



Ian H. Witten & Eibe Frank

# DATA MINING

Practical Machine Learning Tools and Techniques



A close-up photograph of a green chameleon camouflaged among several large, green, fern-like leaves against a black background. The chameleon's body shape and coloration closely match the surrounding foliage.

SECOND EDITION



# Data Mining

**Practical Machine Learning Tools and Techniques**

# The Morgan Kaufmann Series in Data Management Systems

*Series Editor:* Jim Gray, Microsoft Research

*Data Mining: Practical Machine Learning Tools and Techniques, Second Edition*  
Ian H. Witten and Eibe Frank

*Fuzzy Modeling and Genetic Algorithms for Data Mining and Exploration*  
Earl Cox

*Data Modeling Essentials, Third Edition*  
Graeme C. Simsion and Graham C. Witt

*Location-Based Services*  
Jochen Schiller and Agnès Voisard

*Database Modeling with Microsoft® Visio for Enterprise Architects*  
Terry Halpin, Ken Evans, Patrick Hallock, and Bill Maclean

*Designing Data-Intensive Web Applications*  
Stefano Ceri, Piero Fraternali, Aldo Bongio, Marco Brambilla, Sara Comai, and Maristella Matera

*Mining the Web: Discovering Knowledge from Hypertext Data*  
Soumen Chakrabarti

*Advanced SQL: 1999—Understanding Object-Relational and Other Advanced Features*  
Jim Melton

*Database Tuning: Principles, Experiments, and Troubleshooting Techniques*  
Dennis Shasha and Philippe Bonnet

*SQL: 1999—Understanding Relational Language Components*  
Jim Melton and Alan R. Simon

*Information Visualization in Data Mining and Knowledge Discovery*  
Edited by Usama Fayyad, Georges G. Grinstein, and Andreas Wierse

*Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery*  
Gerhard Weikum and Gottfried Vossen

*Spatial Databases: With Application to GIS*  
Philippe Rigaux, Michel Scholl, and Agnès Voisard

*Information Modeling and Relational Databases: From Conceptual Analysis to Logical Design*  
Terry Halpin

*Component Database Systems*  
Edited by Klaus R. Dittrich and Andreas Geppert

*Managing Reference Data in Enterprise Databases: Binding Corporate Data to the Wider World*  
Malcolm Chisholm

*Data Mining: Concepts and Techniques*  
Jiawei Han and Micheline Kamber

*Understanding SQL and Java Together: A Guide to SQLJ, JDBC, and Related Technologies*  
Jim Melton and Andrew Eisenberg

*Database: Principles, Programming, and Performance, Second Edition*  
Patrick O'Neil and Elizabeth O'Neil

*The Object Data Standard: ODMG 3.0*  
Edited by R. G. G. Cattell, Douglas K. Barry, Mark Berler, Jeff Eastman, David Jordan, Craig Russell, Olaf Schadow, Torsten Stanienda, and Fernando Velez

*Data on the Web: From Relations to Semistructured Data and XML*  
Serge Abiteboul, Peter Buneman, and Dan Suciu

*Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*  
Ian H. Witten and Eibe Frank

*Joe Celko's SQL for Smarties: Advanced SQL Programming, Second Edition*  
Joe Celko

*Joe Celko's Data and Databases: Concepts in Practice*  
Joe Celko

*Developing Time-Oriented Database Applications in SQL*  
Richard T. Snodgrass

*Web Farming for the Data Warehouse*  
Richard D. Hackathorn

*Database Modeling & Design, Third Edition*  
Toby J. Teorey

*Management of Heterogeneous and Autonomous Database Systems*  
Edited by Ahmed Elmagarmid, Marek Rusinkiewicz, and Amit Sheth

*Object-Relational DBMSs: Tracking the Next Great Wave, Second Edition*  
Michael Stonebraker and Paul Brown, with Dorothy Moore

*A Complete Guide to DB2 Universal Database*  
Don Chamberlin

*Universal Database Management: A Guide to Object/Relational Technology*  
Cynthia Maro Saracco

*Readings in Database Systems, Third Edition*  
Edited by Michael Stonebraker and Joseph M. Hellerstein

*Understanding SQL's Stored Procedures: A Complete Guide to SQL/PSM*  
Jim Melton

*Principles of Multimedia Database Systems*  
V. S. Subrahmanian

*Principles of Database Query Processing for Advanced Applications*  
Clement T. Yu and Weiyi Meng

*Advanced Database Systems*  
Carlo Zaniolo, Stefano Ceri, Christos Faloutsos, Richard T. Snodgrass, V. S. Subrahmanian, and Roberto Zicari

*Principles of Transaction Processing for the Systems Professional*  
Philip A. Bernstein and Eric Newcomer

*Using the New DB2: IBM's Object-Relational Database System*  
Don Chamberlin

*Distributed Algorithms*  
Nancy A. Lynch

*Active Database Systems: Triggers and Rules For Advanced Database Processing*  
Edited by Jennifer Widom and Stefano Ceri

*Migrating Legacy Systems: Gateways, Interfaces & the Incremental Approach*  
Michael L. Brodie and Michael Stonebraker

*Atomic Transactions*  
Nancy Lynch, Michael Merritt, William Weihl, and Alan Fekete

*Query Processing For Advanced Database Systems*  
Edited by Johann Christoph Freytag, David Maier, and Gottfried Vossen

*Transaction Processing: Concepts and Techniques*  
Jim Gray and Andreas Reuter

*Building an Object-Oriented Database System: The Story of O<sub>2</sub>*  
Edited by François Bancilhon, Claude Delobel, and Paris Kanellakis

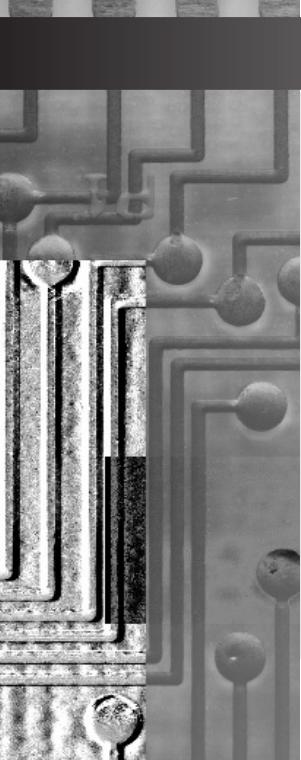
*Database Transaction Models For Advanced Applications*  
Edited by Ahmed K. Elmagarmid

*A Guide to Developing Client/Server SQL Applications*  
Setrag Khoshafian, Arvola Chan, Anna Wong, and Harry K. T. Wong

*The Benchmark Handbook For Database and Transaction Processing Systems, Second Edition*  
Edited by Jim Gray

*Camelot and Avalon: A Distributed Transaction Facility*  
Edited by Jeffrey L. Eppinger, Lily B. Mummert, and Alfred Z. Spector

