

The
Minimum
Description
Length
Principle

The
Minimum
Description
Length
Principle

Peter D. Grünwald

The MIT Press
Cambridge, Massachusetts
London, England

© 2007 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

Typeset in Palatino by the author using $\LaTeX 2_{\epsilon}$ with C. Manning's `fbook.cls` and `statnlpbook.sty` macros.

Printed and bound in the United States of America.

Library of Congress Cataloging-in-Publication Information

Grünwald, Peter D.

The minimum description length principle / Peter D. Grünwald.

p. cm.—(Adaptive computation and machine learning)

Includes bibliographical references and index.

ISBN-13: 978-0-262-07281-6 (alk. paper)

1. Minimum description length (Information theory) I. Title

QA276.9G78 2007

003'.54—dc22

2006046646

10 9 8 7 6 5 4 3 2 1

To my father

Brief Contents

I	Introductory Material	1
1	<i>Learning, Regularity, and Compression</i>	3
2	<i>Probabilistic and Statistical Preliminaries</i>	41
3	<i>Information-Theoretic Preliminaries</i>	79
4	<i>Information-Theoretic Properties of Statistical Models</i>	109
5	<i>Crude Two-Part Code MDL</i>	131
II	Universal Coding	165
6	<i>Universal Coding with Countable Models</i>	171
7	<i>Parametric Models: Normalized Maximum Likelihood</i>	207
8	<i>Parametric Models: Bayes</i>	231
9	<i>Parametric Models: Prequential Plug-in</i>	257
10	<i>Parametric Models: Two-Part</i>	271
11	<i>NML With Infinite Complexity</i>	295
12	<i>Linear Regression</i>	335
13	<i>Beyond Parametrics</i>	369
III	Refined MDL	403
14	<i>MDL Model Selection</i>	409
15	<i>MDL Prediction and Estimation</i>	459
16	<i>MDL Consistency and Convergence</i>	501
17	<i>MDL in Context</i>	523

IV	Additional Background	597
18	<i>The Exponential or "Maximum Entropy" Families</i>	599
19	<i>Information-Theoretic Properties of Exponential Families</i>	623

Contents

List of Figures xix

Series Foreword xxi

Foreword xxiii

Preface xxv

I Introductory Material 1

1 *Learning, Regularity, and Compression* 3

- 1.1 Regularity and Learning 4
- 1.2 Regularity and Compression 4
- 1.3 Solomonoff's Breakthrough – Kolmogorov Complexity 8
- 1.4 Making the Idea Applicable 10
- 1.5 Crude MDL, Refined MDL and Universal Coding 12
 - 1.5.1 From Crude to Refined MDL 14
 - 1.5.2 Universal Coding and Refined MDL 17
 - 1.5.3 Refined MDL for Model Selection 18
 - 1.5.4 Refined MDL for Prediction and Hypothesis Selection 20
- 1.6 Some Remarks on Model Selection 23
 - 1.6.1 Model Selection among Non-Nested Models 23
 - 1.6.2 Goals of Model vs. Point Hypothesis Selection 25
- 1.7 The MDL Philosophy 26
- 1.8 MDL, Occam's Razor, and the "True Model" 29
 - 1.8.1 Answer to Criticism No. 1 30

1.8.2	Answer to Criticism No. 2	32
1.9	History and Forms of MDL	36
1.9.1	What Is MDL?	37
1.9.2	MDL Literature	38
1.10	Summary and Outlook	40
2	<i>Probabilistic and Statistical Preliminaries</i>	41
2.1	General Mathematical Preliminaries	41
2.2	Probabilistic Preliminaries	46
2.2.1	Definitions; Notational Conventions	46
2.2.2	Probabilistic Sources	53
2.2.3	Limit Theorems and Statements	55
2.2.4	Probabilistic Models	57
2.2.5	Probabilistic Model Classes	60
2.3	Kinds of Probabilistic Models*	62
2.4	Terminological Preliminaries	69
2.5	Modeling Preliminaries:	
	Goals and Methods for Inductive Inference	71
2.5.1	Consistency	71
2.5.2	Basic Concepts of Bayesian Statistics	74
2.6	Summary and Outlook	78
3	<i>Information-Theoretic Preliminaries</i>	79
3.1	Coding Preliminaries	79
3.1.1	Restriction to Prefix Coding Systems; Descriptions as Messages	83
3.1.2	Different Kinds of Codes	86
3.1.3	Assessing the Efficiency of Description Methods	90
3.2	The Most Important Section of This Book: Probabilities and Code Lengths	90
3.2.1	The Kraft Inequality	91
3.2.2	Code Lengths “Are” Probabilities	95
3.2.3	Immediate Insights and Consequences	99
3.3	Probabilities and Code Lengths, Part II	101
3.3.1	(Relative) Entropy and the Information Inequality	103
3.3.2	Uniform Codes, Maximum Entropy, and Minimax Codelength	106
3.4	Summary, Outlook, Further Reading	106

4	<i>Information-Theoretic Properties of Statistical Models</i>	109
4.1	Introduction	109
4.2	Likelihood and <i>Observed</i> Fisher Information	111
4.3	KL Divergence and <i>Expected</i> Fisher Information	117
4.4	Maximum Likelihood: Data vs. Parameters	124
4.5	Summary and Outlook	130
5	<i>Crude Two-Part Code MDL</i>	131
5.1	Introduction: Making Two-Part MDL Precise	132
5.2	Two-Part Code MDL for Markov Chain Selection	133
5.2.1	The Code C_2	135
5.2.2	The Code C_1	137
5.2.3	Crude Two-Part Code MDL for Markov Chains	138
5.3	Simplistic Two-Part Code MDL Hypothesis Selection	139
5.4	Two-Part MDL for Tasks Other Than Hypothesis Selection	141
5.5	Behavior of Two-Part Code MDL	142
5.6	Two-Part Code MDL and Maximum Likelihood	144
5.6.1	The Maximum Likelihood <i>Principle</i>	144
5.6.2	MDL vs. ML	147
5.6.3	MDL as a <i>Maximum Probability Principle</i>	148
5.7	Computing and Approximating Two-Part MDL in Practice	150
5.8	Justifying Crude MDL: Consistency and Code Design	152
5.8.1	A General Consistency Result	153
5.8.2	Code Design for Two-Part Code MDL	157
5.9	Summary and Outlook	163
5.A	Appendix: Proof of Theorem 5.1	163
II	Universal Coding	165
6	<i>Universal Coding with Countable Models</i>	171
6.1	Universal Coding: The Basic Idea	172
6.1.1	Two-Part Codes as Simple Universal Codes	174
6.1.2	From Universal Codes to Universal Models	175
6.1.3	Formal Definition of Universality	177
6.2	The Finite Case	178
6.2.1	Minimax Regret and Normalized ML	179
6.2.2	NML vs. Two-Part vs. Bayes	182
6.3	The Countably Infinite Case	184

6.3.1	The Two-Part and Bayesian Codes	184
6.3.2	The NML Code	187
6.4	Prequential Universal Models	190
6.4.1	Distributions as Prediction Strategies	190
6.4.2	Bayes Is Prequential; NML and Two-part Are Not	193
6.4.3	The Prequential Plug-In Model	197
6.5	Individual vs. Stochastic Universality*	199
6.5.1	Stochastic Redundancy	199
6.5.2	Uniformly Universal Models	201
6.6	Summary, Outlook and Further Reading	204
7	<i>Parametric Models: Normalized Maximum Likelihood</i>	207
7.1	Introduction	207
7.1.1	Preliminaries	208
7.2	Asymptotic Expansion of Parametric Complexity	211
7.3	The Meaning of $\int_{\Theta} \sqrt{\det I(\theta)} d\theta$	216
7.3.1	Complexity and Functional Form	217
7.3.2	KL Divergence and Distinguishability	219
7.3.3	Complexity and Volume	222
7.3.4	Complexity and the Number of Distinguishable Distributions*	224
7.4	Explicit and Simplified Computations	226
8	<i>Parametric Models: Bayes</i>	231
8.1	The Bayesian Regret	231
8.1.1	Basic Interpretation of Theorem 8.1	233
8.2	Bayes Meets Minimax – Jeffreys’ Prior	234
8.2.1	Jeffreys’ Prior and the Boundary	237
8.3	How to Prove the Bayesian and NML Regret Theorems	239
8.3.1	Proof Sketch of Theorem 8.1	239
8.3.2	Beyond Exponential Families	241
8.3.3	Proof Sketch of Theorem 7.1	243
8.4	Stochastic Universality*	244
8.A	Appendix: Proofs of Theorem 8.1 and Theorem 8.2	248
9	<i>Parametric Models: Prequential Plug-in</i>	257
9.1	Prequential Plug-in for Exponential Families	257
9.2	The Plug-in vs. the Bayes Universal Model	262
9.3	More Precise Asymptotics	265

9.4	Summary	269
10	<i>Parametric Models: Two-Part</i>	271
10.1	The Ordinary Two-Part Universal Model	271
10.1.1	Derivation of the Two-Part Code Regret	274
10.1.2	Proof Sketch of Theorem 10.1	277
10.1.3	Discussion	282
10.2	The Conditional Two-Part Universal Code*	284
10.2.1	Conditional Two-Part Codes for Discrete Exponential Families	286
10.2.2	Distinguishability and the Phase Transition*	290
10.3	Summary and Outlook	293
11	<i>NML With Infinite Complexity</i>	295
11.1	Introduction	295
11.1.1	Examples of Undefined NML Distribution	298
11.1.2	Examples of Undefined Jeffreys' Prior	299
11.2	Metauniversal Codes	301
11.2.1	Constrained Parametric Complexity	302
11.2.2	Meta-Two-Part Coding	303
11.2.3	Renormalized Maximum Likelihood*	306
11.3	NML with Luckiness	308
11.3.1	Asymptotic Expansion of LNML	312
11.4	Conditional Universal Models	316
11.4.1	Bayesian Approach with Jeffreys' Prior	317
11.4.2	Conditional NML	320
11.4.3	Liang and Barron's Approach	325
11.5	Summary and Remarks	329
11.A	Appendix: Proof of Theorem 11.4	329
12	<i>Linear Regression</i>	335
12.1	Introduction	336
12.1.1	Prelude: The Normal Location Family	338
12.2	Least-Squares Estimation	340
12.2.1	The Normal Equations	342
12.2.2	Composition of Experiments	345
12.2.3	Penalized Least-Squares	346
12.3	The Linear Model	348
12.3.1	Bayesian Linear Model \mathcal{M}^X with Gaussian Prior	354

12.3.2	Bayesian Linear Models $\mathcal{M}^{\mathbf{X}}$ and $\mathcal{S}^{\mathbf{X}}$ with Noninformative Priors	359
12.4	Universal Models for Linear Regression	363
12.4.1	NML	363
12.4.2	Bayes and LNML	364
12.4.3	Bayes-Jeffreys and CNML	365
13	<i>Beyond Parametrics</i>	369
13.1	Introduction	370
13.2	CUP: Unions of Parametric Models	372
13.2.1	CUP vs. Parametric Models	375
13.3	Universal Codes Based on Histograms	376
13.3.1	Redundancy of Universal CUP Histogram Codes	380
13.4	Nonparametric Redundancy	383
13.4.1	Standard CUP Universal Codes	384
13.4.2	Minimax Nonparametric Redundancy	387
13.5	Gaussian Process Regression*	390
13.5.1	Kernelization of Bayesian Linear Regression	390
13.5.2	Gaussian Processes	394
13.5.3	Gaussian Processes as Universal Models	396
13.6	Conclusion and Further Reading	402
III	Refined MDL	403
14	<i>MDL Model Selection</i>	409
14.1	Introduction	409
14.2	Simple Refined MDL Model Selection	411
14.2.1	Compression Interpretation	415
14.2.2	Counting Interpretation	416
14.2.3	Bayesian Interpretation	418
14.2.4	Prequential Interpretation	419
14.3	General Parametric Model Selection	420
14.3.1	Models with Infinite Complexities	420
14.3.2	Comparing Many or Infinitely Many Models	422
14.3.3	The General Picture	425
14.4	Practical Issues in MDL Model Selection	428
14.4.1	Calculating Universal Codelengths	428

