



# LEARNING FROM DATA

**Concepts, Theory, and Methods**

SECOND EDITION



VLADIMIR CHERKASSKY • FILIP MULIER

TN911.7  
C521  
E-2

# LEARNING FROM DATA

---

## Concepts, Theory, and Methods

Second Edition

VLADIMIR CHERKASSKY  
FILIP MULIER



WILEY-INTERSCIENCE  
A JOHN WILEY & SONS, INC., PUBLICATION



E2009000172

Copyright © 2007 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey  
Published simultaneously in Canada

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, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400, fax 978-750-4470, or on the web at [www.copyright.com](http://www.copyright.com). Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, 201-748-6011, fax 201-748-6008, or online at <http://www.wiley.com/go/permission>.

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at 877-762-2974, outside the United States at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at [www.wiley.com](http://www.wiley.com).

Wiley Bicentennial Logo: Richard J. Pacifico

***Library of Congress Cataloging-in-Publication Data:***

Cherkassky, Vladimir S.

Learning from data : concepts, theory, and methods / by Vladimir Cherkassky,

Filip Mulier. – 2nd ed.

p. cm.

ISBN 978-0-471-68182-3 (cloth)

1. Adaptive signal processing. 2. Machine learning. 3. Neural networks (Computer science) 4. Fuzzy systems. I. Mulier, Filip. II. Title.

TK5102.9.C475 2007

006.3'1–dc22

2006038736

Printed in the United States of America  
10 9 8 7 6 5 4 3 2 1

# **LEARNING FROM DATA**



---

## THE WILEY BICENTENNIAL—KNOWLEDGE FOR GENERATIONS

---

Each generation has its unique needs and aspirations. When Charles Wiley first opened his small printing shop in lower Manhattan in 1807, it was a generation of boundless potential searching for an identity. And we were there, helping to define a new American literary tradition. Over half a century later, in the midst of the Second Industrial Revolution, it was a generation focused on building the future. Once again, we were there, supplying the critical scientific, technical, and engineering knowledge that helped frame the world. Throughout the 20th Century, and into the new millennium, nations began to reach out beyond their own borders and a new international community was born. Wiley was there, expanding its operations around the world to enable a global exchange of ideas, opinions, and know-how.

For 200 years, Wiley has been an integral part of each generation's journey, enabling the flow of information and understanding necessary to meet their needs and fulfill their aspirations. Today, bold new technologies are changing the way we live and learn. Wiley will be there, providing you the must-have knowledge you need to imagine new worlds, new possibilities, and new opportunities.

Generations come and go, but you can always count on Wiley to provide you the knowledge you need, when and where you need it!

**WILLIAM J. PESCE**  
PRESIDENT AND CHIEF EXECUTIVE OFFICER

**PETER BOOTH WILEY**  
CHAIRMAN OF THE BOARD

# PREFACE

There are two problems in modern science:

- too many people use different terminology to solve the same problems;
- even more people use the same terminology to address completely different issues.

Anonymous

In recent years, there has been an explosive growth of methods for learning (or estimating dependencies) from data. This is not surprising given the proliferation of

- low-cost computers (for implementing such methods in software)
- low-cost sensors and database technology (for collecting and storing data)
- highly computer-literate application experts (who can pose “interesting” application problems)

A learning method is an algorithm (usually implemented in software) that estimates an unknown mapping (dependency) between a system’s inputs and outputs from the available data, namely from known (input, output) samples. Once such a dependency has been accurately estimated, it can be used for prediction of future system outputs from the known input values. This book provides a unified description of principles and methods for learning dependencies from data.

Methods for estimating dependencies from data have been traditionally explored in diverse fields such as statistics (multivariate regression and classification), engineering (pattern recognition), and computer science (artificial intelligence, machine

learning, and, more recently, data mining). Recent interest in learning from data has resulted in the development of biologically motivated methodologies, such as artificial neural networks, fuzzy systems, and wavelets.

Unfortunately, developments in each field are seldom related to other fields, despite the apparent commonality of issues and methods. The mere fact that hundreds of “new” methods are being proposed each year at various conferences and in numerous journals suggests a certain lack of understanding of the basic issues common to all such methods.

The premise of this book is that there are just a handful of important principles and issues in the field of learning dependencies from data. Any researcher or practitioner in this field needs to be aware of these issues in order to successfully apply a particular methodology, understand a method’s limitations, or develop new techniques.

This book is an attempt to present and discuss such issues and principles (common to all methods) and then describe representative popular methods originating from statistics, neural networks, and pattern recognition. Often methods developed in different fields can be related to a common conceptual framework. This approach enables better understanding of a method’s properties, and it has methodological advantages over traditional “cookbook” descriptions of various learning algorithms.

Many aspects of learning methods can be addressed under a traditional statistical framework. At the same time, many popular learning algorithms and learning methodologies have been developed outside classical statistics. This happened for several reasons:

1. Traditionally, the statistician’s role has been to analyze the inferential limitations of the structural model constructed (proposed) by the application-domain expert. Consequently, the conceptual approach (adopted in statistics) is parameter estimation for model identification. For many real-life problems that require flexible estimation with finite samples, the statistical approach is fundamentally flawed. As shown in this book, learning with finite samples should be based on the framework known as risk minimization, rather than density estimation.
2. Statisticians have been late to recognize and appreciate the importance of computer-intensive approaches to data analysis. The growing use of computers has fundamentally changed the traditional boundaries between a statistician (data modeler) and a user (application expert). Nowadays, engineers and computer scientists successfully use sophisticated empirical data-modeling techniques (i.e., neural nets) to estimate complex nonlinear dependencies from the data.
3. Statistics (being part of mathematics) has developed into a closed discipline, with its own scientific jargon and academic objectives that favor analytic proofs rather than practical methods for learning from data.

