clean”  $y$ .

Even though consuming and producing information may seem to be operations of opposite sign, we can define a distance based on the notion of information flow as the minimal *sum* of amounts of extra information flowing into and out of the computer in the course of the computation transforming  $x$  into  $y$ . This may be viewed as a measure of “work” required during a reversible computation in which the program is not retained.

The resulting distance turns out to be within a logarithmic additive term of the sum of the conditional complexities  $E_3(x, y) = K(y|x) + K(x|y)$ . The reversible distance  $E_2$  defined in the previous section is equal to the length of a “catalytic” program, which allows the interconversion of  $x$  and  $y$  while remaining unchanged itself. Here we consider noncatalytic

reversible computations that consume some information  $p$  besides  $x$  and produce some information  $q$  besides  $y$ .

Consider the enumeration of all reversible Turing machines.

Here we can take either the prefix machines or the nonprefix ones. For the validity of the properties below it doesn't matter whether or not we use prefix machines. Using nonprefix machines gives a smaller sum distance variant though. We therefore use the nonprefix machines.

For a function  $\psi$  computed on a reversible Turing machine, let

$$E_\psi(x, y) = \min\{l(p) + l(q) : \psi(\langle x, p \rangle) = \langle y, q \rangle\}.$$

**Lemma 8.3.2** *There is a universal (nonprefix) reversible Turing machine  $U$  computing the function  $\psi_0$  such that for all functions  $\psi$  computed on a reversible Turing machine, we have*

$$E_{\psi_0}(x, y) \leq E_\psi(x, y) + c_\psi$$

for all  $x$  and  $y$ , where  $c_\psi$  is a constant that depends on  $\psi$  but not on  $x$  or  $y$ .

**Proof.** Use the universal reversible Turing machines in the manner of the proof of Theorem 2.1.1, page 97.  $\square$

**Definition 8.3.4** The *sum distance*  $E_3$  is defined by

$$E_3(x, y) = E_{\psi_0}(x, y).$$

**Theorem 8.3.5**  $E_3(x, y) = K(x|y) + K(y|x) \pm O(\log E_3(x, y))$ .

**Proof.** ( $\geq$ ) We show  $E_3(x, y) \geq C(y|x) + C(x|y)$ . Since  $C(x|y) \geq K(x|y) - 2 \log C(x|y) - O(1)$  for all  $x$  and  $y$  (Example 3.1.3, page 194), this will show the lower bound. To compute  $y$  from  $x$  we must be given a program  $p$  to do so with which to start. By definition,

$$C(y|x) \leq l(p) + O(1).$$

Assume the computation from  $x, p$  ends up with  $y, q$ . Since the computation is reversible we can compute  $x$  from  $y, q$ . Consequently,  $C(x|y) \leq l(q) + O(1)$ .

( $\leq$ ) Assume  $k_1 = K(x|y) \leq k_2 = K(y|x)$  and let  $l = k_2 - k_1$ . According to Theorem 8.3.1, there is a string  $q$  of length  $l + O(\log l)$  such that  $K(qx|y) = k_1 + O(\log k_1)$  and  $K(y|qx) = k_1 + O(\log k_1)$ . We can even assume  $q$  to be self-delimiting: the cost can be included in the  $O(\log l)$

term. According to Theorem 8.3.1 and Theorem 8.3.3, there is a program  $d$  of length  $k_1 + O(\log k_1)$  going reversibly between  $qx$  and  $y$ . Therefore, with a constant extra program  $s$ , the universal reversible machine will go from  $(qd, x)$  to  $(d, y)$ . And by the above estimates

$$l(dq) + l(d) \leq 2k_1 + l + O(\log k_2) = k_1 + k_2 + O(\log k_2).$$

□

Note that all bits supplied in the beginning to the computation, apart from input  $x$ , as well as all bits erased at the end of the computation, are *random* bits. This is because we supply and delete only shortest programs, and a shortest program  $p$  satisfies  $K(p) \geq l(p)$ , that is, it is maximally random.

**Example 8.3.3** It does not matter whether the irreversible providing and erasing operation executions are scattered throughout the computation, or whether all required bits are irreversibly provided at the beginning and all garbage bits are irreversibly erased at the end.

Namely, we can provide the sequence of bits that are irreversibly consumed in the computation as a specially delimited sequence at the beginning of the computation. Whenever the computation requires an irreversibly provided bit, we use the next unused bit from this special string.

Bits that need to be irreversibly erased are put in a special delimited string in the order in which they need to be erased. This special string is then erased at the very end of the computation. ◇

### 8.3.5 Metrics Relations and Density

The metrics we have considered can be arranged in increasing order. Within an additive logarithmic term, we have

$$E_1(x, y) = E_2(x, y) \leq E_3(x, y) \leq 2E_1(x, y),$$

where the various items were defined or proven to be precisely

$$E_1(x, y) = \max\{K(y|x), K(x|y)\},$$

$$E_2(x, y) = KR(y|x) = \min\{l(p) : U(p, x) = y, U(p, y) = x\} + O(1),$$

$$E_3(x, y) = K(x|y) + K(y|x) \pm O(\log E_3(x, y)).$$

The sum distance  $E_3$ , in other words, can be anywhere between the optimum distance  $E_1$  and twice the optimal distance. The former occurs if one of the conditional entropies  $K(y|x)$  and  $K(x|y)$  is zero, the latter if the two conditional entropies are equal.

In a discrete space with a distance function, the rate of growth of the number of elements in balls of size  $d$  can be considered as a kind of “dimension” of the space. The space with distance  $E_1(x, y)$  behaves rather simply from a dimensional point of view.

For a binary string  $x$ , let  $B_1(d, x)$  be the set of strings  $y$  with  $E_1(x, y) \leq d$ . Let the number of elements in  $B_1(d, x)$  be denoted by  $b_1(d, x)$ .

**Theorem 8.3.6** *We have, up to an additive constant,*

$$d - K(d) < \log b_1(d, x) < d - K(d|x).$$

*The same bounds apply to  $B_1(d, x) \cap \{y : l(y) = l(x)\}$ .*

**Proof.** The upper bound is immediate from Exercise 3.3.9, page 205. For the lower bound, let  $i < 2^{d-K(d)}$ , and let  $p_i$  be the  $i$ th binary string of length  $l(x)$ . Let us consider all strings  $y_i = x \oplus p_i$ , where  $\oplus$  means bitwise mod 2 addition. The number of such strings  $y_i$  is  $2^{d-K(d)}$ . We clearly have  $E_1(x, y_i) < K(i) + O(1) < d + O(1)$ .  $\square$

It is interesting that a similar dimension relation holds also for the larger distance  $E_3(x, y)$ , setting it for convenience at  $E_3 = K(y|x) + K(x|y)$ .

**Theorem 8.3.7** *Let  $x$  be a binary string. There is a positive constant  $c$  such that for all sufficiently large  $d$ , the number of binary strings  $y$  with  $E_3(x, y) \leq d$  is at most  $2^d/c$  and at least  $2^d/(d \log^2 d)$ .*

**Proof.** The upper bound follows from the previous theorem since  $E_3 \geq E_1$ . For the lower bound, consider strings  $y$  of the form  $px$  where  $p$  is a self-delimiting program. For all such programs,  $K(x|y) \leq O(1)$ , since  $x$  can be recovered from  $y$ . Therefore  $E_3(x, y) = K(y|x) \pm O(1) = K(p|x)$ . Now, just as in the argument of the previous proof, we obtain the lower bound  $2^{d-K(d)}$  for the number of such strings  $p$  with  $K(p|x) \leq d$ . Substitute  $K(d) \leq \log d + 2 \log \log d$ .  $\square$

For the distance  $E_3$ , for the number of strings of length  $n$  near a random string  $x$  of length  $n$  (that is, a string with  $K(x)$  near  $n$ ), the picture is a little different from that of distance  $E_1$ . In this distance, “tough guys have few neighbors.” In particular, a random string  $x$  of length  $n$  has only about  $2^{d/2}$  strings of length  $n$  within distance  $d$ . The following theorem describes a more general situation.

**Theorem 8.3.8** *Let the binary strings  $x, y$  have length  $n$ . For each  $x$  the number of  $y$ 's such that  $E_3(x, y) \leq d$  is  $2^\alpha$ , with*

$$\alpha = \frac{n + d - K(x)}{2} \pm O(\log n),$$

*while  $n - K(x) \leq d$ . For  $n - K(x) > d$  we have  $\alpha = d \pm O(\log n)$ .*

**Proof.** Let  $K(x) = n - \delta(n)$ . In the remainder of the proof we will assume that all (in)equalities involving  $K$  hold up to an  $O(\log n)$  additive term.

( $\geq$ ) We show that there are at least  $2^{(d+\delta(n))/2}$  elements  $y$  such that  $E_3(x, y) \leq d$  holds. Let  $y = x^*z$  with  $l(z) = \delta(n)$  and let  $x^*$  be the first self-delimiting program for  $x$  of length  $K(x)$  that we find by dovetailing all computations on programs of length less than  $n$ . We can retrieve  $z$  from  $y$  using at most  $O(\log n)$  bits. There are  $2^{\delta(n)}$  different such  $y$ 's. For each such  $y$  we have  $K(x|y) = O(1)$ , since  $x$  can be retrieved from  $y$  using  $x^*$ . Now suppose we further divide  $y = uw$  with  $l(u) = l/2$  for an appropriate  $l$  and choose  $u$  arbitrarily. Then the total number of such  $y$ 's increases to  $2^{\delta(n)+l/2}$ .

These choices of  $y$  must satisfy  $E_3(x, y) \leq d$ . Clearly,  $K(y|x) \leq \delta(n) + l/2$ . Moreover,  $K(x|y) \leq l/2$  since we can retrieve  $x$  by providing  $l/2$  bits. Therefore,

$$K(x|y) + K(y|x) \leq l/2 + \delta(n) + l/2.$$

Since the left-hand side has a value at most  $d$ , the largest  $l$  we can choose is, up to the suppressed additional term  $O(\log n)$ , given by  $l = d - \delta(n)$ .

This puts the number of  $y$ 's such that  $E(x, y) \leq d$  at least at

$$2^{(\delta(n)+d)/2 \pm O(\log n)}.$$

( $\leq$ ) Assume, to the contrary, that there are at least  $2^{(d+\delta(n))/2+c}$  elements  $y$  such that  $E_3(x, y) \leq d$  holds, with  $c$  some large constant. For some  $y$ ,

$$K(y|x) \geq \frac{d + \delta(n)}{2} + c. \quad (8.5)$$

By assumption,

$$K(x) = n - \delta(n), \quad K(y) \leq n.$$

By Theorem 3.9.1, page 232, and substitution, we find

$$n + \frac{d - \delta(n)}{2} + c \leq n + K(x|y).$$

But this means that

$$K(x|y) \geq \frac{d - \delta(n)}{2} + c. \quad (8.6)$$

Together, Equations 8.5 and 8.6 contradict  $K(x|y) + K(y|x) \leq d$ .  $\square$

If  $S$  is a finite set, let us define the prefix Kolmogorov complexity  $K(S)$  for  $S$  to be the length of the shortest binary self-delimiting program that enumerates  $S$  and then halts. It follows from our estimates that in every set of low Kolmogorov complexity almost all elements are far away from each other in terms of distance  $E_3$ .

**Theorem 8.3.9** *For a constant  $c$ , let  $S$  be a set with cardinality  $d(S) = 2^d$  and  $K(S) = c \log d$ . Almost all pairs of elements  $x, y \in S$  have distance  $E_1(x, y) \geq d$ , up to an additive logarithmic term.*

The proof of this theorem is easy and is left as Exercise 8.3.1. A similar statement can be proved for the distance of a string  $x$  (possibly outside  $S$ ) to the majority of elements  $y$  in  $S$ . If  $K(x) \geq n$ , then for almost all  $y \in S$  we have  $E_1(x, y) \geq n + d - O(\log dn)$ .

### 8.3.6 Thermodynamics of Computing

Finally, using the physical theory of reversible computation, the simple difference  $K(x) - K(y)$  turns out to be an appropriate (universal, antisymmetric, and transitive) measure of the amount of thermodynamic work required to transform string  $x$  into string  $y$  by the most efficient process.

Thermodynamics, among other things, deals with the amounts of heat and work ideally required, by the most efficient process, to convert one form of matter to another. For example, at  $0^\circ\text{C}$  and atmospheric pressure, it takes 80 calories of heat and no work to convert a gram of ice into water at the same temperature and pressure.

From an atomic point of view, the conversion of ice to water at  $0^\circ\text{C}$  is a reversible process, in which each melting water molecule gains about 3.8 bits of entropy (representing the approximately  $2^{3.8}$ -fold increased freedom of motion it has in the liquid state), while the environment loses 3.8 bits.

During this ideal melting process, the entropy of the universe remains constant, because the entropy gain by the melting ice is compensated by an equal entropy loss by the environment. Perfect compensation takes place only in the limit of slow melting, with an infinitesimal temperature difference between the ice and the water.

Rapid melting, for example, when ice is dropped into hot water, is thermodynamically irreversible and inefficient, with the environment (the hot water) losing less entropy than the ice gains, resulting in a net and irredeemable entropy increase for the universe as a whole.

Turning again to ideal reversible processes, the entropy change in going from state  $X$  to state  $Y$  is an antisymmetric function of  $X$  and  $Y$ ; thus, when water freezes at  $0^\circ\text{C}$  by the most efficient process, it gives up 3.8 bits of entropy per molecule to the environment. When more than two states are involved, the entropy changes are transitive: thus, the entropy change per molecule of going from ice to water vapor at  $0^\circ\text{C}$  (+32.6 bits) plus that for going from vapor to liquid water (-28.8 bits) sum to the entropy change for going from ice to water directly.

Because of this antisymmetry and transitivity, entropy can be regarded as a thermodynamic potential, or state function: each state has an entropy, and the entropy change in going from state  $X$  to state  $Y$  by the most efficient process is simply the entropy difference between states  $X$  and  $Y$ .

Thermodynamic ideas were first successfully applied to computation by Landauer (Section 8.2.1). According to Landauer's principle, an operation that maps  $n$  states onto a common successor state must be accompanied by an entropy increase of  $\log_2 n$  bits in other, noninformation-bearing degrees of freedom in the computer or its environment. At room temperature, this is equivalent to the production of  $kT \ln 2$  (about  $2.8 \times 10^{-21}$ ) joules of waste heat per bit of information discarded.

Landauer's principle follows from the fact that such a logically irreversible operation would otherwise be able to decrease the thermodynamic entropy of the computer's data without a compensating entropy increase elsewhere in the universe, thereby violating the second law of thermodynamics.

Converse to Landauer's principle is the fact that when a computer takes a physical *randomizing* step, such as tossing a coin, in which a single logical state passes stochastically into one of  $n$  equiprobable successors, that step can, if properly harnessed, be used to remove  $\log_2 n$  bits of entropy from the computer's environment. Models have been constructed (Section 8.2.1) obeying the usual conventions of classical, quantum, and thermodynamic thought-experiments showing both the ability in principle to perform logically reversible computations in a thermodynamically reversible fashion, and the ability to harness entropy increases due to data randomization within a computer to reduce correspondingly the entropy of its environment.

In view of the above considerations, it seems reasonable to assign each string  $x$  an effective thermodynamic entropy equal to its complexity  $K(x)$ . A computation that erases an  $n$ -bit random string would then reduce its entropy by  $n$  bits, requiring an entropy increase in the environment of at least  $n$  bits, in agreement with Landauer's principle.

Conversely, a randomizing computation that starts with a string of  $n$  zeros and produces  $n$  random bits has, as its typical result, an algorithmically random  $n$ -bit string  $x$ , that is, one for which  $K(x) \approx n$ . By the converse of Landauer's principle, this randomizing computation is capable of removing up to  $n$  bits of entropy from the environment, again in agreement with the subsumption of the algorithmic complexity  $K$  in the thermodynamic entropy.

What about computations that start with one random string  $x$  and end with another string  $y$ ? By the transitivity of entropy changes one is led to

say that the thermodynamic cost, that is, the minimal entropy increase in the environment, of a transformation of  $x$  into  $y$ , should be

$$W(y|x) = K(x) - K(y),$$

because the transformation of  $x$  into  $y$  could be thought of as a two-step process in which one first erases  $x$ , and then allows  $y$  to be produced by randomization. By the elementary properties of self-delimiting programs, this cost measure is transitive within an additive constant. A similar antisymmetric measure of the thermodynamic cost of data transformations, such as

$$W'(y|x) = K(x|y) - K(y|x),$$

is slightly nontransitive. There are strings  $x$  with  $K(x^*|x) \approx K(K(x))$ , where  $x^*$  is the minimal program for  $x$  (Theorem 3.8.1 on page 227). According to the  $W'$  measure, erasing such an  $x$  via the intermediate  $x^*$  would generate  $K(K(x))$  less entropy than erasing it directly, while for the  $W$  measure the two costs would be equal within an additive constant. Indeed, erasing in two steps would cost only  $K(x|x^*) - K(x^*|x) + K(x^*|\epsilon) - K(\epsilon|x^*) = K(x) - K(x^*|x)$  while erasing in one step would cost  $K(x|\epsilon) - K(\epsilon|x) = K(x)$ . (Everything up to an additive constant.)

Finally, we consider entropy changes in nonideal computations. Consider the thermodynamics of an intelligent demon or engine that has some capacity to analyze and transform data  $x$  before erasing it. If the demon erases a random-looking string, such as the digits of  $\pi$ , without taking the trouble to understand it, it will commit a thermodynamically irreversible act, in which the entropy of the data is decreased very little, while the entropy of the environment increases by a full  $n$  bits. On the other hand, if the demon recognizes the redundancy in  $\pi$ , it can transform  $\pi$  to a short string by a reversible computation, and thereby accomplish the erasure at very little thermodynamic cost.

More generally, given unlimited time, a demon could approximate the co-enumerable function  $K(x)$  and so compress a string  $x$  to size  $K(x)$  before erasing it. But in limited time, the demon will not be able to compress  $x$  so much and will have to generate more entropy to get rid of it. This tradeoff between speed and thermodynamic efficiency is superficially similar to the tradeoff between speed and efficiency for physical processes such as melting, but the functional form of the tradeoff is very different.

For typical physical state changes such as melting, the excess entropy produced per molecule goes to zero inversely in the time  $t$  allowed for melting to occur. But the time-bounded prefix complexity  $K^t(x)$ , that is, the size of the smallest program to compute  $x$  in time at most  $t$ , in general approaches  $K(x)$  with uncomputable slowness as a function of  $t$  and  $x$ . A formal result along these lines is stated as Exercise 8.3.5.

## Exercises

**8.3.1.** [24] Prove Theorem 8.3.9.

**8.3.2.** [21] Fix a reference universal reversible Turing machine, say  $UR_0$ . Define  $E_3^t(x, y) = \min\{l(p) + l(q) : UR_0(\langle p, x \rangle) = \langle q, y \rangle$  in  $t(n)$  steps of computation, where  $n = l(x)\}$ . Show that for each recursive function  $t$  there is an  $x$  such that  $E_3^t(x, \epsilon) > E_3(x, \epsilon)$ .

**8.3.3.** [32] Prove the upper bound on  $E_3(\cdot, \cdot)$  of Theorem 8.3.5 on page 546 by a direct construction supplying logarithmically small initial programs and ending with logarithmically small garbage.

*Comments.* Source: M. Li and P.M.B. Vitányi, *Proc. Royal Soc. London, Ser. A*, 452(1996), 769–789. This is also the source of the four exercises that follow.

**8.3.4.** [25] We determine the irreversibility cost of effective erasure. Let  $t(|x|) \geq |x|$  be a time bound that is provided at the start of the computation. Show that erasing an  $n$  bit record  $x$  by an otherwise reversible computation can be done in time (number of steps)  $O(2^{|x|} t(|x|))$  at irreversibility cost  $C^t(x) + 2C^t(t|x) + 4 \log C^t(t|x)$  bits. (Typically we consider  $t$  as some standard explicit time bound and the last two terms adding up to  $O(1)$ .)

**8.3.5.** [37] Use the notions of Exercise 8.3.2. Prove a time-energy trade-off hierarchy: For every large enough  $n$  there is a string  $x$  of length  $n$  and a sequence of  $m = \frac{1}{2}\sqrt{n}$  time functions  $t_1(n) < t_2(n) < \dots < t_m(n)$ , such that

$$E_3^{t_1}(x, \epsilon) > E_3^{t_2}(x, \epsilon) > \dots > E_3^{t_m}(x, \epsilon).$$

Improve this result by replacing  $\epsilon$  by an arbitrary string  $y$ .

**8.3.6.** [28] Show there is a recursively enumerable infinite sequence  $\chi$  and some (nonrecursively) large time bound  $T$ , such that for each total recursive time bound  $t$ , for each initial segment  $x$  of  $\chi$

$$E_3^t(x, \epsilon) > c_t 2^{E_3^T(x, \epsilon)/2},$$

where  $c_t > 0$  is a constant depending only on  $t$  and  $\chi$ .

**8.3.7.** [23] Defining a uniform variant  $E_u(\cdot, \cdot)$  of  $E_3(\cdot, \cdot)$  based on the uniform Kolmogorov complexity variant of Exercise 2.3.3 on page 124. Show that the energy-dissipation can be reduced arbitrarily computably far by using enough (that is, a noncomputable amount of) time. That is, there is an infinite sequence  $\omega$  and a (nonrecursively) large time bound  $T$ , such that for each unbounded total recursive function  $f$ , no matter how large, for each total recursive time bound  $t$ , there are infinitely many  $n$  for which

$$E_u^t(\omega_1 \dots \omega_n, \epsilon) > f(E_u^T(\omega_1 \dots \omega_n, \epsilon)).$$

*Comments.* Hint: use Exercise 2.5.13 on page 153 and Exercise 7.1.6 on page 473.

**8.3.8.** [O40] Develop a theory of resource-bounded (polynomial time or logarithmic space) information distance.

## 8.4 Thermo-dynamics

Classical thermodynamics deals with describing the physical properties of substances like gases under variations of macroscopic observables such as volume, temperature, and pressure. Certain regularities are codified on the basis of experience and formulated as axioms of the theory. The two fundamental laws of thermodynamics are

**First Law** The total energy of an isolated system is invariant over time.

**Second Law** No process is possible that has as its only result the transformation of heat into work.

The first law is a statement of the principle of conservation of energy for thermodynamical systems. The variation in energy of a system during any transformation is equal to the amount of energy that the system receives from its environment, minus the amount of energy the system loses to its environment. The first law arose as the result of the impossibility of constructing a machine that would create energy (*perpetuum mobile* of the first kind). It places no limitations on the possibility of transforming energy from one form into another, like the transformation of heat into work or work into heat.

The first law still allows a machine that transforms heat into work by cooling the environment (*perpetuum mobile* of the second kind). The second law rules out this possibility. An equivalent formulation of the second law is: “no process is possible which has as its only result the transfer of heat from a cooler to a warmer body.” To prove that the two formulations are equivalent, it suffices to show that a process can be created that transforms heat from a cooler body to a warmer body if a process is available that transforms heat into work, and conversely; this is not difficult.

### 8.4.1 Classical Entropy

The original definition of entropy was proposed in 1824 by the French engineer N.L. Sadi Carnot in the context of the so-called Carnot cycle. This is a cycle of states of thermodynamic systems such as steam engines. The existence of an entropy can be derived from the second law by a reasoning that depends on the notion of *reversible processes*. The two standard examples are the following:

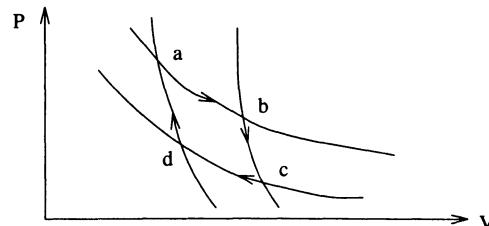
**Example 8.4.1** We have a cylinder filled with an ideal gas, in contact with a heat reservoir at temperature  $T_1$ . Allow the gas to push out a frictionless piston “very slowly.” The gas will expand while staying at the same temperature and meanwhile will take up heat from the reservoir. This is (in the limit) a reversible process: by pushing in the piston again “very slowly,” the whole system can be brought back into exactly the same condition as it was before the expansion. ◇

**Example 8.4.2** Use the same cylinder as before, but now in thermic isolation. Again pull out the piston very slowly. The temperature of the gas will fall. This is again (in the limit) a reversible process. ◇

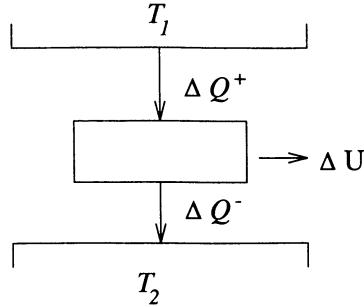
A concatenation of reversible processes is again reversible. By combination of the two types of processes just mentioned, one can obtain a reversible *cycle*: first isothermic expansion, then adiabatic expansion, then isothermic compression, and finally adiabatic compression to the initial state. This is called a *Carnot cycle*. The whole process can be plotted on the  $(V, P)$ -plane, Figure 8.7, where  $V$  denotes “volume” and  $P$  denotes “pressure.” The curves corresponding to isothermic expansion and compression satisfy  $PV$  is constant (actually,  $PV = RT$  where  $T$  is the temperature and  $R > 0$  is a gas-specific constant). The curves corresponding to adiabatic expansion and compression satisfy  $PV^\gamma$  is constant, where  $\gamma > 1$  is another gas-specific constant.

In terms of heat and work, during a Carnot cycle the gas in the cylinder consumes an amount of heat  $\Delta Q^+$  from the reservoir at temperature  $T_1$ , and delivers an amount of heat  $\Delta Q^-$  to another reservoir at a lower temperature  $T_2$ . The total amount of work performed in the cycle is equal to  $\int P dV$ , which is exactly the surface area enclosed in Figure 8.7. A Carnot cycle is therefore a form of a heat engine.

According to the first law of thermodynamics, the amount of work done by a heat engine, Figure 8.8, must equal the difference between the heat extracted from the warm reservoir and the heat delivered to the cold



**FIGURE 8.7.** Carnot cycle



**FIGURE 8.8.** Heat engine

reservoir,

$$\Delta U = \Delta Q^+ + \Delta Q^-,$$

where  $\Delta Q^-$  is negative.

Assuming the process to be reversible, one can derive a further relation. Denote by  $a, b, c, d$  the four intermediate stages of the Carnot cycle, Figure 8.7. Then

$$\Delta Q^+ = \int_a^b P \, dV = \int_a^b RT_1 \frac{dV}{V} = RT_1 \log \frac{V_b}{V_a}.$$

Likewise,

$$\Delta Q^- = RT_2 \log \frac{V_d}{V_c}.$$

The relations  $P_b V_b^\gamma = P_c V_c^\gamma$ ,  $P_b V_b^\gamma = RT_1$  and  $P_c V_c^\gamma = RT_2$  yield

$$T_1 V_b^{\gamma-1} = T_2 V_c^{\gamma-1},$$

and in the same way

$$T_1 V_a^{\gamma-1} = T_2 V_d^{\gamma-1}.$$

Dividing, we obtain  $V_b/V_a = V_c/V_d$  and therefore

$$\frac{\Delta Q^+}{T_1} + \frac{\Delta Q^-}{T_2} = 0.$$

Although this was derived for the Carnot cycle, it follows from the second law that this formula must be universal for reversible processes. Namely, if we had a reversible process that did not satisfy the relation above, then we could couple it to a Carnot cycle in a suitable way such as to fabricate

a process that would violate the second law. For general processes, the law that we have just derived takes the following form:

$$\oint \frac{dQ}{T} = 0, \text{ along reversible paths.}$$

(An alternative form of the integral with the time parameter appearing explicitly is  $\oint Q(t)T^{-1}(t) dt$ .) This property makes it possible to define the entropy, as a function of temperature and volume, as follows:

$$S(T_b, V_b) = S(T_a, V_a) + \int_a^b \frac{dQ}{T},$$

where the integral is taken along a reversible path. It follows from the above that the answer does not depend on the choice of the reversible path. This is the classical thermodynamical definition of entropy. It determines the *entropy S* only up to an additive constant. If the process is irreversible, then

$$\oint dS > 0. \tag{8.7}$$

Thus, the entropy always increases with time, except for reversible systems, where it can stay the same. This is the classic second law.

**Example 8.4.3** We compute the entropy of an ideal gas at a constant temperature  $T$ , so that the entropy will only be a function of volume. Using the ideal gas law  $PV = RT$ , we obtain

$$\begin{aligned} S(V_b) - S(V_a) &= \int_a^b \frac{dQ}{T} = \frac{1}{T} \Delta Q = \frac{1}{T} \int_a^b P dV \\ &= \frac{1}{T} \int_a^b \frac{RT}{V} dV = R(\log V_b - \log V_a). \end{aligned}$$

This result may already be related to the information-theoretic entropy, as follows: Assume that a volume is divided into “cells,” where each cell can contain either one molecule or no molecule. A typical gas molecule may be found with equal probabilities in one of  $N_a$  cells in volume  $V_a$  and in one of  $N_b$  cells in volume  $V_b$ . By Definition 1.11.1, page 67, of Shannon’s entropy, the difference in entropies is proportional to  $\log N_b - \log N_a = \log V_b - \log V_a$ . If the molecules are independent, then we can just multiply by the number of molecules to obtain the difference of the total entropies. In this simple case, the information-theoretic Definition 1.11.1 is in line with the thermodynamic case. ◇

### 8.4.2 Statistical Mechanics and Boltzmann Entropy

The above theory belongs to the “phenomenological” thermodynamics. However, it is desirable to derive all laws in physics from first principles. Eventually, a gas consists of molecules that move about according to the laws of mechanics. This can be classical or quantum. But whether classical or modern, the known equations of physics are reversible in time. This implies that if a system has a certain trajectory of its evolution from state  $a$  to state  $b$ , then there also is a trajectory of evolution from state  $b$  to state  $a$ . In the classical case, it suffices to reverse all final velocities of all particles and the system will trace its trajectory backwards.

This reversibility of equations seems to be in contradiction with the above considered irreversibility of some thermodynamical processes. L. Boltzmann was the first to formulate the idea that entropy measures disorderliness, and to connect this with the origin of the distinction between past and future in a nonrelativistic universe. This distinction is obvious in all our immediate experiences. We have no hesitation to identify a movie that is run backwards by, say, a broken vase becoming whole again. It would be nice if this conviction could be justified from the laws of nature. But the laws of physics (as far as known) tell a different story: if the movie run in one direction is physically appropriate, then so is the movie run in the opposite direction.

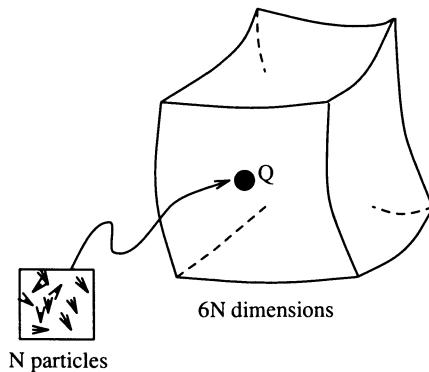
A refinement of the notions of macroscopic state and reversibility is needed. Roughly speaking, the “macro state” of an isolated mechanical system consists of an approximate description of the values of just a few macroscopic observables such as volume, temperature, pressure, and so on.

The “micro state” of such a system determines the behavior of the system completely, whether it is in equilibrium or not. Roughly speaking, a micro state of a physical system consists of an exhaustive description of the values of all the microscopic parameters of the particles. Statistical thermodynamics, also called statistical mechanics, tries to derive the classical laws of thermodynamics from microscopic phenomena.

*Ideal gas* is a convenient tool to study statistical thermodynamics. Consider a (3-dimensional) container with  $N$  identical (ideal) gas molecules. Each molecule is considered as an elastic sphere (also called ball) with no internal freedom. The behavior of the gas is completely determined if we specify the position, the mass, and velocity of each molecule at some time instant. Instead of the velocity, it is more convenient to use the *impulse momentum* which is *mass times velocity*, of each molecule.

The position of a molecule is specified by three space coordinates. The momentum of a molecule is a vector in 3-space and hence is also specified by three coordinates. To specify  $N$  molecules we need  $6N$  parameters.

Since  $N \approx 6 \times 10^{23}$  (Avogadro’s number) for 1 *mole* of gas (1 mole of hydrogen is 2 grams), such a system is hard to analyze. It is convenient to



**FIGURE 8.9.** State space

use a fictional space  $X$ , the so-called *state space* of the system (also called *phase space*), of  $6N$  dimensions. There is a one-to-one correspondence between the dimensions and the coordinates of the molecules.