14.4.2	Computational Efficiency and Practical Quality of Non-NML Universal Codes	429
14.4.3	Model Selection with Conditional NML and Plug-in Codes	431
14.4.4	General Warnings about Model Selection	435
14.5	MDL Model Selection for Linear Regression	438
14.5.1	Rissanen's RNML Approach	439
14.5.2	Hansen and Yu's gMDL Approach	443
14.5.3	Liang and Barron's Approach	446
14.5.4	Discussion	448
14.6	Worst Case vs. Average Case*	451
15	<i>MDL Prediction and Estimation</i>	459
15.1	Introduction	459
15.2	MDL for Prediction and Predictive Estimation	460
15.2.1	Prequential MDL Estimators	461
15.2.2	Prequential MDL Estimators Are Consistent	465
15.2.3	Parametric and Nonparametric Examples	469
15.2.4	Césaro KL consistency vs. KL consistency*	472
15.3	Two-Part Code MDL for Point Hypothesis Selection	476
15.3.1	Discussion of Two-Part Consistency Theorem	478
15.4	MDL Parameter Estimation	483
15.4.1	MDL Estimators vs. Luckiness ML Estimators	487
15.4.2	What Estimator To Use?	491
15.4.3	Comparison to Bayesian Estimators*	493
15.5	Summary and Outlook	498
15.A	Appendix: Proof of Theorem 15.3	499
16	<i>MDL Consistency and Convergence</i>	501
16.1	Introduction	501
16.1.1	The Scenarios Considered	501
16.2	Consistency: Prequential and Two-Part MDL Estimators	502
16.3	Consistency: MDL Model Selection	505
16.3.1	Selection between a Union of Parametric Models	505
16.3.2	Nonparametric Model Selection Based on CUP Model Class	508
16.4	MDL Consistency Peculiarities	511
16.5	Risks and Rates	515

16.5.1	Relations between Divergences and Risk Measures	517
16.5.2	Minimax Rates	519
16.6	MDL Rates of Convergence	520
16.6.1	Prequential and Two-Part MDL Estimators	520
16.6.2	MDL Model Selection	522
17	<i>MDL in Context</i>	523
17.1	MDL and Frequentist Paradigms	524
17.1.1	Sanity Check or Design Principle?	525
17.1.2	The Weak Prequential Principle	528
17.1.3	MDL vs. Frequentist Principles: Remaining Issues	529
17.2	MDL and Bayesian Inference	531
17.2.1	Luckiness Functions vs. Prior Distributions	534
17.2.2	MDL, Bayes, and Occam	539
17.2.3	MDL and Brands of Bayesian Statistics	544
17.2.4	Conclusion: a Common Future after All?	548
17.3	MDL, AIC and BIC	549
17.3.1	BIC	549
17.3.2	AIC	550
17.3.3	Combining the Best of AIC and BIC	552
17.4	MDL and MML	555
17.4.1	Strict Minimum Message Length	556
17.4.2	Comparison to MDL	558
17.4.3	The Wallace-Freeman Estimator	560
17.5	MDL and Prequential Analysis	562
17.6	MDL and Cross-Validation	565
17.7	MDL and Maximum Entropy	567
17.8	Kolmogorov Complexity and Structure Function	570
17.9	MDL and Individual Sequence Prediction	573
17.10	MDL and Statistical Learning Theory	579
17.10.1	Structural Risk Minimization	581
17.10.2	PAC-Bayesian Approaches	585
17.10.3	PAC-Bayes and MDL	588
17.11	The Road Ahead	592
IV	Additional Background	597
18	<i>The Exponential or "Maximum Entropy" Families</i>	599

18.1	Introduction	600
18.2	Definition and Overview	601
18.3	Basic Properties	605
18.4	Mean-Value, Canonical, and Other Parameterizations	609
18.4.1	The Mean Value Parameterization	609
18.4.2	Other Parameterizations	611
18.4.3	Relating Mean-Value and Canonical Parameters**	613
18.5	Exponential Families of General Probabilistic Sources*	617
18.6	Fisher Information Definitions and Characterizations*	619
19	Information-Theoretic Properties of Exponential Families	623
19.1	Introduction	624
19.2	Robustness of Exponential Family Codes	624
19.2.1	If Θ_{mean} Does Not Contain the Mean**	627
19.3	Behavior <i>at</i> the ML Estimate $\hat{\beta}$	629
19.4	Behavior <i>of</i> the ML Estimate $\hat{\beta}$	632
19.4.1	Central Limit Theorem	633
19.4.2	Large Deviations	634
19.5	Maximum Entropy and Minimax Codelength	637
19.5.1	Exponential Families and Maximum Entropy	638
19.5.2	Exponential Families and Minimax Codelength	641
19.5.3	The Compression Game	643
19.6	Likelihood Ratio Families and Rényi Divergences*	645
19.6.1	The Likelihood Ratio Family	647
19.7	Summary	650
	References	651
	List of Symbols	675
	Subject Index	679

List of Figures

1.1	A simple, a complex and a tradeoff (third-degree) polynomial.	13
1.2	Models and Model Classes; (Point) Hypotheses.	15
3.1	Coding systems, codes and description methods as defined in this book. MDL is only concerned with nonsingular codes/coding systems, allowing for lossless coding.	81
3.2	Binary code tree for the Kraft inequality using alphabet $\{a, b, c\}$ and code $C_0(a) = 0; C_0(b) = 10; C_0(c) = 11$.	93
3.3	The most important observation of this book.	96
3.4	The third most important observation of this book.	107
4.1	The horizontal axis represents $\theta - \hat{\theta}(x^n)$ as a function of θ for a particular, fixed, x^n . The vertical axis represents $P_\theta(x^n)$. The function achieves its maximum at $\theta = \hat{\theta}$ and, near the maximum, has the shape of a Gaussian.	114
4.2	$I(\theta)$ as a function of θ for the Bernoulli model.	122
4.3	The horizontal axis represents θ . The vertical axis represents $D(\theta^* \theta)$ (solid thick line), $D(\theta \theta^*)$ (solid thin line), and $0.5(\theta - \theta^*)^2 I(\theta^*)$ (dotted line). In the upper picture, $\theta^* = 0.5$. In the lower picture, $\theta^* = 0.9$.	123

4.4	The top graph shows the negative log-likelihood $-\nu \ln \theta - (1 - \nu) \ln(1 - \theta)$ as a function of θ , where ν represents n_1/n . The graph shows the cases $\nu = 0.5$ (line that is lowest on the left, highest on the right), $\nu = 0.7$ (solid middle line), and $\nu = 0.9$. Note that for $\nu = 0.5$, θ achieves its minimum at 0.5: $\hat{\theta} = \nu$. Similarly for $\nu = 0.7$, $\hat{\theta} = 0.7$, and for $\nu = 0.9$, $\hat{\theta} = 0.9$. Nevertheless, $\theta = 0.7$ assigns a smaller description length to data with $\nu = 0.9$ than to data with $\nu = 0.7$. This is further illustrated in the bottom graph, which shows the negative log-likelihood $-\nu \ln \theta - (1 - \nu) \ln(1 - \theta)$ as a function of ν , for $\theta = 0.5$, $\theta = 0.7$, and $\theta = 0.9$. The corresponding functions are obviously linear. Note that we depict minus log rather than direct likelihoods here, which explains the difference in form between the top figure and the graph in Figure 4.1.	126
7.1	The crazy Bernoulli model.	218
10.1	The structure of the discretization for the case $\Theta \subset \mathbb{R}^2$. The picture shows a single "large" hypercube S containing some "small" hyperrectangles R . The discretized points are the centers of the rectangles, if the rectangles lie completely inside S . Otherwise they are the closest points to the center that still lies within S . For the S that is shown here, the angle between the small and the large grid is 30 degrees; for other "large" S , the angle of the $R \subset S$ will be different. The arrows point in the direction of the eigenvectors of the Fisher information matrix. The length of the arrows is proportional to the square root of the inverse of the eigenvalues.	278
14.1	The Refined MDL Principle for Model Selection	426
14.2	Ignoring codelengths.	439
17.1	Rissanen's MDL, Wallace's MML and Dawid's Prequential Approach	562

Series Foreword

The goal of building systems that can adapt to their environments and learn from their experience has attracted researchers from many fields, including computer science, engineering, mathematics, physics, neuroscience, and cognitive science. Out of this research has come a wide variety of learning techniques that have the potential to transform many scientific and industrial fields. Recently, several research communities have converged on a common set of issues surrounding supervised, unsupervised and reinforcement learning problems. The MIT Press series on Adaptive Computation and Machine Learning seeks to unify the many diverse strands of machine learning research and to foster high-quality research and innovative applications.

Thomas Dietterich

Foreword

This is a splendid account of the latest developments on the minimum description length (MDL) principle and the related theory of stochastic complexity. The MDL principle seeks to place the age-old statistical or inductive inference on a sound foundation. In order to achieve this it requires the drastically different and, for many, unpalatable view that the objective is not to estimate any “true” data-generating mechanism but simply to find a good explanation of data, technically called a model. The author gives an impassionate balanced discussion of the deep philosophical implications of the principle, and he traces the tortuous path from the roots to the current refined stage of the principle, in which the idea of a *universal* model plays a central role. This is a model that allows for an objective comparison of alternative models regardless of their form or number of parameters in case the interest is in model selection. Further, it provides a basis for prediction and classification.

The author describes painstakingly the information- and probability-theoretic notions needed for the reader with a minimum of prerequisites to apply the principle to a variety of statistical problems. This involves an in-depth treatment of the theory of “universal models,” which in its general form is deep and complex. The author’s treatment of it, however, is highly accessible. He achieves this by devoting an extensive section on discussing finite universal models, which are much simpler than the general case but do serve to illustrate the general ideas.

Based on this treatment, he then introduces the MDL principle in its modern, refined form, always emphasizing the ideas that give rise to the actual formulas. He starts out with the simple case of comparing a finite number of parametric models, and gradually builds up the theory to general problems of model selection. He also briefly discusses parameter estimation and

nonparametric inference. For the reader with deeper statistical knowledge, in Chapter 17 he compares MDL to some other more customary statistical techniques.

Jorma Rissanen
Helsinki Institute for Information Technology
Helsinki, Finland
December 2005

Preface

How does one decide among competing explanations of data given limited observations? This is the problem of *model selection*. A central concern in model selection is the danger of *overfitting*: the selection of an overly complex model that, while fitting observed data very well, predicts future data very badly. Overfitting is one of the most important issues in inductive and statistical inference: besides model selection, it also pervades applications such as prediction, pattern classification and parameter estimation.