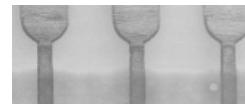
*Readings in Object-Oriented Database Systems*  
Edited by Stanley B. Zdonik and David Maier



# Data Mining

---

**Practical Machine Learning Tools and Techniques,  
Second Edition**



**Ian H. Witten**

Department of Computer Science  
University of Waikato

**Eibe Frank**

Department of Computer Science  
University of Waikato



AMSTERDAM • BOSTON • HEIDELBERG • LONDON  
NEW YORK • OXFORD • PARIS • SAN DIEGO  
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO  
MORGAN KAUFMANN PUBLISHERS IS AN IMPRINT OF ELSEVIER



MORGAN KAUFMANN PUBLISHERS

Publisher:	Diane Cerra
Publishing Services Manager:	Simon Crump
Project Manager:	Brandy Lilly
Editorial Assistant:	Asma Stephan
Cover Design:	Yvo Riezebos Design
Cover Image:	Getty Images
Composition:	SNP Best-set Typesetter Ltd., Hong Kong
Technical Illustration:	Dartmouth Publishing, Inc.
Copyeditor:	Graphic World Inc.
Proofreader:	Graphic World Inc.
Indexer:	Graphic World Inc.
Interior printer:	The Maple-Vail Book Manufacturing Group
Cover printer:	Phoenix Color Corp

Morgan Kaufmann Publishers is an imprint of Elsevier.  
500 Sansome Street, Suite 400, San Francisco, CA 94111

This book is printed on acid-free paper.

© 2005 by Elsevier Inc. All rights reserved.

Designations used by companies to distinguish their products are often claimed as trademarks or registered trademarks. In all instances in which Morgan Kaufmann Publishers is aware of a claim, the product names appear in initial capital or all capital letters. Readers, however, should contact the appropriate companies for more complete information regarding trademarks and registration.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopying, scanning, or otherwise—without prior written permission of the publisher.

Permissions may be sought directly from Elsevier's Science & Technology Rights Department in Oxford, UK: phone: (+44) 1865 843830, fax: (+44) 1865 853333, e-mail: permissions@elsevier.com.uk. You may also complete your request on-line via the Elsevier homepage (<http://elsevier.com>) by selecting "Customer Support" and then "Obtaining Permissions."

#### Library of Congress Cataloging-in-Publication Data

Witten, I. H. (Ian H.)

Data mining : practical machine learning tools and techniques / Ian H. Witten, Eibe

Frank. — 2nd ed.

p. cm. — (Morgan Kaufmann series in data management systems)

Includes bibliographical references and index.

ISBN: 0-12-088407-0

1. Data mining. I. Frank, Eibe. II. Title. III. Series.

QA76.9.D343W58 2005

006.3—dc22

2005043385

For information on all Morgan Kaufmann publications,  
visit our Web site at [www.mkp.com](http://www.mkp.com) or [www.books.elsevier.com](http://www.books.elsevier.com)

Printed in the United States of America  
05 06 07 08 09 5 4 3 2 1

Working together to grow  
libraries in developing countries

[www.elsevier.com](http://www.elsevier.com) | [www.bookaid.org](http://www.bookaid.org) | [www.sabre.org](http://www.sabre.org)

ELSEVIER

BOOK AID  
International

Sabre Foundation

# Foreword

Jim Gray, Series Editor  
Microsoft Research

Technology now allows us to capture and store vast quantities of data. Finding patterns, trends, and anomalies in these datasets, and summarizing them with simple quantitative models, is one of the grand challenges of the information age—turning data into information and turning information into knowledge.

There has been stunning progress in data mining and machine learning. The synthesis of statistics, machine learning, information theory, and computing has created a solid science, with a firm mathematical base, and with very powerful tools. Witten and Frank present much of this progress in this book and in the companion implementation of the key algorithms. As such, this is a milestone in the synthesis of data mining, data analysis, information theory, and machine learning. If you have not been following this field for the last decade, this is a great way to catch up on this exciting progress. If you have, then Witten and Frank’s presentation and the companion open-source workbench, called Weka, will be a useful addition to your toolkit.

They present the basic theory of automatically extracting models from data, and then validating those models. The book does an excellent job of explaining the various models (decision trees, association rules, linear models, clustering, Bayes nets, neural nets) and how to apply them in practice. With this basis, they then walk through the steps and pitfalls of various approaches. They describe how to safely scrub datasets, how to build models, and how to evaluate a model’s predictive quality. Most of the book is tutorial, but Part II broadly describes how commercial systems work and gives a tour of the publicly available data mining workbench that the authors provide through a website. This Weka workbench has a graphical user interface that leads you through data mining tasks and has excellent data visualization tools that help understand the models. It is a great companion to the text and a useful and popular tool in its own right.

This book presents this new discipline in a very accessible form: as a text both to train the next generation of practitioners and researchers and to inform lifelong learners like myself. Witten and Frank have a passion for simple and elegant solutions. They approach each topic with this mindset, grounding all concepts in concrete examples, and urging the reader to consider the simple techniques first, and then progress to the more sophisticated ones if the simple ones prove inadequate.

If you are interested in databases, and have not been following the machine learning field, this book is a great way to catch up on this exciting progress. If you have data that you want to analyze and understand, this book and the associated Weka toolkit are an excellent way to start.

# Contents

**Foreword** v

**Preface** xxiii

*Updated and revised content* xxvii

*Acknowledgments* xxix

## Part I Machine learning tools and techniques 1

### 1 What's it all about? 3

1.1 Data mining and machine learning 4

*Describing structural patterns* 6

*Machine learning* 7

*Data mining* 9

1.2 Simple examples: The weather problem and others 9

*The weather problem* 10

*Contact lenses: An idealized problem* 13

*Irises: A classic numeric dataset* 15

*CPU performance: Introducing numeric prediction* 16

*Labor negotiations: A more realistic example* 17

*Soybean classification: A classic machine learning success* 18

1.3 Fielded applications 22

*Decisions involving judgment* 22

*Screening images* 23

*Load forecasting* 24

*Diagnosis* 25

*Marketing and sales* 26

*Other applications* 28

1.4	Machine learning and statistics	29
1.5	Generalization as search	30
	<i>Enumerating the concept space</i>	31
	<i>Bias</i>	32
1.6	Data mining and ethics	35
1.7	Further reading	37

## **2 Input: Concepts, instances, and attributes 41**

2.1	What's a concept?	42
2.2	What's in an example?	45
2.3	What's in an attribute?	49
2.4	Preparing the input	52
	<i>Gathering the data together</i>	52
	<i>ARFF format</i>	53
	<i>Sparse data</i>	55
	<i>Attribute types</i>	56
	<i>Missing values</i>	58
	<i>Inaccurate values</i>	59
	<i>Getting to know your data</i>	60
2.5	Further reading	60

## **3 Output: Knowledge representation 61**

3.1	Decision tables	62
3.2	Decision trees	62
3.3	Classification rules	65
3.4	Association rules	69
3.5	Rules with exceptions	70
3.6	Rules involving relations	73
3.7	Trees for numeric prediction	76
3.8	Instance-based representation	76
3.9	Clusters	81
3.10	Further reading	82