**Definition 8.4.1** A *micro state* of the entire system of  $N$  molecules corresponds to just one point,

$$Q = (q_1, \dots, q_{3N}, p_1, \dots, p_{3N}),$$

in state space, where the  $q_i$ 's specify the coordinates of the positions of  $N$  molecules and the  $p_i$ 's specify their momenta, Figure 8.9.

Each coordinate is a real number. If  $(\omega^1, \dots, \omega^{6N})$  is a micro state, then it can be encoded as a single infinite binary sequence  $\omega$ .

An encoding like the  $i$ th bit of  $\omega^j$  represented by the  $(j + 6(i - 1)N)$ th bit of  $\omega$  is improper. Since typically the number of molecules is about  $6 \times 10^{23}$ , the first  $10^{23}$  bits of  $\omega$  give little information about the state of the system. Our encoding should have the property that the consecutive bits contain information about the micro state in decreasing order of importance from a macroscopic point of view. For example, the first few bits may describe, to a reasonable degree of precision, the amount of gas in each half of the container, the next few bits may describe the amounts in each quarter, the next few bits may describe the temperature in each half, the next few bits may describe again the amount of gas in each half, but now to more precision, and so on. We can now divide our space into cells of different grain size.

We shall assume that the micro state moves around in a closed volume of the state space, of a standard number of units in each dimension. Then for each dimension we can assume the position of the decimal point

known, and we ignore it in the binary expansion of the coordinates. The set of micro states of the state space is denoted by  $X \subseteq \{0, 1\}^\infty$ .

The total volume of the space need not measure to 1. For example, if  $\lambda$  denotes the Lebesgue measure (or uniform measure), and our state space is a  $6N$ -dimensional hypercube  $R$  with each side 2 units, then  $\lambda(R) = 2^{6N}$ . We restrict ourselves to measures that are finite over the total volume, so that the difference with a probability measure is just a matter of multiplication by a scaling constant.

**Notation 8.4.1** On page 243 we introduced the shorthand notation of  $\mu(x)$  for the  $\mu$ -measure  $\mu(\Gamma_x)$  of a cylinder set  $\Gamma_x$ . In this section we use the common standard notation  $\mu(A)$  for the  $\mu$ -measure of a set  $A$ .

Now suppose there are  $m$  macroscopic parameters  $u_1, \dots, u_m$ . We also will assume that each such parameter takes only a finite number of values: a macroscopic parameter that is a real number will only be taken up to some reasonable precision (four digits are probably more than sufficient). In this way we obtain a partition of the state space  $X$  into cells

$$X = \bigcup_x \Gamma_x,$$

where  $x$  runs over all possible values of parameters  $u_1, \dots, u_m$ . Here each  $\Gamma_x$  consists of the set of all micro states that can occur under macroscopic constraints  $x$ .

**Definition 8.4.2** Let  $x$  be as above. A *macroscopic state* of the system, or *macro state*, is a cell  $\Gamma_x$ . Each cell is identified with the corresponding set of micro states. The partition interpretation of a macro state is called *coarse-graining*.

**Example 8.4.4** Suppose we subdivide a container into  $10 \times 10 \times 10$  parts. Determine the pressure in each part to within 0.1 atmosphere. It is important to note that the state space partition according to the different values of the macro description will have hugely different cell sizes. Therefore, the cell volumes will depend on the partition, but these differences are for “reasonable” partitions negligible in comparison to the cell size differences arising already within a fixed partition. ◇

According to the laws of mechanics, an isolated system undergoes an evolution described by a transformation group  $U^t$ . If at time  $t_0$ , the system was in state  $\omega$ , then at time  $t + t_0$  it will be in state  $U^t\omega$ . The group  $U^t$  is generally given by a system of differential equations that, at least in the example of ideal gas given with coordinates and impulses,

are called the *Hamiltonian equations*. For this case, *Liouville's Theorem* holds, saying that the volume of a domain remains invariant under transformation under  $U^t$ . In most other cases also, a natural measure is found on  $X$  that remains invariant under  $U^t$ , and we will call this measure the *volume*, and denote the volume of a set  $A$  by  $L(A)$ .

The law of energy conservation means that during an evolution of an isolated system, it is confined to a surface of the state space determined by the requirement that the energy is equal to a certain value. Therefore the volume measure to use will be actually obtained by restricting the original volume measure to a thin layer determined by the requirement that the value of energy be in a certain small interval, and normalizing.

When we say that the transformation of a macroscopic description  $x$  to macroscopic description  $y$  is reversible, then this statement refers to the macro states  $\Gamma_x, \Gamma_y$ , and what is meant is that the existence of a reverse transformation follows already from the properties of macro states. The question arises what reversibility means in terms of micro states.

Since macro state  $x$  contains many micro states, the question is whether this means that a reverse transformation exist for all elements of  $\Gamma_y$  or only for some. It is easy to see that it is too much to require that the reverse transformation exists for all micro states in  $\Gamma_y$ , and too little that it exist only for some. The answer is that the transformation is reversible if the reverse transformation exists for most (by volume) micro states in  $\Gamma_y$ . Boltzmann considered as his greatest achievement to have found a microscopic expression for entropy.

**Definition 8.4.3** The *Boltzmann entropy* of a system with macroscopic description  $x$  is proportional to the logarithm of the volume of  $\Gamma_x$ ,

$$S_B(x) = (k \ln 2) \log L(\Gamma_x),$$

where  $k = 1.38 \times 10^{-23}$  joules/°kelvin is called the *Boltzmann constant*, and  $\Gamma_x$  is the cell in state space corresponding to the macroscopic description  $x$  of the system, and  $L(\Gamma_x)$  is its volume.

**Example 8.4.5** For the special case of a discrete state space with  $L$  the counting measure defined as  $L(\Gamma_x) = d(\Gamma_x)$ , that is, the number of elements in the set  $\Gamma_x$ , the Boltzmann entropy has the familiar form

$$S_B(x) = (k \ln 2) \log d(\Gamma_x).$$

◇

**Example 8.4.6** Consider a container of ideal gas consisting of  $n$  identical molecules. Imagine that the container partitioned into  $m$  compartments  $C_1, \dots, C_m$

where  $m$  is much smaller than  $n$ . For simplicity, we ignore the velocities. Let the macro variable  $n_i$  give the number of molecules in compartment  $C_i$ . Let  $n = (n_1, \dots, n_m)$ . A second set of variables, the numbers  $i_1, \dots, i_n$ , indicates that molecule  $j$  is in compartment  $i_j$ .

The description  $i = (i_1, \dots, i_n)$  is, of course, more detailed than the description  $n$ , which is a function  $n(i)$  of  $i$ . Therefore  $\Gamma_n = \bigcup_{n=n(i)} \Gamma_i$ .

Since the molecules are identical, the dynamics, therefore the measure  $L$ , will certainly be invariant with respect to the exchange of molecules. Therefore if  $n(i) = n(j)$ , then  $L(\Gamma_i) = L(\Gamma_j)$ . Assume that molecules are distinguishable. Then  $L(\Gamma_n) = N(n)p_n$ , where  $p_n$  is the volume of each  $\Gamma_i$ , with  $n(i) = n$  and

$$N(n) = \frac{n!}{n_1! \cdots n_m!}.$$

If we also assume that all the compartments  $C_i$  have the same shape and size, and hence have the same volume  $V$ , then  $p_n$  does not depend on  $n$ . Indeed,  $p_n = V^n$ , the  $n$ th power of the volume of a single compartment, since with  $i$  fixed only the exact position of each molecule in its compartment is undetermined. The Boltzmann entropy of state  $n$  is equal to

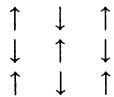
$$S_B = (k \ln 2) \log N(n) \approx (k \ln 2)n \left( - \sum_i f_i \log f_i + \log V \right),$$

where  $f_i = n_i/n$ . If the molecules are indistinguishable, then it is necessary to subtract  $\log n!$ .  $\diamond$

- Example 8.4.7** In high-temperature superconductivity research, a material like CuO<sub>2</sub> loses magnetic moment below a critical temperature. In such a state, the nuclear spins all line up as in Figure 8.10. For Boltzmann entropy, one must agree on some standard partition (coarse-graining) of the space into subsets, and the entropy of  $x$  is the negative logarithm of the measure of the subset that  $x$  is in. There is much arbitrariness in the choice of the partition, but the differences caused by different reasonable choices of partition are negligible compared to the size of the entropy. It is important that the partition is a reasonable one.

$\uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \dots$   
 $\downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \dots$   
 $\uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \dots$   
 $\downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \dots$

**FIGURE 8.10.** Atomic spin in CuO<sub>2</sub> at low temperature



**FIGURE 8.11.** Regular “up” and “down” spins

A *reasonable partition* is a discretization of the state space in which the values of the most important quantities come first. In the case of the ferromagnetic example, the sequence  $x_1, x_2, \dots$  would not be simply the sequence of spins but would begin with various global statistical quantities since presumably these can be measured. One may put first the total number of “up” spins, at least to some precision, then the total number of “up” spins in the left half, and so on. Then, some time later, one may put the total number of neighborhoods of the form as in Figure 8.11.

Let us call this last number  $z$ . If the total number of atoms is  $N$ , then in Figure 8.10 we have approximately  $z = N/2$ . On the other hand, there is only a very small number of configurations with  $z = N/2$ , or even where  $z$  is almost  $N/2$  since they must all essentially look like this. Therefore, the Boltzmann entropy of this configuration is very low.

But suppose that the ups and downs encode the bits of  $3.1415\dots$ . Then in the algorithmic sense the complexity is small, but this will only be seen when the size of the neighborhoods is chosen as large as  $N$ . Considering only local neighborhoods, the statistical properties suggest randomness. In physics, it is assumed that we never look at other than global statistical (“macroscopic”) quantities; hence, when the ferromagnet encodes the digits of  $\pi$  this will forever be its secret. ◇

The most obvious problem with this definition of Boltzmann entropy seems to be its dependence on the particular choice and number of macroscopic variables and the precision with which we want to determine them. Indeed, another digit of precision will decrease the Boltzmann entropy of most states by about  $\log 10$ . The volumes in question are, however, typically so large that such a small difference is negligible (especially if we take into consideration that the actual definition of  $S_B(x)$  involves multiplication by the very small Boltzmann constant  $k$ ).

**Example 8.4.8** Boltzmann entropy possesses the entropy increase properties as required by the second law. This is due to the enormous differences in cell sizes irrespective of a small number of cells. Classical statistical mechanical systems are distinguished by the fact that every reasonable canonical cell division will have this property.

For instance, suppose we are dealing with a vessel containing a mole of gas, page 558, and our macroscopic variable  $x$  is the binary sequence

$x_1 x_2$  giving approximately the relative quantity of gas in the left half of the container. Then the cells  $\Gamma_{00}$  and  $\Gamma_{11}$  are absolutely negligible in volume compared to the cells  $\Gamma_{01}$  and  $\Gamma_{10}$ .

Since entropy measures the amount of state space occupied by a macro state, it sounds plausible that a system tends to be in states with high entropy. According to one formulation of the second law of thermodynamics, entropy in isolated systems cannot decrease. There are some other related entropy increase properties, such as

- an isolated system undergoes an irreversible transformation exactly when entropy actually increases;
- an equilibrium state of an isolated system is a maximum (or at least a strongly pronounced local maximum) of entropy.

It is most important that Boltzmann entropy can be shown to increase in time until equilibrium is reached. Since there are huge differences in the cell volumes  $L(\Gamma_x)$  for different macro states  $\Gamma_x$ , typically  $U^t$  takes a point from a small cell to a bigger cell. Just as the precise notion of reversibility had to be formulated statistically, the same can be expected for these entropy increase properties for Boltzmann entropy.

It turns out that for certain special cases one can prove two properties that together imply the desired entropy increase properties for Boltzmann entropy. The first property states that the union of all “small” cells taken together is small. The second property states that if the system starts from a state in a not very small cell then after a time  $t$ , it is unlikely to end up in any small union of cells. ◇

#### 8.4.3 Gibbs Entropy

A third definition of entropy is due to J.W. Gibbs. Consider some thermodynamic system. Another possible interpretation of a macro state is as a certain distribution over microscopic states. An *ensemble* is some measure with density  $p(\omega)$  over the state space  $X$  interpreted as a distribution of individual points satisfying  $\int p(\omega)L(d\omega) = 1$ . The ensemble  $p_\Gamma(\omega)$  corresponding to a macro state  $\Gamma$  is defined as

$$p_\Gamma(\omega) = \begin{cases} 1/L(\Gamma) & \omega \in \Gamma, \\ 0 & \text{otherwise.} \end{cases}$$

**Definition 8.4.4** The *Gibbs entropy* of a cell  $\Gamma$  is defined as

$$S_G(p) = -(k \ln 2) \int_\Gamma p(\omega) \log p(\omega) L(d\omega). \quad (8.8)$$

**Example 8.4.9** Let  $X$  be the state space of a thermodynamic system, and  $U$  the reversible transformation satisfying Liouville's Theorem, page 560, meaning that  $U$  is volume preserving. The ensemble is imagined to consist of points  $\omega$  moving individually so that point  $\omega$  is transformed to  $U^t\omega$ . For a cell  $\Gamma$  we find  $S_G(p_\Gamma) = (k \ln 2) \log L(\Gamma)$ . That is, the Gibbs entropy is the same as the Boltzmann entropy.

What is the probability density of finding a state  $\omega$  at time  $t+t_0$ ? We call the new ensemble  $p^t$ . We recall Liouville's Theorem, page 560, which says that the volume  $L$  is invariant under the transformation  $U^t$ . This implies that  $p^t(U^t\omega) = p(\omega)$ , from which one can infer that  $S_G(p^t) = S_G(p)$ —the Gibbs entropy of an ensemble does not change at all in an isolated system during evolution. This shows that in the case of the evolution of isolated nonequilibrium systems, the evolution of a Gibbs ensemble does not express adequately what we consider thermodynamic behavior. The problem is that even if at the starting time  $t_0$  the Gibbs ensemble was something simple, it can develop in time  $t$  into a very complicated density function that does not correspond to any reasonable macroscopic description.  $\diamond$

**Example 8.4.10** If the ensemble is a discrete state space  $X$  with a probability density function  $P(x)$ , then the Gibbs entropy reduces to

$$S_G(P) = -(k \ln 2) \sum_{x \in X} P(x) \log P(x).$$

This is proportional to Shannon's information-theoretic entropy  $H(P) = -\sum_{x \in X} P(x) \log P(x)$  of Definition 1.11.1, page 67.  $\diamond$

## 8.5 Entropy Revisited

The previous approaches treated entropy as a probabilistic notion; in particular, each individual micro state of a system has entropy equal to 0. However, it is desirable to find a notion of entropy that assigns a nonzero entropy to each micro state of the system, as a measure of its individual disorder.

We will first give a naive form. Then we give a more sophisticated algorithmic form possessing properties similar to Boltzmann entropy, and in which a version of the second law can be proved.

Assume that the set of micro states of a thermodynamic system is discrete. For convenience we identify them with the natural numbers.

**Definition 8.5.1** The *complexity entropy* of a micro state  $x$  of a system is defined as  $K(x)$ , where  $K$  is the prefix complexity of Chapter 3.

Now, a single regular micro state like the state of CuO<sub>2</sub> in Figure 8.10 has low complexity entropy, as opposed to a complex state. Such a micro state looks “regular” and therefore special, only because it has low complexity. Such nonrandom states have complexity entropy close to zero.

The new definition can even be “justified” to a certain degree by Theorem 8.1.1 on page 523. Let  $X$  be the set of micro states of a thermodynamic system. We prove that a discrete version of Gibbs entropy and the expected complexity entropy are quantitatively about the same under some mild assumptions. Consider a discrete ensemble  $X$ , with  $P(x)$  the probability density function that the system is in micro state  $x \in X$ . Then,  $S_G(P) = -(k \ln 2) \sum_{x \in X} P(x) \log P(x)$  is the Gibbs entropy. If  $P$  is recursive, then by Theorem 8.1.1 on page 523 we have

$$S_G(P) = (k \ln 2) \sum_{x \in X} P(x) K(x) + O(1).$$

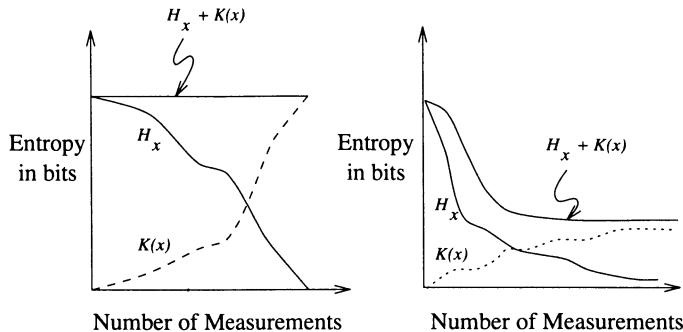
In this sense, the complexity entropy inherits the successes and failures of Gibbs entropy. A serious problem is the choice of precision in describing micro state  $x$ . Our system will certainly have simpler micro states if we have units on the scale of 1 kilometer, rather than of scale  $5 \times 10^{-11}$  meter (the diameter of a hydrogen atom).

### 8.5.1 Algorithmic Entropy

In the previous approaches the macroscopic parameters were not accounted for and appeared out of the blue. Considering a system from an observer’s viewpoint occasions a modified approach. Suppose we have determined the macroscopic parameters describing the classical entropy of a thermodynamical system. Truncate these parameters to the required precision (four decimal places is enough) and encode them in an integer  $x$ .

**Definition 8.5.2** Assume the discussion above, and suppose the system in equilibrium is described by the encoding  $x$  of the approximated macroscopic parameters. The *algorithmic entropy* of the macro state of a system is given by  $K(x) + H_x$ , where  $K(x)$  is the prefix complexity of  $x$ , and  $H_x = S_B(x)/(k \ln 2)$ . Here  $S_B(x)$  is the Boltzmann entropy of the system constrained by macroscopic parameters  $x$ , and  $k$  is Boltzmann’s constant. The physical version of algorithmic entropy is defined as  $S_A(x) = (k \ln 2)(K(x) + H_x)$ .

The second term,  $H_x$ , denotes our “ignorance” about the micro state, given  $x$ . The notion of algorithmic entropy reflects the fact that measurements can increase our knowledge about a system. Initially, we have no knowledge about the state of system, and therefore the algorithmic entropy reduces to Boltzmann entropy.



**FIGURE 8.12.** Algorithmic entropy: left a random micro state, right a regular micro state

If the system is in a nonrandom micro state  $\omega$ , then the  $H_x$  term decreases significantly as we make more and more measurements, while the  $K(x)$  term rises very slowly. Overall, the algorithmic entropy decreases. If a micro state  $\omega$  is random, then we cannot achieve decrease of algorithmic entropy by making more and more measurements. We simply exchange decreased ignorance in  $H_x$  for increased knowledge in  $K(x)$ , for  $x$  is an increasing initial segment of  $\omega$ . The two pictures in Figure 8.12 describe these situations.

**Example 8.5.1 (Maxwell's Demon)** In 1867, J.C. Maxwell described in a letter a thought experiment that has received widespread attention ever since.

"If we conceive a being whose faculties are so sharpened that he can follow every molecule in its course, such a being, whose attributes are still as essentially finite as our own, would be able to do what is at present impossible to us. For we have seen that the molecules in a vessel full of air at uniform temperature are moving with velocities by no means uniform .... Now let us suppose that such a vessel is divided into two portions, A and B, by a division in which there is a small hole, and that a being, who can see the individual molecule, opens and closes this hole, so as to allow only the swifter molecules to pass from A to B and only the slower ones to pass from B to A. He will thus, without expenditure of work, raise the temperature of B and lower that of A, in contradiction to the second law of thermodynamics." J.C. Maxwell, *Theory of Heat*, 1871.

Several solutions to this problem have been proposed and subsequently been disposed of. L. Szilard in 1929 suggested that the required information gathering by the demon saves the second law; and C.H. Bennett in 1982 proposed the current interpretation that it is the destruction of old information that must dissipate energy.

In Figure 8.13 we describe this version of Maxwell's demon: the Szilard engine. The engine runs on a one-molecule gas and is submerged in a

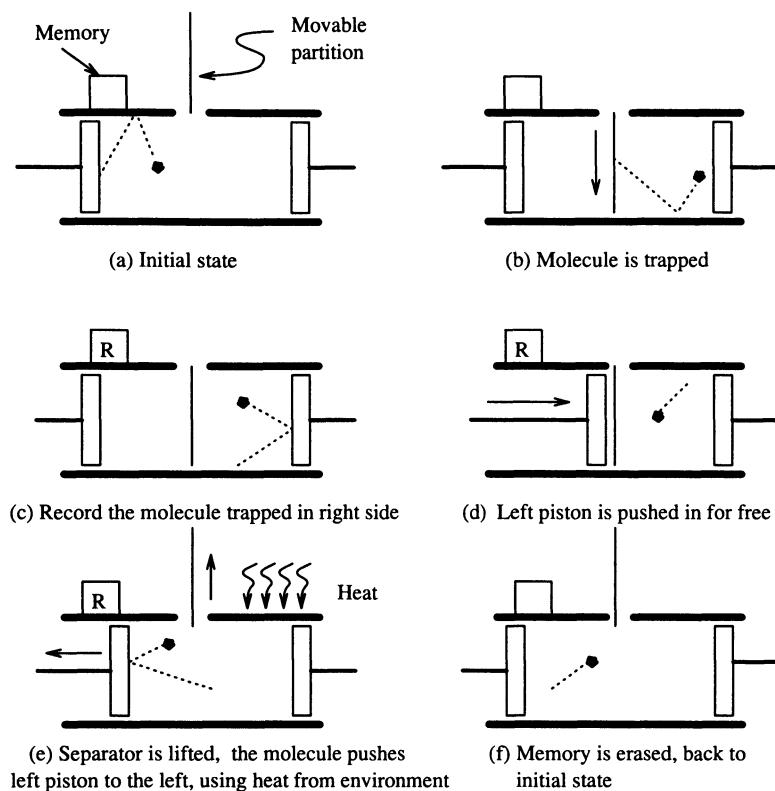


FIGURE 8.13. Szilard engine

heat bath at temperature  $T$ . The initial state of the system is that the molecule is located in either the left or the right chamber of the device and the trap door is up. The demon lowers the trap door, trapping the molecule in either the left or right chamber. The demon records where it is in its memory. It then pushes in the frictionless and massless piston in the chamber that does not contain the molecule. This does not cost work, since there is no pressure (bouncing molecule) to overcome. Since the engine is submerged in a heat bath, the molecule bounces around vigorously. If we raise the trap door, then the molecule will push the piston back to its original position, performing work. At the end of the cycle, we erase memory to resume the initial state.

Common objections against the demon include the idea that opening and closing the door must dissipate energy; the measurement of the molecule's position must dissipate energy (for example, the demon must send at least one photon to the molecule in order to see it), and so on. Such arguments have been found invalid, at least from the theoretical point of view.

However, the demon operating the engine must either store increasingly much information, or it must irreversibly erase some of it. Irreversible erasure of information will dissipate energy by Sections 8.2.1 and 8.3.6. If the demon does not erase its memory each cycle, but keeps on adding new information, then this is not a true cycle. The engine increases the entropy of its memory in order to decrease the entropy of its environment.

The demon can, of course, wait for a while, accumulating information, and then try to guess the shortest program that generates it, and subsequently erase all information by irreversibly erasing only the shortest program for it.

We can obtain some insight into what happens using the algorithmic entropy of Definition 8.5.2. (For a more rigorous treatment we need to develop some theory first; see Theorem 8.5.4 on page 581, Example 8.5.12 on page 581, and Example 8.5.13 on page 581.)

For the sake of the present discussion we make the simplifying assumption that the system consisting of the engine coupled with the computerized demon has entropy  $S_A = (k \ln 2)(K(x) + H_y)$ , where  $k$  is Boltzmann's constant,  $K(x)$  is the complexity of the data in the demon's memory, and  $S_B(y) = (k \ln 2)H_y$  is the Boltzmann entropy of the engine. That is, the demon has no Boltzmann entropy part and the engine has no complexity part.

We make another important assumption by setting the total work performed to the sum of the work needed to change the entropy of the machine and the work needed to change the entropy of the demon. This seems a possibly unrealistic assumption, since it is not obvious that the joint effect cannot be achieved with less work. Think of the triangle inequality. The work to pull something north 1 meter and east 1 meter is more than the work needed to pull it to its final place in a straight line.

**Claim 8.5.1** Assume that the entropy of the system behaves as described above. Then the net heat gained by the engine, coupled with a computerized demon that can perform measurements and control the operation of the engine is no more than

$$\Delta Q = (S_A^f - S_A^i)T,$$

where  $T$  is the temperature in degrees kelvin, and  $S_A^i$  and  $S_A^f$  are the initial and final algorithmic entropies of the system, respectively.

The idea is as follows: The system is operating at temperature  $T$ . Then the heat gained due to the change of Boltzmann entropy, denoted by  $S_B$ , is given by the usual formula

$$\Delta Q^+ = (S_B^f - S_B^i)T.$$

Let  $K(i)$  and  $K(f)$  be the complexity of the demon's initial memory state  $i$  and final memory state  $f$ , respectively. In Section 8.2 we have analyzed the absolute minimum number of bits that need to be erased in a computation from  $i$  to  $f$ , leaving nothing behind except the record  $f$ . This is at least the number of garbage bits left after computing  $f$  from  $i$  by a reversible machine. These garbage bits constitute a program to compute  $f$  from  $i$  (executing the computation from  $i$  to  $f$  in reverse). The length of a self-delimiting program for this purpose is by definition  $K(i|f)$ . By Theorem 3.9.1 on page 232, up to an additive constant,

$$K(i, f) = K(i) + K(f|\langle i, K(i) \rangle) = K(f) + K(i|\langle f, K(f) \rangle).$$

Trivially,  $K(i|f) \geq K(i|\langle f, K(f) \rangle) + O(1)$ . Therefore,  $K(i|f) \geq K(i) - K(f + O(1))$ . (See also Section 8.3.6.) The heat loss to update the demon's memory is the negative of the number of erased bits:

$$\Delta Q^- = (K(f) - K(i))kT \ln 2.$$

Thus, the total heat  $\Delta Q$  gained by the demon is given by

$$\begin{aligned} \Delta Q^+ + \Delta Q^- &= [(S_B^f + (k \ln 2)K(f)) - (S_A^i + (k \ln 2)K(i))]T \\ &= (S_A^f - S_A^i)T. \end{aligned}$$

With  $S_A$  subject to the second law,  $\Delta Q$  in the above claim must be less than or equal to 0.  $\diamond$

### 8.5.2 Algorithmic Entropy and Randomness Tests

Algorithmic entropy can be based on the profound idea of the universal randomness test as developed in Section 2.5 and especially Section 4.5. Again, consider the molecules of an ideal gas as points that move in 3-dimensional space  $\mathcal{R}^3$ , where  $\mathcal{R}$  is the set of real numbers, and use the model of statistical mechanics as described on page 559 and the following pages.

**Definition 8.5.3** Let  $x$  be a binary string of length  $n$ . An  $n$ -cell is the set  $\Gamma_x$  of all infinite binary sequences (micro states) with finite initial segment  $x$ , defined by  $\Gamma_x = \{\omega : \omega = x\zeta, \zeta \in \{0, 1\}^\infty\}$ .

The sets  $\Gamma_x$  are the well-known cylinder sets of measure theory (Section 1.6). The area of state space we consider is divided into  $n$ -cells. For  $n = 1, 2, \dots$  this division becomes progressively finer-grained and discerns more detail. To say that a micro state  $\omega$  is in  $n$ -cell  $\Gamma_x$  means that  $x$  is a prefix of  $\omega$ .

In fact,  $\omega$  is the limit point to which the infinite nested sequence of  $n$  cells, which contain  $\omega$ , converges with increasing  $n$ .

**Definition 8.5.4** Assume the notation above and recall Notation 8.4.1 on page 560. The *algorithmic entropy* of an  $n$ -cell  $\Gamma_x$  with respect to  $\mu$  is defined as

$$H_\mu(\Gamma_x) = K(x|\mu) + \log \mu(\Gamma_x).$$

In terms of Definition 8.5.2, page 566, the algorithmic entropy of an  $n$ -cell  $\Gamma_x$  consists of two parts, where the first part is the prefix complexity of  $x$  (the macroscopic description of the  $n$ -cell) and the second part is the logarithm of the measure of volume  $\Gamma_x$  (our lack of knowledge about a member of the  $n$ -cell). This definition can serve to define successively closer approximations to the entropy of a micro state as the minimum of the entropies of successively smaller  $n$ -cells containing it.

For computable measures  $\mu$ , we have  $K(x|\mu) \geq K(x) - K(\mu)$  for all  $x$ , with  $K(\mu)$  a constant independent of  $x$  that is the length of the shortest program to compute  $\mu$ . Since by definition  $K(x|\mu) \leq K(x) + O(1)$ , we have  $K(x|\mu) = K(x) \pm O(K(\mu))$ .

**Definition 8.5.5** Let  $\mu$  be a recursive measure. The *algorithmic entropy* of  $\omega \in X$  with respect to  $\mu$  is defined by

$$H_\mu(\omega) = \inf_n \{K(\omega_{1:n}|\mu) + \log \mu(\Gamma_{\omega_{1:n}})\}.$$

This definition sets  $H_\mu(\omega)$  equal to  $-\rho_0(\omega|\mu)$ , the universal integral  $\mu$ -test of Corollary 4.5.2, page 291. (Recall Notation 4.2.1.) For the special case where  $\mu$  is the uniform measure see Corollary 3.6.1, page 214.

**Example 8.5.2** (**Generalized Prefix Complexity**) Let  $X = \{0, 1, \dots\}$ , and for every  $x \in X$ , set  $\mu(\{x\}) = 1$ . Then  $H_\mu(x) = K(x) + O(1)$ . Thus, the algorithmic entropy  $H_\mu$ , defined by way of randomness tests, is a *generalization of the prefix complexity*  $K$ . ◇

**Example 8.5.3** How large can  $H_\mu(\omega)$  be? The higher it is, the more random is  $\omega$ . For a finite measure  $\mu$  ( $\mu(X) \leq \infty$ ), the function  $H_\mu$  is bounded above. Such a bound is the ultimate maximum, or equilibrium, of algorithmic entropy in such a system, and it is reached when the system has evolved to total disorder. For infinite measures, it can take arbitrarily large positive values; but it will never be  $+\infty$ . By the theory of testing,  $\omega$  is random in the sense of Martin-Löf (Sections 2.5, 4.5.6) iff  $-H_\mu(\omega) = \rho_0(\omega|\mu) < \infty$ . In other words, a micro state  $\omega$  is random iff  $H_\mu(\omega) > -\infty$ .

To give some idea how  $H_\mu(\omega)$  depends on  $\omega$  and  $\mu$  we give an upper bound. Let  $m(x) = \mu(\Gamma_x)$ . For  $\omega \in \Gamma_x - \Gamma_{x0}$ , we have  $H_\mu(\omega) < \log(m(x) + 1) + 2 \log(\log(m(x) + 1) + 1) + O(1)$ . The right-hand side is an upper bound on the amount by which  $K(x|\mu)$  can exceed  $\log \mu(\Gamma_x)$ . ◇

**Example 8.5.4** The sequence  $\eta = 00\dots$  is nonrandom in the uniform measure  $\lambda$ , and hence (Sections 2.5, 4.5.6)  $\rho_0(\eta|\lambda) = \infty$ . This means that the algorithmic entropy  $H_\lambda(\eta) = -\infty$ . The average  $\omega$  satisfies  $\rho_0(\omega|\lambda) < \infty$ , and therefore  $H_\lambda(\omega) > -\infty$ . (This is the case for all random sequences, and the set of these sequences has  $\lambda$ -measure 1, by Theorem 2.5.3, page 144.)

But some  $\eta \in \{0, 1\}^\infty$  have  $H_\lambda(\eta) = -\infty$ . These are precisely the sequences that are not  $\lambda$ -random.  $\diamond$

**Example 8.5.5** For some recursive measures  $\mu$  there are  $n$ -cells  $\Gamma_x$  such that all sequences in  $\Gamma_x$  are not  $\mu$ -random. Namely, if  $\mu(\Gamma_x) = 0$ , then  $\log \mu(\Gamma_x) = -\infty$ , while  $K(x|\mu)$  is finite. It follows from the definition of algorithmic entropy that then  $H_\mu(\omega) = -\infty$  for all  $\omega \in \Gamma_x$ . (By the way, existence of such a  $\Gamma_x$  means that  $\mu$  is not the uniform measure.)