The minimum description length (MDL) principle is a relatively recent method for inductive inference that provides a generic solution to the model selection problem, and, more generally, to the overfitting problem. MDL is based on the following insight: any regularity in the data can be used to *compress* the data, i.e. to describe it using fewer symbols than the number of symbols needed to describe the data literally. The more regularities there are, the more the data can be compressed. Equating “learning” with “finding regularity,” we can therefore say that the more we are able to compress the data, the more we have *learned* about the data. Formalizing this idea leads to a general theory of inductive inference with several attractive properties:

- 1. Occam’s razor.** MDL chooses a model that trades off goodness-of-fit on the observed data with “complexity” or “richness” of the model. As such, MDL embodies a form of Occam’s razor, a principle that is both intuitively appealing and informally applied throughout all the sciences.
- 2. No overfitting, automatically.** MDL methods *automatically* and *inherently* protect against overfitting and can be used to estimate both the parameters and the structure (e.g., number of parameters) of a model. In contrast, to avoid overfitting when estimating the structure of a model, traditional

methods such as maximum likelihood must be *modified* and extended with additional, typically ad hoc principles.

3. **Bayesian interpretation.** Some (not all) MDL procedures are closely related to Bayesian inference. Yet they avoid some of the interpretation difficulties of the Bayesian approach, especially in the realistic case when it is known a priori to the modeler that none of the models under consideration is true. In fact:
4. **No need for “underlying truth.”** In contrast to other statistical methods, MDL procedures have a clear interpretation independent of whether or not there exists some underlying “true” model.
5. **Predictive interpretation.** Because data compression is formally equivalent to a form of probabilistic prediction, MDL methods can be interpreted as searching for a model with good predictive performance on *unseen* data. This makes MDL related to, yet different from, data-oriented model selection techniques such as cross-validation.

This Book

This book provides an extensive, step-by-step introduction to the MDL principle, with an emphasis on conceptual issues. From the many talks that I have given on the subject, I have noticed that the same questions about MDL pop up over and over again. Often, the corresponding answers can be found only — if at all — in highly technical journal articles. The main aim of this book is to serve as a reference guide, in which such answers can be found in a much more accessible form. There seems to be a real need for such an exposition because, quoting Lanterman (2001), of “the challenging nature of the original works and the preponderance of misinterpretations and misunderstandings in the applied literature.” Correcting such misunderstandings is the second main aim of this book.

First Aim: Accessibility I first learned about MDL in 1993, just before finishing my master’s in computer science. As such, I knew some basic probability theory and linear algebra, but I knew next to nothing about advanced measure-theoretic probability, statistics, and information theory. To my surprise, I found that to access the MDL literature, I needed substantial knowledge about all three subjects! This experience has had a profound influence on this book: in a way, I wanted to write a book which I would have been

able to understand when I was a beginning graduate student. Therefore, since with some difficulty its use can be avoided, there is no measure theory whatsoever in this book. On the other hand, this book is full of statistics and information theory, since these are essential to any understanding of MDL. Still, both subjects are introduced at a very basic level in Part I of the book, which provides an initial introduction to MDL. At least this part of the book should be readable without any prior exposure to statistics or information theory.

If my main aim has succeeded, then this book should be accessible to (a) researchers from the diverse areas dealing with inductive inference, such as statistics, pattern classification, and branches of computer science such as machine learning and data mining; (b) researchers from biology, econometrics, experimental psychology, and other applied sciences that frequently have to deal with inductive inference, especially model selection; and (c) philosophers interested in the foundations of inductive inference. This book should enable such readers to understand what MDL is, how it can be used, and what it does.

Second Aim: A Coherent, Detailed Overview In the year 2000, when I first thought about writing this book, the field had just witnessed a number of advances and breakthroughs, involving the so-called *normalized maximum likelihood code*. These advances had not received much attention outside of a very small research community; most practical applications and assessments of MDL were based on “old” (early 1980s) methods and ideas. At the time, some pervasive myths were that “MDL is just two-part coding”, “MDL is BIC” (an asymptotic Bayesian method for model selection), or “MDL is just Bayes.” This prompted me and several other researchers to write papers and give talks about the new ideas, related to the normalized maximum likelihood. Unfortunately, this may have had somewhat of an adverse effect: I now frequently talk to people who think that MDL is just “normalized maximum likelihood coding.” This is just as much of a myth as the earlier ones! In reality, MDL in its modern form is based on a general notion known in the information-theoretic literature as *universal coding*. There exist many types of universal codes, the main four types being the Bayesian, two-part, normalized maximum likelihood, and prequential plug-in codes. All of these can be used in MDL inference, and which one to use depends on the application at hand. While this emphasis on universal codes is already present in the overview (Barron, Rissanen, and Yu 1998), their paper requires substan-

tial knowledge of information theory and statistics. With this book, I hope to make the universal coding-based MDL theory accessible to a much wider audience.

A Guide for the Reader

This book consists of four parts. Part I is really almost a separate book. It provides a very basic introduction to MDL, as well as an introductory overview of the statistical and information-theoretic concepts needed to understand MDL. Part II is entirely devoted to universal coding, the information-theoretic notion on which MDL is built. Universal coding is really a theory about data compression. It is easiest to introduce without directly connecting it to inductive inference, and this is the way we treat it in Part II. In fact though, there is a very strong relation between universal coding and inductive inference. This connection is formalized in Part III, where we give a detailed treatment of MDL theory as a theory of inductive inference based on universal coding. Part IV can once again be read separately, providing an overview of the statistical theory of *exponential families*. It provides background knowledge needed in the proofs of theorems in Part II.

The Fast Track — How to Avoid Reading Most of This Book I do not suppose that any reader will find the time to read all four parts in detail. Indeed, for readers with prior exposure to MDL, this book may serve more like a reference guide than an introduction in itself. For the benefit of readers with no such prior knowledge, each chapter in Part I and Part II starts with a brief list of its contents as well as a *fast track*-paragraph, which indicates the parts that should definitely be read, and the parts that can be skipped at first reading. This allows a “fast track” through Part I and Part II, so that the reader can quickly reach Part III, which treats state-of-the-art MDL inference. Additionally, some sections are marked with an asterisk (*). Such sections contain advanced material and may certainly be skipped at first reading.

Also, the reader will frequently find paragraphs such as the present one, which are set in smaller font. These provide additional, more detailed discussion of the issues arising in the main text, and may also be skipped at first reading.

Also, at several places, the reader will find boxes like the one below:

Boxes Contain the Most Important Ideas

Each chapter contains several boxes like this one. These contain the most important insights. Together, they form a summary of the chapter.

To further benefit the hurried reader, we now give a brief overview of each part:

Part I Chapter 1 discusses some of the basic ideas underlying MDL in a mostly nonmathematical manner. Chapter 2 briefly reviews general mathematical and probabilistic preliminaries. Chapter 3 gives a detailed discussion of some essential information-theoretic ideas. Chapter 4 applies these notions to statistical models. This chapter gives an extensive analysis of the log-likelihood function and its expectation. It may be of interest for teachers of introductory statistics, since the treatment emphasizes some, in my view, quite important aspects usually not considered in statistics textbooks. For example, we consider in detail what happens if we vary the data, rather than the parameters. Chapter 5 then gives a first mathematically precise implementation of MDL. This is the so-called crude two-part code MDL. I call it “crude” because it is suboptimal, and not explicitly based on universal coding. I included it because it is easy to explain — especially the fact that it has obvious defects raises some serious questions, and thinking about these questions seems the perfect introduction to the “refined” MDL that we introduce in Part III of the book.

Although some basic familiarity with elementary probability theory is assumed throughout the text, all probabilistic concepts needed are briefly reviewed in Chapter 2. They are typically taught in undergraduate courses and can be found in books such as (Ross 1998). Strictly speaking, the text can be read without any prior knowledge of statistics or information theory — all concepts and ideas are introduced in Chapters 3 and 4. Nevertheless, some prior exposure to these subjects is probably needed to fully appreciate the developments in Part II and Part III. More extensive introductions to the statistical concepts needed can be found in, for example (Bain and Engelhardt 1989; Casella and Berger ; Rice 1995).

Part II Part II then treats the general theory of universal coding, with an emphasis on issues that are relevant to MDL. It starts with a brief introduction which gives a high-level overview of the chapters contained in Part II. Its first chapter, Chapter 6, then contains a detailed introduction to the main

ideas, in the restricted context of countable model classes. Each of the four subsequent chapters gives a detailed discussion of one of the four main types of universal codes, in the still restricted context of “parametric models” with (essentially) compact parameter spaces. Chapters 11, 12, and 13 deal with general parametric models — including linear regression models — as well as nonparametric models.

Part III Part III gives a detailed treatment of refined MDL. We call it “refined” so as to mark the contrast with the “crude” form of MDL of Chapter 5. It starts with a brief introduction which gives a high-level overview of refined MDL. Chapter 14 deals with refined MDL for model selection. Chapter 15 is about its other two main applications: hypothesis selection (a basis for parametric and nonparametric density estimation) and prediction. Consistency and rate-of-convergence results for refined MDL are detailed in Chapter 16. Refined MDL is placed in its proper context in Chapter 17, in which we discuss its underlying philosophy and compare it to various other approaches.

Compared to Part I, Part II and Part III contain more advanced material, and some prior exposure to statistics may be needed to fully appreciate the developments. Still, all required information-theoretic concepts — invariably related to *universal coding* — are once again discussed at a very basic level. These parts of the book mainly serve as a reference guide, providing a detailed exposition of the main topics in MDL inference. The discussion of each topic includes details which are often left open in the existing literature, but which are important when devising practical applications of MDL. When pondering these details, I noticed that there are several open questions in MDL theory which previously have not been explicitly posed. We explicitly list and number such open questions in Part II and Part III. These parts also contain several new developments: in order to tell a coherent story about MDL, I provide some new results — not published elsewhere — that connect various notions devised by different authors.

The main innovations are the “distinguishability” interpretation of MDL for finite models in Chapter 6, the “phase transition” view on two-part coding in Chapter 10, the luckiness framework as well as the CNML-1 and CNML-2 extensions of the normalized maximum likelihood code in Chapter 11, and the connections between Césaro and standard KL risk and the use of redundancy rather than resolvability in the convergence theorem for two-part MDL in Chapter 15.

I also found it useful to rephrase and re-prove existing mathematical theorems in a unified way. The many theorems in Part II and Part III usually express results that are similar to existing theorems by various authors, mainly Andrew Barron, Jorma Rissanen, and Bin Yu. Since these theorems were often stated in slightly different contexts, they are hard to compare. In our version, they become easily comparable. Specifically, in Part II, we restrict the treatment to so-called *exponential families* of distributions, which is a weakening of existing results. Yet, the theorems invariably deal with uniform convergence, which is often a strengthening of existing results.

Part IV: Exponential Family Theory The theorems in Part II make heavy use of the general and beautiful theory of *exponential* or, relatedly, *maximum entropy* families of probability distributions. Part IV is an appendix that contains an overview of these families and their mathematical properties. When writing the book, I found that most existing treatments are much too restricted to contain the results that we need in this book. The only general treatments I am aware of (Barndorff-Nielsen 1978; Brown 1986) use measure theory, and give a detailed treatment of behavior at parameters tending to the boundaries of the parameter space. For this reason, they are quite hard to follow. Thus, I decided to write my own overview, which avoids measure theory and boundary issues, but otherwise contains most essential ideas such as sufficiency, mean-value and canonical parameterizations, duality, and maximum entropy interpretations.

Acknowledgments

Tim van Erven, Peter Harremoës, Wouter Koolen, In Jae Myung, Mark Pitt, Teemu Roos, Steven de Rooij, and Tomi Silander read and commented on parts of this text. I would especially like to thank Tim, who provided comments on the entire manuscript.

Mistakes

Of course, the many mistakes which undoubtedly remain in this text are all my (the author's) sole responsibility. I welcome all emails that point out mistakes in the text!