<b>4 Algorithms: The basic methods</b>	<b>83</b>
<b>4.1 Inferring rudimentary rules</b>	<b>84</b>
<i>Missing values and numeric attributes</i>	86
<i>Discussion</i>	88
<b>4.2 Statistical modeling</b>	<b>88</b>
<i>Missing values and numeric attributes</i>	92
<i>Bayesian models for document classification</i>	94
<i>Discussion</i>	96
<b>4.3 Divide-and-conquer: Constructing decision trees</b>	<b>97</b>
<i>Calculating information</i>	100
<i>Highly branching attributes</i>	102
<i>Discussion</i>	105
<b>4.4 Covering algorithms: Constructing rules</b>	<b>105</b>
<i>Rules versus trees</i>	107
<i>A simple covering algorithm</i>	107
<i>Rules versus decision lists</i>	111
<b>4.5 Mining association rules</b>	<b>112</b>
<i>Item sets</i>	113
<i>Association rules</i>	113
<i>Generating rules efficiently</i>	117
<i>Discussion</i>	118
<b>4.6 Linear models</b>	<b>119</b>
<i>Numeric prediction: Linear regression</i>	119
<i>Linear classification: Logistic regression</i>	121
<i>Linear classification using the perceptron</i>	124
<i>Linear classification using Winnow</i>	126
<b>4.7 Instance-based learning</b>	<b>128</b>
<i>The distance function</i>	128
<i>Finding nearest neighbors efficiently</i>	129
<i>Discussion</i>	135
<b>4.8 Clustering</b>	<b>136</b>
<i>Iterative distance-based clustering</i>	137
<i>Faster distance calculations</i>	138
<i>Discussion</i>	139
<b>4.9 Further reading</b>	<b>139</b>

## 5 Credibility: Evaluating what's been learned 143

- 5.1 Training and testing 144
- 5.2 Predicting performance 146
- 5.3 Cross-validation 149
- 5.4 Other estimates 151
  - Leave-one-out* 151
  - The bootstrap* 152
- 5.5 Comparing data mining methods 153
- 5.6 Predicting probabilities 157
  - Quadratic loss function* 158
  - Informational loss function* 159
  - Discussion* 160
- 5.7 Counting the cost 161
  - Cost-sensitive classification* 164
  - Cost-sensitive learning* 165
  - Lift charts* 166
  - ROC curves* 168
  - Recall–precision curves* 171
  - Discussion* 172
  - Cost curves* 173
- 5.8 Evaluating numeric prediction 176
- 5.9 The minimum description length principle 179
- 5.10 Applying the MDL principle to clustering 183
- 5.11 Further reading 184

## 6 Implementations: Real machine learning schemes 187

- 6.1 Decision trees 189
  - Numeric attributes* 189
  - Missing values* 191
  - Pruning* 192
  - Estimating error rates* 193
  - Complexity of decision tree induction* 196
  - From trees to rules* 198
  - C4.5: Choices and options* 198
  - Discussion* 199
- 6.2 Classification rules 200
  - Criteria for choosing tests* 200
  - Missing values, numeric attributes* 201

<i>Generating good rules</i>	202
<i>Using global optimization</i>	205
<i>Obtaining rules from partial decision trees</i>	207
<i>Rules with exceptions</i>	210
<i>Discussion</i>	213
<b>6.3 Extending linear models</b>	<b>214</b>
<i>The maximum margin hyperplane</i>	215
<i>Nonlinear class boundaries</i>	217
<i>Support vector regression</i>	219
<i>The kernel perceptron</i>	222
<i>Multilayer perceptrons</i>	223
<i>Discussion</i>	235
<b>6.4 Instance-based learning</b>	<b>235</b>
<i>Reducing the number of exemplars</i>	236
<i>Pruning noisy exemplars</i>	236
<i>Weighting attributes</i>	237
<i>Generalizing exemplars</i>	238
<i>Distance functions for generalized exemplars</i>	239
<i>Generalized distance functions</i>	241
<i>Discussion</i>	242
<b>6.5 Numeric prediction</b>	<b>243</b>
<i>Model trees</i>	244
<i>Building the tree</i>	245
<i>Pruning the tree</i>	245
<i>Nominal attributes</i>	246
<i>Missing values</i>	246
<i>Pseudocode for model tree induction</i>	247
<i>Rules from model trees</i>	250
<i>Locally weighted linear regression</i>	251
<i>Discussion</i>	253
<b>6.6 Clustering</b>	<b>254</b>
<i>Choosing the number of clusters</i>	254
<i>Incremental clustering</i>	255
<i>Category utility</i>	260
<i>Probability-based clustering</i>	262
<i>The EM algorithm</i>	265
<i>Extending the mixture model</i>	266
<i>Bayesian clustering</i>	268
<i>Discussion</i>	270
<b>6.7 Bayesian networks</b>	<b>271</b>
<i>Making predictions</i>	272
<i>Learning Bayesian networks</i>	276

<i>Specific algorithms</i>	278
<i>Data structures for fast learning</i>	280
<i>Discussion</i>	283

## 7 Transformations: Engineering the input and output 285

<b>7.1 Attribute selection</b>	<b>288</b>
<i>Scheme-independent selection</i>	290
<i>Searching the attribute space</i>	292
<i>Scheme-specific selection</i>	294
<b>7.2 Discretizing numeric attributes</b>	<b>296</b>
<i>Unsupervised discretization</i>	297
<i>Entropy-based discretization</i>	298
<i>Other discretization methods</i>	302
<i>Entropy-based versus error-based discretization</i>	302
<i>Converting discrete to numeric attributes</i>	304
<b>7.3 Some useful transformations</b>	<b>305</b>
<i>Principal components analysis</i>	306
<i>Random projections</i>	309
<i>Text to attribute vectors</i>	309
<i>Time series</i>	311
<b>7.4 Automatic data cleansing</b>	<b>312</b>
<i>Improving decision trees</i>	312
<i>Robust regression</i>	313
<i>Detecting anomalies</i>	314
<b>7.5 Combining multiple models</b>	<b>315</b>
<i>Bagging</i>	316
<i>Bagging with costs</i>	319
<i>Randomization</i>	320
<i>Boosting</i>	321
<i>Additive regression</i>	325
<i>Additive logistic regression</i>	327
<i>Option trees</i>	328
<i>Logistic model trees</i>	331
<i>Stacking</i>	332
<i>Error-correcting output codes</i>	334
<b>7.6 Using unlabeled data</b>	<b>337</b>
<i>Clustering for classification</i>	337
<i>Co-training</i>	339
<i>EM and co-training</i>	340
<b>7.7 Further reading</b>	<b>341</b>