If  $H_\mu(\omega) = -\infty$ , then the system is in an orderly state. For example, all molecules of a gas are in the left half of the containing vessel. The set of all such orderly states has zero probability in the sense of measure  $\mu$ .  $\diamond$

**Example 8.5.6** We show a connection between  $H_\mu$  and the Gibbs entropy  $S_G$ . Consider an ensemble  $X$  with recursive probability density function  $p$  and measure  $\mu$ . Define  $\nu(d\omega) = p(\omega)\mu(d\omega)$ . Then it can be shown that

$$S_G(p)/(k \ln 2) = \int p(\omega) \log p(\omega) \mu(d\omega) = \int H_\mu(\omega) \nu(d\omega) + O(1).$$

This means that the Gibbs entropy is essentially the average algorithmic entropy, while the latter measures the randomness of  $\omega$  with respect to  $\mu$ . The higher it is, the more random is  $\omega$ .  $\diamond$

The  $H_\mu$  measure is important, but it does not satisfy the second law of thermodynamics since we will show in Theorem 8.5.2, page 576, that it has a “no strong increase” property. We solve this problem by considering coarse-grained entropy.

**Definition 8.5.6** Assume the above notation. The  $n$ th approximation of  $H_\mu(\omega)$ , the *coarse-grained algorithmic entropy*, is given by  $H_\mu^n(\omega)$  defined as

$$H_\mu^n(\omega) = \inf_{i \leq n} \{H_\mu(\Gamma_{\omega_{1:i}})\}.$$

**Example 8.5.7** If the notion of temperature is meaningful in our system, then the physical version of algorithmic entropy is to be defined as

$$S_A(\omega) = (k \ln 2) H_\mu(\omega), \quad S_A^n(\omega) = (k \ln 2) H_\mu^n(\omega),$$

where  $k$  is the very small Boltzmann’s constant. For appropriate fixed  $n$ , this new quantity  $(k \ln 2) H_\mu^n(\cdot)$  is our corrected version of coarse-grained Boltzmann entropy  $S_B(\cdot)$ .  $\diamond$

The coarse-grained algorithmic entropy  $H_\mu^n$  is important since we will show in Theorem 8.5.2, page 576, and page 578, that only  $H_\mu^n$  satisfies the second law, by giving the required strong growth of entropy, similar to Boltzmann entropy. The quantity  $H_\mu^n$  is mathematically objective and is related to the observation dependent term  $H_x$  in Definition 8.5.2.

**Notation 8.5.1** The subscript  $\mu$  is usually understood from the context. We dispense with the subscript from now on, using it only when there is risk of confusion.

The relations between the algorithmic entropy of a micro state, the coarse-grained algorithmic entropy, and the  $n$ -cell containing it, are as follows:

**Theorem 8.5.1** *Assume the notation above. Then*

$$H(\omega) = \lim_{n \rightarrow \infty} H^n(\omega) = \inf_{x \in \{0,1\}^*} \{H(\Gamma_x) : \omega \in \Gamma_x\}. \quad (8.9)$$

**Proof.** This follows directly from the definitions.  $\square$

The formula for  $H(\omega)$  says that to determine its value we should consider ever finer-grained  $n$ -cells (test more bits of  $\omega_{1:n}$ ), but only as long as the complexity increase of  $K(\cdot)$  buys us an even greater decrease in  $\log \mu(\Gamma_{\omega_{1:n}})$ . For  $\omega$  that are not  $\mu$ -random, the latter quantity goes down all the way to  $-\infty$  with growing  $n$ .

For  $\mu$ -random  $\omega$ , the value of  $H^n(\omega)$  goes down to at most a finite (possibly negative) value with  $n$ . Hence, there is a value  $n_0$  such that  $H^n(\omega)$  reaches its infimum for  $n = n_0$ . This is the optimal granularity of description for this micro state that yields the final algorithmic entropy. The term  $\log \mu(\Gamma_x)$  represents our a priori uncertainty about  $\omega$  as an element of its  $n$ -cell  $\Gamma_x$ . This corresponds to the second term,  $H_x$ , in Definition 8.5.2.

The following result says that for most elements  $\omega$  of an  $n$ -cell  $\Gamma_x$ , the value of  $H(\omega)$  cannot be much less than  $H(\Gamma_x)$ . That is, if some elements of a cell are (sufficiently) random, then most of them are (sufficiently) random. To appreciate the lemma below, recall that  $K(n) \leq \log n + 2 \log \log n + O(1)$ , for all  $n$ .

**Lemma 8.5.1 (Cell Stability)** *Let  $n = l(x)$ . Then for some constant  $c$ ,*

$$\mu\{\omega \in \Gamma_x : H(\omega) < H(\Gamma_x) - K(n) - m\} < c 2^{-m} \mu(\Gamma_x).$$

**Proof.** Assume the notation above. Let

$$f(x) = 2^{-K(l(x))} \int_{\Gamma_x} 2^{-H(\omega)} \mu(d\omega). \quad (8.10)$$

Clearly,  $f$  is enumerable. We prove  $\sum_x f(x) \leq 1$ , showing that  $f$  is a discrete semimeasure. Since  $\int_X 2^{-H(\omega)} \mu(d\omega) \leq 1$ , we can rewrite  $\sum_x f(x)$  as

$$\begin{aligned} \sum_{n=1}^{\infty} \sum_{l(x)=n} f(x) &= \sum_{n=1}^{\infty} 2^{-K(n)} \int_X 2^{-H(\omega)} \mu(d\omega) \\ &\leq \sum_{n=1}^{\infty} 2^{-K(n)} \leq 1. \end{aligned}$$

Therefore, by Theorem 4.3.1, page 247, and Theorem 4.3.3, page 253,  $f(x) \leq c_1 2^{-K(x)}$  for some constant  $c_1$ . Writing  $f(x)$  according to Equation 8.10,

$$2^{-K(l(x))} \int_{\Gamma_x} 2^{-H(\omega)} \mu(d\omega) \leq c_1 2^{-K(x)}.$$

Substituting according to Definition 8.5.4, page 571, and rearranging yields

$$\mu(\Gamma_x)^{-1} \int_{\Gamma_x} 2^{-H(\omega)} \mu(d\omega) < c_2 2^{-H(\Gamma_x) + K(l(x))}, \quad (8.11)$$

with  $c_2$  a new constant possibly depending on  $\mu$ . Set  $n = l(x)$ . Rewrite the condition

$$H(\omega) < H(\Gamma_x) - K(n) - m$$

as

$$2^{-H(\omega)} > 2^{-H(\Gamma_x) + K(n) + m}.$$

Applying Markov's Inequality ( $\mu(Y > t) \leq \mathbf{E}(Y)/t$  on page 261), we obtain

$$\begin{aligned} \mu\{\omega \in \Gamma_x : 2^{-H(\omega)} > 2^{H(\Gamma_x) - K(n) - m}\} \\ \leq \int_{\Gamma_x} 2^{-H(\omega)} \mu(d\omega) / 2^{-H(\Gamma_x) + K(n) + m}. \end{aligned} \quad (8.12)$$

Combining Equations 8.11 and 8.12, the lemma follows.  $\square$

**Example 8.5.8** For the systems of interest in physics, and for those precisions in the macroscopic parameters for which Boltzmann's entropy is generally considered, we expect the additive term  $K(\omega_{1:n})$ , which is at most  $n+2 \log n$ ,

to be negligible compared to the second term of algorithmic entropy. This log-volume term is usually very large.

For example, consider a mole of gas (about  $6 \times 10^{23}$  molecules, page 558) in a container. Let us be extremely generous and assume that each atom can only assume two states (or each of the  $6N$  dimensions is of size 2 in the continuous case). Then the log-volume is about  $10^{23}$  bits in order to describe the complete state of a few grams of matter. This is the order of magnitude of the total entropy.

The macroscopic information given to us in  $\omega_{1:n}$  is generally just the values of some macroscopic parameters. We may, for example, subdivide the piece of magnet under study into 100 parts and measure the magnetization of each to an unlikely precision of 10 bits: this is still only 1000 bits out of  $10^{23}$  bits. Therefore,  $K(\omega_{1:n})$  is negligible compared to the second log-volume term in such cases. Hence,  $H(\omega)$  is not appreciably larger than its log-volume term when expressed this way, and can be interpreted as an approximation to the Boltzmann entropy.  $\diamond$

### 8.5.3 Entropy Stability and Nondecrease

The following mathematical conditions are satisfied by classic thermodynamic systems. Each point in a state space  $X$  describes a possible micro state of the system. Starting from each position, the point representing the system describes a trajectory in  $X$ . We take time to be discrete and let  $t$  range over the integers.

If the evolution of the system in time can be described by a transformation group  $U^t$  of  $X$  such that  $U^t\omega$  is the trajectory starting from some point  $\omega$  at time 0, then we call the system *dynamical*. Assume that  $U^t\omega$  is computable as a function of  $t$  and  $\omega$  for both positive and negative  $t$ .

#### Definition 8.5.7

Assume that the involved measure  $\mu$  is recursive and satisfies  $\mu(U^tV) = \mu(V)$  for all  $t$  and all measurable sets  $V$  (like the  $\Gamma_x$ 's), as in Liouville's Theorem, as described on page 560 and the following pages. The  $\mu$ -measure of a volume  $V$  of points in state space is invariant over time. A *physical system*  $\mathcal{X}$  consists of the space  $X$ ,  $n$ -cells  $\Gamma_x$ , transformation group  $U^t$ , and the recursive and volume invariant measure  $\mu$ .

#### Example 8.5.9

In general, the partitions of the space given by the  $n$ -cells have no connection with the dynamics of  $U^t$ . In particular, they are not “generated” by  $U^t$  from the initial partition into precisely  $\Gamma_0$  and  $\Gamma_1$ . An exception is the *baker's map*. Let  $X$  be the set of doubly infinite binary sequences

$$\omega = \dots \omega_{-2} \omega_{-1} \omega_0 \omega_1 \omega_2 \dots$$

with the shift transition  $U^t\omega_i = \omega_{i+t}$  over discrete time, and the uniform measure  $\lambda$ . Write  $\omega^n = \omega_{-\lfloor n/2 \rfloor} \dots \omega_{\lceil n/2 \rceil - 1}$ . The  $n$ -cells are  $\Gamma_{\omega^n}$ . They all have the same measure,  $2^{-n}$ .  $\diamond$

**Example 8.5.10** The physical version  $S_A$  of algorithmic entropy differs only by a constant factor  $k \ln 2$  from  $H$ , so we can keep working with  $H$ . Due to the very small size of the Boltzmann constant  $k$ , a difference in  $H$  translates to a much smaller difference in  $S_A$ . Therefore, an assertion of how close  $H$  is to a maximum results in a much stronger such assertion about  $S_A$ . We show that the proportional size of the set of points where the entropy is not close to its maximum is very small. Define the event  $E_m$  by

$$E_m = \{\omega : H(\omega) < -m\}. \quad (8.13)$$

The  $\mu$ -measure of  $E_m$  is

$$\mu\{\omega : \rho_0(\omega|\mu) > m\} \leq 2^{-m}\mu(X).$$

This follows from the definition of the universal integral  $\mu$ -test of Corollary 4.5.2, page 291. That is, the proportional size of the point volume where the algorithmic entropy is far from its maximum is very small.  $\diamond$

We want to analyze the behavior of entropy  $H$  over time. Fix a finite time interval  $\{1, \dots, T\}$  and investigate the dependence of the function  $H(U^t\omega)$  on  $t$ . We assumed that  $U^t\omega$  is computable. The next statement says that the entropy does not decrease (significantly) over time. On page 233 we have defined the notion  $I(x : y) = K(y) - K(y|x)$  as the information that  $x$  carries about  $y$ . We consider a more general form.

**Definition 8.5.8** The *information that  $x$  carries about  $\omega$  with respect to  $\mu$*  is defined as  $I_\mu(x : \omega) = H_\mu(\omega) - H_\mu(\omega|x)$ .

**Theorem 8.5.2** **(Bounded Entropy Change)** *Use the notation above, and let everything be with respect to the measure  $\mu$ . Up to additive constants,*

$$-I(t : \omega) \leq H(U^t\omega) - H(\omega) \leq I(t : U^t\omega).$$

**Proof.** We sketch the argument and omit details. We require the notion of integral  $\mu$ -tests, as treated in Section 4.5.6. The lower bound is argued by showing that  $-H(U^t\omega)$  is an integral  $\mu$ -test parametrized by  $t$ , and that  $-H(\omega) + I(t : \omega)$  is a universal integral  $\mu$ -test parametrized by  $t$ . This shows the lower bound

$$H(U^t\omega) \geq H(\omega) - I(t : \omega).$$

Because  $U$  is reversible, we have  $U^{(-1)t}U^t\omega = \omega$ . Hence, we can repeat the proof above going from  $U^t\omega$  to  $U^{(-1)t}U^t\omega$ . This shows the upper bound.  $\square$

From the proof sketch we see that if  $U^t$  should increase the entropy by a too large amount, then  $U^{-t}$  would not be able to decrease it back to

$H(\omega)$ . The theorem states that the only decrease or increase from  $H(\omega)$  to  $H(U^t\omega)$  we will ever see must be due to the information that  $t$  may have on  $U^t\omega$  or on  $\omega$ . We can bound the entropy increase and decrease more crudely.

**Lemma 8.5.2**  $|H(U^t\omega) - H(\omega)| \leq K(t) + O(1)$ .

**Proof.** By Theorem 8.5.2, we have  $H(U^t\omega) - H(\omega) \leq I(t : U^t\omega) + O(1)$ . Define

$$g(\omega) = \sum_z 2^{-H(\omega|z) - K(z)}.$$

Then  $g$  is enumerable and

$$\int_X g(\omega) \mu(d\omega) \leq \sum_z 2^{-K(z)} \leq 1.$$

Therefore,  $g(\omega) = O(2^{-H(\omega)})$ . The same inequality holds for one term of the sum expression of  $g$ :

$$2^{-H(\omega|z) - K(z)} = O(2^{-H(\omega)}).$$

That is,  $I(t : U^t\omega) \leq K(t) + O(1)$ , and hence we have shown one direction of the lemma. The proof of the other direction is similar.  $\square$

This means that for some simple values of  $t$ , the algorithmic entropy hardly increases at all. For example, if  $t$  is of the form  $2^n$ , then the increase is at most  $2 \log \log t$ . The following result shows that the second law is satisfied in that decreases are very infrequent and get more infrequent the larger they are. Let  $d(A)$  denote the number of elements in a set  $A$ .

**Theorem 8.5.3 (Entropy Nondecrease)** *Assume the notation above. Consider the time interval  $\{1, \dots, T\}$ . Then, for  $1 \leq t \leq T$ ,*

$$d(\{t : H(U^t\omega) < H(\omega) - K(T) - m\}) = O(2^{-m}T).$$

**Proof.** Replace the statement in the theorem by the equivalent statement

$$d(\{t : 2^{-H(U^t\omega)} > 2^{-H(\omega)+K(T)+m}\}) = O(2^{-m}T). \quad (8.14)$$

By Markov's Inequality about expectations, Example 4.3.10, page 261, the left side of Equation 8.14 is bounded above by

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T 2^{-H(U^t\omega)} / 2^{-H(\omega)+K(T)+m} &\leq \sum_{t=1}^T 2^{-H(U^t\omega)} / 2^{-H(\omega)+K(T)+m} \\ &= \sum_{t=1}^T 2^{H(\omega)-H(U^t\omega)} / 2^{K(T)+m}. \end{aligned}$$

Substitute according to the Claim 8.5.2 below and rearrange to obtain the right side of Equation 8.14.

$$\text{Claim 8.5.2 } \sum_{t=1}^T 2^{H(\omega) - H(U^t \omega)} = O(T 2^{K(T)}).$$

**Proof.** For each fixed  $T$ , the function  $f(\omega, T) = T^{-1} \sum_{t=1}^T 2^{-H(U^t \omega)}$  is an enumerable unit integrable function as treated in Section 4.5.6, namely,

$$\int_X f(\omega, T) \mu(d\omega) < \int_X 2^{-H(\omega)} \mu(d\omega) \leq 1.$$

By Theorem 4.5.5 on page 288 (everything parametrized by  $T$ ) the function  $2^{-H(\cdot|T)}$  is a universal enumerable unit integrable function with respect to measure  $\mu$ . Therefore,  $-\log f(\omega, T) \geq H(\omega|T) + O(1)$ . From the definitions we have  $H(\omega|T) + O(1) \geq H(\omega) - K(T) + O(1)$ . Therefore,

$$T^{-1} \sum_{t=1}^T 2^{-H(U^t \omega)} = O(2^{-H(\omega) + K(T)}).$$

Rearrangement yields the claim. □

The theorem says that no matter what interval of time we consider, if the maximum entropy during that interval is  $H$ , then the proportion of time in which the entropy can be  $H - m - K(T)$  will be only about  $2^{-m}$ .

The decrease  $K(T)$  can be made very small by choosing a simple value of  $T$ . If we choose it to be of the form  $T = 2^n$  for a positive or negative integer  $n$ , then  $K(T)$  is at most  $2 \log \log T$ , which is generally negligible even if a unit of  $T$  is of the duration of some atomic event.

We could have chosen an arbitrary interval  $[s, s+T]$  since  $U^{s+t}\omega = U^t(U^s\omega)$ . It is remarkable that the theorem is true on all time scales, not only on the time scale on which the ergodic theorem works (if it is applicable). Thus, if  $\omega$  has an orbit on which the entropy reaches near its upper bound  $\log X$  at some  $t_0$ , then the entropy will stay near this upper bound on the largest proportion of each interval around  $t_0$ . And the decrease of entropy always happens in an “accelerated” fashion. The further the decrease has gone, the more “urgent” it is, since the relative time spent in the lower regions of entropy is decreased by a factor of 2 for every step of decrease, on all time scales.

On the interval  $\{1, \dots, T\}$  the “steepest” (measured from  $T$ ) monotonically increasing function of  $t$  that has the behavior proved here is  $H(U^t \omega) = \log t$ .

## Chaotic Systems

We now argue that the coarse-grained algorithmic entropy  $H^n$  can increase rapidly, like Boltzmann entropy, and in some cases is actually preferable. An example of a chaotic transformation in which Boltzmann

entropy (and the fine-grained algorithmic entropy,  $H$ ) does not increase is the baker's map of Example 8.5.9 on page 575. Since for fixed  $n$  all  $n$ -cells have the same measure no matter what fixed precision  $n$  we choose, the Boltzmann entropy of  $U^t\omega$  does not increase with  $t$ .

The coarse-grained algorithmic entropy  $H^n(U^t\omega)$  can, however, increase rapidly with  $t$ . Let us consider a typical example of an infinite sequence  $\omega$  that has  $\omega^n = 0 \dots 0$  and whose other bits are random. Then  $K(\omega^n) \leq 2 \log n$ ,  $\log \lambda(\Gamma_\omega^n) = -n$ , and therefore,

$$H^n(\omega) \leq 2 \log n - n.$$

Now, for  $t > n$ , the string  $(U^t\omega)^n$  consists of random bits; so essentially,  $H^n(U^t\omega) \geq 0$ . Therefore, between time 0 and  $n$ , this algorithmic coarse-grained entropy increases linearly from  $-n$  to 0. The fine-grained algorithmic entropy  $H(U^t\omega)$  will, on the other hand, only increase slowly since in our choice of the precision  $n$ , we can follow the increase of  $t$ ,

$$\begin{aligned} H(U^t\omega) &\leq H^{t+n}(U^t\omega) + O(1) \leq K(n) + K(t) + t - (t + n) + O(1) \\ &\leq 2 \log n + 2 \log t - n + O(1). \end{aligned}$$

On the right-hand side of the second inequality of the first line, the first two terms upper-bound the complexity of  $n$  and  $t$ ; the third term is the length of the nonzero part of the string  $(U^t\omega)^{t+n}$ . Therefore,  $K((U^t\omega)^{t+n})$  is bounded by the sum of the first three terms, while the last term gives the log-volume of the cylinder. This inequality shows that  $H(U^t\omega)$  can only grow as slowly as  $\log t$ . Theorem 8.5.2 implies that for some simple values of  $t$ , the algorithmic entropy  $H(U^t\omega)$  hardly increases at all. For example, if  $t$  is an integer of the form  $2^n$ , then the increase is at most  $2 \log \log t$ .

The reason that Boltzmann entropy does not increase in the baker's map (pages 575, 578) is that all  $n$ -cells have equal volume. If we use, instead,  $n$ -cells that are extremely different in size, then Boltzmann entropy will probably increase. By "extremely different," we mean that even the sum of the volumes of the small cells is much smaller than the volume of some larger cells.

In fact, when the quantity  $H^n(\omega)$  is high this can be treated as a concise expression of the fact that in every "simple" partition with extremely different cell sizes, the point  $\omega$  would end up in a large cell. Classical statistical mechanical systems are distinguished by the property that every reasonable canonical cell division will have this property, Example 8.4.8.

Perhaps the quantity  $H^n(\omega)$  extends the idea of entropy increase to a wider class of chaotic systems than the baker's map and can also serve as a useful tool for formulating conjectures concerning general notions of chaos.

- Example 8.5.11** Let us consider the much more speculative problem of approach to equilibrium, that is, the argument that  $H^n(U^t\omega)$  must increase fast if it is far from its upper bound  $\log \mu(X)$ . Let  $E_m$  be as in Definition 8.13,

page 576, and let  $E_m^n$  be the union of  $n$ -cells contained in  $E_m$ . The measure of this set is at most that of  $E_m$ , and the set is also “simple” since it is the union of  $n$ -cells. If  $H^n(U^t\omega) < \log \mu(X) - m$ , then  $U^t\omega \in E_m^n$ , that is,  $\omega \in U^{-t}E_m^n$ .

According to standard (heuristic) reasoning on chaotic systems, as  $t$  grows, the measure-preserving transformation  $U^t$  transforms a small and “compact” (since it consists of the union of  $n$ -cells) set into a “thin” one, which has a small intersection with all  $n$ -cells. (How long this takes depends on how narrow the communication channels are among the different portions of the state space for the flow  $U^t$ .)

In other words, after a while, the  $n$ -cells will be “mixed.” The same is true of  $U^{-t}$ ; therefore, the intersection of  $U^{-t}E_m^n$  with the  $n$ -cell containing  $\omega$  at time  $t = 0$  ( $\Gamma_x$ , with  $x = \omega_{1:n}$ ) will be only a small portion of the  $n$ -cell. Therefore, the conditional probability that a micro state, starting from an arbitrary  $n$ -cell, ends up in  $E_m^n$  is small.

The example of the shift transformation of the baker’s map suggests that in some chaotic systems the parameter  $n$  can be made a function of  $t$  as long as it grows slower than linearly with  $t$ . Thus, if  $n(t)$  is a function of  $t$  such that  $\lim_{t \rightarrow \infty} n(t)/t = 0$ , then

$$H^{n(t)}(U^t\omega)$$

will approach  $\log \mu(X)$  almost as fast as if we held  $n$  constant. The speed at which  $n(t)$  can grow relates to the speed of mixing of  $n$ -cells under  $U^t$ . It is a measure of the degree of chaos of  $U^t$ .  $\diamond$

## Joint Systems

Let us consider two systems, say  $\mathcal{X} = (X, \mu, U_X)$  and  $\mathcal{Y} = (Y, \mu, U_Y)$ , where  $\mathcal{Y}$  is considered the environment from which  $\mathcal{X}$  is temporarily isolated and  $\mu$  is the invariant measure. For example, the subsystems could be the heat engine plus computerized demon versus the environment of Example 8.5.1, page 567. Our systems satisfy all common requirements: the transformation groups leave the measure invariant, and so on. Assume that these systems are independent at time 0, when we want to “do something” to  $\mathcal{X}$ .

Consider a third system that is the combination of  $\mathcal{X}$  and  $\mathcal{Y}$ . The impulses and momenta of the joint system are simply the impulses and momenta of the two subsystems; therefore, the measure on  $Z = X \times Y$  is the product of the original measures on  $X$  and  $Y$ . The transformation group in the joint system is  $U^t(\omega, \zeta)$  with  $\omega \in X$  and  $\zeta \in Y$ . Use the notation  $U^t\omega = \omega^t$  and  $(\omega^t, \zeta^t) = U^t(\omega, \zeta)$ .

On page 235 we have given the *mutual information* between two finite strings  $x$  and  $y$  as  $I(x; y) = K(x) + K(y) - K(x, y) + O(1)$ . We consider a more general form.

**Definition 8.5.9** The *mutual information* between infinite sequences  $\omega$  and  $\zeta$  with respect to  $\mu$  is defined as  $I_\mu(\omega; \zeta) = H_\mu(\omega) + H_\mu(\zeta) - H_\mu(\omega, \zeta)$ .

**Theorem 8.5.4** (**Entropy Balance**) *Assume the above notation. Let  $\Delta H(\omega) = H(\omega_t) - H(\omega)$ . Then*

$$\Delta H(\omega) + \Delta H(\zeta) \geq I(\omega^t; \zeta^t) - I(\omega; \zeta) - I(t : (\omega, \zeta)) - O(1).$$

**Proof.** For the joint system, Theorem 8.5.2 implies  $\Delta H(\omega, \zeta) \geq -I(t : (\omega, \zeta)) - O(1)$ . Note that this formula cannot be applied now to the parts of the system since they do not have their own transformations now. Using this formula, we have (up to additive  $O(1)$  terms)

$$\begin{aligned} H(\omega^t) + H(\zeta^t) &= H(\omega^t, \zeta^t) + I(\omega^t; \zeta^t) \\ &\geq H(\omega, \zeta) - I(t : (\omega, \zeta)) + I(\omega^t; \zeta^t) \\ &= H(\omega) + H(\zeta) + I(\omega^t; \zeta^t) - I(\omega; \zeta) - I(t : (\omega, \zeta)), \end{aligned}$$

which gives the statement by rearrangement.  $\square$

**Example 8.5.12** If the two systems were originally independent (this means  $I(\omega; \zeta) \approx 0$ ), then apart from those rare times  $t$  that contain information about  $(\omega, \zeta)$ , a decrease in the algorithmic entropy of  $\omega$  must be accompanied by an increase in the algorithmic entropy of  $\zeta$ .

The entropy balance theorem is not new, of course, for Boltzmann entropy. But there is a situation that is apparently not easily treated by Boltzmann entropy and may be better understood in terms of  $H(\omega)$ . For example, let  $\zeta$  be Maxwell's demon, Example 8.5.1 on page 567, and let  $\omega$  be a micro state of the gas whose entropy it is trying to decrease. Then Theorem 8.5.4 provides an explanation, in a neat inequality, of why it is that the demon will, after a while, have trouble concentrating.

Landauer (Section 8.2.1) claims that irreversibly erasing one bit of information from a system requires “heat dissipation” of size  $kT \ln 2$ . If we interpret “erasing a bit” as decreasing the algorithmic entropy of  $\omega$  by 1 and instead of “dissipating heat by  $kT \ln 2$ ” we write “increasing the algorithmic entropy of  $\zeta$ , the environment of  $\omega$ , by 1,” then our algorithmic entropy balance theorem confirms the thesis.  $\diamond$

**Example 8.5.13** Often, the first term  $K(\cdot | \mu)$  of the algorithmic entropy  $\inf_n \{K(\omega_{1:n} | \mu) + \log \mu(\Gamma_{\omega_{1:n}})\}$  is insignificant compared to the second term, which measures the Boltzmann entropy. An example of a system for which also the complexity term of algorithmic entropy is needed is a computer memory.

The memory cells are not individual atoms but they are rather tiny. The total number of bits in a modern computer is already measured in gigabytes. It is natural to view all the information in the computer memory

as macroscopic when we describe this system. However, as the memory cells get smaller, the boundary between macroscopic and microscopic information becomes somewhat blurred, and when describing the whole system one will probably want to use the whole formula. All memory states of the computer would have about the same classical Boltzmann entropy, but more complex memory states would have larger algorithmic entropy. Thus, the tendency of noise to increase the complexity of memory could be considered the same phenomenon as the entropy increase in classical physical systems.

Such computer systems are typical of the systems considered by Landauer. Interpret “erasing a bit,” done by irreversible computation as in Section 8.2, as decreasing the entropy by 1. Suppose that all of the memory together contains the string  $x$ . Then we may say that  $H(\omega) = K(x) + \log \mu(\Gamma_x)$ , where  $\omega \in \Gamma_x$ , and  $\Gamma_x$  is a subset of the memory states that can be described by saying what bits are in the memory. This means that  $K(\Gamma_x) = K(x)$ .

Erasing the information means putting 0’s everywhere into the memory. Denote the micro state of the environment by  $\zeta$ . Then the development of the system consisting of memory and environment in time is denoted by  $U^t(\omega, \zeta) = (\omega^t, \zeta^t)$ . Arguably, we will have  $H(\omega^t) = K(y) + \log \mu(\Gamma_y)$ , where  $\omega^t \in \Gamma_y$ . Here,  $\Gamma_y$  just says that the memory contains all 0’s; hence  $K(\Gamma_y) = O(1)$ .

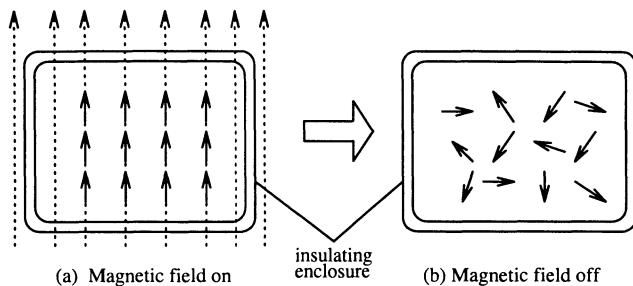
If we require that the memory have approximately the same amount of Boltzmann entropy  $\log \mu(\Gamma_x)$  no matter what bits it holds, then this means  $\log \mu(\Gamma_x) = \log \mu(\Gamma_y)$ . In this interpretation, the erasure indeed means  $\Delta H(\omega) = K(x)$ , which will be balanced by a similar increase in  $H(\zeta)$ .

In order to argue, on the basis of Boltzmann entropy alone, that turning all these bits to 0 results in heat dissipation, we need to consider this operation of turning all bits to 0 as some kind of general operation that does something similar with all possible memory contents, that is, decreases the volume of the whole state space. This artificiality is avoided here. ◇

**Example 8.5.14** Let us go back to Example 8.4.7 on page 562 again. The state of CuO<sub>2</sub> in Figure 8.10 can be described by a program of a few bits:

```
repeat forever : print ↑; print ↓.
```

Thus, the complexity term in the algorithmic entropy almost vanishes. Similarly, the sequence of arrows that encodes 3.1415… has a very low complexity term, and corresponding low algorithmic entropy. ◇



**FIGURE 8.14.** Adiabatic demagnetization to achieve low temperature

**Example 8.5.15** Adiabatic demagnetization is an important technique that has been used with great success to achieve record low temperatures (near zero kelvin). A sample such as a chrome-alum salt (whose molecules may be considered as tiny magnets) is placed in a thermally insulating (adiabatic) enclosure at the lowest attainable temperature during cooling. A strong magnetic field is applied by an external magnet so that the tiny atomic magnets (spins) line up, forming a very ordered state, as Figure 8.14(a) shows. Then the magnet is removed so that its field is no longer present.

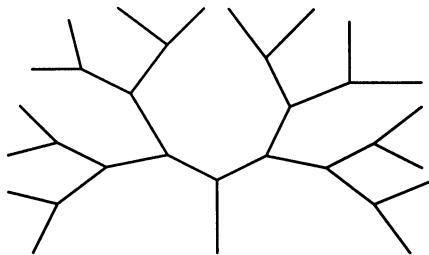
Redistributing the energy evenly among all the degrees of freedom lowers the temperature of the specimen, since the temperature is derived only from some of these degrees of freedom (the translational velocities). But the state of having the same amount of energy in each degree of freedom has the highest entropy. This clearly includes a significant increase of Kolmogorov complexity of the spins. ◇

## 8.6 Compression in Nature

Learning, in general, appears to involve compression of observed data or the results of experiments. It may be assumed that this holds both for the organisms that learn, as well as for the phenomena that are being learned. If the learner can't compress the data, he or she doesn't learn. If the data are not regular enough to be compressed, the phenomenon that produces the data is random, and its description does not add value to the data of generalizing or predictive nature.