Among those who have helped shape my views on statistical inference, two people stand out: Phil Dawid and Jorma Rissanen. Other people who have

strongly influenced my thinking on these matters are Vijay Balasubramanian, Andrew Barron, Richard Gill, Teemu Roos, Paul Vitányi, Volodya Vovk, and Eric-Jan Wagenmakers. My wife Louise de Rooij made a very visible and colourful contribution. Among the many other people who in some way or other had an impact on this book I should mention Petri Myllymäki, Henry Tirri, Richard Shiffrin, Johan van Benthem, and, last but not least, Herbert, Christa and Wiske Grünwald. As leaders of our research group at CWI (the National Research Institute for Mathematics and Computer Science in the Netherlands), Harry Buhrman and Paul Vitányi provided the pleasant working environment in which this book could be written. The initial parts of this book were written in 2001, while I was visiting the University of California at Santa Cruz. I would like to thank Manfred Warmuth and David Draper for hosting me. Finally and most importantly, I would like to thank my lovely wife Louise for putting up with my foolishness for so long.

PART I

Introductory Material

1 *Learning, Regularity, and Compression*

Overview The task of inductive inference is to find laws or regularities underlying some given set of data. These laws are then used to gain insight into the data or to classify or predict future data. The minimum description length (MDL) principle is a general method for inductive inference, based on the idea that the more we are able to *compress* (describe in a compact manner) a set of data, the more regularities we have found in it and therefore, the more we have *learned* from the data. In this chapter we give a first, preliminary and informal introduction to this principle.

Contents In Sections 1.1 and 1.2 we discuss some of the fundamental ideas relating description length and regularity. In Section 1.3 we describe what was historically the first attempt to formalize these ideas. In Section 1.4 we explain the problems with using the original formalization in practice, and indicate what must be done to make the ideas practicable. Section 1.5 introduces the practical forms of MDL we deal with in this book, as well as the crucial concept of “universal coding.” Section 1.6 deals with some issues concerning *model selection*, which is one of the main MDL applications. The philosophy underlying MDL is discussed in Section 1.7. Section 1.8 shows how the ideas behind MDL are related to “Occam’s razor.” We end in Section 1.9 with a brief historical overview of the field and its literature.

Fast Track This chapter discusses, in an informal manner, several of the complicated issues we will deal with in this book. It is therefore essential for readers without prior exposure to MDL. Readers who are familiar with the basic ideas behind MDL may just want to look at the boxes.

1.1 Regularity and Learning

We are interested in developing a method for *learning* the laws and regularities in data. The following example will illustrate what we mean by this and give a first idea of how it can be related to descriptions of data.

Example 1.1 We start by considering binary data. Consider the following three sequences. We assume that each sequence is 10000 bits long, and we just list the beginning and the end of each sequence.

00010001000100010001...000100010001000100010001 (1.1)

01110100110100100110...1010111010111011000101100010 (1.2)

00011000001010100000...0010001000010000001000110000 (1.3)

The first of these three sequences is a 2500-fold repetition of 0001. Intuitively, the sequence looks regular; there seems to be a simple “law” underlying it; it might make sense to conjecture that future data will also be subject to this law, and to predict that future data will behave according to this law. The second sequence has been generated by tosses of a fair coin. It is, intuitively speaking, as “random as possible,” and in this sense there is no regularity underlying it.¹ Indeed, we cannot seem to find such a regularity either when we look at the data. The third sequence contains exactly four times as many 0s as 1s. It looks less regular, more random than the first; but it looks less random than the second. There is still some discernible regularity in these data, but of a statistical rather than of a deterministic kind. Again, noticing that such a regularity is there and predicting that future data will behave according to the same regularity seems sensible.

1.2 Regularity and Compression

What do we mean by a “regularity”? The fundamental idea behind the MDL principle is the following insight: every regularity in the data can be used to *compress* the data, i.e. to describe it using fewer symbols than the number of symbols needed to describe the data literally. Such a description should always uniquely specify the data it describes - hence given a description or

1. Unless we call “generated by a fair coin toss” a “regularity” too. There is nothing wrong with that view - the point is that, the *more* we can compress a sequence, the *more* regularity we have found. One can avoid all terminological confusion about the concept of “regularity” by making it *relative* to something called a “base measure,” but that is beyond the scope of this book (Li and Vitányi 1997).

encoding D' of a particular sequence of data D , we should always be able to fully reconstruct D using D' .

For example, sequence (1.1) above can be described using only a few words; we have actually done so already: we have not given the complete sequence — which would have taken about the whole page — but rather just a one-sentence description of it that nevertheless allows you to reproduce the complete sequence if necessary. Of course, the description was done using natural language and we may want to do it in some more formal manner.

If we want to identify regularity with compressibility, then it should also be the case that nonregular sequences can *not* be compressed. Since sequence (1.2) has been generated by fair coin tosses, it should not be compressible. As we will show below, we can indeed prove that *whatever* description method C one uses, the length of the description of a sequence like (1.2) will, with overwhelming probability, be not much shorter than sequence (1.2) itself.

Note that the description of sequence (1.3) that we gave above does not uniquely define sequence (1.3). Therefore, it does not count as a “real” description: one cannot regenerate the whole sequence if one has the description. A unique description that still takes only a few words may look like this: “Sequence (1.3) is one of those sequences of 10000 bits in which there are four times as many 0s as there are 1s. In the lexicographical ordering of those sequences, it is number i .” Here i is some large number that is explicitly spelled out in the description. In general, there are 2^n binary sequences of length n , while there are only $\binom{n}{\nu n}$ sequences of length n with a fraction of ν 1s. For every rational number ν except $\nu = 1/2$, the ratio of $\binom{n}{\nu n}$ to 2^n goes to 0 exponentially fast as n increases (this is shown formally in Chapter 4; see Equation (4.36) on page 129 and the text thereunder; by the method used there one can also show that for $\nu = 1/2$, it goes to 0 as $O(1/\sqrt{n})$). It follows that compared to the total number of binary sequences of length 10000, the number of sequences of length 10000 with four times as many 0s as 1s is vanishingly small. Direct computation shows it is smaller than 2^{7213} , so that the ratio between the number of sequences with four times as many 0s than 1s and the total number of sequences is smaller than 2^{-2787} . Thus, $i < 2^{7213} \ll 2^{10000}$ and to write down i in binary we need approximately $(\log_2 i) < 7213 \ll 10000$ bits.

Example 1.2 [Compressing Various Regular Sequences] The regularities underlying sequences (1) and (3) were of a very particular kind. To illustrate that *any* type of regularity in a sequence may be exploited to compress that sequence, we give a few more examples:

The Number π Evidently, there exists a computer program for generating the first n digits of π – such a program could be based, for example, on an infinite series expansion of π . This computer program has constant size, except for the specification of n which takes no more than $O(\log n)$ bits. Thus, when n is very large, the size of the program generating the first n digits of π will be very small compared to n : the π -digit sequence is deterministic, and therefore extremely regular.

Physics Data Consider a two-column table where the first column contains numbers representing various heights from which an object was dropped. The second column contains the corresponding times it took for the object to reach the ground. Assume both heights and times are recorded to some finite precision. In Section 1.5 we illustrate that such a table can be substantially compressed by first describing the coefficients of the second-degree polynomial H that expresses Newton’s law; then describing the heights; and then describing the deviation of the time points from the numbers predicted by H .

Natural Language Most sequences of words are not valid sentences according to the English language. This fact can be exploited to substantially compress English text, as long as it is syntactically mostly correct: by first describing a grammar for English, and then describing an English text D with the help of that grammar (Grünwald 1996), D can be described using much less bits than are needed without the assumption that word order is constrained.

Description Methods In order to formalize our idea, we have to replace the part of the descriptions above that made use of natural language by some formal language. For this, we need to fix a *description method* that maps sequences of data to their descriptions. Each such sequence will be encoded as another sequence of symbols coming from some finite or countably infinite *coding alphabet*. An *alphabet* is simply a countable set of distinct symbols. An example of an alphabet is the binary alphabet $\mathbb{B} = \{0, 1\}$; the three data sequences above are sequences over the binary alphabet. A sequence over a binary alphabet will also be called a *binary string*. Sometimes our data will consist of real numbers rather than binary strings. In practice, however, such numbers are always truncated to some finite precision. We can then again model them as symbols coming from a finite data alphabet.

More precisely, we are given a *sample* or equivalently *data sequence* $D = (x_1, \dots, x_n)$ where each x_i is a member of some set \mathcal{X} , called the *space of observations* or the *sample space for one observation*. The set of all potential samples of length n is denoted \mathcal{X}^n and is called the *sample space*. We call

x_i a single *observation* or, equivalently, a *data item*. For a general note about how our terminology relates to the usual terminology in statistics, machine learning and pattern recognition, we refer to the box on page 72.

Without any loss of generality we may describe our data sequences as binary strings (this is explained in Chapter 3, Section 3.2.2). Hence all the description methods we consider map data sequences to sequences of bits. All description methods considered in MDL satisfy the *unique decodability property*: given a description D' , there is at most one (“unique”) D that is encoded as D' . Therefore, given any description D' , one should be able to fully reconstruct the original sequence D . Semiformally:

Description Methods

Definition 1.1 A *description method* is a *one-many* relation from the sample space to the set of binary strings of arbitrary length.

A truly formal definition will be given in Chapter 3, Section 3.1. There we also explain how our notion of “description method” relates to the more common and closely related notion of a “code.” Until then, the distinction between codes and description methods is not that important, and we use the symbol C to denote both concepts.

Compression and Small Subsets We are now in a position to show that strings which are “intuitively” random cannot be substantially compressed. We equate intuitively random with “having been generated by independent tosses of a fair coin.” We therefore have to prove that it is virtually impossible to substantially compress sequences that have been generated by fair coin tosses. By “it is virtually impossible” we mean “it happens with vanishing probability.” Let us take some arbitrary but fixed description method C over the data alphabet consisting of the set of all binary sequences of length ≥ 1 . Such a code maps binary strings to binary strings. Suppose we are given a data sequence of length n (in Example 1.1, $n = 10000$). Clearly, there are 2^n possible data sequences of length n . We see that only two of these can be mapped to a description of length 1 (since there are only two binary strings of length 1: 0 and 1). Similarly, only a subset of at most 2^m sequences can have a description of length m . This means that at most $\sum_{i=1}^m 2^i < 2^{m+1}$ data sequences can have a description length $\leq m$. The fraction of data sequences of length n that can be compressed by more than k bits is therefore at

most 2^{-k} and as such decreases exponentially in k . If data are generated by n tosses of a fair coin, then all 2^n possibilities for the data are equally probable, so the probability that we can compress the data by more than k bits is smaller than 2^{-k} . For example, the probability that we can compress the data by more than 20 bits is smaller than one in a million.

Most Data Sets Are Incompressible

Suppose our goal is to encode a binary sequence of length n . Then

- No matter what description method we use, only a fraction of at most 2^{-k} sequences can be compressed by more than k bits.
- Thus, if data are generated by fair coin tosses, then no matter what code we use, the probability that we can compress a sequence by more than k bits is at most 2^{-k} .
- This observation will be generalized to data generated by an arbitrary distribution in Chapter 3. We then call it the *no-hypercompression inequality*. It can be found in the box on page 103.

Seen in this light, having a short description length for the data is equivalent to identifying the data as belonging to a tiny, very *special* subset out of all a priori possible data sequences; see also the box on page 31.