**8 Moving on: Extensions and applications 345**

- 8.1 Learning from massive datasets 346
- 8.2 Incorporating domain knowledge 349
- 8.3 Text and Web mining 351
- 8.4 Adversarial situations 356
- 8.5 Ubiquitous data mining 358
- 8.6 Further reading 361

**Part II The Weka machine learning workbench 363****9 Introduction to Weka 365**

- 9.1 What's in Weka? 366
- 9.2 How do you use it? 367
- 9.3 What else can you do? 368
- 9.4 How do you get it? 368

**10 The Explorer 369**

- 10.1 Getting started 369
  - Preparing the data* 370
  - Loading the data into the Explorer* 370
  - Building a decision tree* 373
  - Examining the output* 373
  - Doing it again* 377
  - Working with models* 377
  - When things go wrong* 378
- 10.2 Exploring the Explorer 380
  - Loading and filtering files* 380
  - Training and testing learning schemes* 384
  - Do it yourself: The User Classifier* 388
  - Using a metalearner* 389
  - Clustering and association rules* 391
  - Attribute selection* 392
  - Visualization* 393
- 10.3 Filtering algorithms 393
  - Unsupervised attribute filters* 395
  - Unsupervised instance filters* 400
  - Supervised filters* 401

<b>10.4</b>	<b>Learning algorithms</b>	<b>403</b>
	<i>Bayesian classifiers</i>	403
	<i>Trees</i>	406
	<i>Rules</i>	408
	<i>Functions</i>	409
	<i>Lazy classifiers</i>	413
	<i>Miscellaneous classifiers</i>	414
<b>10.5</b>	<b>Metalearning algorithms</b>	<b>414</b>
	<i>Bagging and randomization</i>	414
	<i>Boosting</i>	416
	<i>Combining classifiers</i>	417
	<i>Cost-sensitive learning</i>	417
	<i>Optimizing performance</i>	417
	<i>Retargeting classifiers for different tasks</i>	418
<b>10.6</b>	<b>Clustering algorithms</b>	<b>418</b>
<b>10.7</b>	<b>Association-rule learners</b>	<b>419</b>
<b>10.8</b>	<b>Attribute selection</b>	<b>420</b>
	<i>Attribute subset evaluators</i>	422
	<i>Single-attribute evaluators</i>	422
	<i>Search methods</i>	423

## **11 The Knowledge Flow interface** **427**

<b>11.1</b>	<b>Getting started</b>	<b>427</b>
<b>11.2</b>	<b>The Knowledge Flow components</b>	<b>430</b>
<b>11.3</b>	<b>Configuring and connecting the components</b>	<b>431</b>
<b>11.4</b>	<b>Incremental learning</b>	<b>433</b>

## **12 The Experimenter** **437**

<b>12.1</b>	<b>Getting started</b>	<b>438</b>
	<i>Running an experiment</i>	439
	<i>Analyzing the results</i>	440
<b>12.2</b>	<b>Simple setup</b>	<b>441</b>
<b>12.3</b>	<b>Advanced setup</b>	<b>442</b>
<b>12.4</b>	<b>The Analyze panel</b>	<b>443</b>
<b>12.5</b>	<b>Distributing processing over several machines</b>	<b>445</b>

**13 The command-line interface 449**

- 13.1 Getting started 449
- 13.2 The structure of Weka 450
  - Classes, instances, and packages* 450
  - The weka.core package* 451
  - The weka.classifiers package* 453
  - Other packages* 455
  - Javadoc indices* 456
- 13.3 Command-line options 456
  - Generic options* 456
  - Scheme-specific options* 458

**14 Embedded machine learning 461**

- 14.1 A simple data mining application 461
- 14.2 Going through the code 462
  - main()* 462
  - MessageClassifier()* 462
  - updateData()* 468
  - classifyMessage()* 468

**15 Writing new learning schemes 471**

- 15.1 An example classifier 471
  - buildClassifier()* 472
  - makeTree()* 472
  - computeInfoGain()* 480
  - classifyInstance()* 480
  - main()* 481
- 15.2 Conventions for implementing classifiers 483

**References 485****Index 505****About the authors 525**



# List of Figures

Figure 1.1	Rules for the contact lens data.	13
Figure 1.2	Decision tree for the contact lens data.	14
Figure 1.3	Decision trees for the labor negotiations data.	19
Figure 2.1	A family tree and two ways of expressing the sister-of relation.	46
Figure 2.2	ARFF file for the weather data.	54
Figure 3.1	Constructing a decision tree interactively: (a) creating a rectangular test involving <i>petallength</i> and <i>petalwidth</i> and (b) the resulting (unfinished) decision tree.	64
Figure 3.2	Decision tree for a simple disjunction.	66
Figure 3.3	The exclusive-or problem.	67
Figure 3.4	Decision tree with a replicated subtree.	68
Figure 3.5	Rules for the Iris data.	72
Figure 3.6	The shapes problem.	73
Figure 3.7	Models for the CPU performance data: (a) linear regression, (b) regression tree, and (c) model tree.	77
Figure 3.8	Different ways of partitioning the instance space.	79
Figure 3.9	Different ways of representing clusters.	81
Figure 4.1	Pseudocode for 1R.	85
Figure 4.2	Tree stumps for the weather data.	98
Figure 4.3	Expanded tree stumps for the weather data.	100
Figure 4.4	Decision tree for the weather data.	101
Figure 4.5	Tree stump for the <i>ID code</i> attribute.	103
Figure 4.6	Covering algorithm: (a) covering the instances and (b) the decision tree for the same problem.	106
Figure 4.7	The instance space during operation of a covering algorithm.	108
Figure 4.8	Pseudocode for a basic rule learner.	111
Figure 4.9	Logistic regression: (a) the logit transform and (b) an example logistic regression function.	122

- Figure 4.10 The perceptron: (a) learning rule and (b) representation as a neural network. 125
- Figure 4.11 The Winnow algorithm: (a) the unbalanced version and (b) the balanced version. 127
- Figure 4.12 A *k*D-tree for four training instances: (a) the tree and (b) instances and splits. 130
- Figure 4.13 Using a *k*D-tree to find the nearest neighbor of the star. 131
- Figure 4.14 Ball tree for 16 training instances: (a) instances and balls and (b) the tree. 134
- Figure 4.15 Ruling out an entire ball (gray) based on a target point (star) and its current nearest neighbor. 135
- Figure 4.16 A ball tree: (a) two cluster centers and their dividing line and (b) the corresponding tree. 140
- Figure 5.1 A hypothetical lift chart. 168
- Figure 5.2 A sample ROC curve. 169
- Figure 5.3 ROC curves for two learning methods. 170
- Figure 5.4 Effects of varying the probability threshold: (a) the error curve and (b) the cost curve. 174
- Figure 6.1 Example of subtree raising, where node C is “raised” to subsume node B. 194
- Figure 6.2 Pruning the labor negotiations decision tree. 196
- Figure 6.3 Algorithm for forming rules by incremental reduced-error pruning. 205
- Figure 6.4 RIPPER: (a) algorithm for rule learning and (b) meaning of symbols. 206
- Figure 6.5 Algorithm for expanding examples into a partial tree. 208
- Figure 6.6 Example of building a partial tree. 209
- Figure 6.7 Rules with exceptions for the iris data. 211
- Figure 6.8 A maximum margin hyperplane. 216
- Figure 6.9 Support vector regression: (a)  $\epsilon = 1$ , (b)  $\epsilon = 2$ , and (c)  $\epsilon = 0.5$ . 221
- Figure 6.10 Example datasets and corresponding perceptrons. 225
- Figure 6.11 Step versus sigmoid: (a) step function and (b) sigmoid function. 228
- Figure 6.12 Gradient descent using the error function  $x^2 + 1$ . 229
- Figure 6.13 Multilayer perceptron with a hidden layer. 231
- Figure 6.14 A boundary between two rectangular classes. 240
- Figure 6.15 Pseudocode for model tree induction. 248
- Figure 6.16 Model tree for a dataset with nominal attributes. 250
- Figure 6.17 Clustering the weather data. 256