### 8.6.1 Compression by Ants

In everyday life, we continuously compress information that is presented to us by the environment. Perhaps animals do this as well, as the following experiment reported by Zh.I. Reznikova and B.Ya. Ryabko [*Problems Inform. Transmission*, 22:3(1986), 245–249] suggests. The authors claim that the transmission of information by ants using tactile code is a well-established fact. This led the researchers to probe both the information



**FIGURE 8.15.** The maze: a binary tree constructed from matches

transmission rate and message-compressing capabilities of ants. We describe the latter issue.

Colonies of ants (*Formica sanguinea*) were put in artificial nests on laboratory tables. The ants were fed only once every three days, but only at some exit of a maze. This maze was a binary tree, as shown in Figure 8.15, with the root connecting to the nest and the food source at a leaf.

Each leaf contained a feeder. Only one feeder contained sugar syrup. All other feeders were empty. Since the tree was constructed by matches floating on water, the ants could not take short cuts. To reach the maze, an ant had to cross a bridge. This allowed hiding the maze in order that the remaining ants could not see the maze from their nest. To prevent marking of the trail by odorous substances, matches crossed by the ants were periodically replaced with fresh ones. During the experiments, the ants were marked with individual and group labels. Experiments with different numbers of branches were conducted in several sessions. In each case, the number of correct and wrong turns made by the ants was recorded, and also the total duration of tactile contacts between scouts and foragers was measured.

First, the scout ants would venture out to look for food. After a scout located the syrup in the maze, it returned to the nest to communicate the way to the syrup to the forager ants. Subsequently, the scouts were isolated and the foragers went in search of the syrup with, presumably, second-hand information. The simplest maze used was a two-leaf tree. In this setting, the scouts only had to communicate one bit. In the course of the experiment the depth of the tree was increased to 6. Since the scouts still managed to transmit the location of the syrup, the number of messages the ants could communicate is at least  $2^7 - 1 = 127$ .

Let “L” mean a left turn and “R” mean a right turn. It was found that the ants could communicate simple roads, like LLLLLL or LRLRLR, faster than more random roads. This seems to indicate that the ants compress the information before transmitting it.

No.	SEQUENCE OF TURNS TO SYRUP	MEAN TIME SEC.	SAMPLE STANDARD DEVIATION	NUMBER OF TESTS
1	LLL	72	8	18
2	RRR	75	5	15
3	LLLLL	84	6	9
4	RRRRR	78	8	10
5	LLLLL	90	9	8
6	RRRRRR	88	9	5
7	LRLRLR	130	11	4
8	RLRLRL	135	9	8
9	LLR	69	4	12
10	LRLL	100	11	10
11	RLLLRL	120	9	6
12	RRLRL	150	16	8
13	RLRRRL	180	20	6
14	RRLRRR	220	15	7
15	LRLLRL	200	18	5

**FIGURE 8.16.** Time required for *Formica sanguinea* scouts to transmit information about the direction to the syrup to the forager ants

The table in Figure 8.16 contains the results of the experiments. Apparently, it takes a longer time for the scouts to communicate “random” sequences to the foragers than to communicate “regular” sequences.

### 8.6.2 Compression by Science

Science may be regarded as the art of data compression. Compression of a great number of experimental data into the form of a short natural law yields an increase in predictive power, as shown in Chapter 5.

Whether science can exist would seem to depend on the question of whether the mass of experimentally obtained data form a compressible sequence. A randomly generated string (or universe) is with overwhelming probability not algorithmically compressible.

In Section 2.6 it was shown that maximally random sequences must contain logarithmically long very regular sequences. It may be the case that our part of the universe is an oasis of regularity in a maximally random universe. In practice, discovery of scientific laws with great predictive powers and many applications progresses spectacularly. This evidences, even if it doesn’t prove, inherent order in our universe. According to Francis Bacon (1561–1626), “The eye of the understanding is like the eye of the sense: for as you may see great objects through small crannies or levels, so you may see great axioms of nature through small and contemptible instances.”

## 8.7 History and References

---

Section 8.1 is partially based on T.M. Cover, P. Gács and R.M. Gray, *Ann. Probab.*, 17:3(1989), 840–865. Relations between physical entropy and expected Kolmogorov complexity are studied by W.H. Zurek in [*Phys. Rev. A*, 40:8(1989), 4731–4751; pp. 73–89 in: *Complexity, Entropy and the Physics of Information*, W.H. Zurek (Ed.), Addison-Wesley, 1991; *Nature*, 341(September 1989), 119–124]; by C.M. Caves in [pp. 91–115 in: *Complexity, Entropy and the Physics of Information*, W.H. Zurek (Ed.), Addison-Wesley, 1991; pp. 47–89 in: *Physical Origins of Time Asymmetry*, J.J. Halliwell, J. Pérez-Mercader and W.H. Zurek, Eds., Cambridge Univ. Press, 1994]; and by R. Schack ['Algorithmic information and simplicity in statistical physics', *Phys. Rev. E*, Submitted].

The relation between information and energy was derived from thermal noise considerations by J. H. Felker [*Proc. IRE* 40(1952), 728–729; *Ibid.* 42(1954), 1191]. There, a derivation based on statistical thermodynamic considerations is mentioned and attributed to J.R. Pierce. See also [J. R. Pierce and C. C. Cutler, pp. 55–109 in: F. I. Ordway, III, Ed., *Advances in Space Science*, Vol. 1, Academic Press, Inc., New York, 1959]. A further attribution to J. von Neumann appears in the edited volume [*Theory of Self-Reproducing Automata*, A.W. Burks, Ed., Univ. Illinois Press, Urbana, 1966]. (However, this relation is not mentioned in von Neumann's own writing.) The thermodynamic cost of erasing information was first studied by R. Landauer in [*IBM J. Res. Develop.*, 5(1961), 183–191]. Logical reversibility of computation was first studied by Y. Lecerf, *Comptes Rendus*, 257(1963), 2597–2600; C.H. Bennett, *IBM J. Res. Develop.*, 17(1973), 525–532. Section 8.2 is to a large extent based on [C.H. Bennett, *Scientific American*, (Nov. 1987), 108–117; C.H. Bennett and R. Landauer, *Scientific American*, (July 1985), 48–56; C.H. Bennett, *Int. J. Theoret. Phys.*, 21:12(1982), 905–940; see also *Phys. Rev. Lett.*, 53:12(1984), 1202–1206]. These papers discuss thermodynamics of computation in relation to reversible computing. The ballistic (molecular) computer in Section 8.2.3 and its reversibility were investigated in [T. Toffoli, *Proc. 7th Int. Colloq. Automata, Languages and Prog.*, Lect. Notes Comp. Sci., vol. 85, Springer-Verlag, 1980, pp. 632–644; *Math. Systems Theory*, 14(1981), 13–23; and E. Fredkin and T. Toffoli, *Int. J. Theoret. Phys.*, 21(1982), 219–253]. The issue *Int. J. Theoret. Phys.*, 21(1982) is dedicated to the physics of computing, containing papers presented in a conference on this topic at MIT in 1981. The second and third conferences on these and related issues took place in Dallas in 1992 and 1994 [*Proc. IEEE Workshop on Physics and Computation*, Dallas, 1992 and 1994 (also: Preliminary Proceedings 1992)]. C.H. Bennett [*Int. J. Theoret. Phys.*, 21:12(1982), 905–940] suggested a Brownian computer as a more stable alternative to the ballistic computer. Other stable alternatives in terms of quantum mechanical models

were devised by P.A. Benioff [*Int. J. Theoret. Physics*, 21(1982), 177–202; *Ann. New York Acad. Sci.*, 480(1986), 475–486]; and R.P. Feynman [*Optics News*, 11(1985), 11]. Related work about time/space tradeoffs in reversible computing is [C.H. Bennett, *SIAM J. Comput.*, 18(1989), 766–776; R.Y. Levine and A.T. Sherman, *SIAM J. Comput.*, 19(1990), 673–677; M. Li and P.M.B. Vitányi, *Proc. Royal Soc. London, Ser. A*, 452(1996), 769–789]. A comprehensive study in time/energy tradeoffs in reversible computation and space/energy tradeoffs in reversible simulation of irreversible computation is [M. Li and P.M.B. Vitányi, *Proc. Royal Soc. London, Ser. A*, 452(1996), 769–789]. Related material is contained in [C.H. Bennett, P. Gács, M. Li, P.M.B. Vitányi, and W.H. Zurek, *Proc. 25th ACM Symp. Theory Comput.*, 1993, pp. 21–30].

Other proposals to implement real physical reversible computing devices that dissipate no energy (or almost no energy) are bistable magnetic devices for reversible canceling/copying of records in [R. Landauer, *IBM J. Res. Develop.*, 5(1961), 183–191] and Brownian computers [R.W. Keyes and R. Landauer, *IBM J. Res. Develop.*, 14(1970), 152–157]; for Turing machines and Brownian enzymatic computers [C.H. Bennett, *IBM J. Res. Develop.*, 17(1973), 525–532; *BioSystems*, 11(1979), 85–90; C.H. Bennett and R. Landauer, *Scientific American*, (July 1985), 48–56]; with respect to reversible Boolean circuits in [E. Fredkin and T. Toffoli, *Int. J. Theoret. Phys.*, 21(1982), 219–253]; Brownian computing using Josephson devices in [K. Likharev, *Int. J. Theoret. Phys.*, 21(1982), 311–326].

All these models seem mutually simulatable. For more background information, see [C.H. Bennett, *Int. J. Theoret. Phys.*, 21(1982), 905–940]. Related material on reversible computing is [C.H. Bennett, *IBM J. Res. Develop.*, 32(1988), 16–23; R.W. Keyes and R. Landauer, *IBM J. Res. Develop.*, 14(1970), 152–157; R. Landauer, *Int. J. Theoret. Phys.*, 21(1982), 283; *Nature*, 335(October 1988), 779–784; R.W. Keyes, *Science*, 230(October 1985), 138–144; *Rev. Mod. Phys.*, 61(1989), 279–287; *IBM J. Res. Develop.*, 32(1988), 24–28; *Phys. Today*, 45(August 1992), 42–48]. Proposals for implementation of reversible computation in existing electronic circuit technology (CMOS, nMOS) are given in [R.C. Merkle, *Nanotechnology*, 4(1993), 21–40]; see also [*Proc. IEEE Workshop on Physics and Computation*, Dallas, 1992, 227–228, 237–247, 267–270; K.E. Drexler, *Nanosystems: Molecular Machinery, Manufacturing, and Computation*, Wiley, 1992].

Section 8.3 on “information distances” is based on [M. Li and P.M.B. Vitányi, *Proc. IEEE Workshop on Physics and Computation*, Dallas, 1992, 42–46 (Complete version in Preliminary Proceedings 1992); C.H. Bennett, P. Gács, M. Li, P.M.B. Vitányi, and W.H. Zurek, *Proc. 25th ACM Symp. Theory Comput.*, 1993, pp. 21–30]. The current proof of Theorem 8.3.1 is due to J. Tromp. Related work on “information measures” is [M. Gell-Mann and S. Lloyd, *Complexity*, 2:1(1996), 44–52].

For classical thermodynamics and entropy see [E. Fermi, *Thermodynamics*, Dover, 1956]. We partly used the survey [J. Schumacher, *CWI Quarterly*, 6:2(1993), 97–120]. The study of statistical thermodynamics using Kolmogorov complexity was started by C.H. Bennett, [*Int. J. Theoret. Phys.*, 21:12(1982), 905–940; *Scientific American*, 255:11(1987), 108–117] and W.H. Zurek. Initially, the algorithmic entropy of Section 8.5.1 was proposed by W.H. Zurek in [*Phys. Rev. A*, 40:8(1989), 4731–4751; *Complexity, entropy and the physics of information*, W.H. Zurek (Ed.), Addison-Wesley, 1991, pp. 73–89; *Nature*, 341(September 1989), 119–124] under the name “physical entropy.”

The initial discussion of Maxwell’s demon accounting for the thermodynamic cost of irreversible information erasure as in Example 8.5.1 is due to C.H. Bennett [*Int. J. Theoret. Phys.*, 21:12 (1982), 905–940; *Scientific American*, 255:11(1987), 108–117]. These papers contain an excellent exposition on Maxwell’s demon, including a construction of a device for measuring the position of a molecule. This type of solution was further explored by W.H. Zurek using “physical entropy” arguments in [*Nature*, 341(September 1989), 119–124; *Phys. Rev. A*, 40:8(1989), 4731–4751; *Complexity, entropy and the physics of information*, W.H. Zurek (Ed.), Addison-Wesley, 1991, pp. 73–89]. The discussion on Maxwell’s demon in Example 8.5.1, and Claim 8.5.1 follow by and large Zurek’s discussion. The Szilard engine was described by L. Szilard in a paper entitled “On the decrease of entropy in a thermodynamic system by the intervention of intelligent beings” [*Z. Phys.*, 53(1929), 840–856]. Szilard in this paper in fact also discovered the relationship between entropy and information. However, this was not generally accepted until this relation was rediscovered by C.E. Shannon in the 1940s. A collection of these and other key papers on this topic is [H.S. Leff and A.F. Rex, Eds., *Maxwell’s Demon: Entropy, Information, Computing*, Princeton University Press, 1990]. See also [R.G. Brewer and E.L. Hahn, *Scientific American*, December 1984, 42–49].

An initial approach to an “algorithmic thermodynamics” was partially published in [M. Li and P.M.B. Vitányi, *Proc. 19th Int. Colloq. Automata, Languages and Prog., Lect. Notes Comp. Sci.*, vol. 623, Springer-Verlag, Heidelberg, 1992, 1–16]. P. Gács removed scale-dependence of our approach by adding coarse-graining and randomness tests, obtaining the “physical entropy” of W.H. Zurek, *Ibid.* from mathematical first principles. The definition of algorithmic entropy in Section 8.5.2 and most results and examples are based on the treatment using continuous time by P. Gács, *Proc. 2nd IEEE Workshop on Physics and Computation*, 1994, 209–216. We have tried to simplify the discussion by discretizing time. We have considered a state space  $\{0, 1\}^\infty$ . An analogous treatment can be given with the state space  $\mathcal{R}$ , the real numbers. Since this space has a different metric, we need to define appropriate notions of conti-

nuity, recursiveness, and enumerability of functions over it. Let us look at some inherent distinctions between the  $\mathcal{R}$  and  $\{0, 1\}^\infty$  spaces. Our space is  $\{0, 1\}^\infty$  with the  $\mu$ -metric. In this space, the function  $f(\omega) = 0$  for  $\omega \in \Gamma_0$  and  $f(\omega) = 1$  for  $\omega \in \Gamma_1$  is a continuous recursive function. This is different from the situation on the set of real numbers  $\mathcal{R}$ . If we take the space to be  $\mathcal{R}$  with the  $\mu$ -metric defined by the measure of all nonzero intervals (instead of just the cylinders) then the analogous function  $f(r) = 0$  for  $r < \frac{1}{2}$ , and  $f(r) = 1$  for  $r \geq \frac{1}{2}$  is not continuous and not recursive. Namely, however close we approximate  $r = \frac{1}{2}$ , we may never know whether  $f(r) = 0$  or  $f(r) = 1$ .

Separating the information in an object into a part (or model) accounting for the *useful* information, the regularities, and a part describing the remaining *random* information was first proposed by Kolmogorov as the “Kolmogorov minimal sufficient statistic” treated in Section 2.2.2. Related ideas are “sophistication” on page 114, and in general the idea of two-part codes (Section 2.1.1) and the “minimum description length (MDL)” principle of Section 5.5. The latter approach is a statistical inference method to obtain the right hypothesis, or model, for a given data sample. Here “right” means capturing the regular, or useful, aspects of the data. Similar ideas were recently proposed in [M. Gell-Mann, *The Quark and the Jaguar*, W.H. Freeman, New York, 1994; M. Gell-Mann, *Complexity*, 1:1(1995), 16–19; M. Gell-Mann and S. Lloyd, *Complexity*, 2:1(1996), 44–52] and applied to quantum-mechanics theory in [M. Gell-Mann and J.B. Hartle, *Proc. 4th Drexel Symp. Quantum Non-Intergrability—The Quantum-Classical Approach*, D.H. Feng, Ed., To appear] and to the theory of adaptation and control in [S. Lloyd and J.J. Slotine, *Int. J. Adapt. Contr. Sign. Proc.*, To appear]. The sum of the useful information and the random information is called “total information.” It is another version of W.H. Zurek’s [*Ibid.*] “physical entropy,” and is closely related to both the MDL principle (Section 5.5), and algorithmic entropy of Section 8.5.2.

[J. Ford, *Phys. Today*, (April 1983), 40–47; ‘Chaos: Solving the unsolvable, predicting the unpredictable,’ in *Chaotic Dynamics and Fractals*, M.F. Barnsley and S.G. Demko (Ed.), Academic Press, 1986] may be the earliest papers applying Kolmogorov complexity to chaos. Use of Kolmogorov complexity and information theory to analyze quantum chaos related to issues in Section 8.5.2 like the “baker’s map” are [R. Schack, G.M. D’Ariano, and C.M. Caves, *Phys. Rev. E*, 50(1994); 972]; and R. Schack and C. Caves [*Phys. Rev. E*, 53:4(1996), 3257–3270; *Phys. Rev. E*, 53:4(1996), 3387–3401]. See [D. Ruelle, *Chance and Chaos*, Penguin, 1993] for an informal and witty discussion of some topics in this chapter. The experiment with ants is reported by Zh.I. Reznikova and B.Ya. Ryabko [*Problems Inform. Transmission*, 22(1986), 245–249]. Francis Bacon is quoted from *Sylva Sylvarum*, 337, 1627.

---

# References

- [1] A.M. Abramov. Kolmogorov's pedagogic legacy. *Russian Math. Surveys*, 43(6):45–88, 1988.
- [2] Y.S. Abu-Mostafa. The complexity of information extraction. *IEEE IEEE Trans. Inform. Theory*, IT-32(4):513–525, 1986.
- [3] L. Adleman. Time, space, and randomness. Technical Report TM-131, MIT, Lab. Comput. Sci., March 1979.
- [4] V.N. Agafonov. Normal sequences and finite automata. *Soviet Math. Dokl.*, 9:324–325, 1968.
- [5] V.N. Agafonov. *On algorithms, frequency and randomness*. PhD thesis, University of Novosibirsk, Novosibirsk, 1970.
- [6] P.S. Aleksandrov. A few words on A.N. Kolmogorov. *Russian Math. Surveys*, 38(4):5–7, 1983.
- [7] V.M. Alekseev and M.V. Yakobson. Symbolic dynamics and hyperbolic dynamical systems. *Physics Reports*, 75:287–325, 1981.
- [8] E. Allender. *Invertible functions*. PhD thesis, Georgia Instit. Tech., 1985.
- [9] E. Allender. Some consequences of the existence of pseudorandom generators. *J. Comput. System Sci.*, 39:101–124, 1989.
- [10] E. Allender. Applications of time-bounded Kolmogorov complexity in complexity theory. In O. Watanabe, editor, *Kolmogorov Complexity and Computational Complexity*, pages 4–22. Springer-Verlag, 1992.
- [11] E. Allender and R.S. Rubinstein. P-printable sets. *SIAM J. Comput.*, 17:1193–1202, 1988.
- [12] E. Allender and O. Watanabe. Kolmogorov complexity and degrees of tally sets. *Inform. Comput.*, 86:160–178, 1990.
- [13] A. Ambainis. Application of Kolmogorov complexity to inductive inference with limited memory. In *Proc. 6th Int'l Workshop Algorithmic Learning Theory*, volume 997 of *Lect. Notes Artificial Intelligence*, pages 313–318. Springer-Verlag, 1995.
- [14] D. Angluin. Algorithmic theory of information and randomness. Lecture notes postgraduate course, Edinburgh University, 1977/1978.

- [15] V.I. Arnol'd. A few words on Andrei Nikolaevich Kolmogorov. *Russian Math. Surveys*, 43(6):43–44, 1988.
- [16] E.A. Asarin. Individual random continuous functions. In Yu.V. Prokhorov, editor, *Summaries of Reports of the First All-World Congress of the Bernouilli Society of Mathematical Statistics and Probability Theory*, volume 1, page 450. Nauka, 1986. In Russian.
- [17] E.A. Asarin. Some properties of Kolmogorov  $\delta$ -random finite sequences. *SIAM Theory Probab. Appl.*, 32:507–508, 1987.
- [18] E.A. Asarin. On some properties of finite objects random in the algorithmic sense. *Soviet Math. Dokl.*, 36:109–112, 1988.
- [19] E.A. Asarin and A.V. Pokrovskii. Application of Kolmogorov complexity to analyzing the dynamics of controlled systems. *Automat. and Telemekh.*, 1:25–33, 1986. In Russian.
- [20] J. Balcazar, H. Buhrman, and M. Hermo. Learnability of Kolmogorov-easy circuit expressions via queries. In P.M.B. Vitányi, editor, *Computational Learning Theory; Proc. 2nd European Conf.*, volume 904 of *Lect. Notes Artificial Intelligence*, pages 112–124. Springer-Verlag, 1995.
- [21] J.L. Balcázar and R.V. Book. Sets with small generalized Kolmogorov complexity. *Acta Informatica*, 23:679–688, 1986.
- [22] J.L. Balcázar, J. Diaz, and J. Gabarró. *Structural Complexity*. Springer-Verlag, 1988.
- [23] J.L. Balcázar and U. Schöning. Logarithmic advice classes. *Theoret. Comput. Sci.*, 99:279–290, 1992.
- [24] A.R. Barron and T.M. Cover. Minimum complexity density estimation. *IEEE Trans. Inform. Theory*, IT-37:1034–1054, 1991.
- [25] J.M. Barzdins. Complexity of programs to determine whether natural numbers not greater than  $n$  belong to a recursively enumerable set. *Soviet Math. Dokl.*, 9:1251–1254, 1968.
- [26] J.M. Barzdins. On computability by probabilistic machines. *Soviet Math. Dokl.*, 10:1464–1467, 1969.
- [27] J.M. Barzdins. On the relative frequency of solution of algorithmically unsolvable mass problems. *Soviet Math. Dokl.*, 11:459–462, 1970.
- [28] J.M. Barzdins. Algorithmic information theory. In *Encyclopaedia of Mathematics, volume 1*, pages 140–142. D. Reidel (Kluwer Academic Publishers), 1988. Updated and annotated translation of the *Soviet Mathematical Encyclopaedia*.
- [29] J.M. Barzdins and R.V. Freivalds. On the prediction of general recursive functions. *Soviet Math. Dokl.*, 13:1251–1254 (1224–1228), 1972.
- [30] Th. Bayes. An essay towards solving a problem in the doctrine of chances. *Philos. Trans. Royal Soc.*, 53:376–398, 1763. Ibid., 54:298–310, 1764, R. Price, Ed.
- [31] T. C. Bell, J. G. Cleary, and I. H. Witten. *Text Compression*. Prentice Hall, 1990.
- [32] A.M. Ben-Amram. Average-case complexity of sorting on a RAM. Report 95/23, DIKU, the University of Copenhagen, 1995.
- [33] A.M. Ben-Amram and Z. Galil. On pointers versus addresses. *J. Assoc. Comput. Mach.*, 39(3):617–648, 1992.

- [34] P.A. Benioff. Models of Zermelo Frankel set theory as carriers for the mathematics or physics. i, ii. *J. Math. Phys.*, 17(5):618–628,629–640, 1976.
- [35] C.H. Bennett. The thermodynamics of computation—a review. *Int'l J. Theoret. Physics*, 21:905–940, 1982.
- [36] C.H. Bennett. Demons, engines and the second law. *Scientific American*, pages 108–116, Nov. 1987.
- [37] C.H. Bennett. Dissipation, information, computational complexity and the definition of organization. In D. Pines, editor, *Emerging Syntheses in Science*, volume 1 of *Santa Fe Institute Studies in the Science of Complexity*, pages 297–313. Addison-Wesley, 1987.
- [38] C.H. Bennett. Logical depth and physical complexity. In R. Herken, editor, *The Universal Turing Machine; A Half-Century Survey*, pages 227–258. Oxford University Press, 1988. In Germany: Kammerer & Unverzagt.
- [39] C.H. Bennett. How to define complexity in physics, and why. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 137–148. Addison-Wesley, 1991.
- [40] C.H. Bennett, P. Gács, M. Li, P.M.B. Vitányi, and W. Zurek. Thermodynamics of computation and information distance. In *Proc. 25th ACM Symp. Theory Comput.*, pages 21–30, 1993.
- [41] C.H. Bennett and M. Gardner. The random number omega bids fair to hold the mysteries of the universe. *Scientific American*, 241:20–34, May 1979.
- [42] C.H. Bennett and R. Landauer. The fundamental physical limits of computation. *Scientific American*, pages 48–56, July 1985.
- [43] M. Blum. On the size of machines. *Inform. Contr.*, 11:257–265, 1967.
- [44] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *J. Assoc. Comput. Mach.*, 36(4):929–965, 1989.
- [45] D.E. Boekee, R.A. Kraak, and E. Backer. On complexity and syntactic information. *IEEE Trans. Systems Man Cybernet.*, SMC-12:71–79, 1982.
- [46] N.N. Bogolyubov, B.V. Gnedenko, and S.L. Sobolev. Andrei Nikolae-vich Kolmogorov (on his eightieth birthday). *Russian Math. Surveys*, 38(4):9–27, 1983.
- [47] R.V. Book. On sets with small information content. In O. Watanabe, editor, *Kolmogorov Complexity and Computational Complexity*, pages 23–42. Springer-Verlag, 1992.
- [48] R.V. Book. On languages reducible to algorithmically random languages. *SIAM J. Comput.*, 23:1275–1282, 1994.
- [49] R.V. Book. Relativizations of the  $p=?np$  and other problems: developments in structural complexity theory. *SIAM Review*, 36:157–175, 1994.
- [50] R.V. Book and J.H. Lutz. On languages with very high space-bounded Kolmogorov complexity. *SIAM J. Comput.*, 22(2):395–402, 1993.
- [51] R.V. Book, J.H. Lutz, and K. Wagner. An observation on probability versus randomness with applications to complexity classes. *Math. Systems Theory*, 27:201–209, 1994.

- [52] R.V. Book, P. Orponen, D. Russo, and O. Watanabe. Lowness properties of sets in the exponential-time hierarchy. *SIAM J. Comput.*, 17:504–516, 1988.
- [53] R.V. Book and O. Watanabe. On random hard sets for NP. *Inform. Comput.*, 125:70–76, 1996.
- [54] A.A. Brudno. On the complexity of paths of dynamic systems. *Uspekhi Mat. Nauk*, 33:207–208, 1978.
- [55] A.A. Brudno. Entropy and the complexity of trajectories of a dynamical system. *Trans. Mosc. Math. Soc.*, 44:127–151, 1983.
- [56] H. Buhrman, E. Hemaspaandra, and L. Longpré. SPARSE reduces conjunctively to TALLY. *SIAM J. Comput.*, 1996/1997. To appear. Preliminary version in *Proc. 8th IEEE Conf. Structure in Complexity Theory*, 1993, pp. 208–214.
- [57] H. Buhrman, J.H. Hoepman, and P.M.B. Vitányi. Optimal routing tables. In *Proc. 15th ACM Symp. Principles Distribut. Comput.*, pages 134–142. ACM Press, 1996.
- [58] H. Buhrman, M. Li, and P.M.B. Vitányi. Kolmogorov random graphs. Manuscript, CWI, Amsterdam, November 1995.
- [59] H. Buhrman and L. Longpré. Compressibility and resource bounded measure. In *Proc. 13th Symp. Theoret. Aspects Comput. Sci.*, volume 1046 of *Lect. Notes Comput. Sci.*, pages 13–24. Springer-Verlag, 1996.
- [60] H. Buhrman and E. Mayordomo. An excursion to the kolmogorov random strings. In *Proc. 10th IEEE Conf. Structure in Complexity Theory*, pages 197–203, 1995. To appear in *J. Comput. System Sci.*
- [61] H. Buhrman and P. Orponen. Random strings make hard instances. *J. Comput. System Sci.*, 53(2):261–266, 1996.
- [62] H.M. Buhrman and L. Fortnow. Resource-bounded Kolmogorov complexity revisited. In *Proc. 14th Symp. Theoret. Aspects Comput. Sci.*, Lect. Notes Comput. Sci. Springer-Verlag, 1997.
- [63] M.S. Burgin. Generalized Kolmogorov complexity and duality in computational theory. *Soviet Math. Dokl.*, 25(3):559–564, 1982.
- [64] J.-Y. Cai and J. Hartmanis. On Hausdorff and topological dimensions of the Kolmogorov complexity of the real line. *J. Comput. System Sci.*, 49(3):605–619, 1994.
- [65] J.-Y. Cai and L. Hemachandra. A note on enumerative counting. *Inform. Process. Lett.*, 38:215–219, 1991.
- [66] C. Calude. *Theories of Computational Complexity*, chapter 4. North-Holland, 1988.
- [67] C. Calude. *Information and Randomness: An Algorithmic Perspective*. Springer-Verlag, 1994.
- [68] C. Calude, I. Chitescu, and L. Staiger. P. Martin-Löf tests: representability and embedability. *Revue Roumaine Math. Pures Appl.*, 30:719–732, 1985.
- [69] C. Calude and H. Jürgensen. Randomness as an invariant for number representations. In H. Maurer, J. Karhumaki, and G. Rozenberg, editors, *Results and Trends in Theoretical Computer Science*, pages 44–66. Springer-Verlag, 1994.
- [70] R. Carnap. *Logical Foundations of Probability*. Univ. Chicago Press, 1950.

- [71] J. Castro and J.L. Balcázar. Simple pac learning of simple decision lists. In *Proceedings 6th Int'l Workshop on Algorithmic Learning Theory*, volume 997 of *Lect. Notes Artificial Intelligence*, pages 239–248. Springer-Verlag, 1995.
- [72] C.M. Caves. Entropy and information: How much information is needed to assign a probability? In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 91–115. Addison-Wesley, 1991.
- [73] C.M. Caves. Information, entropy and chaos. In J.J. Halliwell, J. Pérez-Mercader, and W.H. Zurek, editors, *Physical Origins of Time Asymmetry*, pages 47–89. Cambridge Univ. Press, 1994.
- [74] G.J. Chaitin. On the length of programs for computing finite binary sequences. *J. Assoc. Comput. Mach.*, 13:547–569, 1966.
- [75] G.J. Chaitin. On the length of programs for computing finite binary sequences: statistical considerations. *J. Assoc. Comput. Mach.*, 16:145–159, 1969.
- [76] G.J. Chaitin. On the simplicity and speed of programs for computing infinite sets of natural numbers. *J. Assoc. Comput. Mach.*, 16:407–422, 1969.
- [77] G.J. Chaitin. On the difficulty of computations. *IEEE Trans. Inform. Theory*, IT-16:5–9, 1970.
- [78] G.J. Chaitin. Towards a mathematical definition of ‘life’. *SIGACT News*, 4:12–18, January 1970.
- [79] G.J. Chaitin. Computational complexity and Gödel’s incompleteness theorem. *SIGACT News*, 9:11–12, 1971.
- [80] G.J. Chaitin. Information-theoretic computational complexity. *IEEE Trans. Inform. Theory*, IT-20:10–15, 1974. Reprinted in T. Tymoczko, *New Directions in the Philosophy of Mathematics*, Birkhauser, 1986.
- [81] G.J. Chaitin. Information-theoretic limitations of formal systems. *J. Assoc. Comput. Mach.*, 21:403–424, 1974.
- [82] G.J. Chaitin. Randomness and mathematical proof. *Scientific American*, 232:47–52, May 1975.
- [83] G.J. Chaitin. A theory of program size formally identical to information theory. *J. Assoc. Comput. Mach.*, 22:329–340, 1975.
- [84] G.J. Chaitin. Algorithmic entropy of sets. *Comput. Math. Appl.*, 2:233–245, 1976.
- [85] G.J. Chaitin. Information-theoretic characterizations of recursive infinite strings. *Theoret. Comput. Sci.*, 2:45–48, 1976.
- [86] G.J. Chaitin. Algorithmic information theory. *IBM J. Res. Develop.*, 21:350–359, 1977.
- [87] G.J. Chaitin. Program size, oracles, and the jump operation. *Osaka J. Math.*, 14:139–149, 1977.
- [88] G.J. Chaitin. Toward a mathematical definition of “life”. In R.D. Levine and M. Tribus, editors, *The Maximal Entropy Formalism*, pages 477–498. MIT Press, 1979.
- [89] G.J. Chaitin. Algorithmic information theory. In *Encyclopedia of Statistical Sciences, volume 1*, pages 38–41. Wiley, 1982.
- [90] G.J. Chaitin. Gödel’s theorem and information. *Int'l J. Theoret. Physics*, 22:941–954, 1982. Reprinted in T. Tymoczko, *New Directions in the Philosophy of Mathematics*, Birkhauser, Boston, 1986.