1.3 Solomonoff's Breakthrough – Kolmogorov Complexity

It seems that what data are compressible and what are not is extremely dependent on the specific description method used. In 1964 – in a pioneering paper that may be regarded as the starting point of all MDL-related research (Solomonoff 1964) – Ray Solomonoff suggested the use of a *universal computer language* as a description method. By a universal language we mean a computer language in which a universal Turing machine can be implemented. All commonly used computer languages, like Pascal, LISP, C, are “universal.” Every data sequence D can be encoded by a computer program P that prints D and then halts. We can define a description method that maps each data sequence D to the *shortest program* that prints D and then

halts.² Clearly, this is a description method in our sense of the word in that it defines a 1-many (even 1-1) mapping from sequences over the data alphabet to a subset of the binary sequences.

The shortest program for a sequence D is then interpreted as the *optimal hypothesis* for D . Let us see how this works for sequence (1.1) above. Using a language similar to C, we can write a program

```
for i = 1 to 2500; do {print '0001'}; halt
```

which prints sequence (1.1) but is clearly a lot shorter than it. If we want to make a fair comparison, we should rewrite this program in a binary alphabet; the resulting number of bits is still much smaller than 10000. The shortest program printing sequence (1.1) is at least as short as the program above, which means that sequence (1.1) is indeed highly compressible using Solomonoff's code. By the arguments of the previous section we see that, given an arbitrary description method C , sequences like (1.2) that have been generated by tosses of a fair coin are very likely not substantially compressible using C . In other words, the shortest program for sequence (1.1) is, with extremely high probability, not much shorter than the following:

```
print '01110100110100001010.....10111011000101100010'; halt
```

This program has size about equal to the length of the sequence. Clearly, it is nothing more than a repetition of the sequence.

Kolmogorov Complexity We define the *Kolmogorov complexity* of a sequence as the length of the shortest program that prints the sequence and then halts. Kolmogorov complexity has become a large subject in its own right; see (Li and Vitányi 1997) for a comprehensive introduction.

The lower the Kolmogorov complexity of a sequence, the *more regular* or equivalently, the *less random*, or, yet equivalently, the *simpler* it is. Measuring regularity in this way confronts us with a problem, since it depends on the particular programming language used. However, in his 1964 paper, Ray Solomonoff (Solomonoff 1964) showed that *asymptotically* it does not matter what programming language one uses, as long as it is universal: for every sequence of data $D = (x_1, \dots, x_n)$, let us denote by $L_{UL}(D)$ the length of the shortest program for D using universal language UL. We can show that for

2. If there exists more than one shortest program, we pick the one that comes first in enumeration order.

every two universal languages UL_1 and UL_2 , the difference between the two lengths $L_{UL_1}(D) - L_{UL_2}(D)$ is bounded by a constant that depends on UL_1 and UL_2 but not on the length n of the data sequence D . This implies that if we have a lot of data (n is large), then the difference in the two description lengths is negligible compared to the size of the data sequence. This result is known as the *invariance theorem* and was proved independently in (Solomonoff 1964), (Kolmogorov 1965) (hence the name Kolmogorov complexity), and (Chaitin 1969). The proof is based on the fact that one can write a compiler for every universal language UL_1 in every other universal language UL_2 . Such a compiler is a computer program with length $L_{1 \rightarrow 2}$. For example, we can write a program in Pascal that translates every C program into an equivalent Pascal program. The length (in bits) of this program would then be $L_{C \rightarrow \text{Pascal}}$. We can simulate each program P_1 written in language UL_1 by program P_2 written in UL_2 as follows: P_2 consists of the compiler from UL_1 to UL_2 , followed by P_1 . The length of program P_2 is bounded by the length of P_1 plus $L_{1 \rightarrow 2}$. Hence for all data D , the maximal difference between $L_{UL_1}(D)$ and $L_{UL_2}(D)$ is bounded by $\max\{L_{1 \rightarrow 2}, L_{2 \rightarrow 1}\}$, a constant which only depends on UL_1 and UL_2 but not on D .

1.4 Making the Idea Applicable

Problems There are two major problems with applying Kolmogorov complexity to practical learning problems:

1. **Uncomputability.** The Kolmogorov complexity cannot be computed in general;
2. **Large constants.** The description length of any sequence of data involves a constant depending on the description method used.

By “Kolmogorov complexity cannot be computed” we mean the following: there is no computer program that, for every sequence of data D , when given D as input, returns the shortest program that prints D and halts. Neither can there be a program, that for every data D returns only the *length* of the shortest program that prints D and then halts. Assuming such a program exists leads to a contradiction (Li and Vitányi 1997).

The second problem relates to the fact that in many realistic settings, we are confronted with very small data sequences for which the invariance theorem is not very relevant since the length of D is small compared to the constant $L_{1 \rightarrow 2}$.

“Idealized” or “Algorithmic” MDL If we ignore these problems, we may use Kolmogorov complexity as our fundamental concept and build a theory of idealized inductive inference on top of it. This road has been taken by Solomonoff (1964, 1978), starting with the 1964 paper in which he introduced Kolmogorov complexity, and by Kolmogorov, when he introduced the *Kolmogorov minimum sufficient statistic* (Li and Vitányi 1997; Cover and Thomas 1991). Both Solomonoff’s and Kolmogorov’s ideas have been substantially refined by several authors. We mention here P. Vitányi (Li and Vitányi 1997; Gács, Tromp, and Vitányi 2001; Vereshchagin and Vitányi 2002; Vereshchagin and Vitányi 2004; Vitányi 2005), who concentrated on Kolmogorov’s ideas, and M. Hutter (2004), who concentrated on Solomonoff’s ideas. Different authors have used different names for this area of research: “ideal MDL,” “idealized MDL,” or “algorithmic statistics.” It is closely related to the celebrated theory of *random sequences* due to P. Martin-Löf and Kolmogorov (Li and Vitányi 1997). We briefly return to idealized MDL in Chapter 17, Section 17.8.

Practical MDL Like most authors in the field, we concentrate here on non-idealized, practical versions of MDL that explicitly deal with the two problems mentioned above. The basic idea is to scale down Solomonoff’s approach so that it does become applicable. This is achieved by using description methods that are less expressive than general-purpose computer languages. Such description methods C should be restrictive enough so that for any data sequence D , we can always compute the length of the shortest description of D that is attainable using method C ; but they should be general enough to allow us to compress many of the intuitively “regular” sequences. The price we pay is that, using the “practical” MDL principle, there will always be some regular sequences which we will not be able to compress. But we already know that there can be *no* method for inductive inference at all which will always give us all the regularity there is — simply because there can be no automated method which for any sequence D finds the shortest computer program that prints D and then halts. Moreover, it will often be possible to guide a suitable choice of C by a priori knowledge we have about our problem domain. For example, below we consider a description method C that is based on the class of all polynomials, such that with the help of C we can compress all data sets which can meaningfully be seen as points on some polynomial.

1.5 Crude MDL, Refined MDL and Universal Coding

Let us recapitulate our main insights so far:

MDL: The Basic Idea

The goal of statistical inference may be cast as trying to find regularity in the data. “Regularity” may be identified with “ability to compress.” MDL combines these two insights by *viewing learning as data compression*: it tells us that, for a given set of hypotheses \mathcal{H} and data set D , we should try to find the hypothesis or combination of hypotheses in \mathcal{H} that compresses D most.

This idea can be applied to all sorts of inductive inference problems, but it turns out to be most fruitful in (and its development has mostly concentrated on) problems of *model selection* and, more generally, dealing with *overfitting*. Here is a standard example (we explain the difference between “model” and “hypothesis” after the example).

Example 1.3 [Model Selection and Overfitting] Consider the points in Figure 1.1. We would like to learn how the y -values depend on the x -values. To this end, we may want to fit a polynomial to the points. Straightforward linear regression will give us the leftmost polynomial - a straight line that seems overly simple: it does not capture the regularities in the data well. Since for any set of n points there exists a polynomial of the $(n - 1)$ st degree that goes exactly through all these points, simply looking for the polynomial with the least error will give us a polynomial like the one in the second picture. This polynomial seems overly complex: it reflects the random fluctuations in the data rather than the general pattern underlying it. Instead of picking the overly simple or the overly complex polynomial, it seems more reasonable to prefer a relatively simple polynomial with small but nonzero error, as in the rightmost picture. This intuition is confirmed by numerous experiments on real-world data from a broad variety of sources (Rissanen 1989; Vapnik 1998; Ripley 1996): if one naively fits a high-degree polynomial to a small sample (set of data points), then one obtains a very good fit to the data. Yet if one *tests* the inferred polynomial on a second set of data coming from the same source, it typically fits this test data very badly in the sense that there is a large distance between the polynomial and the new data points. We say that the polynomial *overfits* the data. Indeed, all model selection methods that are used in practice either implicitly or explicitly choose

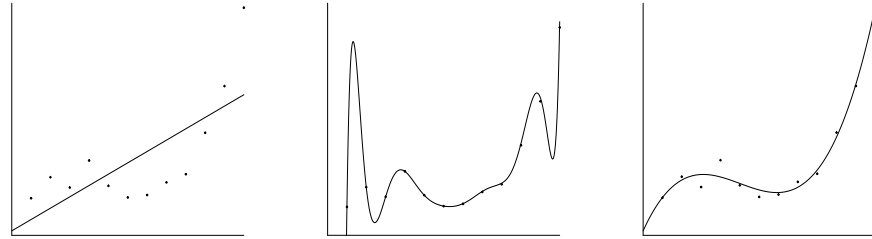


Figure 1.1 A simple, a complex and a tradeoff (third-degree) polynomial.

a tradeoff between goodness-of-fit and complexity of the models involved. In practice, such tradeoffs lead to much better predictions of test data than one would get by adopting the “simplest” (one degree) or most “complex”³ ($n - 1$ -degree) polynomial. MDL provides one particular means of achieving such a tradeoff.

It will be useful to distinguish between “model”, “model class” and “(point) hypothesis.” This terminology is explained in the box on page 15, and will be discussed in more detail in Section 2.4, page 69. In our terminology, the problem described in Example 1.3 is a “point hypothesis selection problem” if we are interested in selecting both the degree of a polynomial and the corresponding parameters; it is a “model selection problem” if we are mainly interested in selecting the degree.

To apply MDL to polynomial or other types of hypothesis and model selection, we have to make precise the somewhat vague insight “learning may be viewed as data compression.” This can be done in various ways. We first explain the earliest and simplest implementation of the idea. This is the so-called *two-part code* version of MDL:

3. Strictly speaking, in our context it is not very accurate to speak of “simple” or “complex” polynomials; instead we should call the *set* of first degree polynomials “simple,” and the *set* of 100th-degree polynomials “complex.”

Crude Two-Part Version of MDL Principle (Informally Stated)

Let $\mathcal{H}_1, \mathcal{H}_2, \dots$ be a list of candidate models (e.g., \mathcal{H}_γ is the set of γ th degree polynomials), each containing a set of point hypotheses (e.g., individual polynomials). The best point hypothesis $H \in \mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$ to explain the data D is the one which minimizes the sum $L(H) + L(D|H)$, where

- $L(H)$ is the length, in bits, of the description of the hypothesis; and
- $L(D|H)$ is the length, in bits, of the description of the data when encoded with the help of the hypothesis.

The best *model* to explain D is the smallest model containing the selected H .

The terminology “crude MDL” is explained in the next subsection. It is not standard, and it is introduced here for pedagogical reasons.