- Figure 6.18 Hierarchical clusterings of the iris data. 259  
Figure 6.19 A two-class mixture model. 264  
Figure 6.20 A simple Bayesian network for the weather data. 273  
Figure 6.21 Another Bayesian network for the weather data. 274  
Figure 6.22 The weather data: (a) reduced version and (b) corresponding AD tree. 281  
Figure 7.1 Attribute space for the weather dataset. 293  
Figure 7.2 Discretizing the *temperature* attribute using the entropy method. 299  
Figure 7.3 The result of discretizing the *temperature* attribute. 300  
Figure 7.4 Class distribution for a two-class, two-attribute problem. 303  
Figure 7.5 Principal components transform of a dataset: (a) variance of each component and (b) variance plot. 308  
Figure 7.6 Number of international phone calls from Belgium, 1950–1973. 314  
Figure 7.7 Algorithm for bagging. 319  
Figure 7.8 Algorithm for boosting. 322  
Figure 7.9 Algorithm for additive logistic regression. 327  
Figure 7.10 Simple option tree for the weather data. 329  
Figure 7.11 Alternating decision tree for the weather data. 330  
Figure 10.1 The Explorer interface. 370  
Figure 10.2 Weather data: (a) spreadsheet, (b) CSV format, and (c) ARFF. 371  
Figure 10.3 The Weka Explorer: (a) choosing the Explorer interface and (b) reading in the weather data. 372  
Figure 10.4 Using J4.8: (a) finding it in the classifiers list and (b) the *Classify* tab. 374  
Figure 10.5 Output from the J4.8 decision tree learner. 375  
Figure 10.6 Visualizing the result of J4.8 on the iris dataset: (a) the tree and (b) the classifier errors. 379  
Figure 10.7 Generic object editor: (a) the editor, (b) more information (click *More*), and (c) choosing a converter (click *Choose*). 381  
Figure 10.8 Choosing a filter: (a) the *filters* menu, (b) an object editor, and (c) more information (click *More*). 383  
Figure 10.9 The weather data with two attributes removed. 384  
Figure 10.10 Processing the CPU performance data with M5'. 385  
Figure 10.11 Output from the M5' program for numeric prediction. 386  
Figure 10.12 Visualizing the errors: (a) from M5' and (b) from linear regression. 388

- Figure 10.13 Working on the segmentation data with the User Classifier:  
    (a) the data visualizer and (b) the tree visualizer. 390
- Figure 10.14 Configuring a metalearner for boosting decision  
    stumps. 391
- Figure 10.15 Output from the Apriori program for association rules. 392
- Figure 10.16 Visualizing the Iris dataset. 394
- Figure 10.17 Using Weka's metalearner for discretization: (a) configuring  
    *FilteredClassifier*, and (b) the menu of filters. 402
- Figure 10.18 Visualizing a Bayesian network for the weather data (nominal  
    version): (a) default output, (b) a version with the  
    maximum number of parents set to 3 in the search  
    algorithm, and (c) probability distribution table for the  
    *windy* node in (b). 406
- Figure 10.19 Changing the parameters for J4.8. 407
- Figure 10.20 Using Weka's neural-network graphical user  
    interface. 411
- Figure 10.21 Attribute selection: specifying an evaluator and a search  
    method. 420
- Figure 11.1 The Knowledge Flow interface. 428
- Figure 11.2 Configuring a data source: (a) the right-click menu and  
    (b) the file browser obtained from the *Configure* menu  
    item. 429
- Figure 11.3 Operations on the Knowledge Flow components. 432
- Figure 11.4 A Knowledge Flow that operates incrementally: (a) the  
    configuration and (b) the strip chart output. 434
- Figure 12.1 An experiment: (a) setting it up, (b) the results file, and  
    (c) a spreadsheet with the results. 438
- Figure 12.2 Statistical test results for the experiment in  
    Figure 12.1. 440
- Figure 12.3 Setting up an experiment in advanced mode. 442
- Figure 12.4 Rows and columns of Figure 12.2: (a) row field, (b) column  
    field, (c) result of swapping the row and column selections,  
    and (d) substituting *Run* for *Dataset* as rows. 444
- Figure 13.1 Using Javadoc: (a) the front page and (b) the *weka.core*  
    package. 452
- Figure 13.2 *DecisionStump*: A class of the *weka.classifiers.trees*  
    package. 454
- Figure 14.1 Source code for the message classifier. 463
- Figure 15.1 Source code for the ID3 decision tree learner. 473

# List of Tables

Table 1.1	The contact lens data.	6
Table 1.2	The weather data.	11
Table 1.3	Weather data with some numeric attributes.	12
Table 1.4	The iris data.	15
Table 1.5	The CPU performance data.	16
Table 1.6	The labor negotiations data.	18
Table 1.7	The soybean data.	21
Table 2.1	Iris data as a clustering problem.	44
Table 2.2	Weather data with a numeric class.	44
Table 2.3	Family tree represented as a table.	47
Table 2.4	The sister-of relation represented in a table.	47
Table 2.5	Another relation represented as a table.	49
Table 3.1	A new iris flower.	70
Table 3.2	Training data for the shapes problem.	74
Table 4.1	Evaluating the attributes in the weather data.	85
Table 4.2	The weather data with counts and probabilities.	89
Table 4.3	A new day.	89
Table 4.4	The numeric weather data with summary statistics.	93
Table 4.5	Another new day.	94
Table 4.6	The weather data with identification codes.	103
Table 4.7	Gain ratio calculations for the tree stumps of Figure 4.2.	104
Table 4.8	Part of the contact lens data for which <i>astigmatism</i> = <i>yes</i> .	109
Table 4.9	Part of the contact lens data for which <i>astigmatism</i> = <i>yes</i> and <i>tear production rate</i> = <i>normal</i> .	110
Table 4.10	Item sets for the weather data with coverage 2 or greater.	114
Table 4.11	Association rules for the weather data.	116
Table 5.1	Confidence limits for the normal distribution.	148