- [91] G.J. Chaitin. Randomness and Gödel's theorem. *Mondes en Developpement*, 14, No. 54-55:125–128, 356., 1986.
- [92] G.J. Chaitin. *Algorithmic Information Theory*. Cambridge University Press, 1987.
- [93] G.J. Chaitin. Beyond Gödel's proof. *IBM Res. Mag.*, 25:12–15, Fall 1987.
- [94] G.J. Chaitin. Computing the busy beaver function. In T.M. Cover and B. Gopinath, editors, *Open Problems in Communication and Computation*, pages 108–112. Springer-Verlag, 1987.
- [95] G.J. Chaitin. Incompleteness theorems for random reals. *Advances in Applied Math.*, 8:119–146, 1987.
- [96] G.J. Chaitin. *Information, Randomness and Incompleteness—Papers on Algorithmic Information Theory*. World Scientific, 1987.
- [97] G.J. Chaitin. An algebraic equation for the halting probability. In R. Herken, editor, *The Universal Turing Machine; A Half-Century Survey*, pages 279–284. Oxford University Press, 1988. In Germany: Kammerer & Unverzagt.
- [98] G.J. Chaitin. *Information-Theoretic Incompleteness*. World Scientific, 1992.
- [99] G.J. Chaitin. A new version of algorithmic information theory. *Complexity*, 1(4):55–59, 1995/1996.
- [100] G.J. Chaitin and J.T. Schwartz. A note on Monte-Carlo primality tests and algorithmic information theory. *Comm. Pure Applied Math.*, 31:521–527, 1978.
- [101] D.G. Champernowne. The construction of decimals normal in the scale of ten. *J. London Math. Soc.*, 8:254–260, 1933.
- [102] N. Chater. Reconciling simplicity and likelihood principles in perceptual organization. *Psychological Review*, 103:566–581, 1996.
- [103] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sums of observations. *Ann. Math. Stat.*, 23:493–509, 1952.
- [104] M. Chrobak and M. Li.  $k + 1$  heads are better than  $k$  for PDAs. *J. Comput. System Sci.*, 37:144–155, 1988.
- [105] F.R.K. Chung, R.E. Tarjan, W.J. Paul, and R. Reischuk. Coding strings by pairs of strings. *SIAM J. Algebra Disc. Math.*, 6(3):445–461, 1985.
- [106] A. Church. On the concept of a random sequence. *Bull. Amer. Math. Soc.*, 46:130–135, 1940.
- [107] J.T. Coffey and R.M. Goodman. Any code of which we cannot think is good. *IEEE Trans. Inform. Theory*, IT-36:1453–1461, 1990.
- [108] J.D. Collier. Two faces of Maxwell's demon reveal the nature of irreversibility. *Stud. Hist. Phil. Sci.*, 21(2):257–268, 1990.
- [109] T.M. Cover. Enumerative source encoding. *IEEE Trans. Inform. Theory*, IT-19:73–77, 1973.
- [110] T.M. Cover. Generalization on patterns using Kolmogorov complexity. In *Proc. 1st Int'l Conf. Pattern Recognition*, pages 551–553, 1973.
- [111] T.M. Cover. On the determination of the irrationality of the mean of a random variable. *Ann. Statist.*, 1:862–871, 1973.
- [112] T.M. Cover. Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin. Technical Report 12, Statistics Department, Stanford University, October 1974.

- [113] T.M. Cover. Kolmogorov complexity, data compression, and inference. In J.K. Skwirzynski, editor, *The Impact of Processing Techniques on Communications*, pages 23–33. Martinus Nijhoff Publishers, 1985.
- [114] T.M. Cover, P. Gács, and R.M. Gray. Kolmogorov’s contributions to information theory and algorithmic complexity. *Ann. Probab.*, 17:840–865, 1989.
- [115] T.M. Cover and R.C. King. A convergent gambling estimate of the entropy of English. *IEEE Trans. Inform. Theory*, IT-24:413–421, 1978.
- [116] T.M. Cover and J.A. Thomas. *Elements of Information Theory*. Wiley & Sons, 1991.
- [117] J.P. Crutchfield and K. Young. Computation at the onset of chaos. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 223–269. Addison-Wesley, 1991.
- [118] R.R. Cuykendall. *Kolmogorov information and VLSI lower bounds*. PhD thesis, University of California, Los Angeles, Dec. 1984.
- [119] R.P. Daley. Complexity and randomness. In R. Rustin, editor, *Computational Complexity; Courant Comput. Sci. Symp. 7*, pages 113–122. Algorithmics Press, 1971.
- [120] R.P. Daley. An example of information and computation resource trade-off. *J. Assoc. Comput. Mach.*, 20(4):687–695, 1973.
- [121] R.P. Daley. Minimal-program complexity of sequences with restricted resources. *Inform. Contr.*, 23:301–312, 1973.
- [122] R.P. Daley. The extent and density of sequences within the minimum-program complexity hierarchies. *J. Comput. System Sci.*, 9:151–163, 1974.
- [123] R.P. Daley. Minimal-program complexity of pseudo-recursive and pseudo-random sequences. *Math. Systems Theory*, 9:83–94, 1975.
- [124] R.P. Daley. Noncomplex sequences: characterizations and examples. *J. Symbolic Logic*, 41:626–638, 1976.
- [125] R.P. Daley. On the inference of optimal descriptions. *Theoret. Comput. Sci.*, 4:301–319, 1977.
- [126] R.P. Daley. Quantitative and qualitative information in computation. *Inform. Contr.*, 45:236–244, 1980.
- [127] R.P. Daley. The process complexity and the understanding of sequences. In *Proc. Symp. Summer School MFCS*, High Tatras, September 1973.
- [128] P.C.W. Davies. Why is the physical world so comprehensible? In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 61–70. Addison-Wesley, 1991.
- [129] L.D. Davisson. Universal noiseless encoding. *IEEE Trans. Inform. Theory*, IT-19:783–795, 1973.
- [130] B. de Finetti. *Probability, Induction, and Statistics*. Wiley, 1972.
- [131] A. DeSantis, G. Markowsky, and M.N. Wegman. Learning probabilistic prediction functions. In *Proc. 29th IEEE Symp. Found. Comput. Sci.*, pages 110–119, 1988.
- [132] J.-E. Dies. Information et complexité. *Ann. Inst. Henri Poincaré, B*, 12:365–390, 1976. *Ibid.*, 14:113–118, 1978.
- [133] M. Dietzfelbinger. *Lower bounds on computation time for various models in computational complexity theory*. PhD thesis, Dept. Comput. Sci., Univ. Illinois at Chicago, 1987.

- [134] M. Dietzfelbinger. The speed of copying on one-tape off-line Turing machines. *Inform. Process. Lett.*, 33:83–90, 1989/1990.
- [135] M. Dietzfelbinger and W. Maass. The complexity of matrix transposition on one-tape off-line Turing machines with output tape. *Theoret. Comput. Sci.*, 108(2):271–290, 1993.
- [136] M. Dietzfelbinger, W. Maass, and G. Schnitger. The complexity of matrix transposition on one-tape off-line Turing machines. *Theoret. Comput. Sci.*, 82(1):113–129, 1991.
- [137] J.L. Doob. Kolmogorov’s early work on convergence theory and foundations. *Ann. Probab.*, 17:815–821, 1989.
- [138] P. Duris, Z. Galil, W.J. Paul, and R. Reischuk. Two nonlinear lower bounds for on-line computations. *Inform. Contr.*, 60:1–11, 1984.
- [139] E.B. Dynkin. Kolmogorov and the theory of Markov processes. *Ann. Probab.*, 17:822–832, 1989.
- [140] P. Erdős and J. Spencer. *Probabilistic Methods in Combinatorics*. Academic Press, 1974.
- [141] M. Feder. Maximal entropy as special case of the minimum length description criterion. *IEEE Trans. Inform. Theory*, IT-32:847–849, 1986.
- [142] J. H. Felker. A link between information and energy. *Proc. IRE*, 40:728–729, 1952. Discussion, *Ibid.*, 52(1954), 1191.
- [143] W. Feller. *An Introduction to Probability Theory and Its Applications*, volume 1. Wiley, 1968. Third edition.
- [144] T.L. Fine. On the apparent convergence of relative frequencies and its implications. *IEEE Trans. Inform. Theory*, IT-16:251–257, 1970.
- [145] T.L. Fine. *Theories of Probability*. Academic Press, 1973.
- [146] T.L. Fine. Uniformly reasonable source encoding is often practically impossible. *IEEE Trans. Inform. Theory*, IT-21:368–373, 1975.
- [147] R.A. Fisher. On the mathematical foundations of theoretical statistics. *Philos. Trans. Royal Soc. London, Ser. A*, 222:309–368, 1925.
- [148] J. Ford. How random is a random coin toss. *Physics Today*, 36:40–47, April 1983.
- [149] J. Ford. Chaos: solving the unsolvable, predicting the unpredictable. In M.F. Barnsley and S.G. Demko, editors, *Chaotic Dynamics and Fractals*. Academic Press, 1986.
- [150] L. Fortnow and M. Kummer. Resource-bounded instance complexity. *Theoret. Comput. Sci. A*, 161:123–140, 1996.
- [151] L. Fortnow and S. Laplante. Circuit complexity à la Kolmogorov. *Inform. Comput.*, 123:121–126, 1995.
- [152] W.L. Fouché. Identifying randomness given by high descriptive complexity. *Acta Applicandae Mathematicae*, 34:313–328, 1994.
- [153] W.L. Fouché. Descriptive complexity and reflective properties of combinatorial configurations. *J. London Math. Soc.*, 54(2):199–208, 1996.
- [154] R.V. Freivalds. On the running time of deterministic and nondeterministic Turing machines. *Latv. Mat. Ezhegodnik*, 23:158–165, 1979. In Russian.
- [155] B. Fu. With quasi-linear queries, EXP is not polynomial-time Turing reducible to sparse sets. *SIAM J. Comput.*, 24(5):1082–1090, 1995.

- [156] P.H. Fuchs. Statistical characterization of learnable sequences. In *Automata Theory and Formal Languages, Proc. 2nd GI Conference*, volume 33 of *Lect. Notes Comput. Sci.*, pages 52–57. Springer-Verlag, 1975.
- [157] P. Gács. On the symmetry of algorithmic information. *Soviet Math. Dokl.*, 15:1477–1480, 1974. Correction, *Ibid.*, 15:1480, 1974.
- [158] P. Gács. *Komplexität und Zufälligkeit*. PhD thesis, Fachbereich Mathematik, J.W. Goethe Universität, Frankfurt am Main, 1978.
- [159] P. Gács. Exact expressions for some randomness tests. *Z. Math. Logik Grundl. Math.*, 26:385–394, 1980.
- [160] P. Gács. On the relation between descriptional complexity and algorithmic probability. *Theoret. Comput. Sci.*, 22:71–93, 1983.
- [161] P. Gács. Every sequence is reducible to a random sequence. *Inform. Contr.*, 70:186–192, 1986. See also: A. Kučera, Measure,  $\Pi_1^0$ -classes and complete extensions of PA, pp. 245–259 in: H.-D. Ebbinghaus, G. H. Müller en G. E. Sacks, eds., *Recursion Theory Week*, Lect. Notes in Math., Vol. 1141, Springer Verlag, Heidelberg, 1985.
- [162] P. Gács. Randomness and probability - complexity of description. In Kotz-Johnson, editor, *Encyclopedia of Statistical Sciences*, volume 7, pages 551–555. Wiley, 1986.
- [163] P. Gács. Lecture notes on descriptional complexity and randomness. Technical report, Comput. Sci. Dept., Boston Univ., 1988. Prelim. version 1987.
- [164] P. Gács. Gregory J. Chaitin, Algorithmic Information Theory. *J. Symbolic Logic*, 54:624–627, 1989. Book review.
- [165] P. Gács. Self-correcting two-dimensional arrays. In Silvio Micali, editor, *Randomness in Computation*, volume 5 of *Advances in Computing Research (a scientific annual)*, pages 223–326. JAI Press, 1989.
- [166] P. Gács. The Boltzmann entropy and randomness tests. In *Proc. 2nd IEEE Workshop on Physics and Computation*, pages 209–216, 1994.
- [167] P. Gács and J. Körner. Common information is far less than mutual information. *Problems of Control and Inform. Theory*, 2:149–162, 1973.
- [168] H. Gaifman and M. Snir. Probabilities over rich languages, randomness and testing. *J. Symbolic Logic*, 47:495–548, 1982.
- [169] Z. Galil, R. Kannan, and E. Szemerédi. On 3-pushdown graphs with large separators. *Combinatorica*, 9:9–19, 1989.
- [170] Z. Galil, R. Kannan, and E. Szemerédi. On nontrivial separators for k-page graphs and simulations by nondeterministic one-tape Turing machines. *J. Comput. System Sci.*, 38:134–149, 1989.
- [171] Z. Galil and J. Seiferas. Time-space-optimal string matching. *J. Comput. System Sci.*, 26:3:280–294, 1983.
- [172] R.G. Gallager. *Information Theory and Reliable Communication*. Wiley, 1968.
- [173] Q. Gao and M. Li. An application of minimum description length principle to online recognition of handprinted alphanumerals. In *11th Int'l Joint Conf. Artificial Intelligence*, pages 843–848. Morgan Kaufmann, 1989.
- [174] M.R. Garey and D.S. Johnson. *Computers and Intractability, a guide to the theory of NP-completeness*. W.H. Freeman and Co., 1979.

- [175] P. Gaspard and X.-J. Wang. Sporadicity: Between periodic and chaotic dynamical behaviors. *Proc. Nat'l Acad. Sci. USA*, 85:4591–4595, 1988.
- [176] R. Gavaldá. *Kolmogorov Randomness and its Applications to Structural Complexity Theory*. PhD thesis, Universitat Politècnica de Catalunya, Barcelona, 1992.
- [177] R. Gavaldá, L. Torenvliet, O. Watanabe, and J.L. Balcázar. Generalized Kolmogorov complexity in relativized separations. In *Proc. 15th Conf. Math. Found. Comput. Sci.*, volume 452 of *Lect. Notes Comput. Sci.*, pages 269–276. Springer-Verlag, 1991.
- [178] R. Gavaldá and O. Watanabe. On the computational complexity of small descriptions. *SIAM J. Comput.*, 22(6):1257–1275, 1993.
- [179] M. Gell-Mann. *The Quark and the Jaguar*. W.H. Freeman, 1994.
- [180] M. Gell-Mann. Remarks on simplicity and complexity. *Complexity*, 1(1):16–19, 1995.
- [181] M. Gell-Mann and J.B. Hartle. Strong decoherence. In D.H. Feng, editor, *Proc. 4th Drexel Symp. Quantum Non-Integrability—The Quantum-Classical Approach*. To appear.
- [182] M. Gell-Mann and S. Lloyd. Information measures, effective complexity, and total information. *Complexity*, September/October 1996.
- [183] M. Geréb-Graus and M. Li. Three one-way heads cannot do string matching. *J. Comput. System Sci.*, 48:1–8, 1994.
- [184] B.V. Gnedenko. Andrei Nikolaevich Kolmogorov (on the occasion of his seventieth birthday). *Russian Math. Surveys*, 28(5):5–16, 1973.
- [185] E.M. Gold. Language identification in the limit. *Inform. Contr.*, 10:447–474, 1967.
- [186] A. Goldberg and M. Sipser. Compression and ranking. *SIAM J. Comput.*, 20:524–536, 1991.
- [187] O. Goldreich, S. Goldwasser, and S. Micali. How to construct random functions. *J. Assoc. Comput. Mach.*, 33:792–807, 1986.
- [188] P. Grassberger. Towards a quantitative theory of self-generated complexity. *Int'l J. Theoret. Physics*, 25(9):907–938, 1986.
- [189] R.I. Grigorchuk. A connection between algorithmic problems and entropy characteristics of groups. *Soviet Math. Dokl.*, 32:356–360, 1985.
- [190] Yu. Gurevich. The logic of computer science column. *EATCS Bulletin*, 35:71–82, June 1988.
- [191] U. Hahn and N. Chater. Concepts and similarity. In K. Lamberts and D. Shanks, editors, *Knowledge, Concepts and Categories*. UCL Press, In Press.
- [192] D. Hammer and A.Kh. Shen'. A strange application of Kolmogorov complexity. *Math. Systems Theory*, 1996. To appear.
- [193] J. Hartmanis. Generalized Kolmogorov complexity and the structure of feasible computations. In *Proc. 24th IEEE Symp. Found. Comput. Sci.*, pages 439–445, 1983.
- [194] J. Hartmanis and L. Hemachandra. On sparse oracles separating feasible complexity classes. *Inform. Process. Lett.*, 28:291–295, 1988.
- [195] J. Hartmanis and R.E. Stearns. On the computational complexity of algorithms. *Trans. Amer. Math. Soc.*, 117:285–306, 1969.
- [196] R. Heim. On the algorithmic foundation of information theory. *IEEE Trans. Inform. Theory*, IT-25:557–566, 1979.

- [197] L. Hemachandra. On ranking. In *Proc. 2nd IEEE Conf. Structure in Complexity Theory*, pages 103–117, 1987.
- [198] L. Hemachandra and S. Rudich. On the complexity of ranking. *J. Comput. System Sci.*, 41:2:251–271, 1990.
- [199] L. Hemachandra and G. Wechsung. Use randomness to characterize the complexity of computation. *Theoret. Comput. Sci.*, 83:313–322, 1991.
- [200] M. Hermo and E. Mayordomo. A note on polynomial size circuits with low resource-bounded Kolmogorov complexity. *Math. Systems Theory*, 27:247–356, 1994.
- [201] G. Hotz. Komplexität als Kriterium in der Theorienbildung. *Akademie der Wissenschaften und der Literatur (Mainz)/Abhandlungen Mathematisch-Naturwissenschaftliche Klasse*, 1, 1988. Steiner-Verlag, Wiesbaden.
- [202] T. Housel and V.A. Kanevsky. Re-engineering business processes: a complexity theory approach. *Inform. Syst. Operations Res.*, 33(4), 1995.
- [203] D.A. Huffman. A method for construction of minimum-redundancy codes. *Proceedings IRE*, 40:1098–1101, 1952.
- [204] M. Hühne. Linear speed-up does not hold on Turing machines with tree storages. Forschungsbericht nr. 462, Universität Dortmund, Fachbereich Informatik, 1993.
- [205] M. Hühne. On the power of several queues. *Theoret. Comput. Sci.*, 113(1):75–91, 1993.
- [206] D.T. Huynh. Resource-bounded Kolmogorov complexity of hard languages. In *Proc. 1st IEEE Conf. Structure in Complexity Theory*, volume 223 of *Lect. Notes Comput. Sci.*, pages 184–195. Springer-Verlag, 1986.
- [207] D.T. Huynh. Effective entropies and data compression. *Inform. Comput.*, 90(1):67–85, 1991.
- [208] D.T. Huynh. The effective entropies of some extensions of context-free languages. *Inform. Process. Lett.*, 37:165–169, 1991.
- [209] D.T. Huynh. Non-uniform complexity and the randomness of certain complete languages. *Theoret. Comput. Sci.*, 96:305–324, 1992.
- [210] K. Jacobs. Turingmaschinen und zufällige 0-1-Folgen. *Selecta Mathematica*, 2:141–167, 1970.
- [211] A.K. Jagota and K.W. Regan. Performance of MAX-CLIQUE approximation heuristics under description-length weighed distributions. Technical Report TR 92-24, Comput. Sci. Dept, SUNY at Buffalo, 1992.
- [212] A. Jakoby, R. Reischuk, and C. Schindelhauer. Malign distributions for average case circuit complexity. In *Proc. 12th Symp. Theoret. Aspects Comput. Sci.*, volume 900 of *Lect. Notes Comput. Sci.*, pages 628–639. Springer-Verlag, 1995.
- [213] E.T. Jaynes. Prior probabilities. *IEEE Trans. Systems Man Cybernet.*, SMC-4:227–241, 1968.
- [214] E.T. Jaynes. On the rationale of maximum entropy methods. *Proceedings of the IEEE*, 70:939–952, 1982.
- [215] E.T. Jaynes. *Papers on Probability, Statistics, and Statistical Physics*. Kluwer Academic Publishers, 1989. Second edition.
- [216] T. Jiang and M. Li. k one-way heads cannot do string-matching. In *Proc. 25th ACM Symp. Theory Comput.*, pages 62–70, 1993.

- [217] T. Jiang and M. Li. On the approximation of shortest common supersequences and longest common subsequences. *SIAM J. Comput.*, 24(5):1122–1139, 1995.
- [218] T. Jiang, J.I. Seiferas, and P.M.B. Vitányi. Two heads are better than two tapes. *J. Assoc. Comput. Mach.* To appear; preliminary version in *Proc. 26th ACM Symp. Theory Comput.*, 1994, pp. 668–675.
- [219] D. Joseph and M. Sitharam. Kolmogorov complexity, restricted nondeterminism and generalized spectra. Technical report, University of Wisconsin, December, 1990.
- [220] D.J. Juedes, J.I. Lathrop, and J.H. Lutz. Computational depth and reducibility. In *Int'l Colloq. Automata, Languages, Programming*, Lect. Notes Comput. Sci. Springer-Verlag, 1993. To appear in *Theoret. Comput. Sci.*
- [221] D.W. Juedes and J.H. Lutz. Kolmogorov complexity, complexity cores, and the distribution of hardness. In O. Watanabe, editor, *Kolmogorov Complexity and Computational Complexity*, pages 43–65. Springer-Verlag, 1992.
- [222] H. Jürgensen and L. Staiger. Local Hausdorff dimension and quadtrees. Technical Report 259, Comput. Sci. Dept., Univ. Western Ontario, March 1991.
- [223] B. Kalyanasundaram and G. Schnitger. The probabilistic communication complexity of set intersection. *SIAM J. Discrete Math.*, 5(4):545–557, 1992.
- [224] T. Kamae. On Kolmogorov's complexity and information. *Osaka J. Math.*, 10:305–307, 1973.
- [225] T. Kamae. Subsequences of normal sequences. *Israel J. Math.*, 16:121–149, 1973.
- [226] T. Kamae and B. Weiss. Normal numbers and selection rules. *Israel J. Math.*, 21:101–110, 1975.
- [227] M.I. Kanovich. On the decision complexity of algorithms. *Soviet Math. Dokl.*, 10:700–701, 1969.
- [228] M.I. Kanovich. Complexity of resolution of a recursively enumerable set as a criterion of its universality. *Soviet Math. Dokl.*, 11:1224–1228, 1970.
- [229] M.I. Kanovich. On the complexity of enumeration and decision of predicates. *Soviet Math. Dokl.*, 11:17–20, 1970.
- [230] M.I. Kanovich. On the decision complexity of recursively enumerable sets. *Soviet Math. Dokl.*, 11:704–706, 1970.
- [231] M.I. Kanovich. On the complexity of Boolean function minimization. *Soviet Math. Dokl.*, 12(3):720–724, 1971.
- [232] M.I. Kanovich. On the precision of a complexity criterion for nonrecursiveness and universality. *Soviet Math. Dokl.*, 18:232–236, 1977.
- [233] M.I. Kanovich. An estimate of the complexity of arithmetic incompleteness. *Soviet Math. Dokl.*, 19:206–210, 1978.
- [234] M.I. Kanovich and N.V. Petri. Some theorems on the complexity of normal algorithms and their computations. *Soviet Math. Dokl.*, 10(1):233–234, 1969.
- [235] F. Kaspar and H.G. Schuster. The complexity of spatiotemporal patterns. *Phys. Review A*, 36(2):842–848, July 1987.

- [236] H.P. Katseff. Complexity dips in infinite binary sequences. *Inform. Contr.*, 38:258–263, 1978.
- [237] H.P. Katseff and M. Sipser. Several results in program size complexity. *Theoret. Comput. Sci.*, 15:291–309, 1981.
- [238] J. Kececioglu, M. Li, and J.T. Tromp. Reconstructing a DNA sequence from erroneous copies. *Theoret. Comput. Sci.*, 1996/1997. To appear.
- [239] J.G. Kemeny. The use of simplicity in induction. *Philos. Rev.*, 62:391–408, 1953.
- [240] H.A. Keuzenkamp and M. McAleer. Simplicity, scientific inference and econometric modelling. *The Economic Journal*, 105:1–21, 1995.
- [241] A.I. Khinchin. *Mathematical Foundations of Information Theory*. Dover, 1957.
- [242] J.C. Kieffer and E.-H. Yang. Sequential codes, lossless compression of individual sequences, and Kolmogorov complexity. *IEEE Trans. Inform. Theory*, IT-42(1):29–39, 1996.
- [243] W.W. Kirchherr. Kolmogorov complexity and random graphs. *Inform. Process. Lett.*, 41:125–130, 1992.
- [244] D.E. Knuth. *The Art of Computer Programming. Volume 1: Fundamental Algorithms*. Addison-Wesley, 1973. Second edition.
- [245] D.E. Knuth. *The Art of Computer Programming. Volume 2: Seminumerical Algorithms*, pages 163–166. Addison-Wesley, 1981. Second edition.
- [246] Ker-I Ko. On the definition of infinite pseudo-random sequences. *Theoret. Comput. Sci.*, 48:9–34, 1986.
- [247] K. Kobayashi. On the structure of one-tape nondeterministic Turing machine time hierarchy. *Theoret. Comput. Sci.*, 40:175–193, 1985.
- [248] K. Kobayashi.  $\Sigma_n^0$ -complete properties of programs and Martin-Löf randomness. *Inform. Process. Lett.*, 46:37–42, 1993.
- [249] K. Kobayashi. On malign input distributions for algorithms. *IEICE Trans. Inform. and Syst.*, E76-D(6):634–640, 1993.
- [250] A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*. Springer-Verlag, 1933. English translation (by N. Morrison): *Foundations of the Theory of Probability*, Chelsea, 1956; 2nd Russian edition: *Osnovnye Poniatija Teorii Verojatnostej*, Nauka, 1974.
- [251] A.N. Kolmogorov. On tables of random numbers. *Sankhyā, The Indian Journal of Statistics, Ser. A*, 25:369–376, 1963.
- [252] A.N. Kolmogorov. Three approaches to the quantitative definition of information. *Problems Inform. Transmission*, 1(1):1–7, 1965.
- [253] A.N. Kolmogorov. Logical basis for information theory and probability theory. *IEEE Trans. Inform. Theory*, IT-14(5):662–664, 1968.
- [254] A.N. Kolmogorov. Some theorems on algorithmic entropy and the algorithmic quantity of information. *Uspekhi Mat. Nauk*, 23(2):201, 1968. Meeting of the Moscow Mathematical Society.
- [255] A.N. Kolmogorov. On the logical foundations of information theory and probability theory. *Problems Inform. Transmission*, 5:1–4, 1969.
- [256] A.N. Kolmogorov. Combinatorial foundations of information theory and the calculus of probabilities. *Russian Math. Surveys*, 38(4):29–40, 1983.
- [257] A.N. Kolmogorov. On logical foundations of probability theory. In K. Itô and Yu.V. Prokhorov, editors, *Probability Theory and Mathematical*

- Statistics*, volume 1021 of *Lect. Notes Math.*, pages 1–5. Springer-Verlag, 1983.
- [258] A.N. Kolmogorov. Memories of P.S. Aleksandrov. *Russian Math. Surveys*, 41(6):225–246, 1986.
  - [259] A.N. Kolmogorov. *Information Theory and Theory of Algorithms, Selected Works*, volume 3. Nauka, 1987. Edited by Yu.V. Prokhorov and A.N. Shiryaev. In Russian.
  - [260] A.N. Kolmogorov. Letters of A.N. Kolmogorov to A. Heyting. *Russian Math. Surveys*, 43(6):89–93, 1988.
  - [261] A.N. Kolmogorov and V.A. Uspensky. On the definition of an algorithm. *Uspekhi Mat. Nauk*, 13(4):3–28, 1958. In Russian. English translation: *Amer. Math. Soc. Translat.*, 29:2(1963), 217–245.
  - [262] A.N. Kolmogorov and V.A. Uspensky. Algorithms and randomness. *SIAM J. Theory Probab. Appl.*, 32:389–412, 1987. Without annoying translation errors pp. 3–53 in: Yu.V. Prokhorov and V.V. Sazonov, Eds., *Proc. 1st World Congress of the Bernoulli Society (Tashkent 1986)*, Vol. 1: *Probab. Theory and Appl.*, VNU Science Press, Utrecht, 1987.
  - [263] D.K. Kondepudi. Non-equilibrium polymers, entropy, and algorithmic information. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 199–206. Addison-Wesley, 1991.
  - [264] M. Koppel. Complexity, depth, and sophistication. *Complex Systems*, 1:1087–1091, 1987.
  - [265] M. Koppel. Structure. In R. Herken, editor, *The Universal Turing Machine; A Half-Century Survey*, pages 435–452. Oxford University Press, 1988. In Germany: Kammerer & Unverzagt, Hamburg.
  - [266] L.G. Kraft. A device for quantizing, grouping and coding amplitude modulated pulses. Master’s thesis, Dept. of Electrical Engineering, M.I.T., Cambridge, Mass., 1949.
  - [267] E. Kranakis and D. Krizanc. Lower bounds for compact routing. In C. Puech and R. Reischuk, editors, *Proc. 13th Symp. Theoret. Aspects Comput. Sci.*, volume 1046 of *Lect. Notes Comput. Sci.*, pages 529–540. Springer-Verlag, 1996.
  - [268] E. Kranakis, D. Krizanc, and F. Luccio. String recognition on anonymous rings. In Peter Hajek, editor, *Proc. 13th Symp. Math. Found. Comput. Sci.*, volume 969 of *Lect. Notes Comput. Sci.*, pages 392–401. Springer-Verlag, 1995.
  - [269] V. Kreinovich and L. Longpré. Why are symmetries a universal language of physics? (on the unreasonable effectiveness of symmetries in physics). *Int'l J. Theoret. Phys.*, 1996. To appear.
  - [270] V. Kreinovich and R. Watson. How difficult is it to invent a nontrivial game. *Cybernetics and Systems: An Int'l Journal*, 25:629–640, 1994.
  - [271] R.E. Krichevskii. Universal encoding and Kolmogorov complexity. In *Proc. 5th Int'l Symp. Inform. Theory, Part 1, Abstracts of Papers*, pages 22–25, Moscow-Tbilisi, 1979. In Russian.
  - [272] R.E. Krichevskii and V.K. Trofimov. The performance of universal encoding. *IEEE Trans. Inform. Theory*, IT-27:199–207, 1983.
  - [273] M. Kummer. Kolmogorov complexity and instance complexity of recursively enumerable sets. *SIAM J. Comput.*, 25(6):1123–1143, 1996.

- [274] M. Kummer. On the complexity of random strings. In *Proc. 13th Symp. Theoret. Aspects Comput. Sci.*, volume 1046 of *Lect. Notes Comput. Sci.*, pages 25–36. Springer-Verlag, 1996.
- [275] R. Landauer. Irreversibility and heat generation in the computing process. *IBM J. Res. Develop.*, 5:183–191, 1961.
- [276] P.S. Laplace. *A Philosophical Essay on Probabilities*, pages 16–17. Dover, 1952. Originally published in 1819. Translated from 6th French edition.
- [277] A. Lempel and J. Ziv. On the complexity of finite sequences. *IEEE Trans. Inform. Theory*, IT-22:75–81, 1976.
- [278] S.K. Leung-Yan-Cheong and T.M. Cover. Some equivalences between Shannon entropy and Kolmogorov complexity. *IEEE Trans. Inform. Theory*, IT-24:331–339, 1978.
- [279] L.A. Levin. *Some theorems on the algorithmic approach to probability theory and information theory*. PhD thesis, Moscow University, 1971. In Russian.
- [280] L.A. Levin. On storage capacity for algorithms. *Soviet Math. Dokl.*, 14:1464–1466, 1973.
- [281] L.A. Levin. On the notion of a random sequence. *Soviet Math. Dokl.*, 14:1413–1416, 1973.
- [282] L.A. Levin. Universal search problems. *Problems Inform. Transmission*, 9:265–266, 1973.
- [283] L.A. Levin. Laws of information conservation (non-growth) and aspects of the foundation of probability theory. *Problems Inform. Transmission*, 10:206–210, 1974.
- [284] L.A. Levin. On the principle of conservation of information in intuitionistic mathematics. *Soviet Math. Dokl.*, 17:601–605, 1976.
- [285] L.A. Levin. Uniform tests of randomness. *Soviet Math. Dokl.*, 17:337, 1976.
- [286] L.A. Levin. Various measures of complexity for finite objects (axiomatic description). *Soviet Math. Dokl.*, 17:522–526, 1976.
- [287] L.A. Levin. Randomness conservation inequalities; information and independence in mathematical theories. *Inform. Contr.*, 61:15–37, 1984.
- [288] L.A. Levin. One-way functions and pseudorandom generators. *Combinatorica*, 7:357–363, 1987.
- [289] L.A. Levin and V.V. V'yugin. Invariant properties of information bulks. In *Lect. Notes Comput. Sci.*, volume 53, pages 359–364. Springer-Verlag, 1977.
- [290] M. Li, L. Longpré, and P.M.B. Vitányi. The power of the queue. *SIAM J. Comput.*, 21(4):697–712, 1992.
- [291] M. Li and P.M.B. Vitányi. Kolmogorovskaya slozhnost' dvadsat' let sputstia. *Uspekhi Mat. Nauk*, 43(6):129–166, 1988. In Russian.
- [292] M. Li and P.M.B. Vitányi. Tape versus queue and stacks: The lower bounds. *Inform. Comput.*, 78:56–85, 1988.
- [293] M. Li and P.M.B. Vitányi. Two decades of applied Kolmogorov complexity: In memoriam A.N. Kolmogorov 1903–1987. In *Proc. 3rd IEEE Conf. Structure in Complexity Theory*, pages 80–101, 1988.
- [294] M. Li and P.M.B. Vitányi. Applications of Kolmogorov complexity in the theory of computation. In A.L. Selman, editor, *Complexity Theory Retrospective*, pages 147–203. Springer-Verlag, 1990.