Example 1.4 [Polynomials, cont.] In our previous example, the candidate hypotheses were polynomials. We can describe a polynomial by describing its coefficients at a certain precision (number of bits per parameter). Thus, the higher the degree of a polynomial or the precision, the more bits we need to describe it and the more “complex” it becomes. A description of the data “with the help of” a hypothesis means that the better the hypothesis fits the data, the shorter the description will be. A hypothesis that fits the data well gives us a lot of *information* about the data. Such information can always be used to compress the data. Intuitively, this is because we only have to code the *errors* the hypothesis makes on the data rather than the full data. In our polynomial example, the better a polynomial H fits D , the fewer bits we need to encode the discrepancies between the actual y -values y_i and the predicted y -values $H(x_i)$. We can typically find a very complex point hypothesis (large $L(H)$) with a very good fit (small $L(D|H)$). We can also typically find a very simple point hypothesis (small $L(H)$) with a rather bad fit (large $L(D|H)$). The sum of the two description lengths will be minimized at a hypothesis that is quite (but not too) “simple,” with a good (but not perfect) fit.

1.5.1 From Crude to Refined MDL

Crude MDL picks the H minimizing the sum $L(H) + L(D|H)$. To make this procedure well defined, we need to agree on precise definitions for the

Models and Model Classes; (Point) Hypotheses

We use the word *model* to refer to a *set* of probability distributions or functions of the same functional form. E.g., the “first-order Markov model” is the set of all probability distributions that are first-order Markov chains. The “model of k th degree polynomials” is the set of all k th degree polynomials for some fixed k .

We use the word *model class* to refer to a family (set) of models, e.g. “the model class of all polynomials” or “the model class of all Markov chains of each order.” The definitions of “model” and “model class” are chosen so that they agree with how these words are used in statistical practice. Therefore they are intentionally left somewhat imprecise.

We use the word *hypothesis* to refer to an *arbitrary* set of probability distributions or functions. We use the word *point hypothesis* to refer to a *single* probability distribution (e.g. a Markov chain with all parameter values specified) or function (e.g. a particular polynomial). In parametric inference (Chapter 2), a point hypothesis corresponds to a particular parameter value. A point hypothesis may also be viewed as an *instantiation* of a model.

What we call “point hypothesis” is called “*simple hypothesis*” in the statistics literature; our use of the word “model (selection)” coincides with its use in much of the statistics literature; see Section 2.3, page 62 where we give several examples to clarify our terminology.

Figure 1.2 Models and Model Classes; (Point) Hypotheses.

codes (description methods) giving rise to lengths $L(D|H)$ and $L(H)$. We now discuss these codes in more detail. We will see that the definition of $L(H)$ is problematic, indicating that we somehow need to “refine” our crude MDL principle.

Definition of $L(D|H)$ Consider a two-part code as described above, and assume for the time being that all H under consideration define probability distributions. If H is a polynomial, we can turn it into a distribution by mak-

ing the additional assumption that the Y -values are given by $Y = H(X) + Z$, where Z is a normally distributed noise term with mean 0.

For each H we need to define a code with length $L(\cdot | H)$ such that $L(D|H)$ can be interpreted as “the codelength of D when encoded with the help of H .” It turns out that for probabilistic hypotheses, there is only one reasonable choice for this code; this is explained at length in Chapter 5. It is the so-called *Shannon-Fano code*, satisfying, for all data sequences D , $L(D|H) = -\log P(D|H)$, where $P(D|H)$ is the probability mass or density of D according to H . Such a code always exists, as we explain in Chapter 3, in the box on page 96.

Definition of $L(H)$: A Problem for Crude MDL It is more problematic to find a good code for hypotheses H . Some authors have simply used “intuitively reasonable” codes in the past, but this is not satisfactory: since the description length $L(H)$ of any fixed point hypothesis H can be very large under one code, but quite short under another, our procedure is in danger of becoming arbitrary. Instead, *we need some additional principle for designing a code for \mathcal{H} .*

In the first publications on MDL (Rissanen 1978; Rissanen 1983), it was implicitly advocated to choose some sort of *minimax code* for each \mathcal{H}_γ , minimizing the shortest worst-case total description length $L(H) + L(D|H)$, where the worst-case is over all possible data sequences. Thus, the MDL principle is employed at a “meta-level” to choose a code for \mathcal{H}_γ . This idea, already implicit in Rissanen’s early work about perhaps for the first time stated and formalized in a completely precise way Barron and Cover (1991), is the first step towards “refined” MDL.

More Problems for Crude MDL We can use crude MDL to code any sequence of data D with a total description length $L(D) := \min_H \{L(D|H) + L(H)\}$. But it turns out that this code is *incomplete*: one can show that there exist other codes L' which for some D achieve strictly smaller codelength ($L'(D) < L(D)$), and for no D achieve larger codelength (Chapter 6, Example 6.4). It seems strange that our “minimum description length” principle should be based on codes which are incomplete (inefficient) in this sense. Another, less fundamental problem with two-part codes is that, if designed in a minimax way as indicated above, they require a cumbersome discretization of the model space \mathcal{H} , which is not always feasible in practice. The final problem we mention is that, while it is clear how to use crude two-part codes for

point hypothesis and model selection, it is not immediately clear how they can be used for *prediction*.

Later, Rissanen (1984) realized that these problems could be side-stepped by using *one-part* rather than *two-part codes*. As we explain below, it depends on the situation at hand whether a one-part or a two-part code should be used. Combining the idea of designing codes so as to achieve essentially minimax optimal codelengths with the combined use of one-part and two-part codes (whichever is appropriate for the situation at hand) has culminated in a theory of inductive inference that we call *refined MDL*. We discuss it in more detail in the next subsection.

Crude Two-Part MDL (Part I, Chapter 5 of this book)

In this book, we use the term “crude MDL” to refer to applications of MDL for model and hypothesis selection of the type described in the box on page 14, as long as the hypotheses $H \in \mathcal{H}$ are encoded in “intuitively reasonable” but ad-hoc ways.

Refined MDL is sometimes based on one-part codes, sometimes on two-part codes, and sometimes on a combination of these, but, in contrast to crude MDL, the codes are invariably designed according to some minimax principles. If there is a choice, one should always prefer refined MDL, but in some exotic modeling situations, the use of crude MDL is inevitable.

Part I of this book first discusses all probabilistic, statistical and information-theoretic preliminaries (Chapters 2–4) and culminates in a description of crude two-part MDL (Chapter 5). Refined MDL is described only in Part III.

1.5.2 Universal Coding and Refined MDL

In refined MDL, we associate a code for encoding D not with a single $H \in \mathcal{H}$, but with the full model \mathcal{H} . Thus, given model \mathcal{H} , we encode data not in two parts but we design a single *one-part code* with lengths $\bar{L}(D|\mathcal{H})$. This code is designed such that *whenever there is a member of (parameter in) \mathcal{H} that fits the data well, in the sense that $L(D | H)$ is small, then the codelength $\bar{L}(D|\mathcal{H})$ will also be small*. Codes with this property are called *universal codes* in the information-theoretic literature (Barron, Rissanen, and Yu 1998):

Universal Coding (Part II of This Book)

There exist at least four types of universal codes:

1. The normalized maximum likelihood (NML) code and its variations.
2. The Bayesian mixture code and its variations.
3. The prequential plug-in code
4. The two-part code

These codes are all based on entirely different coding schemes, but in practice, lead to very similar codelengths $\bar{L}(D|\mathcal{H})$. Part II of this book is entirely devoted to universal coding. The four types of codes are introduced in Chapter 6. This is followed by a separate chapter for each code.

For each model \mathcal{H} , there are many different universal codes we can associate with \mathcal{H} . When applying MDL, we have a preference for the one that is *minimax optimal* in a sense made precise in Chapter 6. For example, the set \mathcal{H}_3 of third-degree polynomials is associated with a code with lengths $\bar{L}(\cdot | \mathcal{H}_3)$ such that, the better the data D are fit by the best-fitting third-degree polynomial, the shorter the codelength $\bar{L}(D | \mathcal{H})$. $\bar{L}(D | \mathcal{H})$ is called the *stochastic complexity* of the data given the model.

Refined MDL is a general theory of inductive inference based on universal codes that are designed to be minimax, or close to minimax optimal. It has mostly been developed for model selection, estimation and prediction. To give a first flavor, we initially discuss model selection, where, arguably, it has the most new insights to offer:

1.5.3 Refined MDL for Model Selection

Parametric Complexity A fundamental concept of refined MDL for model selection is the *parametric complexity* of a parametric model \mathcal{H} which we denote by $\text{COMP}(\mathcal{H})$. This is a measure of the “richness” of model \mathcal{H} , indicating its ability to fit random data. This complexity is related to the number of degrees-of-freedom (parameters) in \mathcal{H} , but also to the geometrical structure of \mathcal{H} ; see Example 1.5. To see how it relates to stochastic complexity, let, for given data D , \hat{H} denote the distribution in \mathcal{H} which maximizes the probability, and hence minimizes the codelength $L(D | \hat{H})$ of D . It turns out

that

$$\bar{L}(D | \mathcal{H}) = \text{stochastic complexity of } D \text{ given } \mathcal{H} = L(D | \hat{H}) + \text{COMP}(\mathcal{H}).$$

Refined MDL model selection between two parametric models \mathcal{H}_1 and \mathcal{H}_2 (such as the models of first and second degree polynomials) now proceeds as follows. We encode data D in two stages. In the first stage, we encode a number $j \in \{1, 2\}$. In the second stage, we encode the data using the universal code with lengths $\bar{L}(D | \mathcal{H}_j)$. As in the two-part code principle, we then select the \mathcal{M}_j achieving the minimum total two-part codelength,

$$\min_{j \in \{1, 2\}} \{L(j) + \bar{L}(D | \mathcal{H}_j)\} = \min_{j \in \{1, 2\}} \{L(j) + L(D | \hat{H}) + \text{COMP}(\mathcal{H})\}. \quad (1.4)$$

Since the worst-case optimal code to encode j needs only 1 bit to encode either $j = 1$ or $j = 2$, we use a code for the first-part such that $L(1) = L(2) = 1$. But this means that $L(j)$ plays no role in the minimization, and we are effectively selecting the model such that the stochastic complexity of the given data D is smallest.⁴ Thus, in the end we select the model *minimizing the one-part codelength of the data*. Nevertheless, refined MDL model selection involves a tradeoff between two terms: a goodness-of-fit term $L(D | \hat{H})$ and a complexity term $\text{COMP}(\mathcal{H})$. However, because we do not explicitly encode hypotheses H anymore, there is no potential for arbitrary codelengths anymore. The resulting procedure can be interpreted in several different ways, some of which provide us with rationales for MDL model selection beyond the pure coding interpretation (Chapter 14):

1. **Counting/differential geometric interpretation** The parametric complexity of a model is the logarithm of the number of *essentially different, distinguishable* point hypotheses within the model.
2. **Two-part code interpretation** For large samples, the stochastic complexity can be interpreted as a two-part codelength of the data after all, where hypotheses H are encoded with a special code that works by first discretizing the model space \mathcal{H} into a set of “maximally distinguishable hypotheses,” and then assigning equal codelength to each of these.
3. **Bayesian interpretation** In many cases, refined MDL model selection coincides with Bayes factor model selection based on a *noninformative prior* such as *Jeffreys’ prior* (Bernardo and Smith 1994).

4. The reason we include $L(j)$ at all in (1.4) is to maintain consistency with the case where we need to select between an infinite number of models. In that case, it is necessary to include $L(j)$.

4. Prequential interpretation MDL model selection can be interpreted as selecting the model with the best predictive performance when sequentially predicting *unseen* test data, in the sense described in Chapter 6, Section 6.4 and Chapter 9. This makes it an instance of Dawid’s (1984) *prequential* model validation and also relates it to *cross-validation* methods; see Chapter 17, Sections 17.5 and 17.6.

In Section 1.6.1 we show that refined MDL allows us to compare models of different functional form. It even accounts for the phenomenon that different models with the same number of parameters may not be equally “complex.”