Table 5.2	Confidence limits for Student's distribution with 9 degrees of freedom.	155
Table 5.3	Different outcomes of a two-class prediction.	162
Table 5.4	Different outcomes of a three-class prediction: (a) actual and (b) expected.	163
Table 5.5	Default cost matrixes: (a) a two-class case and (b) a three-class case.	164
Table 5.6	Data for a lift chart.	167
Table 5.7	Different measures used to evaluate the false positive versus the false negative tradeoff.	172
Table 5.8	Performance measures for numeric prediction.	178
Table 5.9	Performance measures for four numeric prediction models.	179
Table 6.1	Linear models in the model tree.	250
Table 7.1	Transforming a multiclass problem into a two-class one: (a) standard method and (b) error-correcting code.	335
Table 10.1	Unsupervised attribute filters.	396
Table 10.2	Unsupervised instance filters.	400
Table 10.3	Supervised attribute filters.	402
Table 10.4	Supervised instance filters.	402
Table 10.5	Classifier algorithms in Weka.	404
Table 10.6	Metalearning algorithms in Weka.	415
Table 10.7	Clustering algorithms.	419
Table 10.8	Association-rule learners.	419
Table 10.9	Attribute evaluation methods for attribute selection.	421
Table 10.10	Search methods for attribute selection.	421
Table 11.1	Visualization and evaluation components.	430
Table 13.1	Generic options for learning schemes in Weka.	457
Table 13.2	Scheme-specific options for the J4.8 decision tree learner.	458
Table 15.1	Simple learning schemes in Weka.	472

# Preface

The convergence of computing and communication has produced a society that feeds on information. Yet most of the information is in its raw form: data. If *data* is characterized as recorded facts, then *information* is the set of patterns, or expectations, that underlie the data. There is a huge amount of information locked up in databases—information that is potentially important but has not yet been discovered or articulated. Our mission is to bring it forth.

Data mining is the extraction of implicit, previously unknown, and potentially useful information from data. The idea is to build computer programs that sift through databases automatically, seeking regularities or patterns. Strong patterns, if found, will likely generalize to make accurate predictions on future data. Of course, there will be problems. Many patterns will be banal and uninteresting. Others will be spurious, contingent on accidental coincidences in the particular dataset used. In addition real data is imperfect: Some parts will be garbled, and some will be missing. Anything discovered will be inexact: There will be exceptions to every rule and cases not covered by any rule. Algorithms need to be robust enough to cope with imperfect data and to extract regularities that are inexact but useful.

Machine learning provides the technical basis of data mining. It is used to extract information from the raw data in databases—information that is expressed in a comprehensible form and can be used for a variety of purposes. The process is one of abstraction: taking the data, warts and all, and inferring whatever structure underlies it. This book is about the tools and techniques of machine learning used in practical data mining for finding, and describing, structural patterns in data.

As with any burgeoning new technology that enjoys intense commercial attention, the use of data mining is surrounded by a great deal of hype in the technical—and sometimes the popular—press. Exaggerated reports appear of the secrets that can be uncovered by setting learning algorithms loose on oceans of data. But there is no magic in machine learning, no hidden power, no

alchemy. Instead, there is an identifiable body of simple and practical techniques that can often extract useful information from raw data. This book describes these techniques and shows how they work.

We interpret machine learning as the acquisition of structural descriptions from examples. The kind of descriptions found can be used for prediction, explanation, and understanding. Some data mining applications focus on prediction: forecasting what will happen in new situations from data that describe what happened in the past, often by guessing the classification of new examples. But we are equally—perhaps more—interested in applications in which the result of “learning” is an actual description of a structure that can be used to classify examples. This structural description supports explanation, understanding, and prediction. In our experience, insights gained by the applications’ users are of most interest in the majority of practical data mining applications; indeed, this is one of machine learning’s major advantages over classical statistical modeling.

The book explains a variety of machine learning methods. Some are pedagogically motivated: simple schemes designed to explain clearly how the basic ideas work. Others are practical: real systems used in applications today. Many are contemporary and have been developed only in the last few years.

A comprehensive software resource, written in the Java language, has been created to illustrate the ideas in the book. Called the Waikato Environment for Knowledge Analysis, or Weka<sup>1</sup> for short, it is available as source code on the World Wide Web at <http://www.cs.waikato.ac.nz/ml/weka>. It is a full, industrial-strength implementation of essentially all the techniques covered in this book. It includes illustrative code and working implementations of machine learning methods. It offers clean, spare implementations of the simplest techniques, designed to aid understanding of the mechanisms involved. It also provides a workbench that includes full, working, state-of-the-art implementations of many popular learning schemes that can be used for practical data mining or for research. Finally, it contains a framework, in the form of a Java class library, that supports applications that use embedded machine learning and even the implementation of new learning schemes.

The objective of this book is to introduce the tools and techniques for machine learning that are used in data mining. After reading it, you will understand what these techniques are and appreciate their strengths and applicability. If you wish to experiment with your own data, you will be able to do this easily with the Weka software.

---

<sup>1</sup> Found only on the islands of New Zealand, the *weka* (pronounced to rhyme with *Mecca*) is a flightless bird with an inquisitive nature.

The book spans the gulf between the intensely practical approach taken by trade books that provide case studies on data mining and the more theoretical, principle-driven exposition found in current textbooks on machine learning. (A brief description of these books appears in the *Further reading* section at the end of Chapter 1.) This gulf is rather wide. To apply machine learning techniques productively, you need to understand something about how they work; this is not a technology that you can apply blindly and expect to get good results. Different problems yield to different techniques, but it is rarely obvious which techniques are suitable for a given situation: you need to know something about the range of possible solutions. We cover an extremely wide range of techniques. We can do this because, unlike many trade books, this volume does not promote any particular commercial software or approach. We include a large number of examples, but they use illustrative datasets that are small enough to allow you to follow what is going on. Real datasets are far too large to show this (and in any case are usually company confidential). Our datasets are chosen not to illustrate actual large-scale practical problems but to help you understand what the different techniques do, how they work, and what their range of application is.

The book is aimed at the technically aware general reader interested in the principles and ideas underlying the current practice of data mining. It will also be of interest to information professionals who need to become acquainted with this new technology and to all those who wish to gain a detailed technical understanding of what machine learning involves. It is written for an eclectic audience of information systems practitioners, programmers, consultants, developers, information technology managers, specification writers, patent examiners, and curious laypeople—as well as students and professors—who need an easy-to-read book with lots of illustrations that describes what the major machine learning techniques are, what they do, how they are used, and how they work. It is practically oriented, with a strong “how to” flavor, and includes algorithms, code, and implementations. All those involved in practical data mining will benefit directly from the techniques described. The book is aimed at people who want to cut through to the reality that underlies the hype about machine learning and who seek a practical, nonacademic, unpretentious approach. We have avoided requiring any specific theoretical or mathematical knowledge except in some sections marked by a light gray bar in the margin. These contain optional material, often for the more technical or theoretically inclined reader, and may be skipped without loss of continuity.