- [295] M. Li and P.M.B. Vitányi. Kolmogorov complexity and its applications. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, chapter 4, pages 187–254. Elsevier and MIT Press, 1990.
- [296] M. Li and P.M.B. Vitányi. Learning simple concepts under simple distributions. *SIAM J. Comput.*, 20(5):915–935, 1991.
- [297] M. Li and P.M.B. Vitányi. Inductive reasoning and Kolmogorov complexity. *J. Comput. System Sci.*, 44(2):343–384, 1992.
- [298] M. Li and P.M.B. Vitányi. Worst case complexity is equal to average case complexity under the universal distribution. *Inform. Process. Lett.*, 42:145–149, 1992.
- [299] M. Li and P.M.B. Vitányi. Mathematical theory of thermodynamics of computation. In *IEEE Proc. Workshop on Physics and Computation*, pages 42–46, 1993. Complete version in prelim. proceedings, Oct., 1992.
- [300] M. Li and P.M.B. Vitányi. Inductive reasoning. In E.S. Ristad, editor, *Language Computations; Proc. DIMACS Workshop on Human Language*, volume 17 of *DIMACS Series in Discr. Math. Theoret. Comput. Sci.*, pages 127–148. American Math. Society, 1994.
- [301] M. Li and P.M.B. Vitányi. Kolmogorov complexity arguments in combinatorics. *J. Comb. Theory, Ser. A*, 66(2):226–236, 1994. Erratum, *Ibid.*, 69(1995), 183.
- [302] M. Li and P.M.B. Vitányi. Statistical properties of finite sequences with high Kolmogorov complexity. *Math. Systems Theory*, 27:365–376, 1994.
- [303] M. Li and P.M.B. Vitányi. Computational machine learning in theory and praxis. In J. van Leeuwen, editor, *Computer Science Today, Recent Trends and Developments*, volume 1000 of *Lect. Notes Comput. Sci.*, pages 518–535. Springer-Verlag, 1995.
- [304] M. Li and P.M.B. Vitányi. A new approach to formal language theory by Kolmogorov complexity. *SIAM J. Comput.*, 24(2):398–410, 1995.
- [305] M. Li and P.M.B. Vitányi. Reversibility and adiabatic computation: trading time and space for energy. *Proc. Royal Soc. London, Ser. A*, 452:769–789, 1996.
- [306] M. Li and Y. Yesha. String-matching cannot be done by 2-head 1-way deterministic finite automata. *Inform. Process. Lett.*, 22:231–235, 1986.
- [307] M. Li and Y. Yesha. New lower bounds for parallel computation. *J. Assoc. Comput. Mach.*, 36:671–680, 1989.
- [308] S. Lloyd and J.J. Slotine. Algorithmic Lyapunov functions for stable adaptation and control. *Int'l J. Adapt. Contr. Sign. Proc.*, 1996.
- [309] L. Löfgren. Explicability of sets and transfinite automata. In E. Caianiello, editor, *Automata Theory*, pages 251–268. Academic Press, 1966.
- [310] L. Löfgren. Recognition of order and evolutionary systems. In J. Tou, editor, *Computer and Information Sciences II*, pages 165–175. Academic Press, 1967.
- [311] L. Löfgren. Complexity of descriptions of systems. *Int'l J. General Systems*, 3:197–214, 1977.
- [312] L. Longpré. *Resource bounded Kolmogorov complexity, a link between computational complexity and information theory*. PhD thesis, Comput. Sci. Dept., Cornell Univ., 1986.

- [313] L. Longpré. Resource bounded Kolmogorov complexity and statistical tests. In O. Watanabe, editor, *Kolmogorov Complexity and Computational Complexity*, pages 66–84. Springer-Verlag, 1992.
- [314] L. Longpré and V. Kreinovich. Randomness as incompressibility: a non-algorithmic analogue. Manuscript, Comput. Sci. Dept., Univ. Texas at El Paso, 1996.
- [315] L. Longpré and S. Mocas. Symmetry of information and one-way functions. *Inform. Proc. Lett.*, 46(2):95–100, 1993.
- [316] L. Longpré and O. Watanabe. On symmetry of information and polynomial time invertibility. *Inform. Comput.*, 121(1):14–22, 1995.
- [317] A. López-Ortiz. New lower bounds for element distinctness on a one-tape Turing machine. *Inform. Proc. Lett.*, 51(6):311–314, 1994.
- [318] M.C. Loui. Optimal dynamic embedding of trees into arrays. *SIAM J. Comput.*, 12:463–472, 1983.
- [319] M.C. Loui. Minimizing access pointers into trees and arrays. *J. Comput. System Sci.*, 28:359–378, 1984.
- [320] M.C. Loui and D.R. Luginbuhl. The complexity of simulations between multidimensional Turing machines and random access machines. *Math. Systems Theory*, 25(4):293–308, 1992.
- [321] M.C. Loui and D.R. Luginbuhl. Optimal on-line simulations of tree machines by random access machines. *SIAM J. Comput.*, 21(5):959–971, 1992.
- [322] D.W. Loveland. The Kleene hierarchy classification of recursively random sequences. *Trans. Amer. Math. Soc.*, 125:497–510, 1966.
- [323] D.W. Loveland. A new interpretation of von Mises' concept of a random sequence. *Z. Math. Logik und Grundlagen Math.*, 12:279–294, 1966.
- [324] D.W. Loveland. On minimal-program complexity measures. In *Proc. (1st) ACM Symp. Theory Comput.*, pages 61–66, 1969.
- [325] D.W. Loveland. A variant of the Kolmogorov concept of complexity. *Inform. Contr.*, 15:510–526, 1969.
- [326] D.R. Luginbuhl. *Computational complexity of random access models*. PhD thesis, University of Illinois at Urbana-Champaign, 1990.
- [327] J.H. Lutz. Category and measure in complexity classes. *SIAM. J. Comput.*, 19:6:1100–1131, 1990.
- [328] J.H. Lutz. An upward measure separation theorem. *Theoret. Comput. Sci.*, 81:127–135, 1991.
- [329] J.H. Lutz. Almost everywhere high nonuniform complexity. *J. Comput. System Sci.*, 44:220–258, 1992.
- [330] J.H. Lutz. The quantitative structure of exponential time. In *Proc. 8th IEEE Conf. Structure in Complexity Theory*, pages 158–175, 1993.
- [331] W. Maass. Combinatorial lower bound arguments for deterministic and nondeterministic Turing machines. *Trans. Amer. Math. Soc.*, 292:675–693, 1985.
- [332] W. Maass, E. Szemerédi, G. Schnitger, and G. Turan. Two tapes versus one for off-line Turing machines. *Computational Complexity*, 3:392–401, 1993.
- [333] H.G. Mairson. The program complexity of searching a table. In *Proc. 24th IEEE Symp. Found. Comput. Sci.*, pages 40–47, 1983.

- [334] G.B. Marandzhyan. On certain properties of asymptotically optimal recursive function. *Izv. Akad. Nauk Armyan. SSSR*, 4:3–22, 1969.
- [335] A.A. Markov. On normal algorithms which compute Boolean functions. *Soviet Math. Dokl.*, 5:922–924, 1964.
- [336] A.A. Markov. On normal algorithms associated with the computation of Boolean functions and predicates. *Izv. Akad. Nauk USSR Ser. Mat.*, 31:161–208, 1967.
- [337] P. Martin-Löf. Algorithmen und zufällige Folgen. Lecture notes, University of Erlangen, 1966.
- [338] P. Martin-Löf. The definition of random sequences. *Inform. Contr.*, 9:602–619, 1966.
- [339] P. Martin-Löf. On the concept of a random sequence. *Theory Probability Appl.*, 11:177–179, 1966.
- [340] P. Martin-Löf. Algorithms and randomness. *Rev. Int. Statist. Inst.*, 37:265–272, 1969.
- [341] P. Martin-Löf. The literature on von Mises' Kollektivs revisited. *Theoria*, 35(1):12, 1969.
- [342] P. Martin-Löf. *Notes on Constructive Mathematics*. Almqvist and Wiksell, 1970.
- [343] P. Martin-Löf. On the notion of randomness. In A. Kino et al., editor, *Intuitionism and Proof Theory*, pages 73–78. North-Holland, 1970.
- [344] P. Martin-Löf. Complexity oscillations in infinite binary sequences. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 19:225–230, 1971.
- [345] P. Martin-Löf. The notion of redundancy and its use as a quantitative measure of the discrepancy between a statistical hypothesis and a set of observational data. *Scand. J. Stat.*, 1:3–18, 1974.
- [346] P. Martin-Löf. Reply to Sverdrup's polemical article “Tests without power”. *Scand. J. Stat.*, 2:161–165, 1975.
- [347] K. Mehlhorn. On the program-size of perfect and universal hash functions. In *Proc. 23rd IEEE Symp. Found. Comput. Sci.*, pages 170–175, 1982.
- [348] N.C. Metropolis, G. Reitweiser, and J. von Neumann. Statistical treatment of values of the first 2,000 decimal digits of  $e$  and  $\pi$  calculated on the ENIAC. In A.H. Traub, editor, *John von Neumann, Collected Works, Vol. V*. Macmillan, 1963.
- [349] P.B. Miltersen. The complexity of malign ensembles. *SIAM J. Comput.*, 22(1):147–156, 1993.
- [350] M.L. Minsky. Problems of formulation for artificial intelligence. In R.E. Bellman, editor, *Mathematical Problems in the Biological Sciences*, Proc. Symposia in Applied Mathematics XIV, page 43. American Mathematical Society, 1962.
- [351] M.L. Minsky. Steps towards artificial intelligence. *Proceedings I.R.E.*, pages 8–30, January, 1961.
- [352] An.A. Muchnik. Lower limits on frequencies in computable sequences and relativized a priori probability. *SIAM Theory Probab. Appl.*, 32:513–514, 1987.
- [353] D.W. Müller. Randomness and extrapolation. In *Sixth Berkeley Symposium*, pages 1–31, 1970.

- [354] A. Naik, K.W. Regan, and D. Sivakumar. On quasilinear time complexity theory. *Theoret. Comput. Sci.*, 148(2):325–349, 1995.
- [355] S.M. Nikol'skii. Aleksandrov and Kolmogorov in Dnjepropetrovsk. *Russian Math. Surveys*, 38(4):41–55, 1983.
- [356] S.P. Novikov. Memories of A.N. Kolmogorov. *Russian Math. Surveys*, 43(6):40–42, 1988.
- [357] Obituary. Mr. Andrei Kolmogorov—Giant of mathematics. *Times*, October 26 1987.
- [358] Obituary. Andrei Nikolaevich Kolmogorov. *Bull. London Math. Soc.*, 22(1):31–100, 1990.
- [359] P. Odifreddi. *Classical Recursion Theory*. North-Holland, 1989.
- [360] P. Orponen, Ker-I Ko, U. Schöning, and O. Watanabe. Instance complexity. *J. Assoc. Comput. Mach.*, 41:96–121, 1994.
- [361] D. Pager. On the problem of finding minimal programs for tables. *Inform. Contr.*, 14:550–554, 1969.
- [362] R. Paturi. *Study of certain probabilistic models of information transfer and on-line computation*. PhD thesis, Penn State University, 1985.
- [363] R. Paturi and J. Simon. Lower bounds on the time of probabilistic on-line simulations. In *Proc. 24th IEEE Symp. Found. Comput. Sci.*, pages 343–350, 1983.
- [364] R. Paturi, J. Simon, R.E. Newman-Wolfe, and J. Seiferas. Milking the Aanderaa argument. *Inform. Comput.*, 88:88–104, 1990.
- [365] W.J. Paul. Kolmogorov's complexity and lower bounds. In L. Budach, editor, *Proc. 2nd Int'l Conf. Fundamentals of Computation Theory*, pages 325–334. Akademie Verlag, 1979.
- [366] W.J. Paul. On-line simulation of  $k+1$  tapes by  $k$  tapes requires nonlinear time. *Inform. Contr.*, 53:1–8, 1982.
- [367] W.J. Paul. On heads versus tapes. *Theoret. Comput. Sci.*, 28:1–12, 1984.
- [368] W.J. Paul, J.I. Seiferas, and J. Simon. An information theoretic approach to time bounds for on-line computation. *J. Comput. System Sci.*, 23(2):108–126, 1981.
- [369] W.J. Paul and R.J. Solomonoff. Autonomous theory building systems. In P. Bock, M. Loew, and M. Richter, editors, *Neural Networks and Adaptive Learning*, Knowledge Processing and Its Applications Series. Elsevier Science Publishers, 1992.
- [370] J. Pearl. On the connection between the complexity and credibility of inferred models. *Int'l J. Gen. Syst.*, 4:255–264, 1978.
- [371] E.P.D. Pednault. Some experiments in applying inductive inference principles to surface reconstruction. In *11th Int'l Joint Conf. Artificial Intelligence*, pages 1603–1609. Morgan-Kaufmann, 1989.
- [372] G. Peterson. Succinct representations, random strings and complexity classes. In *Proc. 21st IEEE Symp. Found. Comput. Sci.*, pages 86–95, 1980.
- [373] N.V. Petri. Algorithms connected with predicates and Boolean functions. *Soviet Math. Dokl.*, 10:294–297, 1969.
- [374] N.V. Petri. The complexity of algorithms and their operating time. *Soviet Math. Dokl.*, 10:547–549, 1969.

- [375] J.R. Pierce and C.C. Cutler. Interplanetary communications. In F.I. Ordway, III, editor, *Advances in Space Science, Vol. 1*, pages 55–109, New York, 1959. Academic Press, Inc.
- [376] N. Pippenger. An information-theoretic method in combinatorial theory. *J. Comb. Theory, Ser. A*, 23:99–104, 1977.
- [377] K.R. Popper. *The Logic of Scientific Discovery*. University of Toronto Press, 1959.
- [378] M.B. Pour-El and J.I. Richards. *Computability in Analysis and Physics*. Springer-Verlag, 1989.
- [379] J. Quinlan and R. Rivest. Inferring decision trees using the minimum description length principle. *Inform. Comput.*, 80:227–248, 1989.
- [380] K.W. Regan. On superlinear lower bounds in complexity theory. In *Proc. 10th IEEE Conf. Structure in Complexity Theory*, pages 50–64, 1995.
- [381] K.W. Regan and J. Wang. The quasilinear isomorphism challenge. *SIGACT News*, 25:106–113, September 1994.
- [382] S. Reisch and G. Schnitger. Three applications of Kolmogorov complexity. In *Proc. 23rd IEEE Symp. Found. Comput. Sci.*, pages 45–52, 1982.
- [383] Zh.I. Reznikova and B.Ya. Ryabko. Analysis of the language of ants by information-theoretical methods. *Problems Inform. Transmission*, 22:245–249, 1986.
- [384] J.J. Rissanen. Modeling by the shortest data description. *Automatica*, 14:465–471, 1978.
- [385] J.J. Rissanen. *Stochastical Complexity and Statistical Inquiry*. World Scientific, 1989.
- [386] J.J. Rissanen. Complexity of models. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 117–125. Addison-Wesley, 1991.
- [387] J.J. Rissanen. Stochastic complexity in learning. In P.M.B. Vitányi, editor, *Computational Learning Theory, Proc. 2nd European Conf.*, volume 904 of *Lect. Notes in Artificial Intelligence*, pages 196–210. Springer-Verlag, 1995. *J. Comput. System Sci.*, To appear.
- [388] J.J. Rissanen. Fisher information and stochastic complexity. *IEEE Trans. Inform. Theory*, IT-42(1):40–47, 1996.
- [389] R. Rivest. Lecture notes in machine learning. Manuscript, 1987.
- [390] H. Rogers, Jr. *Theory of Recursive Functions and Effective Computability*. McGraw-Hill, 1967.
- [391] R. Rubinstein. *Structural complexity classes of sparse sets: intractability, data compression and printability*. PhD thesis, Northeastern Univ., 1988.
- [392] B.Ya. Ryabko. Encoding of combinatorial sources and Hausdorff dimension. *Dokl. Akad. Nauk SSSR*, 27:1066–1070, 1984.
- [393] B.Ya. Ryabko. Noiseless coding of combinatorial sources, Hausdorff dimension, and Kolmogorov complexity. *Problems Inform. Transmission*, 22:170–179, 1986.
- [394] B.Ya. Ryabko. The complexity and effectiveness of prediction algorithms. *J. Complexity*, 10:281–295, 1994.
- [395] R. Schack. Algorithmic information and simplicity in statistical physics. *Physical Review E*. Manuscript, submitted.
- [396] R. Schack and C.M. Caves. Chaos for Liouville probability densities. *Physical Review E*, 53(4):3387–3401, 1996.

- [397] R. Schack and C.M. Caves. Information-theoretic characterization of quantum chaos. *Physical Review E*, 53(4):3257–3270, 1996.
- [398] R. Schack, G.M. D’Ariano, and C.M. Caves. Hypersensitivity to perturbation in a quantum kicked top. *Physical Review E*, 50:972, 1994.
- [399] J. Schmidhuber. Discovering solutions with low Kolmogorov complexity and high generalization capability. In *Proc. 12th Intern. Conf. on Machine Learning*, pages 488–496. Morgan Kaufmann Publishers, 1995.
- [400] C.P. Schnorr. Eine Bemerkung zum Begriff der zufälligen Folge. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 14:27–35, 1969.
- [401] C.P. Schnorr. Klassifikation der Zufallsgesetze nach Komplexität und Ordnung. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 16:1–21, 1970.
- [402] C.P. Schnorr. Über die Definition von effektiven Zufallstests, i-ii. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 15:297–312, 313–328, 1970.
- [403] C.P. Schnorr. Optimal Gödel numberings. In *Proc. 1971 IFIP Congress, TA-2*, pages 12–14, Ljubljana, Yugoslavia, 1971.
- [404] C.P. Schnorr. A unified approach to the definition of random sequences. *Math. Systems Theory*, 5:246–258, 1971.
- [405] C.P. Schnorr. *Zufälligkeit und Wahrscheinlichkeit; Eine algorithmische Begründung der Wahrscheinlichkeitstheorie*, volume 218 of *Lect. Notes Math.* Springer-Verlag, 1971.
- [406] C.P. Schnorr. Process complexity and effective random tests. *J. Comput. System Sci.*, 7:376–388, 1973.
- [407] C.P. Schnorr. *Rekursive Funktionen und ihre Komplexität*. Teubner, 1974.
- [408] C.P. Schnorr. A survey of the theory of random sequences. In R.E. Butts and J. Hintikka, editors, *Basic Problems in Methodology and Linguistics*, pages 193–210. D. Reidel, 1977.
- [409] C.P. Schnorr. A review of the theory of random sequences. In *Proc. 5th Int'l Congr. Logic, Meth. Phil. of Sci.*, London, Ontario, August 1975.
- [410] C.P. Schnorr and P. Fuchs. General random sequences and learnable sequences. *J. Symbolic Logic*, 42:329–340, 1977.
- [411] C.P. Schnorr and H. Stimm. Endliche Automaten und Zufallsfolgen. *Acta Informatica*, 1:345–359, 1972.
- [412] C.P. Schnorr and G. Stumpe. A characterization of complexity sequences. *Z. Math. Logik und Grudl. Math.*, 21:47–56, 1975.
- [413] J. Seiferas. The symmetry of information, and An application of the symmetry of information. Unpublished memo, Comput. Sci. Dept, University of Rochester, August 1985.
- [414] J. Seiferas. A simplified lower bound for context-free-language recognition. *Inform. Contr.*, 69:255–260, 1986.
- [415] J. Shallit and Y. Breitbart. Automaticity: Properties of a measure of descriptional complexity. *J. Comput. System Sci.*, 53(1):10–25, 1996.
- [416] C.E. Shannon. The mathematical theory of communication. *Bell System Tech. J.*, 27:379–423, 623–656, 1948.
- [417] C.E. Shannon. A universal Turing machine with two internal states. In C.E. Shannon and J. McCarthy, editors, *Automata Studies*, pages 129–153. Princeton University Press, 1956.
- [418] A.Kh. Shen’. The frequency approach to defining a random sequence. *Semiotika i Informatika*, 19:14–42, 1982. In Russian.

- [419] A.Kh. Shen'. The concept of Kolmogorov  $(\alpha, \beta)$ -stochasticity and its properties. *Soviet Math. Dokl.*, 28:295–299, 1983.
- [420] A.Kh. Shen'. Algorithmic variants of the notion of entropy. *Soviet Math. Dokl.*, 29(3):569–573, 1984.
- [421] A.Kh. Shen'. Connections between different algorithmic definitions of randomness. *Soviet Math. Dokl.*, 38(2):316–319, 1989.
- [422] A.N. Shiryaev. A.N. Kolmogorov: Life and creative activities. *Ann. Probab.*, 17:866–944, 1989. Publications of Kolmogorov: pages 945–964.
- [423] M. Sipser. A complexity theoretic approach to randomness. In *Proc. 15th ACM Symp. Theory Comput.*, pages 330–335, 1983.
- [424] S.S. Skiena. Further evidence for randomness in  $\pi$ . *Complex Systems*, 1:361–366, 1987.
- [425] D. Sleator, R. Tarjan, and W. Thurston. Short encodings of evolving structures. *SIAM J. Discr. Math.*, 5:428–450, 1992.
- [426] R.J. Solomonoff. The mechanization of linguistic learning. In *2nd Int'l Congress on Cybernetics*, pages 180–193, 1958.
- [427] R.J. Solomonoff. A new method for discovering the grammars of phrase structure languages. In *Information Processing*, pages 285–290, Paris, 1959. Unesco.
- [428] R.J. Solomonoff. A preliminary report on a general theory of inductive inference. Technical Report ZTB-138, Zator Company, Cambridge, Mass., November 1960.
- [429] R.J. Solomonoff. An inductive inference code employing definitions. Technical Report ZTB-141, Rockford Research, Cambridge, Mass., April 1962.
- [430] R.J. Solomonoff. A formal theory of inductive inference, part 1 and part 2. *Inform. Contr.*, 7:1–22, 224–254, 1964.
- [431] R.J. Solomonoff. Inductive inference research status, spring 1967. Technical report, Rockford Research Inst., July 1967. Distributed by Clearinghouse, US Dept. of Commerce.
- [432] R.J. Solomonoff. Inductive inference theory—a unified approach to problems in pattern recognition and artificial intelligence. In *4th Int'l Conf. Artificial Intelligence*, pages 274–280, Tbilisi, Georgia, USSR, 1975.
- [433] R.J. Solomonoff. Complexity based induction systems: comparisons and convergence theorems. Technical Report RR-329, Rockford Research, Cambridge, Mass., August 1976.
- [434] R.J. Solomonoff. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Inform. Theory*, IT-24:422–432, 1978.
- [435] R.J. Solomonoff. Perfect training sequences and the costs of corruption—a progress report on inductive inference research. Memorandum, Oxbridge Research, P.O. box 559, Cambridge, Mass. 02238, August 1982.
- [436] R.J. Solomonoff. Optimum sequential search. Memorandum, Oxbridge Research, P.O. box 559, Cambridge, Mass. 02238, June 1984.
- [437] R.J. Solomonoff. An application of algorithmic probability to problems in artificial intelligence. In L.N. Kanal and J.F. Lemmer, editors, *Uncertainty in Artificial Intelligence*, pages 473–491. North-Holland, 1986.

- [438] R.J. Solomonoff. The application of algorithmic probability to machine learning. Grant proposal manuscript, Oxbridge Research, P.O. box 559, Cambridge, Mass. 02238, September 1988.
- [439] R.J. Solomonoff. A system for machine learning based on algorithmic probability. In *Proc. 6th Israeli Conf. on AI and Computer Vision*, 1989.
- [440] R.J. Solomonoff. The discovery of algorithmic probability: A guide for the programming of true creativity. In P.M.B. Vitányi, editor, *Computational Learning Theory; Proc. 2nd European Conf.*, volume 904 of *Lect. Notes Artificial Intelligence*, pages 1–18. Springer-Verlag, 1995. *J. Comput. System Sci.*, To appear.
- [441] R.M. Solovay. Lecture notes on algorithmic complexity. Unpublished, UCLA, 1975.
- [442] R.M. Solovay. On random r.e. sets. In A.I. Arruda et al., editor, *Non-Classical Logic, Model Theory and Computability*, pages 283–307. North-Holland, 1977.
- [443] L. Staiger. Complexity and entropy. In *Proc. Math. Found. Comput. Sci.*, volume 118 of *Lect. Notes Comput. Sci.*, pages 508–514. Springer-Verlag, 1981.
- [444] L. Staiger. Representable P. Martin-Löf tests. *Kybernetika*, 21:235–243, 1985.
- [445] L. Staiger. Kolmogorov complexity, channel capacity and Hausdorff dimension. In *Proc. 4th Joint Swedish-USSR Workshop on Information Theory*, pages 132–137. Lund University, Student-litteratur, 1989.
- [446] L. Staiger. Kolmogorov complexity and Hausdorff dimension. *Inform. Comput.*, 120(2):159–194, 1993.
- [447] L. Staiger. A tight upper bound on Kolmogorov complexity by Hausdorff dimension and uniformly optimal prediction. *Math. Systems Theory*, 1996/1997. To appear.
- [448] J. Storer. *Data Compression: Method and Theory*. Computer Science Press, 1988.
- [449] E. Sverdrup. Tests without power. *Scand. J. Stat.*, 2:158–160, 1975.
- [450] L. Szilard. On the decrease of entropy in a thermodynamic system by the intervention of intelligent beings. *Z. Phys.*, 53:840–856, 1929.
- [451] V.M. Tikhomirov. The life and work of Andrei Nikolaevich Kolmogorov. *Russian Math. Surveys*, 43(6):1–39, 1988.
- [452] B.A. Trakhtenbrot. A survey of Russian approaches to perebor (brute-force-search) algorithms. *Ann. Hist. Comput.*, 6:384–400, 1984.
- [453] J.T. Tromp. A cl-based definition of kolmogorov complexity. Manuscript, CWI, Amsterdam, November 1996.
- [454] A.M. Turing. On computable numbers with an application to the Entscheidungsproblem. *Proc. London Math. Soc.*, Ser. 2, 42:230–265, 1936. Correction, *Ibid*, 43:544–546, 1937.
- [455] J. Tyszkiewicz. On the Kolmogorov expressive power of Boolean query languages. In G. Gottlob and M.Y. Vardi, editors, *Proc. 8th Intern. Conf. Database Theory*, volume 893 of *Lect. Notes Comput. Sci.*, pages 97–110. Springer-Verlag, 1995. To appear in *Theoret. Comput. Sci.*
- [456] J. Tyszkiewicz. The Kolmogorov expression complexity of logics. Preliminary version in *Proc. Comput. Sci. Logic Conf.*, 1995; submitted to *Inform. Comput.*, 1996.

- [457] V.A. Uspensky. Complexity and entropy: an introduction to the theory of Kolmogorov complexity. In O. Watanabe, editor, *Kolmogorov Complexity and Computational Complexity*, pages 85–102. Springer-Verlag, 1992.
- [458] V.A. Uspensky. Kolmogorov and mathematical logic. *J. Symb. Logic*, 57(2):385–412, 1992.
- [459] V.A. Uspensky and A.L. Semenov. *Algorithms: Main Ideas and Applications*. Kluwer Academic Publishers, 1993. Also: in Lect. Notes Comput. Sci., vol. 122, A.P. Ershov and D.E. Knuth, Eds., Springer-Verlag, 1981, pp. 100–234.
- [460] V.A. Uspensky, A.L. Semenov, and A.Kh. Shen'. Can an individual sequence of zeros and ones be random? *Russian Math. Surveys*, 45(1):121–189, 1990.
- [461] V.A. Uspensky and A.Kh. Shen'. Relations between varieties of Kolmogorov complexities. *Math. Systems Theory*, 29:271–292, 1996.
- [462] L.G. Valiant. A theory of the learnable. *Comm. Assoc. Comput. Mach.*, 27:1134–1142, 1984.
- [463] P.J. van Heerden. A general theory of prediction. Technical report, Polaroid Corporation, Cambridge, Mass., 1963.
- [464] M. van Lambalgen. *Random Sequences*. PhD thesis, Universiteit van Amsterdam, Amsterdam, 1987.
- [465] M. van Lambalgen. Von Mises' definition of random sequences reconsidered. *J. Symbolic Logic*, 52:725–755, 1987.
- [466] M. van Lambalgen. Algorithmic Information Theory. *J. Symbolic Logic*, 54:1389–1400, 1989.
- [467] V.N. Vapnik and A.Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and Its Applications*, 16(2):264–280, 1971.
- [468] J. Ville. *Étude Critique de la Notion de Collectif*. Gauthier-Villars, 1939.
- [469] P.M.B. Vitányi. Randomness. In *Matematica, Logica, Informatica*, Storia del XX Secolo. Instituto della Enciclopedia Italiana, Milan, Italy. To appear.
- [470] P.M.B. Vitányi. On the simulation of many storage heads by one. *Theoret. Comput. Sci.*, 34:157–168, 1984.
- [471] P.M.B. Vitányi. On two-tape real-time computation and queues. *J. Comput. System Sci.*, 29:303–311, 1984.
- [472] P.M.B. Vitányi. An  $n^{1.618}$  lower bound on the time to simulate one queue or two pushdown stores by one tape. *Inform. Process. Lett.*, 21:147–152, 1985.
- [473] P.M.B. Vitányi. Square time is optimal for the simulation of a pushdown store by an oblivious one-head tape unit. *Inform. Process. Lett.*, 21:87–91, 1985.
- [474] P.M.B. Vitányi. Andrei Nikolaevich Kolmogorov. *CWI Quarterly*, 1(2):3–18, June 1988.
- [475] P.M.B. Vitányi. Multiprocessor architectures and physical law. In *Proc. 2nd IEEE Workshop on Physics and Computation*, pages 24–29. IEEE Comput. Soc. Press, 1994.