1.5.4 General Refined MDL: Prediction and Hypothesis Selection

Model selection is just one application of refined MDL. The two other main applications are *point hypothesis selection* and *prediction*. These applications can also be interpreted as methods for parametric and nonparametric *estimation*. In fact, it turns out that large parts of MDL theory can be reinterpreted as a theory about *sequential prediction of future data given previously seen data*. This “prequential” interpretation of MDL (Chapter 15) is at least as important as the coding interpretation. It is based on the fundamental correspondence between probability distributions and codes via the Shannon-Fano code that we alluded to before, when explaining the code with lengths $L(D | H)$; see the box on page 96. This correspondence allows us to view any universal code $\bar{L}(\cdot | \mathcal{H})$ as a strategy for sequentially predicting data, such that the better \mathcal{H} is suited as a model for the data, the better the predictions will be.

MDL prediction and hypothesis selection are mathematically cleaner than MDL model selection: in Chapter 15, we provide theorems (Theorem 15.1 and Theorem 15.3) which, in the respective contexts of prediction and hypothesis selection, express that, in full generality, *good data compression implies fast learning*, where “learning” is defined as “finding a hypothesis that is in some sense close to an imagined “true state of the world.” There are similar theorems for model selection, but these lack some of the simplicity and elegance of Theorem 15.1 and Theorem 15.3.

Probabilistic vs. Nonprobabilistic MDL Like most other authors on MDL, in this book we confine ourselves to *probabilistic hypotheses*, also known as *probabilistic sources*. These are hypotheses that take the form of *probability distributions* over the space of possible data sequences. The examples we give in this chapter (Examples 1.3 and 1.5) involve hypotheses H that are functions

from some space \mathcal{X} to another space \mathcal{Y} ; at first sight, these are not “probabilistic.” We will usually assume that for any given x , we have $y = H(x) + Z$ where Z is a *noise term* with a known distribution. Typically, the noise Z will be assumed to be Gaussian (normally) distributed. With such an additional assumption, we may view “functional” hypotheses $H : \mathcal{X} \rightarrow \mathcal{Y}$ as “probabilistic” after all. Such a technique of turning functions into probability distributions is customary in statistics, and we will use it throughout large parts of this book. Whenever we refer to MDL, we implicitly assume that we deal with probabilistic models. We should note though that there exists variations of MDL that *directly* work with universal codes relative to functional hypotheses such as polynomials (see Section 1.9.1, and Chapter 17, Section 17.10).

Fixing Notation

We use the symbol H for general point hypotheses, that may either represent a probabilistic source or a deterministic function. We use \mathcal{H} for sets of such general point hypotheses. We reserve the symbol \mathcal{M} for probabilistic models and model classes. We denote probabilistic point hypotheses by P , and point hypotheses that are deterministic functions by h .

Individual-Sequence vs. Expectation-based MDL Refined MDL is based on minimax optimal universal codes. Broadly speaking, there are two different ways to define what we mean by minimax optimality. One is to look at the worst-case codelength over *all possible sequences*. We call this *individual-sequence MDL*. An alternative is to look at *expected codelength*, where the expectation is taken over some probability distribution, usually but not always assumed to be a member of the model class \mathcal{M} under consideration. We call this *expectation-based MDL*. We discuss the distinction in detail in Part III of the book; see also the box on page 407. The individual-sequence approach is the one taken by Rissanen, the main originator of MDL, and we will mostly follow it throughout this book.

The Luckiness Principle In the individual-sequence approach, the minimax optimal universal code is given by the normalized maximum likelihood (NML) code that we mentioned above. A problem is that for many (in fact, most) practically interesting models, the NML code is not well defined. In

such cases, a minimax optimal code does not exist. As we explain in Chapter 11, in some cases one can get around this problem using so-called “conditional NML” codes, but in general, one needs to use codes based on a modified minimax principle, which we call the *luckiness principle*. Although it has been implicitly used in MDL since its inception, I am the first to use the term “luckiness principle” in an MDL context; see the box on page 92, Chapter 3; the developments in Chapter 11, Section 11.3, where we introduce the concept of a *luckiness function*; and the discussion in Chapter 17, Section 17.2.1.

The luckiness principle reintroduces some subjectivity in MDL code design. This seems to bring us back to the ad-hoc codes used in crude two-part MDL. The difference however is that with luckiness functions, we can precisely quantify the effects of this subjectivity: for each possible data sample D that we may observe, we can indicate how “lucky” we are on the sample, i.e. how many extra bits we need compared to encode D compared to the best hypothesis that we have available for D . This idea significantly extends the applicability of refined MDL methods.

MDL is a Principle Contrary to what is often thought, MDL, and even, “modern, refined MDL” is *not* a unique, single method of inductive inference. Rather, it represents a general *principle* for doing inductive inference. The principle may (and will) be formulated precisely enough to allow us to establish, for many given methods (procedures, learning algorithms) “this method is an instance of MDL” or “this is *not* an instance of MDL. But nevertheless:

MDL Is a Principle, Not a Unique Method

Being a *principle*, MDL gives rise to several *methods* of inductive inference. There is no single “uniquely optimal MDL method/procedure/algorithm.” Nevertheless, in *some special situations* (e.g. simple parametric statistical models), one can clearly distinguish between good and not so good versions of MDL, and something close to “an optimal MDL method” exists.

Summary: Refined MDL (Part III of This Book)

Refined MDL is a method of inductive inference based on *universal codes* which are designed to have some *minimax optimality properties*. Each model \mathcal{H} under consideration is associated with a corresponding universal code. In this book we restrict ourselves to probabilistic \mathcal{H} . Refined MDL has mainly been developed for model selection, point hypothesis selection and prediction.

Refined MDL comes in two versions: individual-sequence and expectation-based refined MDL, depending on whether the universal codes are designed to be optimal in an individual-sequence or in an expected sense. If the minimax optimal code relative to a model \mathcal{M} is not defined, some element of subjectivity is introduced into the coding using a *luckiness function*. A more precise overview is given in the box on page 406.

In the remainder of this chapter we will mostly concentrate on MDL for model selection.

1.6 Some Remarks on Model Selection

Model selection is a controversial topic in statistics. Although most people agree that it is important, many say it can only be done on external grounds, and never by merely looking at the data. Still, a plethora of automatic model selection methods has been suggested in the literature. These can give wildly different results on the same data, one of the main reasons being that they have often been designed with different goals in mind. This section starts with a further example that motivates the need for model selection, and it then discusses several goals that one may have in mind when doing model selection. These issues are discussed in a lot more detail in Chapter 14. See also Chapter 17, especially Section 17.3, where we compare MDL model selection to the standard model selection methods AIC and BIC.

1.6.1 Model Selection among Non-Nested Models

Model selection is often used in the following context: two researchers or research groups A and B propose entirely different models \mathcal{M}_A and \mathcal{M}_B as an explanation for the same data D . This situation occurs all the time in applied sciences like econometrics, biology, experimental psychology, etc. For

example, group A may have some general theory about the phenomenon at hand which prescribes that the trend in data D is given by some polynomial. Group B may think that the trend is better described by some neural network; a concrete case will be given in Example 1.3 below. A and B would like to have some way of deciding which of their two models is better suited for the data at hand. If they simply decide on the model containing the hypothesis (parameter instantiation) that best fits the data, they once again run the risk of overfitting: if model \mathcal{M}_A has more degrees of freedom (parameters) than model \mathcal{M}_B , it will typically be able to better fit random noise in the data. It may then be selected even if \mathcal{M}_B actually better captures the underlying trend (regularity) in the data. Therefore, just as in the hypothesis selection example, deciding whether \mathcal{M}_A or \mathcal{M}_B is a better explanation for the data should somehow depend on how well \mathcal{M}_A and \mathcal{M}_B fit the data and on the respective “complexities” of \mathcal{M}_A and \mathcal{M}_B .

In the polynomial case discussed before, there was a countably infinite number of “nested” \mathcal{M}_γ (i.e. $\mathcal{M}_\gamma \subset \mathcal{M}_{\gamma+1}$). In contrast, we now deal with a finite number of entirely unrelated models \mathcal{M}_γ . But there is nothing that stops us from using MDL model selection as “defined” above.

Example 1.5 [Selecting Between Models of Different Functional Form]

Consider two models from psychophysics describing the relationship between physical dimensions (e.g., light intensity) and their psychological counterparts (e.g. brightness) (Myung, Balasubramanian, and Pitt 2000): $y = ax^b + Z$ (Stevens’s model) and $y = a \ln(x + b) + Z$ (Fechner’s model) where Z is a normally distributed noise term. Both models have two free parameters; nevertheless, according to the refined version of MDL model selection to be introduced in Part III, Chapter 14 of this book, Stevens’s model is in a sense “more complex” than Fechner’s (see page 417). Roughly speaking, this means there are a lot more data patterns that can be *explained* by Stevens’s model than can be explained by Fechner’s model. Somewhat more precisely, the number of data patterns (sequences of data) of a given length that can be fit well by Stevens’s model is much larger than the number of data patterns of the same length that can be fit well by Fechner’s model. Therefore, using Stevens’s model we run a larger risk of “overfitting.”

In the example above, the goal was to select between a power law and a logarithmic relationship. In general, we may of course come across model selection problems involving neural networks, polynomials, Fourier or wavelet expansions, exponential functions - anything may be proposed and tested.

could have tried to learn g using a model class \mathcal{H} containing the function $y = \exp(x)$. But in general, both our imagination and our computational resources are limited, and we may be forced to use imperfect models.

If, based on a small sample, we choose the best-fitting polynomial \hat{h} within the set of *all* polynomials, then, even though \hat{h} will fit the data very well, it is likely to be quite unrelated to the “true” g , and \hat{h} may lead to disastrous predictions of future data. The reason is that, for small samples, the set of all polynomials is very large compared to the set of possible data patterns that we might have observed. Therefore, any particular data pattern can only give us very limited information about which high-degree polynomial best approximates g . On the other hand, if we choose the best-fitting \hat{h}° in some much smaller set such as the set of second-degree polynomials, then it is highly probable that the prediction quality (mean squared error) of \hat{h}° on future data is about the same as its mean squared error on the data we observed: the size (complexity) of the contemplated model is relatively small compared to the set of possible data patterns that we might have observed. Therefore, the particular pattern that we do observe gives us a lot of information on what second-degree polynomial best approximates g .

Thus, (a) \hat{h}° typically leads to better predictions of future data than \hat{h} ; and (b) unlike \hat{h} , \hat{h}° is *reliable* in that it gives a correct impression of how good it will predict future data *even if the “true” g is “infinitely” complex*. This idea does not just appear in MDL, but is also the basis of the structural risk minimization approach (Vapnik 1998) and many standard statistical methods for nonparametric inference; see Chapter 17, Section 17.10. In such approaches one acknowledges that the data-generating machinery can be infinitely complex (e.g., not describable by a finite degree polynomial). Nevertheless, it is still a good strategy to approximate it by simple hypotheses (low-degree polynomials) as long as the sample size is small. Summarizing:

The Inherent Difference between Under- and Overfitting

If we choose an overly simple model for our data, then the best-fitting point hypothesis within the model is likely to be almost the best predictor, within the simple model, of future data coming from the same source. If we overfit (choose a very complex model) and there is noise in our data, then, *even if the complex model contains the “true” point hypothesis*, the best-fitting point hypothesis within the model may lead to very bad predictions of future data coming from the same source.

This statement is very imprecise and is meant more to convey the general idea than to be completely true. The fundamental consistency theorems for MDL prediction and hypothesis selection (Chapter 15, Theorem 15.1 and Theorem 15.3), as well as their extension to model selection (Chapter 16), are essentially just variations of this statement that are provably true.