The book is organized in layers that make the ideas accessible to readers who are interested in grasping the basics and to those who would like more depth of treatment, along with full details on the techniques covered. We believe that consumers of machine learning need to have some idea of how the algorithms they use work. It is often observed that data models are only as good as the person

who interprets them, and that person needs to know something about how the models are produced to appreciate the strengths, and limitations, of the technology. However, it is not necessary for all data model users to have a deep understanding of the finer details of the algorithms.

We address this situation by describing machine learning methods at successive levels of detail. You will learn the basic ideas, the topmost level, by reading the first three chapters. Chapter 1 describes, through examples, what machine learning is and where it can be used; it also provides actual practical applications. Chapters 2 and 3 cover the kinds of input and output—or *knowledge representation*—involved. Different kinds of output dictate different styles of algorithm, and at the next level Chapter 4 describes the basic methods of machine learning, simplified to make them easy to comprehend. Here the principles involved are conveyed in a variety of algorithms without getting into intricate details or tricky implementation issues. To make progress in the application of machine learning techniques to particular data mining problems, it is essential to be able to measure how well you are doing. Chapter 5, which can be read out of sequence, equips you to evaluate the results obtained from machine learning, addressing the sometimes complex issues involved in performance evaluation.

At the lowest and most detailed level, Chapter 6 exposes in naked detail the nitty-gritty issues of implementing a spectrum of machine learning algorithms, including the complexities necessary for them to work well in practice. Although many readers may want to ignore this detailed information, it is at this level that the full, working, tested implementations of machine learning schemes in Weka are written. Chapter 7 describes practical topics involved with engineering the input to machine learning—for example, selecting and discretizing attributes—and covers several more advanced techniques for refining and combining the output from different learning techniques. The final chapter of Part I looks to the future.

The book describes most methods used in practical machine learning. However, it does not cover reinforcement learning, because it is rarely applied in practical data mining; genetic algorithm approaches, because these are just an optimization technique; or relational learning and inductive logic programming, because they are rarely used in mainstream data mining applications.

The data mining system that illustrates the ideas in the book is described in Part II to clearly separate conceptual material from the practical aspects of how to use it. You can skip to Part II directly from Chapter 4 if you are in a hurry to analyze your data and don't want to be bothered with the technical details.

Java has been chosen for the implementations of machine learning techniques that accompany this book because, as an object-oriented programming language, it allows a uniform interface to learning schemes and methods for pre- and postprocessing. We have chosen Java instead of C++, Smalltalk, or other

object-oriented languages because programs written in Java can be run on almost any computer without having to be recompiled, having to undergo complicated installation procedures, or—worst of all—having to change the code. A Java program is compiled into byte-code that can be executed on any computer equipped with an appropriate interpreter. This interpreter is called the *Java virtual machine*. Java virtual machines—and, for that matter, Java compilers—are freely available for all important platforms.

Like all widely used programming languages, Java has received its share of criticism. Although this is not the place to elaborate on such issues, in several cases the critics are clearly right. However, of all currently available programming languages that are widely supported, standardized, and extensively documented, Java seems to be the best choice for the purpose of this book. Its main disadvantage is speed of execution—or lack of it. Executing a Java program is several times slower than running a corresponding program written in C language because the virtual machine has to translate the byte-code into machine code before it can be executed. In our experience the difference is a factor of three to five if the virtual machine uses a just-in-time compiler. Instead of translating each byte-code individually, a *just-in-time compiler* translates whole chunks of byte-code into machine code, thereby achieving significant speedup. However, if this is still too slow for your application, there are compilers that translate Java programs directly into machine code, bypassing the byte-code step. This code cannot be executed on other platforms, thereby sacrificing one of Java's most important advantages.

## Updated and revised content

We finished writing the first edition of this book in 1999 and now, in April 2005, are just polishing this second edition. The areas of data mining and machine learning have matured in the intervening years. Although the core of material in this edition remains the same, we have made the most of our opportunity to update it to reflect the changes that have taken place over 5 years. There have been errors to fix, errors that we had accumulated in our publicly available errata file. Surprisingly few were found, and we hope there are even fewer in this second edition. (The errata for the second edition may be found through the book's home page at <http://www.cs.waikato.ac.nz/ml/weka/book.html>.) We have thoroughly edited the material and brought it up to date, and we practically doubled the number of references. The most enjoyable part has been adding new material. Here are the highlights.

Bowing to popular demand, we have added comprehensive information on neural networks: the perceptron and closely related Winnow algorithm in Section 4.6 and the multilayer perceptron and backpropagation algorithm

in Section 6.3. We have included more recent material on implementing nonlinear decision boundaries using both the kernel perceptron and radial basis function networks. There is a new section on Bayesian networks, again in response to readers' requests, with a description of how to learn classifiers based on these networks and how to implement them efficiently using all-dimensions trees.

The Weka machine learning workbench that accompanies the book, a widely used and popular feature of the first edition, has acquired a radical new look in the form of an interactive interface—or rather, three separate interactive interfaces—that make it far easier to use. The primary one is the Explorer, which gives access to all of Weka's facilities using menu selection and form filling. The others are the Knowledge Flow interface, which allows you to design configurations for streamed data processing, and the Experimenter, with which you set up automated experiments that run selected machine learning algorithms with different parameter settings on a corpus of datasets, collect performance statistics, and perform significance tests on the results. These interfaces lower the bar for becoming a practicing data miner, and we include a full description of how to use them. However, the book continues to stand alone, independent of Weka, and to underline this we have moved all material on the workbench into a separate Part II at the end of the book.

In addition to becoming far easier to use, Weka has grown over the last 5 years and matured enormously in its data mining capabilities. It now includes an unparalleled range of machine learning algorithms and related techniques. The growth has been partly stimulated by recent developments in the field and partly led by Weka users and driven by demand. This puts us in a position in which we know a great deal about what actual users of data mining want, and we have capitalized on this experience when deciding what to include in this new edition.

The earlier chapters, containing more general and foundational material, have suffered relatively little change. We have added more examples of fielded applications to Chapter 1, a new subsection on sparse data and a little on string attributes and date attributes to Chapter 2, and a description of interactive decision tree construction, a useful and revealing technique to help you grapple with your data using manually built decision trees, to Chapter 3.

In addition to introducing linear decision boundaries for classification, the infrastructure for neural networks, Chapter 4 includes new material on multinomial Bayes models for document classification and on logistic regression. The last 5 years have seen great interest in data mining for text, and this is reflected in our introduction to string attributes in Chapter 2, multinomial Bayes for document classification in Chapter 4, and text transformations in Chapter 7. Chapter 4 includes a great deal of new material on efficient data structures for searching the instance space: *kD*-trees and the recently invented ball trees. These

are used to find nearest neighbors efficiently and to accelerate distance-based clustering.