- [476] P.M.B. Vitányi. Physics and the new computation. In *Proc. 20th Int'l Symp. Math. Found. Comput. Sci.*, volume 969 of *Lect. Notes Comput. Sci.*, pages 106–128. Springer-Verlag, 1995.
- [477] P.M.B. Vitányi and M. Li. Algorithmic arguments in physics of computation. In *Proc. 4th Workshop on Algorithms and Data Structures*, volume 955 of *Lect. Notes Comput. Sci.*, pages 315–333. Springer-Verlag, 1995.
- [478] P.M.B. Vitányi and M. Li. Ideal MDL and its relation to Bayesianism. In D. Dowe, K. Korb, and J. Oliver, editors, *Proc. ISIS: Information, Statistics and Induction in Science Conference*, pages 282–291. World Scientific, 1996.
- [479] P.M.B. Vitányi and M. Li. Minimum description length induction, Bayesianism, and Kolmogorov complexity. Manuscript, CWI, Amsterdam, 1996.
- [480] R. von Mises. Grundlagen der Wahrscheinlichkeitsrechnung. *Mathematische Zeitsch.*, 5:52–99, 1919.
- [481] R. von Mises. *Probability, Statistics and Truth*. Macmillan, 1939. Reprint: Dover, 1981.
- [482] J. von Neumann. Various techniques used in connection with random digits. In A.H. Traub, editor, *John von Neumann, Collected Works, volume V*. Macmillan, 1963.
- [483] V.G. Vovk. Algorithmic information theory and prediction problems. In M.I. Kanovich et al., editor, *Complexity Problems of Mathematical Logic*, pages 21–24. Kalininsk. Gos. Univ., Kalinin, 1985. In Russian.
- [484] V.G. Vovk. The law of the iterated logarithm for random Kolmogorov, or chaotic, sequences. *SIAM Theory Probab. Appl.*, 32(3):413–425, 1987.
- [485] V.G. Vovk. On a randomness criterion. *Soviet Math. Dokl.*, 35:656–660, 1987.
- [486] V.G. Vovk. Prediction of stochastic sequences. *Problems Inform. Transmission*, 25:285–296, 1989.
- [487] V.G. Vovk. Universal forecasting algorithms. *Inform. Comput.*, 96:245–277, 1992.
- [488] V.V. V'yugin. Algorithmic entropy (complexity) of finite objects, and its application to defining randomness and quantity of information. *Semiotika and Informatika*, 16:14–43, 1981. In Russian. Translated into English in: *Selecta Mathematica formerly Sovietica*, 13:4(1994), 357–389.
- [489] V.V. V'yugin. The algebra of invariant properties of binary sequences. *Problems Inform. Transmission*, 18:147–161, 1982.
- [490] V.V. V'yugin. On nonstochastic objects. *Problems Inform. Transmission*, 21:3–9, 1985.
- [491] V.V. V'yugin. On the defect of randomness of a finite object with respect to measures with given complexity bounds. *SIAM Theory Probab. Appl.*, 32:508–512, 1987.
- [492] V.V. V'yugin. Bayesianism: an algorithmic analysis. Manuscript, Independent University of Moscow, Russia, 1994.
- [493] V.V. V'yugin. Ergodic theorem for individual random sequences. Manuscript, Institute for Information Transmission Problems, Russian Academy of Sciences, Russia, 1996.

- [494] A. Wald. Sur la notion de collectif dans la calcul des probabilités. *Comptes Rendus des Séances de l'Académie des Sciences*, 202:1080–1083, 1936.
- [495] A. Wald. Die Widerspruchsfreiheit des Kollektivbegriffes der Wahrscheinlichkeitsrechnung. *Ergebnisse eines mathematischen Kolloquiums*, 8:38–72, 1937.
- [496] C.S. Wallace and D.M. Boulton. An information measure for classification. *Comput. Journal*, 11:185–195, 1968.
- [497] C.S. Wallace and P.R. Freeman. Estimation and inference by compact coding. *J. Royal Stat. Soc.*, 49:240–251, 1987. Discussion: pages 252–265.
- [498] X.-J. Wang. Intermittent fluctuations and complexity. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 319–330. Addison-Wesley, 1991.
- [499] Y. Wang. The law of the iterated logarithm for  $p$ -random sequences. In *Proc. 11th IEEE Conf. Structure in Complexity Theory*, pages 180–189, 1996.
- [500] Y. Wang. Randomness and complexity. Ph.D. Thesis, Heidelberg, 1996.
- [501] O. Watanabe. Comparison of polynomial time completeness notions. *Theoret. Comput. Sci.*, 53:249–265, 1987.
- [502] O. Watanabe, editor. *Kolmogorov Complexity and Computational Complexity*. Springer-Verlag, 1992.
- [503] M. Wax and I. Ziskind. Detection of the number of coherent signals by the MDL principle. *IEEE Trans. Acoust. Speech Signal Process.*, ASSP-37(8):1190–1196, 1989.
- [504] H.S. White. Algorithmic complexity of points in dynamical systems. *Ergodic Theory and Dynamical Systems*, 13:807–830, 1993.
- [505] M. Wiering and J. Schmidhuber. Solving POMDPs with Levin search and EIRA. In *Proc. 13th Intern. Conf. on Machine Learning*, page To appear. Morgan Kaufmann Publishers, 1996.
- [506] D.G. Willis. Computational complexity and probability constructions. *J. Assoc. Comput. Mach.*, 17:241–259, 1970.
- [507] C.H. Woo. Laws and boundary conditions. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 127–135. Addison-Wesley, 1991.
- [508] E.-H. Yang. Universal almost sure data compression for abstract alphabets and arbitrary fidelity criterions. *Probl. Contr. Inform. Theory*, 20(6):397–408, 1991.
- [509] B. Yu. Minimum description length principle: a review. Preprint, Dept. of Stat., Univ. of Calif. Berkeley, 1996.
- [510] I.G. Zhurbenko. *The Spectral Analysis of Time Series*, pages 231–236. Series in Statistics and Probability. North-Holland, 1986. Appendix II: Kolmogorov's algorithm of the Random Number Generator.
- [511] M. Zimand. On the topological size of sets of random strings. *Zeitschr. f. math. Logik und Grundlagen d. Math.*, 32:81–88, 1986.
- [512] M. Zimand. A high-low Kolmogorov complexity law equivalent to the 0-1 law. *Inform. Process. Lett.*, 57:59–64, 1996.
- [513] J. Ziv. On the complexity of an individual sequence. Technical Report 146, Faculty of Electrical Engineering, Technion, Israel, June 1971.

- [514] J. Ziv and A. Lempel. Compression of individual sequences via variable-rate encoding. *IEEE Trans. Inform. Theory*, IT-24:530–536, 1978.
- [515] W.H. Zurek. Algorithmic randomness and physical entropy. *Physical Review, Ser. A*, 40(8):4731–4751, 1989.
- [516] W.H. Zurek. Thermodynamic cost of computation, algorithmic complexity and the information metric. *Nature*, 341:119–124, 1989.
- [517] W.H. Zurek. Algorithmic information content, Church-Turing thesis, physical entropy, and Maxwell’s demon. In W.H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, pages 73–89. Addison-Wesley, 1991.
- [518] W.H. Zurek, editor. *Complexity, Entropy and the Physics of Information*. Addison-Wesley, 1991.
- [519] A.K. Zvonkin and L.A. Levin. The complexity of finite objects and the development of the concepts of information and randomness by means of the theory of algorithms. *Russian Math. Surveys*, 25(6):83–124, 1970.

---

# Index

- $(k, \delta)$ -stochasticity, *see* string,  
 $(k, \delta)$ -stochastic  
 $(n)_k$ : number of variations, 8  
 $0'$ : Turing degree halting problem,  
 220  
 $A^*$ : set of all finite sequences of  
 elements of set  $A$ , 12  
 $A^{=n}$ : set of all words of length  $n$   
 in  $A$ , 476  
 $A^\infty$ : set of one-way infinite  
 sequences over set  $A$ , 13  
 $A^{\leq n}$ : set of words of length  $\leq n$   
 in  $A$ , 112  
 $C$ : complexity, 98  
 $C(x|l(x))$ : length-conditional  $C$ ,  
 111  
 $C(x; n)$ : uniform complexity, 124  
 $C[f(n), t(n), s(n)]$ , 469  
 $C^+$ : monotone upper bound on  
 $C$ -complexity, 207  
 $C^s$ : space-bounded version of  
 $C^{t,s}$ , 462  
 $C^t$ : time-bounded version of  $C^{t,s}$ ,  
 462  
 $C^{t,s}$ : time-space-bounded  $C$   
 complexity, 460  
 $C_r$ : complexity of  $r$ -ary strings,  
 107  
 $C_\phi$ : complexity with respect to  $\phi$ ,  
 97
- $C_n$ : normalized complexity for  
 reals, 125  
 $H$ : entropy stochastic source, 67  
 $I_C(x:y)$ : algorithmic information  
 in  $x$  about  $y$ , 179  
 $K$ : diagonal halting set, 34  
 $K(K(x)|x)$ : complexity of the  
 complexity function,  
 227  
 $K(x|l(x))$ : length-conditional  $K$ ,  
 195  
 $K(x)$ : prefix complexity, 194  
 $K^+$ : monotone upper bound on  
 $K$ -complexity, 224  
 $K^{t,s}$ : time-space-bounded  $K$   
 complexity, 463  
 $K_0$ : halting set, 34  
 $L(x)$ : uniform discrete distribution  
 on  $\mathcal{N}$ , 23  
 $O(f(x))$ : at most of order of  
 magnitude  $f(x)$ , 15  
 $\Delta_i^p$ : class in polynomial hierarchy,  
 40  
 $\Delta_n^0$ : class in arithmetic hierarchy,  
 46  
 $\Gamma_x$ : cylinder generated by  $x$ , 14  
 $l^*(x)$ : optimal universal code-word  
 length, 79  
 $\ell^*(x)$ : lower bound on  $l^*(x)$ , 80

- $\Omega(f(x))$ : at least of order of magnitude  $f(x)$ , 15  
 $\Omega$ : halting probability, 217, 427  
 $\Pi_i^p$ : class in polynomial hierarchy, 40  
 $\Pi_n^0$ : class in arithmetic hierarchy, 46  
 $\Sigma_i^p$ : class in polynomial hierarchy, 40  
 $\Sigma_n^0$ : class in arithmetic hierarchy, 46  
 $\Theta(f(x))$ : of order of magnitude  $f(x)$ , 16  
 $\bar{x}$ : prefix-code  $1^{l(x)}0x$  for  $x$ , 13  
 $\emptyset$ : empty set, 6  
 $\exists$ : there exists, 8  
 $\exists^\infty$ : there exist infinitely many, 8  
 $\forall$ : for all, 8  
 $\forall^\infty$ : for all but finitely many, 8  
 $\lambda(\omega)$ : uniform continuous distribution on  $[0, 1)$ , 23  
 $\langle \cdot \rangle$ : pairing function, 7  
 $\lceil \cdot \rceil$ : ceiling of a number, 8  
 $\lfloor \cdot \rfloor$ : floor of a number, 8  
 $\ln$ : natural logarithm, 8  
log-DNF, 345, 349  
log-decision list, 350  
log: binary logarithm, 8  
 $\log^* x$ : number of terms in  $l^*(x)$ , 79  
 $\omega$ : infinite sequence of elements of  $\mathcal{B}$ , 13  
 $\phi(x) < \infty$ :  $\phi(x)$  converges, 7  
 $\phi(x) = \infty$ :  $\phi(x)$  diverges, 7  
 $\rho_0$ : universal integral test, 214  
 $\sigma$ -algebra, 20  
 $d(A)$ : cardinality of set  $A$ , 6  
 $l(x)$ : length of string  $x$ , 13  
 $n$ -cell, 570  
 $n$ -string, 112, 116, 123–125, 152, 186  
 $n(T)$ : index of  $T$ , 30  
 $o(f(x))$ : asymptotically less than  $f(x)$ , 15  
 $x$ : finite sequence of elements of  $\mathcal{B}$ , 12  
 $x^*$ : first shortest program for  $x$  in enumeration order, 102  
 $x^R$ : reverse of string  $x$ , 12  
 $x_i$ :  $i$ -th letter of  $x$ , 12  
 $x_{1:n}$ : first  $n$  letters of  $x$ , 12  
M: universal enumerable continuous semimeasure, 272  
 $M_{norm}$ : Solomonoff measure, 281, 300, 301  
 $Mc$ : universal enumerable extension semimeasure, 302  
m: universal enumerable discrete semimeasure, 247, 248  
 $\mathcal{B}$ : basic elements, 12, 242  
 $\mathcal{N}$ : the nonnegative integers, 6  
 $\mathcal{Q}$ : the rational numbers, 6  
 $\mathcal{R}$ : the real numbers, 6  
 $\mathcal{Z}$ : the integers, 6  
 $CD[f(n), t(n), s(n)]$ , 471  
 $CD^s$ : space-bounded version of  $CD^{t,s}$ , 462  
 $CD^t$ : time-bounded version of  $CD^{t,s}$ , 462, 463  
 $CD^{t,s}$ : time-space-bounded accepting complexity, 461  
 $CU[f(n), t(n), s(n)]$ , 473  
 $KD^{t,s}$ :  $K$  version of  $CD^{t,s}$ , 463  
 $KM$ : negative logarithm of  $M(x)$ , 282  
Kc: Chaitin's conditional prefix complexity, 235  
Km: monotone complexity, 282  
Kt: Levin-complexity, 502–505  
ic<sup>t</sup>: instance complexity, 495  
 $\binom{n}{k}$ : number of combinations, 9  
 $|\cdot|$ : absolute value of a number, 8  
Aanderaa, S.O., 452, 455  
Abel, N.H., 83  
acceptable numbering, 41, 104  
accepting a language, 37  
Ackermann, W., 45  
Adleman, L., 516, 519  
Agafonov, V.N., 58, 178  
Aleksandrov, P.S., 90  
algorithmic complexity, *see* complexity  
algorithmic complexity theory, vii, 93–238

- algorithmic entropy, *see* entropy, algorithmic  
 algorithmic information theory,  
*see* information theory, algorithmic  
 algorithmic probability theory,  
 239–314  
 Allender, E., 494, 505, 517, 519,  
 520  
 Allison, L., 376  
 Alon, N., 394, 453, 454  
 Andrews, J., 177  
 Angluin, D., 87, 373, 374  
 Anthony, M., 374  
 ants, 583  
 Archimedes, 48  
 arithmetic hierarchy, 46  
 Asarin, E.A., 166  
 Asmis, E., 372  
 asymptotic notation, 15–17  
 Atlan, H., 100, 185  
 average-case  
     adder design, 382–383  
     complexity, 268–272, 307,  
         310, 382–383, 385–389,  
         396–420  
     Heapsort, 412–417  
     longest common subsequence,  
         417–419  
     routing in networks, 408–409  
     shortest common  
         supersequence, 420  
 Avogadro’s number, 558  
 Bachman, P., 15  
 Bacon, F., 585, 589  
 baker’s map, 575  
 Balcázar, J.L., 86, 350, 494, 517,  
 519  
 Baranyai, Zs., 403  
 Barendregt, H.P., 198  
 Barron, A.R., 376  
 Barzdins’s Lemma, 171, 173, 174,  
 187, 224, 230, 426, 427,  
 463, 466, 514  
 Barzdins, J.M., 108, 171, 173,  
 178, 187, 373, 427, 466,  
 472, 516  
 basic element, 12, 242, 272  
 Bayes’s Rule, 20, 19–20, 59, 62,  
 63, 87, 89, 300, 308,  
 309, 319–324, 333, 335,  
 353, 370, 372, 375  
 Bayes, T., 59, 319, 323, 372  
 Bayesian reasoning, 319–323  
 Beigel, R., 452  
 Ben-Amram, A.M., 457  
 Benedek, G., 348, 374  
 Benioff, P.A., 157, 532  
 Bennett, C.H., 176, 187, 218, 238,  
 510, 515, 520, 532, 536,  
 537, 567, 586–588  
 Berman, P., 5, 84, 493  
 Berman-Hartmann conjecture,  
 489  
 Bernoulli process, 59, 63–65, 184,  
 263, 300  
 Bernoulli, J., 59, 63  
 Berry, G.G., 170  
 betting, 262–265  
 Biggs, N., 374  
 binary interval, 254, 283  
 binomial coefficient, 9  
 bit: binary digit, v  
 Blum, M., 476  
 Blumer, A., 349, 374  
 Bogolyubov, N.N., 90  
 Bollobás, B., 454  
 Boltzmann constant, 528, 561  
 Boltzmann, L., 558  
 Bolyai, J., 89  
 Bolyai, W., 89  
 Book, R.V., 457, 494, 517, 518,  
 520  
 Boolean formula, 344, 492  
 Boolean matrix rank, 383  
 Boppana, R., 454  
 Borel, E., 20, 158, 223  
 Borel-Cantelli Lemmas, 64, 151  
 Boswell, J., 239, 307  
 Boulton, D.M., 375  
 Brady, A.H., 46  
 Brebner, G., 396  
 Breitbart, Y., 457  
 Brewer, R.G., 588  
 Briley, B.E., 453  
 Brouwer, L.E.J., 156  
 Brownian computer, 532

- Bruijn sequence, 455  
 Bruijn, N.G. de, 455  
 Buhrman, H.M., 117, 403, 405,  
     409, 453, 454, 485, 494,  
     501, 517, 518  
 Burks, A.W., 382, 452, 586  
  
 Cai, J.-Y., 126, 487  
 Calude, C., 150, 151, 219  
 Cantelli, F.P., 64  
 Cantor, D.G., 392  
 Cantor, G., 40  
 Cardano, 23  
 cardinality, 6, 13  
 Carnap, R., 89, 308, 323  
 Carnot cycle, 555  
 Carnot, N.L.S., 554  
 Cartesian product, 6  
 Castro, J., 350  
 Cauchy-Schwartz Inequality, 457  
 Caves, C.M., 527, 586, 589  
 Chaitin, G.J., 84, 89, 92, 96, 125,  
     152, 179, 185–187, 197,  
     198, 204, 212, 214, 219,  
     220, 222–225, 237, 238,  
     302, 311, 375, 424, 520  
 Champernowne's number,  
     *see* sequence,  
     Champernowne's  
 Champernowne, D.G., 53, 87  
 chaos, 579, 580  
 characteristic function, *see*  
     function, characteristic  
 characteristic sequence, 119, 171,  
     464  
     of  $K_0$ , 173, 230  
     of a language, 423, 426, 427,  
         487, 520  
     of high Turing degree set,  
         176  
     of hyperimmune set, 178  
     of immune set, 177  
     of nonrecursively enumerable  
         set, 172, 174  
     of recursively enumerable set,  
         153, 171, 173, 176, 224,  
         230, 463, 466, 472, 514  
     of semirecursive set, 176  
     random, 224, 494  
  
 Chater, N., 377  
 Cheeseman, P., 376  
 Chernoff bounds, 61, 159, 322,  
     397, 406  
 Chernoff, H., 87  
 Chervonenkis, A.Ya., 374  
 Chitescu, I., 150, 151  
 Chomsky hierarchy, 420  
 Chomsky, N., 309  
 Chrobak, M., 431  
 Church random sequence,  
     *see* sequence,  
     Mises-Wald-Church  
     random  
 Church's Thesis, 24, 29, 53  
 Church, A., 24, 35, 41, 51, 53, 87,  
     149  
 Chvátal, V., 374, 420  
 clause, 344  
 clique, 392, 394  
 CNF formula, 344  
 coarse-graining, 560  
 code  
     additively optimal universal,  
         236  
     ASCII, 72  
     asymptotically optimal  
         universal, 78, 79  
     average word length, 75, 88,  
         191  
     fixed-length, 72  
     Hamming, 117  
     instantaneous, *see* code,  
         prefix-  
     Morse, 66, 70  
     optimal prefix-, 75  
     prefix-, 5, 13, 15, 68, 72,  
         71–84, 88, 191  
     self-delimiting, 76  
     Shannon-Fano, 67, 76, 80,  
         88, 253, 259, 303, 353,  
         375, 513, 523, 527  
     two-part, 99–100, 589  
     uniquely decodable, 72, 75,  
         81, 82  
     universal, 78, 77–80, 82, 88,  
         257  
     variable-length, 72  
 code sequence, 72

- code word, **13**, 71  
 Cohen, P., 92  
 coin-weighing problem, 392  
 collective, 50, **51**, 53, 55, 57, 87,  
     136, 148, 149, 155  
 combination, 8  
 combinatorics, 8–11, 86, 389–396  
 combinatory logic, 198, 237  
 complexity  
     **C**-, **98**, 95–188, 197  
     **C<sup>s</sup>**, **462**, 468  
     **C<sup>t</sup>**, 462–468, 472, 476–488  
     **C<sup>t,s</sup>**, **460**, 460–463, 469  
     **K**-, see complexity, prefix  
     **K<sup>t,s</sup>**, 463, 472  
     **K<sub>μ</sub>-**, 303  
     r-ary strings **C**-, 107  
     **CD<sup>s</sup>**, 462  
     **CD<sup>t</sup>**, **462**, 476–488, 496  
     **CD<sup>t,s</sup>**, **461**, 461–463, 472  
     **Ct**, 505  
     **KD<sup>s</sup>**, 472  
     **KD<sup>t,s</sup>**, 463, 472  
     **KM**-, **282**, 303  
     **Kc**-, **235**, 236, 237  
     **Km**-, see complexity,  
         monotone  
     **Kt**, 502–505  
     additivity of, **101**, 111, 184,  
         189, 194, 229, 231, 233  
     additivity of **C**-, 188  
     additivity of **K**-, **233**, 234  
     additivity of **Kc**-, 235  
     algorithmic, 1, 65, 90, 180  
     alternative approach to  
         define it, 196  
     approximation of, 118, 121  
     average-case, 268–272, 307,  
         310  
     conditional **C**-, **98**, 111  
     conditional **K**-, 194  
     continuity of, 116, 122  
     expected **C**-, 116, **181**,  
         522–528  
     expected **K**-, 181, **231**,  
         522–528  
     extension-, 206  
     fluctuation of, 122  
     instance-, **495**, 495–502  
     length-conditional, 111, 116,  
         120, 123, 153, 154, 186  
     length-conditional **K**-, **195**,  
         204, 207  
     lower bound on **C**-, 119  
     lower bound on **C(x|l(x))**-,  
         **120**, 123  
     lower bound on **K**-, 206  
     majorant of, **241**, 463  
     monotone, 197, 212, 216,  
         **282**, 282–285, 303, 306,  
         311, 312  
     monotonic upper bound on  
         **K**-, 207, **224**, 225, 226  
     monotonicity on prefixes,  
         111, 189, 191, 211  
     noncomputability of, 121,  
         185  
     normalized for real numbers,  
         125  
     number of states-, 177  
     of complexity function, 175,  
         **227**, 226–230, 237, 238  
     of function, 108, **227**  
     prefix, 96, 125, **194**, 189–238,  
         310, 372  
     relation between **C** and **K**,  
         205  
     resource bound hierarchies,  
         469–472  
     resource-bounded, viii, 90,  
         459–520  
     space, 37  
     state-symbol product, **84**,  
         84–85, 89, 92  
     stochastic, vi, 375  
     time, 37  
     time-bounded uniform, 473  
     time-space-bounded, 460–463  
     uniform, **124**, 152, 154, 186,  
         189, 197, 285, 473, 516  
     worst-case space-, 269, 307  
     worst-case time-, 269, 307  
 complexity class  
     **Δ<sub>2</sub><sup>E</sup>**, 493  
     **Δ<sub>i</sub><sup>P</sup>**, 40  
     **Π<sub>i</sub><sup>P</sup>**, **40**, 517  
     **Σ<sub>i</sub><sup>P</sup>**, **40**, 517  
     **#P**, 486

- BPP, 480, 517, 518
- DSPACE, 38, 472, 476, 517
- DTIME, 37, 474, 487, 489, 490, 494, 518
- E, 490, 491, 493, 495, 518
- ESPACE, 517
- EXPTIME, 494
- IC[log, poly], 497
- NE, 493
- NP, 38, 462, 481, 486, 489, 493, 498, 505, 518
- NSPACE, 38, 494
- NTIME, 37
- P, 38, 486, 489, 493, 498, 505, 518
- P/log, 500
- P/poly, 485, 500
- PSPACE, 38, 486, 494
- R, 481
- complexity oscillation, 92, 136, 136–140, 148, 151, 152, 186, 187, 190, 208, 211, 215, 223
  - of  $K$ -, 216, 219–221
  - of  $Km$ -, 216, 311
- computability theory, 24–47, 86
- computable majorants, 463–468
- computational complexity, viii, 36–40, 488–502, 517
- computational learning theory, vii, 6, 339–350
- concatenation, 12
- context-free language, 505
- convergence
  - apparent, of relative frequency, 135
  - recursive, of series, 145, 148, 151, 154, 219
  - regulator of, 152
- Cook, S.A., 518
- Copolov, D.L., 376
- counting method, 390
- Cover, T.M., 88, 89, 91, 105, 134, 185, 204, 206, 208, 237, 300, 302, 314, 372, 373, 376, 526, 527, 586
- crossing sequence, 380, 447
- Csiszar, I., 117
- Culik II, K., 416
- Cutler, C.C., 586
- cylinder, 14, 21, 55, 137, 141, 349
- D’Ariano, G.M., 589
- Dúriš, P., 442, 455
- Daley, R.P., 149, 153, 154, 176, 472–474, 516
- decision list, 349–350
- decision tree, 364–368
- degree of unsolvability, 44
- Dekker, J.C.E., 43, 44
- DeSantis, A., 373
- descriptive complexity, *see* complexity
- Dewdney, A.K., 46
- DFA, *see* finite automaton
- diagonalization method, 34
- Diaz, J., 86
- Dietzfelbinger, M., 440, 441
- dimension
  - Hausdorff, 126
  - topological, 126
  - Vapnik-Chervonenkis, 349
- Diophantine equation, 172, 173, 224, 238
- distance
  - max, 539
  - picture, 541
  - reversible, 544
  - sum, 546
- distribution
  - binomial, 61, 322
  - Bose-Einstein, 11, 260
  - computable universal, 506–509
  - conditional universal, 255
  - Fermi-Dirac, 11, 260
  - malign, 509
  - Maxwell-Boltzmann, 11
  - normal, 362, 364
  - of description length, 202, 205, 238, 257, 266
  - simple, 343
  - uniform, 21, 68, 76, 129, 131
  - uniform continuous, 23
  - uniform discrete, 23, 262
  - universal, 6, 253, 246–280, 307, 319, 523

- universal time-limited, 506–509  
 distribution-free learning, 339–350  
 Dix, T.I., 376  
 DNF formula, 344  
 Doob, J.L., 298, 313  
 Dowe, D.L., 376, 377  
 Drexler, K.E., 587  
 Duns Scotus, John, 62  
 Edgoose, T., 376  
 Edmonds, J.E., 39  
 effective enumeration, 29  
 Ehrenfeucht, A., 349, 374, 420  
 element distinctness problem, 438  
 Elias, P., 56, 82, 88  
 ensemble, 65, 564  
 entropy, 65, 67, 75, 78, 80, 149,  
     180, 181, 184, 187, 190,  
     191, 231, 522–528  
*n*-cell algorithmic, 571  
 algorithmic, 571, 565–583,  
     588  
 Boltzmann, 578  
 classical, 554  
 coarse-grained algorithmic,  
     572, 573, 578  
 complexity, 565, 566  
 conditional, 69, 231  
 Gibbs, 564, 566, 572  
 joint, 231  
 of English, 105  
 of Russian, 88  
 physical, 588  
 relation with complexity,  
     522–528  
 Epicurus, 315, 317, 319, 323, 372  
 Erdős, P., 87, 390, 392, 394, 453,  
     454  
 event, 18  
     certain, 18, 21  
     impossible, 18, 21  
     mutually independent, 20  
     probability of, 18  
 Fano, R.M., 88  
 Feder, M., 377  
 Felker, J.H., 586  
 Feller, W., 10, 11, 22, 23, 63–65,  
     86, 308  
 Ferguson, T.S., 56  
 Fermat, P. de, 23, 172  
 Fermi, E., 89, 588  
 Feynman, R.P., 532  
 Fich, F., 447  
 field, 19  
     Borel, 20  
     Borel extension, 21  
     probability, 19  
 Fine, T.L., 87, 90, 135, 372  
 Finetti, B. de, 372  
 finite automaton  
     *k*-head DFA, 430  
     *k*-pass DFA, 431  
     deterministic (DFA), 317,  
         385, 431  
     nondeterministic (NFA), 385  
     sweeping two-way DFA, 431  
 Fisher, R.A., 369, 370, 377  
 Floyd, R.W., 412, 454  
 Ford, J., 589  
 Fortnow, L., 117, 456, 485,  
     499–501, 517, 518  
 Foulser, D., 420  
 fractal, 125, 126  
 Fredkin gate, 530, 532  
 Fredkin, E., 586, 587  
 Freeman, P.R., 375  
 Freivalds, R.V., 373, 385, 452  
 frequency, 66  
     a priori, 268  
     lower, 267  
     relative, 50  
 frequency interpretation of  
     probability, 50, 87  
 Friedberg, R.A., 42, 44  
 Fu, B., 495  
 function  
     Ackermann, 45, 82, 285  
     additively optimal, *see*  
         function, universal  
     additively optimal partial  
         recursive prefix, 193  
 Busy Beaver, 45, 123, 178,  
     301  
 characteristic, 7, 32, 33, 339,  
     495

- co-enumerable, 35, 303  
 complexity-, 196  
 composition of, 7  
 computable, *see* function,  
     partial recursive  
 consistent, 495  
 convergence of, 7  
 decoding, 13, 71  
 distribution, 22  
 divergence of, 7  
 encoding, 71  
 enumerable, 35, 35–36, 86,  
     108, 240–242, 308  
 factorial, 8, 17  
 generalized exponential, 45  
 honest, 488, 489  
 inverse of, 7  
 many-to-one, 7  
 monotone, 277, 278, 279  
 noncomputable, *see*  
     nonrecursive  
 nonrecursive, 45, 167, 178,  
     226  
 one-to-one, 7  
 pairing, 7  
 parity, 451  
 partial, 7  
 partial recursive, 29  
 partial recursive prefix, 192  
 payoff, 263, 262–265, 296  
 predicate, 29  
 predictor, 58  
 primitive recursive, 83  
 probability, 22, 303  
 probability density, 22  
 ranking, 491  
 real-valued enumerable, 36  
 real-valued recursive, 36  
 recursive, 29, 35–36, 40, 46,  
     53, 108, 126, 240–242,  
     308  
 recursive real, 35  
 regular, 277  
 semicomputable, *see*  
     function, enumerable  
 shape, 126  
 successor, 29  
 total, 7  
 total recursive, *see* function,  
     recursive  
 universal, 95, 96, 97  
 universal co-enumerable, 241  
 universal enumerable, 241  
 universal partial recursive, 31  
 Fundamental Inequality, 356–359  
 Gödel number, 30  
 Gödel numbering, *see* numbering,  
     acceptable  
 Gödel, K., 3, 33, 34, 89, 168, 170,  
     187  
 Gács, P., 91, 107, 108, 134, 166,  
     175, 176, 178, 184, 186,  
     187, 197, 204, 205, 208,  
     210, 223, 226, 230, 235,  
     237, 238, 266, 267, 281,  
     284, 300, 303, 311, 313,  
     314, 373, 517, 519, 520,  
     526, 527, 586–588  
 Gabarró, J., 86  
 Gaifman, H., 157  
 Galil, Z., 430, 439, 442, 455–457  
 Gallager, R.G., 82, 88  
 Gallaire, H., 455  
 Gao, Q., 376  
 garbage bits, 530  
 Gardner, M., 218, 238  
 Garey, M.R., 86  
 Gasarch, W., 175, 176, 452, 453,  
     486  
 Gavaldà, R., 374, 494  
 Gavoille, C., 410, 411  
 Gell-Mann, M., 587, 589  
 generalized Kolmogorov  
     complexity, 459–520  
 genericity, 92  
 Geréb-Graus, M., 431  
 Gibbs, J.W., 564  
 Gnedenko, B.V., 90  
 Gold, E.M., 335, 373  
 Goldbach's Conjecture, 218  
 Goldberg, A., 482, 486, 517  
 Goldstine, H.H., 382, 452  
 Good, I.J., 372  
 Graham, R., 17  
 Gray, R.M., 91, 134, 314, 526,  
     527, 586