The Future and The Past Our analysis depends on the data items (x_i, y_i) to be probabilistically independent. While this assumption may be substantially weakened, we can justify the use of MDL and other forms of Occam's razor *only* if we are willing to adopt some (possibly very weak) assumption of the sort "training data and future data are from the same source": future data should (at least with high probability) be subject to some of the same regularities as training data. Otherwise, D and D' may be completely unrelated and *no* method of inductive inference can be expected to work well. This is indirectly related to the *grue*-paradox (Goodman 1955).

MDL and Occam's Razor

While MDL does have a built-in preference for selecting "simple" models (with small "parametric complexity"), this does *not at all* mean that applying MDL only makes sense in situations where simpler models are more likely to be true. MDL is a *methodology for inferring models from data, not a statement about how the world works!* For small sample sizes, it prefers simple models. It does so not because these are "more likely to be true" (they often are not). Instead, it does so because this tends to select the model that leads to the best predictions of future data from the same source. For small sample sizes this may be a model much simpler than the model containing the "truth" (assuming for the time being that such a model containing the "truth" exists in the first place).

In fact, some of MDL's most useful and successful applications are in nonparametric statistics where the "truth" underlying data is typically assumed to be "infinitely" complex (see Chapter 13 and Chapter 15).

1.9 History and Forms of MDL

The practical MDL principle that we discuss in this book has mainly been developed by J. Rissanen in a series of papers starting with (Rissanen 1978). It has its roots in the theory of Kolmogorov complexity (Li and Vitányi 1997), developed in the 1960s by Solomonoff (1964), Kolmogorov (1965) and Chaitin (1966, 1969). Among these authors, Solomonoff (a former student of the famous philosopher of science, Rudolf Carnap) was explicitly interested in inductive inference. The 1964 paper contains explicit suggestions on how the underlying ideas could be made practical, thereby foreshadowing some of the later work on two-part MDL. While Rissanen was not aware of Solomonoff's work at the time, Kolmogorov's [1965] paper did serve as an inspiration for Rissanen's (1978) development of MDL. Still, Rissanen's practical MDL is quite different from the idealized forms of MDL that have been directly based on Kolmogorov complexity, which we discussed in Section 1.4.

Another important inspiration for Rissanen was Akaike's AIC method for model selection (Chapter 17, Section 17.3), essentially the first model selection method based on information-theoretic ideas (Akaike 1973). Even though Rissanen was inspired by AIC, both the actual method and the underlying philosophy are substantially different from MDL.

Minimum Message Length MDL is much closer related to the *Minimum Message Length (MML) Principle* (Wallace 2005), developed by Wallace and his coworkers in a series of papers starting with the groundbreaking (Wallace and Boulton 1968); other milestones are (Wallace and Boulton 1975) and (Wallace and Freeman 1987). Remarkably, Wallace developed his ideas without being aware of the notion of Kolmogorov complexity. Although Rissanen became aware of Wallace's work before the publication of (Rissanen 1978), he developed his ideas mostly independently, being influenced rather by Akaike and Kolmogorov. Indeed, despite the close resemblance of both methods in practice, the underlying philosophy is very different - see Chapter 17, Section 17.4.

Refined MDL The first publications on MDL only mention two-part codes. Important progress was made by Rissanen (1984), in which prequential codes are employed for the first time and Rissanen (1987), who introduced the Bayesian mixture codes into MDL. This led to the development of the notion of stochastic complexity as the shortest codelength of the data given a model

(Rissanen 1986c; Rissanen 1987). However, the connection to Shtarkov's *normalized maximum likelihood code* was not made until 1996, and this prevented the full development of the notion of "parametric complexity." In the mean time, in his impressive Ph.D. thesis, Barron (1985) showed how a specific version of the two-part code criterion has excellent frequentist statistical consistency properties. This was extended by Barron and Cover (1991) who achieved a breakthrough for two-part codes: they gave clear prescriptions on how to design codes for hypotheses, relating codes with good minimax code-length properties to rates of convergence in statistical consistency theorems. Some of the ideas of Rissanen (1987) and Barron and Cover (1991) were, as it were, unified when Rissanen (1996) introduced the normalized maximum likelihood code. The resulting theory was summarized for the first time by Barron, Rissanen, and Yu (1998), and is the subject of this book. Whenever we need to distinguish it from other forms of MDL, we call it "refined MDL."

1.9.1 What Is MDL?

"MDL" is used by different authors in somewhat different meanings, and it may be useful to review these. Some authors use MDL as a broad umbrella term for all types of inductive inference based on finding a short codelength for the data. This would, for example, include the "idealized" versions of MDL based on Kolmogorov complexity (page 11) and Wallaces's MML principle (see above). Some authors take an even broader view and include all inductive inference that is based on data compression, even if it cannot be directly interpreted in terms of codelength minimization. This includes, for example the work on similarity analysis and clustering based on the *normalized compression distance* (Cilibrasi and Vitányi 2005).

On the other extreme, for historical reasons, some authors use the *MDL Criterion* to describe a very specific (and often not very successful) model selection criterion equivalent to BIC (see Chapter 17, Section 17.3).

As already indicated, we adopt the meaning of the term that is embraced in the survey (Barron, Rissanen, and Yu 1998), written by arguably the three most important contributors to the field: we use MDL for general *inference based on universal models*. Although we concentrate on hypothesis selection, model selection and prediction, this idea can be further extended to many other types of inductive inference. These include *denoising* (Rissanen 2000; Hansen and Yu 2000; Roos, Myllymäki, and Tirri 2005), *similarity analysis* and *clustering* (Kontkanen, Myllymäki, Buntine, Rissanen, and Tirri 2005), *outlier detection* and *transduction* (as defined in (Vapnik 1998)), and many others. In

such areas there has been less research and a “definitive” universal-model based MDL approach has not yet been formulated. We do expect, however, that such research will take place in the future: one of the main strengths of “MDL” in this broad sense is that it can be applied to ever more exotic modeling situations, in which the models do not resemble anything that is usually encountered in statistical practice. An example is the model of context-free grammars, already considered by Solomonoff (1964).

Another application of universal-model based MDL is the type of problem usually studied in *statistical learning theory* (Vapnik 1998); see also Chapter 17, Section 17.10. Here the goal is to directly learn functions (such as polynomials) to predict Y given X , without making any specific probabilistic assumptions about the noise. MDL has been developed in some detail for such problems, most notably *classification* problems, where Y takes its values in a finite set – spam filtering is a prototypical example; here X stands for an email message, and Y encodes whether or not it is spam. An example is the application of MDL to decision tree learning (Quinlan and Rivest 1989; Wallace and Patrick 1993; Mehta, Rissanen, and Agrawal 1995). Some MDL theory for such cases has been developed (Meir and Merhav 1995; Yamanishi 1998; Grünwald 1998), but the existing MDL methods in this area can behave suboptimally. This is explained in Chapter 17, Section 17.10.2. Although we certainly consider it a part of “refined” MDL, we do not consider this “nonprobabilistic” MDL further in this book, except in Section 17.10.2.

1.9.2 MDL Literature

Theoretical Contributions There have been numerous contributors to refined MDL theory, but there are three researchers that I should mention explicitly: J. Rissanen, B. Yu and A. Barron, who jointly wrote (Barron, Rissanen, and Yu 1998). For example, most of the results that connect MDL to traditional statistics (including Theorem 15.1 and Theorem 15.3 in Chapter 15) are due to A. Barron. This book contains numerous references to their work.

There is a close connection between MDL theory and work in *universal coding* ((Merhav and Feder 1998); see also Chapter 6) and *universal prediction* ((Cesa-Bianchi and Lugosi 2006); see also Chapter 17, Section 17.9).

Practical Contributions There have been numerous practical applications of MDL. The only three applications we describe in detail are a crude MDL method for learning Markov chains (Chapter 5); a refined MDL method for

learning densities based on histograms (Chapter 13 and Chapter 15); and MDL regression (Chapter 12 and Chapter 14). Below we give a few representative examples of other applications and experimental results that have appeared in the literature. We warn the reader that this list is by no means complete! Hansen and Yu (2001) apply MDL to a variety of practical problems involving regression, clustering analysis, and time series analysis. In (Tabus, Rissanen, and Astola 2002; Tabus, Rissanen, and Astola 2003), MDL is used for classification problems arising in genomics. Lee (2002a,b) describes additive clustering with MDL. use MDL for image denoising and apply MDL to decision tree learning. use MDL for sequential prediction. In (Myung, Pitt, Zhang, and Balasubramanian 2000; Myung, Balasubramanian, and Pitt 2000), MDL is applied to a variety of model selection problems arising in cognitive psychology. All these authors apply modern, “refined” versions of MDL. Some references to older work, in which “crude” (but often quite sensible) ad-hoc codes are used, are (Friedman, Geiger, and Goldszmidt 1997; Allen and Greiner 2000; Allen, Madani, and Greiner 2003; Rissanen and Ristad 1994; Quinlan and Rivest 1989; Nowak and Figueiredo 2000; Liu and Moulin 1998; Ndili, Nowak, and Figueiredo 2001; Figueiredo, J. Leitão, and A.K.Jain 2000; Gao and Li 1989). In these papers, MDL is applied to learning Bayesian networks, grammar inference and language acquisition, learning decision trees, analysis of Poisson point processes (for biomedical imaging applications), image denoising, image segmentation, contour estimation, and Chinese handwritten character recognition respectively. MDL has also been extensively studied in time-series analysis, both in theory (Hannan and Rissanen 1982; Gerencsér 1987; Wax 1988; Hannan, McDougall, and Poskitt 1989; Hemerly and Davis 1989b; Hemerly and Davis 1989a; Gerencsér 1994) and practice (Wei 1992; Wagenmakers, Grünwald, and Steyvers 2006).

Finally, we should note that there have been a number of applications, especially in *natural language learning*, which, although practically viable, have been primarily inspired by “idealized MDL” and Kolmogorov complexity, rather than by the Rissanen-Barron-Yu style of MDL that we consider here. These include (Adriaans and Jacobs 2006; Osborne 1999; Starkie 2001) and my own (Grünwald 1996).

Other Tutorials, Introductions and Overviews The reader who prefers a shorter introduction to MDL than the present one may want to have a look at (Barron, Rissanen, and Yu 1998) (very theoretical and very comprehensive; presumes knowledge of information theory), (Hansen and Yu 2001)

(presumes knowledge of statistics; describes several practical applications), (Lantermann 2001) (about comparing MDL, MML and asymptotic Bayesian approaches to model selection), or perhaps my own (Grünwald 2005), which is part of (Grünwald, Myung, and Pitt 2005), a “source book” for MDL theory and applications that contains chapters by most of the main contributors to the field.

Rissanen (1989,2007) has written two books on MDL. While outdated as an introduction to MDL, the “little green book” (Rissanen 1989) is still very much worth reading for its clear exposition of the philosophy underlying MDL. (Rissanen 2007) contains a brief general introduction and then focuses on some recent research of Rissanen’s, applying the renormalized maximum likelihood (RNML) distribution (Chapter 11) in regression and denoising, and formalizing the connection between MDL and Kolmogorov’s structure function. In contrast to myself, Rissanen writes in accord with his own principle: while containing a lot of information, both texts are quite short.

1.10 Summary and Outlook

We have discussed the relationship between compression, regularity, and learning. We have given a first idea of what the MDL principle is all about, and of the kind of problems we can apply it to. In the next chapters, we present the mathematical background needed to describe such applications in detail.