Chapter 5 describes the principles of statistical evaluation of machine learning, which have not changed. The main addition, apart from a note on the Kappa statistic for measuring the success of a predictor, is a more detailed treatment of cost-sensitive learning. We describe how to use a classifier, built without taking costs into consideration, to make predictions that are sensitive to cost; alternatively, we explain how to take costs into account during the training process to build a cost-sensitive model. We also cover the popular new technique of cost curves.

There are several additions to Chapter 6, apart from the previously mentioned material on neural networks and Bayesian network classifiers. More details—gory details—are given of the heuristics used in the successful RIPPER rule learner. We describe how to use model trees to generate rules for numeric prediction. We show how to apply locally weighted regression to classification problems. Finally, we describe the X-means clustering algorithm, which is a big improvement on traditional  $k$ -means.

Chapter 7 on engineering the input and output has changed most, because this is where recent developments in practical machine learning have been concentrated. We describe new attribute selection schemes such as race search and the use of support vector machines and new methods for combining models such as additive regression, additive logistic regression, logistic model trees, and option trees. We give a full account of LogitBoost (which was mentioned in the first edition but not described). There is a new section on useful transformations, including principal components analysis and transformations for text mining and time series. We also cover recent developments in using unlabeled data to improve classification, including the co-training and co-EM methods.

The final chapter of Part I on new directions and different perspectives has been reworked to keep up with the times and now includes contemporary challenges such as adversarial learning and ubiquitous data mining.

## Acknowledgments

Writing the acknowledgments is always the nicest part! A lot of people have helped us, and we relish this opportunity to thank them. This book has arisen out of the machine learning research project in the Computer Science Department at the University of Waikato, New Zealand. We have received generous encouragement and assistance from the academic staff members on that project: John Cleary, Sally Jo Cunningham, Matt Humphrey, Lyn Hunt, Bob McQueen, Lloyd Smith, and Tony Smith. Special thanks go to Mark Hall, Bernhard Pfahringer, and above all Geoff Holmes, the project leader and source of inspi-

ration. All who have worked on the machine learning project here have contributed to our thinking: we would particularly like to mention Steve Garner, Stuart Inglis, and Craig Nevill-Manning for helping us to get the project off the ground in the beginning when success was less certain and things were more difficult.

The Weka system that illustrates the ideas in this book forms a crucial component of it. It was conceived by the authors and designed and implemented by Eibe Frank, along with Len Trigg and Mark Hall. Many people in the machine learning laboratory at Waikato made significant contributions. Since the first edition of the book the Weka team has expanded considerably: so many people have contributed that it is impossible to acknowledge everyone properly. We are grateful to Remco Bouckaert for his implementation of Bayesian networks, Dale Fletcher for many database-related aspects, Ashraf Kibriya and Richard Kirkby for contributions far too numerous to list, Niels Landwehr for logistic model trees, Abdelaziz Mahoui for the implementation of  $K^*$ , Stefan Mutter for association rule mining, Gabi Schmidberger and Malcolm Ware for numerous miscellaneous contributions, Tony Voyle for least-median-of-squares regression, Yong Wang for Pace regression and the implementation of  $M5'$ , and Xin Xu for *JRip*, logistic regression, and many other contributions. Our sincere thanks go to all these people for their dedicated work and to the many contributors to Weka from outside our group at Waikato.

Tucked away as we are in a remote (but very pretty) corner of the Southern Hemisphere, we greatly appreciate the visitors to our department who play a crucial role in acting as sounding boards and helping us to develop our thinking. We would like to mention in particular Rob Holte, Carl Gutwin, and Russell Beale, each of whom visited us for several months; David Aha, who although he only came for a few days did so at an early and fragile stage of the project and performed a great service by his enthusiasm and encouragement; and Kai Ming Ting, who worked with us for 2 years on many of the topics described in Chapter 7 and helped to bring us into the mainstream of machine learning.

Students at Waikato have played a significant role in the development of the project. Jamie Littin worked on ripple-down rules and relational learning. Brent Martin explored instance-based learning and nested instance-based representations. Murray Fife slaved over relational learning, and Nadeeka Madapathage investigated the use of functional languages for expressing machine learning algorithms. Other graduate students have influenced us in numerous ways, particularly Gordon Paynter, YingYing Wen, and Zane Bray, who have worked with us on text mining. Colleagues Steve Jones and Malika Mahoui have also made far-reaching contributions to these and other machine learning projects. More recently we have learned much from our many visiting students from Freiburg, including Peter Reutemann and Nils Weidmann.

Ian Witten would like to acknowledge the formative role of his former students at Calgary, particularly Brent Krawchuk, Dave Maulsby, Thong Phan, and Tanja Mitrovic, all of whom helped him develop his early ideas in machine learning, as did faculty members Bruce MacDonald, Brian Gaines, and David Hill at Calgary and John Andreea at the University of Canterbury.

Eibe Frank is indebted to his former supervisor at the University of Karlsruhe, Klaus-Peter Huber (now with SAS Institute), who infected him with the fascination of machines that learn. On his travels Eibe has benefited from interactions with Peter Turney, Joel Martin, and Berry de Bruijn in Canada and with Luc de Raedt, Christoph Helma, Kristian Kersting, Stefan Kramer, Ulrich Rückert, and Ashwin Srinivasan in Germany.

Diane Cerra and Asma Stephan of Morgan Kaufmann have worked hard to shape this book, and Lisa Royse, our production editor, has made the process go smoothly. Bronwyn Webster has provided excellent support at the Waikato end.

We gratefully acknowledge the unsung efforts of the anonymous reviewers, one of whom in particular made a great number of pertinent and constructive comments that helped us to improve this book significantly. In addition, we would like to thank the librarians of the Repository of Machine Learning Databases at the University of California, Irvine, whose carefully collected datasets have been invaluable in our research.

Our research has been funded by the New Zealand Foundation for Research, Science and Technology and the Royal Society of New Zealand Marsden Fund. The Department of Computer Science at the University of Waikato has generously supported us in all sorts of ways, and we owe a particular debt of gratitude to Mark Apperley for his enlightened leadership and warm encouragement. Part of the first edition was written while both authors were visiting the University of Calgary, Canada, and the support of the Computer Science department there is gratefully acknowledged—as well as the positive and helpful attitude of the long-suffering students in the machine learning course on whom we experimented.

In producing the second edition Ian was generously supported by Canada's Informatics Circle of Research Excellence and by the University of Lethbridge in southern Alberta, which gave him what all authors yearn for—a quiet space in pleasant and convivial surroundings in which to work.

Last, and most of all, we are grateful to our families and partners. Pam, Anna, and Nikki were all too well aware of the implications of having an author in the house (“not again!”) but let Ian go ahead and write the book anyway. Julie was always supportive, even when Eibe had to burn the midnight oil in the machine learning lab, and Immo and Ollig provided exciting diversions. Between us we hail from Canada, England, Germany, Ireland, and Samoa: New Zealand has brought us together and provided an ideal, even idyllic, place to do this work.