- Grossman, J.W., 45  
 Gurevitch, Yu., 519
- Hühne, M., 440  
 Hahn, E.L., 588  
 halting probability, 217, 216–219,  
     221, 223, 237, 238, 242,  
     252, 514  
 halting problem, 33, 33–35, 42,  
     178, 217  
 Hamiltonian equations, 561  
 Hammer, D., 108, 457  
 Hancock, T., 368  
 Hanson, N.R., 318  
 Hardy, G.H., 16  
 Harrison, M.A., 426, 431, 454  
 Hartle, J.B., 589  
 Hartmanis, J., 86, 126, 438, 455,  
     476, 493, 516, 517  
 Haussler, D., 349, 374  
 Heapsort, 412–417, 454  
 Heim, R., 135, 313  
 Hellinger distance, 303  
 Hemachandra, L., 486, 487, 517,  
     519  
 Hemaspaandra, E., 494  
 Hennie, F.C., 452, 455, 460  
 Hermo, M., 519  
 Heyting, A., 156  
 Hilbert’s Tenth Problem, 173  
 Hilbert, D., 45, 172  
 Hoeffding, W., 56  
 Hoepman, J.H., 405, 409, 454  
 Homer, 93  
 Hopcroft, J.E., 454  
 Hume, D., 323  
 Hunter, L., 376  
 Hurewicz, W., 126  
 Huynh, D.T., 486, 488  
 Hwang, K., 453  
 Hästad, J., 456
- Ibarra, O.H., 430, 431  
 immune  
     bi-, 487  
     P/poly, 487  
 incompressibility method, 379–457  
 induction, 315  
     in recursion theory, 335–338
- inductive inference, 315  
 Gold paradigm, 335, **336**,  
     373  
 inductive reasoning, vii, 6, 59, 89,  
     90, 308, 315–372  
     using M, 326–332  
 information, 65  
     algorithmic, 179  
     dispersal of, 118  
     Fisher, 305  
     in  $x$  about  $y$ , 231, 233, 576  
     mutual, 70, 236, 267, 580  
     mutual using  $K$ , 235  
     symmetry of, *see* symmetry  
         of information  
     symmetry of algorithmic,  
         182, 188, 229, 234, 235,  
         238  
 information theory, 48, 65–84, 88,  
     179, 180  
     algorithmic, 179–185, 191,  
         229–237, 522–528  
 instance complexity, *see*  
     complexity, instance-  
 instance complexity conjecture,  
     496  
 invariance theorems, *see* Theorem,  
     Invariance  
 irreversible computation, 528, 529  
 Israeli, A., 456  
 Itai, A., 348, 374  
 Itoh, S., 376
- Jürgenson, H., 219  
 Jaffe, J., 420  
 Jagota, A.K., 310, 520  
 Jakoby, A., 310, 509, 520  
 Jaynes, E.T., 370, 377  
 Jiang, T., 368, 374, 395, 420, 431,  
     443, 454  
 Jockusch, C.G., 176  
 Johnson, D.S., 86, 374  
 Johnson, Dr. Samuel, 239, 307  
 Jones, J.P., 224, 225  
 Joseph, D., 500  
 Juedes, D.W., 518, 520  
 Jurka, J., 377
- König’s Infinity Lemma, 126

- Kahn, J., 453  
 Kalyanasundaram, B., 396  
 Kamae, T., 58, 174, 204  
 Kanefsky, F., 376  
 Kannan, R., 439, 456  
 Kanovich, M.I., 177  
 Karp, R.M., 486  
 Kasami, T., 455  
 Katseff, H.P., 126, 151–153, 186  
 Kearns, M., 374  
 Kemeny, J.G., 372  
 Keuzenkamp, H.A., 377  
 Keyes, R.W., 587  
 Keynes, J.M., 55  
 Khintchin, A.I., 65, 81  
 Kieffer, J.C., 528  
 Kim, C.E., 430  
 Kim, J., 453  
 Kirchherr, W.W., 404  
 Kleene, S.C., 35, 41  
 Knopp, K., 83  
 Knuth, D.E., ix, 16, 17, 86, 187  
 Ko, K.-I., 474, 499, 501, 502, 516, 518  
 Kobayashi, K., 310, 444, 509, 520  
 Kolmogorov Axioms, 18  
 Kolmogorov complexity, *see*  
     complexity, algorithmic  
 Kolmogorov minimal sufficient  
     statistic, 115, 114–115,  
     119, 358  
 Kolmogorov random graphs,  
     396–404  
 Kolmogorov, A.N., 18, 49, 50, 52,  
     55, 56, 65, 70, 86–92,  
     95, 96, 103, 118, 136,  
     149, 150, 166, 185–188,  
     212, 238, 267, 268, 302,  
     306–308, 312, 459, 516,  
     518, 589  
 Kolmogorov-Chaitin ran-  
     domness, *see*  
     complexity, algorithmic  
 Koppel, M., 100, 185  
 Korb, K.B., 377  
 Kraft Inequality, 74, 74–76, 80,  
     82, 83, 88, 191, 195, 202,  
     213, 214, 220, 232, 254  
 Kraft, L.G., 74, 88  
 Kranakis, E., 410, 411, 454  
 Kreinovich, V., 157  
 Krizanc, D., 410, 411, 454  
 Kullback divergence, 328  
 Kummer, M., 117, 174, 485, 500,  
     502, 518  
 Löfgren, L., 185  
 López-Ortiz, A., 438  
 Lambalgen, M. van, 87, 187, 216,  
     220, 221, 311, 453  
 Landauer, R., 528, 532, 586, 587  
 language compression, 477,  
     476–488  
     optimal, 477, 484–485  
     P-rankable, 486  
     probabilistic, 481–484  
     ranking, 477, 484–485  
     with respect to  $C^p$ , 481–484  
     with respect to  $CD^p$ ,  
         477–481  
 Laplace, P.S., 20, 49, 59, 63, 239,  
     300, 307, 372  
 Laplante, S., 456  
 Lathrop, J.I., 520  
 Law  
     0-1, 457  
     Complete Probabilities, 19  
     High-Low Kolmogorov  
         Complexity, 457  
     Infinite Recurrence, 57, 149,  
         306  
     Inverse Weak Law of Large  
         Numbers, 64  
     Iterated Logarithm, 54, 57,  
         65, 87, 140, 149, 223,  
         265, 306  
     of Large Numbers, 54, 63,  
         140, 155, 263, 265  
     of Probability, 263, 265, 295  
     of Randomness, 54, 57, 140  
     of Succession, 63, 300  
     Slow Growth, 514, 516  
     Strong Law of Large  
         Numbers, 64, 306  
     Weak Law of Large Numbers,  
         59, 63, 64  
 learning  
     log-DNF, 345–347

- log-DNF formula, 344
- log-decision list, 350
- by enumeration, 336
- decision list, 349–350
- decision tree, 364–368
- distribution-free, 342
- in the limit, 336
- monotone  $k$ -term DNF, 350
- simple DNF, 349
- under  $M$ , 347–349
- under  $\mathbf{m}$ , 344–347
- under computable distributions, 342–350
- Lecerf, Y., 586
- Leeuw, K. de, 178
- Lemma
  - Barzdins's, *see* Barzdins's Lemma
  - Coding, 479
  - Honesty, 488, 490
  - Jamming, 433
  - KC-Regularity, 422
  - Pumping, 420
  - Switching, 449
- Leung-Yan-Cheong, S.K., 89, 204, 206, 237
- Levin, L.A., 86, 125, 156, 178, 184–188, 196, 197, 209, 212, 230, 235, 238, 281, 283, 300, 301, 303, 310–314, 472, 504, 517–519
- Levine, R.Y., 537, 587
- Levy, P., 311, 313
- Lewis II, P., 476
- Li, M., 187, 272, 310, 344, 348–350, 368, 373–376, 395, 403, 416, 420, 425, 426, 431, 438–440, 447, 452–454, 456, 509, 519, 537, 553, 587, 588
- Likharev, K., 587
- Lindström, B., 392
- Lipton, R.J., 448, 486
- List, B., 526
- literal, 344
- Littlestone, N., 373
- Littlewood, J.E., 16, 87
- Lloyd, S., 587, 589
- logical depth, 510–516
- ( $d, b$ )-deep, 512, 513
- machine independence, 516
- of  $\Omega$ , 514
- shallow, 514
- stability, 516
- longest common subsequence, 417, 417–419
- Longpré, L., 157, 440, 475, 476, 495, 517, 518, 520
- Loui, M.C., 443, 444
- Lovász, L., 374
- Loveland, D.W., 117, 124, 125, 149, 153, 154, 186, 424, 427, 516
- Low, L.H., 376
- lower bounds, 404–457
  - $k$ -PDA, 431
  - $k$ -head automaton, 430
  - $k$ -pass DFA, 431
  - Boolean matrix rank, 383
  - circuit depth, 448–452
  - converting NFA to DFA, 385
  - for Turing machines, 432–444
  - in formal language theory, 420–427
  - multihead automata, 430–431
  - one-tape Turing machine, 380
  - online CFL recognition, 427–430
  - parallel computation, 445–448
  - Ramsey theory, 392
  - routing in networks, 407–408
  - string-matching, 430, 431
  - sweeping two-way DFA, 431
  - VLSI, 447–448
- Luccio, F.L., 411, 454
- Lucretius, 316
- Luginbuhl, D.R., 444
- Luo, Z.Q., 395
- Lutz, J.H., 494, 517–519
- Méré, Chevalier de, 23
- Maass, W., 438–442, 455, 456
- machine
  - $k$ -pushdown store, 438

- k*-queue, 439  
*k*-stack, 438, 442  
*k*-tape, 442  
 deterministic Turing, 28, 326  
 Kolmogorov-Uspensky, 519  
 monotone, 276, 280, 283,  
     309–312  
 nondeterministic 1-tape  
     Turing, 439  
 nondeterministic Turing, 37  
 offline Turing, 440  
 one-way Turing, 432–444  
 online Turing, 432, 432–444  
 oracle Turing, 38  
 PRAM, 445  
 prefix, 193, 310  
 probabilistic Turing, 177,  
     177, 385, 480  
 reference monotone, 280  
 reference prefix, 194  
 reference Turing, 98, 99  
 reversible Turing, 534  
 Turing, 27, 37, 40, 84–86,  
     380  
 two-dimensional tape, 442  
 universal Turing, 30, 84, 185  
 macro state, 558, 560  
 Mahaney, S., 493  
 Mairson, H.G., 457  
 majorant of complexity, *see*  
     complexity, majorant of  
 malign, *see* distribution, malign  
 Mamitsuka, H., 376  
 Mandelbrot, B., 126  
 Margolus, N., 532  
 Markov process, 20, 326  
 Markov’s Inequality, 134, 261,  
     265, 271  
 Markov, A.A., 41, 185  
 Markowsky, G., 373  
 Martin-Löf, P., 54, 92, 113, 127,  
     136, 149, 151, 152, 154,  
     156, 157, 186, 210, 212,  
     313  
 martingale, 296, 311, 313, 518  
 Marxen, H., 46  
 matching, 540  
 Matijasevich, Yu.V., 173, 224, 225  
 Matthew effect, 90  
 Maxwell’s demon, 567–570  
 Maxwell, J.C., 567  
 Mayordomo, E., 519  
 McAleer, M., 377  
 McCarthy, J., 308  
 McGorry, P.D., 376  
 McKenzie, D.P., 376  
 McMillan, B., 82  
 measure, 19, 21, 324  
     constructive, 186  
     continuous, 21  
     countable, 21  
     defective, 308  
     discrete, 21  
     discrete enumerable, 230  
     discrete recursive, 230  
     Lebesgue, *see* measure,  
         uniform  
     of random sequences, 146,  
         220  
     probability, 243  
     recursive, 36, 244, 277, 300  
     recursive continuous,  
         272–280, 282, 304, 332  
     recursive discrete, 245–268  
     resource-bounded, 517  
     simple, 347  
     Solomonoff, 281, 300, 301  
     uniform, 21, 244, 278, 347,  
         468, 517  
     universal enumerable, 347  
     universal recursive, 265  
 Meertens, L., 17, 86  
 Mehlhorn, K., 457  
 Merkle, R.C., 532, 587  
 Meyer auf der Heide, F., 456  
 Meyer, A.R., 125, 424  
 micro state, 558, 559  
 Mills, W.H., 392  
 Milosavljević, A., 377  
 Miltersen, P.B., 310, 509, 520  
 Minsky, M., 31, 85, 90, 308, 309  
 Mises, R. von, 20, 49–51, 53,  
     55–57, 59, 87, 89, 91,  
     134, 136, 148, 151,  
     155–157, 295, 313, 372  
 Mises-Wald-Church random  
     sequence, *see* sequence,

- Mises-Wald-Church random
- Miyano, S., 431
- Mocas, S., 475
- monomial, 344, 349
- monotone  $k$ -term DNF, 350
- Mooers, C., 309
- Moran, S., 411, 456
- Moser, L., 392
- Muchnik, A.A., 44, 150
- Muchnik, An.A., 150, 268
- multinomial coefficient, 10
- Multiplication Rule, 19
- multiplicative domination, 246
- Munro, I., 412, 454
- Naik, A., 518
- Natarajan, B.K., 374, 476
- Nelson, C.G., 430
- Neumann, J. von, ix, 49, 56, 86, 382, 452, 528, 586
- Newman, I., 453
- Newman-Wolfe, R., 442, 453
- Newton, I., v, 9, 23, 317
- NFA, *see* finite automaton
- normalization of semimeasure, *see* semimeasure, normalization of
- NP-complete, 39, 489
- NP-hard, 39
- null set, 140
- $\Pi_n^0$ , 156
  - $\mu$ , 142
  - constructive, 142, 144, 145, 156
  - total recursive, 156
- number
- $\Omega$ -like real, 222, 242
  - arithmetically random real, 222
  - enumerable, *see* sequence, enumerable
  - enumerable real, 300
  - nonrecursive real, 219
  - normal, *see* sequence, normal of Wisdom, 218
  - prime, 4, 17, 32, 83
  - random real, 218, 222, 300
- recursive real, *see* sequence, recursive
- transcendental, 218, 219
- numbering, acceptable, 104
- O'Connor, M.G., 58
- Oberschelp, A., 46
- Occam algorithm, 340
- Occam's razor, v, 62, 240, 252, 299, 317, 318, 331, 340
- occupancy number, 10
- Ockham, William of, 62, 240, 317, 319, 323
- Odifreddi, P., 41, 42, 45, 86, 175
- Oliver, J.J., 376
- oracle, 38, 40, 462, 481–483, 488–492, 494
- Baker-Gill-Solovay, 488
- order
- partial, 7
  - total, 7
- order of magnitude symbols, 17
- Orponen, P., 499, 501, 502, 517, 518
- outcome of experiment, 18, 256, 257, 263, 331
- Ozhegov, S.I., 88
- P-isomorphic, 489
- P-printability, 491
- Pérennès, S., 410, 411
- Péter, R., 45
- Pólya, G., 64
- pac-learning, 339, 339–350, 374
- simple, 339, 342–350, 374
- Paradox
- Bertrand, 316
  - Richard-Berry, 1, 170
  - Russell, 170
- Parberry, I., 456
- Parikh, R., 420
- partition, 10
- Pascal, B., 23
- Patashnik, O., 17
- Patrick, J.D., 376
- Paturi, R., 442, 453
- Paul, W.J., 442, 443, 452, 455
- PDA, *see* pushdown automaton
- Peano Arithmetic, 35

- Pearl, J., 374  
 Pednault, E.P.D., 376  
 Penrose, R., 198, 202  
 Pepys, S., 23  
 permutation, 8  
 perpetuum mobile  
     of the first kind, 554  
     of the second kind, 554  
 Peterson, G., 519  
 Petri, N.V., 178, 472  
 phase space, 559  
 Pierce, J.R., 586  
 Pippenger, N., 393, 395, 444, 453  
 Pitt, L., 350  
 place-selection rule, 87, 153  
     according to Kolmogorov, 56  
     according to Kolmogorov-Loveland, 149  
     according to Mises-Wald-Church, 53, 149, 154  
     according to von Mises, 51, 149  
     finite-state, 58  
     total recursive function, 154  
     Turing machine, 166  
 polynomial complexity core, 497  
 polynomial hierarchy, 40, 494  
 polynomial many-to-one  
     reduction, *see*  
     reducibility, polynomial many-to-one  
 polynomial Turing reduction, *see*  
     reducibility, polynomial Turing  
 Popper, K., 319, 323  
 Post's Problem, 44  
 Post, E.L., 41, 43, 44, 86, 169  
 Pour-El, M.B., 86  
 Pratt, V., 177  
 predicate, 29, 461  
 prediction error, 5, 55, 59, 327, 328, 331, 333, 373  
     expected using M, 304  
 prefix-code, *see* code, prefix-  
 Price, R., 372  
 principle  
     excluded gambling system, 52  
     indifference, 62, 316  
     insufficient reason, 316  
     maximum entropy, 370, 370–372, 377  
     maximum likelihood, 369, 377  
     minimum description length, vi, 90, 115, 351, 351–369, 371, 376, 377, 589  
     minimum message length, 375  
     multiple explanations, 315, 317  
     pigeon-hole, 379  
     simplicity, 62  
 probabilistic communication complexity, 396  
 probabilistic method, 379, 388, 390–392  
 probability  
     a priori, *see* probability, prior algorithmic, 252, 253  
     conditional, 19, 69  
     inferred, *see* probability, posterior  
     inverse, 59  
     joint, 68  
     posterior, 20, 59, 63, 64, 321, 322  
     prior, 20, 59, 63, 319, 321, 322, 325, 333, 370–372  
     prior continuous, 276–280  
     uniform, 268  
     universal prior, 62, 63, 87, 90, 185, 190, 237, 252, 253, 255, 262, 275, 280, 308–311, 319, 335, 511  
 probability theory, 18–23, 86  
 program-size complexity, *see* complexity  
 queue, 439  
 Quicksort, 268  
 Quinlan, J.R., 368, 376  
 Rényi, A., 392, 453  
 Rabin, M.O., 118, 452, 455  
 Rado, T., 45

- Ragde, P., 447, 453  
 Ramsey number, 392  
 random variable, 22  
     (in)dependent, 22  
     continuous, 22  
     discrete, 22  
 randomness deficiency, 103, 113,  
     112–115, 118, 119, 130,  
     132, 185, 210, 259, 267,  
     294, 304, 305, 331,  
     386–388, 397  
 Ratner, M., 88  
 Razborov, A., 456  
 recursion theory, *see*  
     computability theory  
 recursively uniform limit, 122  
 reducibility  
     many-to-one, 43, 169, 177,  
     495  
     one-to-one, 43  
     polynomial many-to-one, 38,  
     493, 518  
     polynomial truth-table, 495  
     polynomial Turing, 38, 492  
     self, 499  
     truth-table, 495  
     Turing, 43, 169, 175, 176,  
     495  
     weak truth-table, 176  
 Regan, K.W., 310, 457, 518, 520  
 regular language, 421  
 Reisch, S., 396, 452  
 Reischuk, R., 310, 442, 443, 455,  
     509, 520  
 relation  
     *n*-ary, 7  
     binary, 7  
     encoding, 71  
 relative frequency stabilization,  
     135  
 resource-bounded Kolmogorov  
     complexity, 459–520  
 reversible  
     ballistic computer, 531–533,  
     586  
     Boolean gate, 530, 536, 587  
     circuit, 530–531, 586  
     computation, 528–537, 586  
 Turing machine, 533–537,  
     586  
 Reznikova, Zh.I., 589  
 Richards, J.I., 86  
 Riemann's Hypothesis, 218  
 Rissanen, J.J., 82, 83, 89, 351,  
     361, 362, 375–377  
 Rivest, R., 350, 368, 374, 376, 430  
 Robinson, R.M., 45  
 Rockford Research, 309  
 Rogers, H., Jr., 41–44, 46, 47, 86,  
     104  
 Rosenberg, A., 430  
 routing in networks, 404–411  
 routing table, 405  
 Rozenberg, G., 420  
 Rubinstein, R., 517  
 Rudich, S., 486  
 Ruelle, D., 589  
 run of zeros, 110  
 Russell, B., 170, 317  
 Russo, D., 517  
 Ryabko, B.Ya., 589  
 Sakoda, W.J., 431  
 sample space, 5, 18  
     continuous, 18, 20  
     discrete, 18  
 Sankoff, D., 420  
 SAT, 39, 486, 489, 490, 492, 493,  
     497, 498, 500, 503  
 satisfiable, 39  
 Savitch, W.J., 494  
 Schöning, U., 499, 501, 502, 518,  
     519  
 Schack, R., 527, 586, 589  
 Schaffer, R., 454  
 Schay, G., 453  
 Schindelhauer, C., 310, 509, 520  
 Schmidhuber, J., 519  
 Schmidt-Goettsch, K., 46  
 Schnitger, G., 396, 440, 442, 452  
 Schnorr's Thesis, 156  
 Schnorr, C.P., 58, 108, 154–156,  
     187, 197, 212, 222, 238,  
     311–313  
 Schumacher, J., 588  
 Sedgewick, R., 448, 454

- Seiferas, J.I., 430, 442, 443, 452, 453, 455  
 self-delimiting code, *see* code, self-delimiting  
 Semenov, A.L., 150, 187, 238  
 semimeasure, **244**, 242–245, 308, 311  
     conditional, 326  
     discrete, **245**, 245–268  
     enumerable, 244  
     enumerable continuous, 272–280, 305  
     enumerable discrete, 245–268  
     extension-, 302  
     maximal, *see* semimeasure, universal  
     maximal enumerable, *see* semimeasure, universal enumerable  
     normalization of, **280**, 300, 301, 310, 311  
     recursive, 244  
     recursive continuous, 304  
     relative enumerable, 268  
     Solomonoff normalization of, **280**, 300, 328  
     universal, **246**, 272  
     universal enumerable, 237, 300, 310, 311  
     universal enumerable conditional, 255  
     universal enumerable continuous, **272**, 272–276, 280, 301, 303  
     universal enumerable discrete, **247**, 248, 253, 255, 265, 353  
     universal relative enumerable, 268  
 sequence  
      $\Delta_2^0$ -definable, 216, 221, **221**  
      $\Pi_n^0$ -random, 156  
      $\infty$ -distributed, *see* sequence, normal  
      $\mu$ -random, *see* sequence, random  
      $k$ -distributed, 58  
     Bernoulli, **59**, 135, 155, 331  
 set  
      $C[f(n), t(n), s(n)]$ , **469**, 469–476, 488–495  
      $Q$ -immune, 471

- m*-complete, 169, 170, 177  
*r*-complete ( $r = 1, m, T$ ), 44  
*r*-hard ( $r = 1, m, T$ ), 44  
 $CD[f(n), t(n), s(n)]$ , 471  
 $CU[f(n), t(n), s(n)]$ , 473  
arithmetic, 222  
Borel, 157, 222  
complete, 230  
continuous, 7  
countable, 7  
cylinder, *see* cylinder  
diagonal halting, 34, 36, 44, 81, 175  
effectively immune, 168, 175  
effectively simple, 168, 175  
empty, 6  
enumerable without repetitions, 43  
exponentially low, 490  
fractal, 126  
halting, 34, 44, 169, 170, 173, 177  
hyperarithmetic, 157  
hyperimmune, 178  
immune, 43, 174, 177  
intuitionistic measure zero, 156  
Kolmogorov, 175  
meager, 112, 186  
P-printable, 491, 505  
RAND, 168, 174  
recursive, 32  
recursively enumerable, 32, 171–173  
recursively enumerable-complete, 489  
relative enumerable, 268  
semirecursive, 176  
simple, 43, 44, 167, 169, 174, 177  
sparse, 186, 472, 482, 485, 491–493  
tally, 491  
Turing complete, 175, 177  
weak truth-table complete, 176  
Shakespeare, W., 105  
Shallit, J., 457  
Shamir, A., 118, 487  
Shannon, C.E., 48, 65, 70, 75, 81, 84, 85, 88, 89, 92, 93, 185, 191, 308, 588  
Shen', A.Kh., 108, 119, 150, 152, 185, 187, 197, 238, 313, 457  
Sherman, A.T., 537, 587  
Shiryaev, A.N., 91, 92  
shortest common supersequence, 417, 420  
shortest program, 102, 110, 204, 205, 208, 235, 236, 256, 282  
prefix complexity of, 204  
Simon, J., 442, 452, 453  
Simons, G., 56  
simple DNF, 349  
Singh, B.S., 376  
Sipser, M., 126, 151–153, 431, 482, 486, 516, 517  
Sivakumar, D., 518  
Skiena, S.S., 86  
Sleator, D., 416  
Slotine, J.J., 589  
Smith, C., 373  
Smoluchowski, M. von, 57  
Snell, J.L., 313  
Snir, M., 157  
Sobolev, S.L., 90  
Solomonoff measure, *see* measure, Solomonoff  
Solomonoff normalization, *see* semimeasure, Solomonoff normalization of  
Solomonoff's induction theory, 324–335  
Solomonoff's Inductive Formula, 331, 372  
Solomonoff, R.J., 59, 62, 87, 89–92, 96, 185, 186, 191, 237, 280, 281, 301–303, 308–311, 323, 324, 335, 372, 519  
Solovay, R.M., 92, 125, 152, 174, 177, 195, 206, 211, 212, 220–222, 226, 238, 267, 301, 311  
sophistication, 100, 185, 589

- source sequence, 72  
 source word, **13**, 71  
 space-energy tradeoff, 537  
 Spencer, J.H., 87, 390, 394, 453,  
     454  
 stack, *see* pushdown store  
 Stanat, D., 420  
 state space, 559  
 statistical properties  
     of graphs, 400–403  
     of sequences, 158–166, 187  
 Stearns, P., 390  
 Stearns, R., 438, 455, 460, 476  
 Stimm, H., 58  
 Stirling's formula, **17**, 67, 134,  
     180, 403  
 Stirling, J., 17  
 stochastic complexity, vi, 375  
 stochastic source, 67  
 string, 12  
     ( $k, \delta$ )-stochastic, 114, 185  
     C-random, 210  
     K-random, 210  
      $\delta$ -random, 118, 166  
     c-incompressible, **109**, 133  
     c-random, 133  
     n-string, *see* n-string  
     absolutely nonrandom, 119  
     binary, 12–15  
     compressibility of, 108  
     cyclic shift, 154, 204  
     empty, 12  
     incompressibility of  
         substring, 110  
     incompressible, 117, 127,  
         135, 203, 208  
     incompressible w.r.t. K, **203**,  
         209  
     infinite, *see* sequence  
     length of, 13  
     random, 90, 113, **133**,  
         127–136, 186, 208–211  
     reverse of, **12**, 101  
     self-delimiting, 13  
     typical, 119  
 submartingale, 296  
     universal enumerable, **296**,  
         295–297  
 Sudborough, H., 430  
 Svetlova, N.D., 88  
 symmetry of information, 233, 385  
 algorithmic, 182, **182**, 229,  
     238  
 algorithmic C-, 188  
 algorithmic K-, **233**, 229–237  
 algorithmic  $K_C$ -, 236  
 for  $K_\mu$ , 303  
 resource-bounded, 474  
 stochastic, **70**, 303  
 Szemerédi, E., 439, 442, 456  
 Szilard engine, 567  
 Szilard, L., 567, 588  
 Tang, S., 494  
 Tarjan, R., 416  
 test, 129, **129**, 134, 186, 311  
     P-, 129  
     Bernoulli, 135, 136  
     Bose-Einstein distribution,  
         260  
     Cauchy's condensation test,  
         83  
     confidence interval of, 128  
     critical region of, **127**, 130  
     Fermi-Dirac distribution,  
         260–261  
     in statistics, 86, **127**, 156  
     integral, 214, **288**  
     Levin's, 223  
     Martin-Löf, 129  
     martingale, 293  
     pseudo randomness, 49  
     ptime-pseudorandom, 468  
     randomness, 295  
     reference universal for  
         uniform distribution,  
             133  
     sequential, **141**, 186, 468  
     sequential  $\mu$ -, 141  
     sequential Bernoulli, 155  
     sequential for uniform  
         distribution, 141  
     sequential Martin-Löf, 141  
     sequential martingale, 294  
     sequential ptime-, 468  
     significance level of, 128  
     Solovay randomness, 221  
     statistical, 263

- sum, 257
- testing for randomness, 128
- universal, 130, 131, 186, 209, 210, 259, 262, 263, 311
- universal Bernoulli, 135
- universal for arbitrary computable  $P$ , 210
- universal for uniform distribution, 132, 210, 259
- universal integral, 214
- universal martingale, 294, 292–295, 297, 305
- universal sequential, 142, 186, 210, 214, 304
- universal sequential Bernoulli, 155
- universal sequential for the uniform measure, 144
- universal sequential martingale, 294
- universal sum, 257, 257–262
- universal uniform, 156
- Thackeray, W.M., 295
- Theorem
  - Basic of RE Sets, 42
  - Binomial, 9
  - Blum Speed-up, 476
  - Chinese Remainder, 118
  - Coding, 253, 253, 256, 257, 266, 267, 279, 282–284, 310
  - Coding, Continuous Version, 282
  - Conditional Coding, 255
  - Entropy Uniqueness, 81
  - Enumeration, 31, 42
  - Equality Stochastic Entropy and Algorithmic Complexity, 180
  - Fermat's Last, 172, 218
  - Fine, 135
  - Hierarchy, 47
  - Incompleteness, 3, 34, 35, 170, 187
  - Incompressibility, 109
  - Invariance, 90–92, 96, 97, 185, 190, 193, 235
  - Invariance for  $K$ , 185
  - Invariance Instance Complexity, 496
  - Invariance Uniform Complexity, 124
  - Invariance, for  $C^{t,s}$ , 460
  - Jones-Matijasevich, 224, 225
  - Kamae, 119, 174, 204
  - KC-Characterization, 423
  - Liouville, 561, 575
  - McMillan-Kraft, 82, 196
  - Myhill-Nerode, 421
  - Noiseless Coding, 75, 77, 78, 88, 191, 257, 259, 355, 523
  - Noncomputability, 120
  - Occam's Razor, 340, 341, 347, 349, 350, 374
  - Prime Number, 17, 83, 478
  - Recursion, 46, 126
  - s-m-n, 41, 41, 42
  - Savitch's, 494
  - Submartingale Convergence, 298, 332
  - Symmetry of Information ( $C$ -version), 182
  - Symmetry of Information ( $K$ -version), 232, 233
  - theory
    - axiomatizable, 34, 168–170, 178
    - consistent, 34
    - decidable, 34
    - sound, 34, 168–170, 178
  - thermodynamics
    - first law, 554
    - of computation, 528–552
    - second law, 554
    - statistical, 558
  - Thomas, J.A., 88, 185, 300, 314, 373
  - Thompson, C.D., 448
  - Thurston, W., 417
  - Tikhomirov, V.M., 90
  - time-energy tradeoff, 553
  - Todt, G., 46
  - Toffoli, T., 586
  - Torenvliet, L., 494
  - tournament, 389, 389–392
  - transitive, 389

- Tromp, J.T., 5, 84, 218, 237, 266, 368, 416, 502, 587  
 Turan, G., 442  
 Turing degree, 176  
 Turing machine, *see* machine, Turing  
 Turing's Thesis, 24  
 Turing, A.M., 24, 33, 41, 86, 185  
 Tyszkiewicz, J., 274, 457  
 Ullman, J.D., 454  
 undecidable statement, 35, 168–170  
 uniform limit, 123  
 upper bounds  
     carry sequence, 382–383  
     combinatorics, 389–396  
     covering families, 395  
     routing in networks, 405–407  
     tournaments, 389  
 Uspensky, V.A., 91, 92, 118, 150, 185, 187, 197, 212, 238, 268, 312, 313, 520  
 USSR, research in former, 186  
 V'yugin, V.V., 119, 185, 186, 212, 267, 306, 308, 312, 453  
 Valiant learning model, 339–350  
 Valiant, L.G., 87, 339, 350, 374, 396, 485, 486  
 Vapnik, V.N., 374  
 variation, 8  
 Vazirani, U., 374, 457  
 Vazirani, V., 457, 485  
 Ville, J., 53, 57, 87, 311, 313  
 Vitányi, P.M.B., 17, 86, 117, 187, 272, 310, 344, 348–350, 373–375, 395, 403, 405, 409, 425, 426, 438, 440, 443, 453–456, 509, 519, 537, 553, 587, 588  
 VLSI complexity, 447  
 Vovk, V.G., 244, 304, 306, 368, 373, 375  
 Wagner, K.W., 494  
 Wald, A., 53, 87, 157  
 Wallace, C.S., 351, 375–377  
 Wallman, H., 126  
 Wang, J., 518  
 Wang, W.G., 520  
 Wang, Y., 516  
 Warmuth, M., 349, 373, 374, 411  
 Watanabe, O., 197, 475, 494, 495, 499, 501, 502, 517, 518, 520  
 Wax, M., 376  
 Wechsung, G., 519  
 Wegman, M., 373  
 Weiss, B., 58  
 Weiss, S., 420  
 Whitehead, A.N., 170  
 Wiering, M., 519  
 Wigderson, A., 447, 453, 456  
 Wiles, A., 172  
 Williams, J.W.J., 412, 454  
 Willis, D.G., 303, 308  
 Wood, D., 416  
 word  
     finite, *see* string  
     infinite, *see* sequence  
 Yamanishi, K., 368, 375, 376  
 Yang, E., 528  
 Yang, Q., 420  
 Yao, A., 430, 456  
 Yee, C.N., 376  
 Yesha, Y., 431, 447, 453, 456  
 Young, P., 500  
 Yu, B., 376  
 Yu, S., 426  
 Zambella, D., 221  
 Zator Company, 89, 309  
 Zeitman, R.Z., 45  
 Zhang, L., 416, 452  
 Zimand, M., 457  
 Zurek, W.H., 527, 532, 586–589  
 Zvonkin, A.K., 86, 125, 178, 184, 186–188, 237, 300, 308, 311, 313, 472, 516