Historically, we can identify three stages in the development of predictive learning methods. First, in 1985–1992 classical statistics gave way to neural networks (and other empirical methods, such as fuzzy systems) due to an early enthusiasm and naive claims that biologically inspired methods (i.e., neural nets) can achieve model-free learning not subject to statistical limitations. Even though such claims later proved to be false, this stage had a positive impact by showing the power and usefulness of flexible nonlinear modeling based on the risk minimization approach. Then in 1992–1996 came the return of statistics as the researchers and practitioners of neural networks became aware of their statistical limitations, initiating a trend toward interpretation of learning methods using a classical statistical framework. Finally, the third stage, from 1997 to present, is dominated by the wide popularity of support vector machines (SVMs) and similar margin-based approaches (such as boosting), and the growing interest in the Vapnik–Chervonenkis (VC) theoretical framework for predictive learning.

This book is intended for readers with varying interests, including researchers/practitioners in data modeling with a classical statistics background, researchers/practitioners in data modeling with a neural network background, and graduate students in engineering or computer science.

The presentation does not assume a special math background beyond a good working knowledge of probability, linear algebra, and calculus on an undergraduate level. Useful background material on optimization and linear algebra is included in Appendixes A and B, respectively. We do not provide mathematical proofs, but, whenever possible, in place of proofs we provide intuitive explanations and arguments. Likewise, mathematical formulation and discussion of the major concepts and results are provided as needed. The goal is to provide a unified treatment of diverse methodologies (i.e., statistics and neural networks), and to that end we carefully define the terminology used throughout the book. This book is not easy reading because it describes fairly complex concepts and mathematical models for solving inherently difficult (ill-posed) problems of learning with finite data. To aid the reader, each chapter starts with a brief overview of its contents. Also, each chapter is concluded with a summary containing an overview of open research issues and pointers to other (relevant) chapters.

Book chapters are conceptually organized into three parts:

- *Part I: Concepts and Theory* (Chapters 1–4). Following an introduction and motivation given in Chapter 1, we present formal specification of the inductive learning problem in Chapter 2 that also introduces major concepts and issues in learning from data. In particular, it describes an important concept called an *inductive principle*. Chapter 3 describes the regularization (or penalization) framework adopted in statistics. Chapter 4 describes Vapnik’s statistical learning theory (SLT), which provides the theoretical basis for predictive learning with finite data. SLT, aka VC theory, is important for understanding various learning methods developed in neural networks, statistics, and pattern recognition, and for developing new approaches, such as SVMs



(described in Chapter 9) and noninductive learning settings (described in Chapter 10).

- *Part II: Constructive Learning Methods* (Chapters 5–8). This part describes learning methods for regression, classification, and density approximation problems. The objective is to show conceptual similarity of methods originating from statistics, neural networks, and signal processing and to discuss their relative advantages and limitations. Whenever possible, we relate constructive learning methods to the conceptual framework of Part I. Chapter 5 describes nonlinear optimization strategies commonly used in various methods. Chapter 6 describes methods for density approximation, which include statistical, neural network, and signal processing techniques for data reduction and dimensionality reduction. Chapter 7 provides descriptions of statistical and neural network methods for regression. Chapter 8 describes methods for classification.
- *Part III: VC-Based Learning Methodologies* (Chapters 9 and 10). Here we describe constructive learning approaches that originate in VC theory. These include SVMs (or margin-based methods) for several inductive learning problems (in Chapter 9) and various noninductive learning formulations (described in Chapter 10).

The chapters should be followed in a sequential order, as the description of constructive learning methods is related to the conceptual framework developed in the first part of the book. A shortened sequence of Chapters 1–3 followed by Chapters 5, 6, 7 and 8 is recommended for the beginning readers who are interested only in the description of statistical and neural network methods. This sequence omits the mathematically and conceptually challenging Chapters 4 and 9. Alternatively, more advanced readers who are primarily interested in SLT and SVM methodology may adopt the sequence of Chapters 2, 3, 4, 9, and 10.

In the course of writing this book, our understanding of the field has changed. We started with the currently prevailing view of learning methods as a collection of tricks. Statisticians have their own bag of tricks (and terminology), neural networks have a different set of tricks, and so on. However, in the process of writing this book, we realized that it is possible to understand the various heuristic methods (tricks) by a sound general conceptual framework. Such a framework is provided by SLT developed mainly by Vapnik over the past 35 years. This theory combines fundamental concepts and principles related to learning with finite data, well-defined problem formulations, and rigorous mathematical theory. Although SLT is well known for its *mathematical* aspects, its *conceptual* contributions are not fully appreciated. As shown in our book, the conceptual framework provided by SLT can be used for improved understanding of various learning methods even where its mathematical results cannot be directly applied. Modern learning methods (i.e., flexible approaches using finite data) have slowly drifted away from the original problem statements posed in classical statistical decision and estimation theory. A major conceptual contribution of SLT is in revisiting the problem

statement appropriate for modern data mining applications. On the very basic level, SLT makes a clear distinction between the problem formulation and a solution approach (aka inductive principle) used to solve a problem. Although this distinction appears trivial on the surface, it leads to a fundamentally new understanding of the learning problem not explained by classical theory. Although it is tempting to skip directly to constructive solutions, this book devotes enough attention to the learning problem formulation and important concepts *before* describing actual learning methods.

Over the past 10 years (since the first edition of this book), we have witnessed considerable growth of interest in SVM-related methods. Nowadays, SVM (aka kernel) methods are commonly used in data mining, statistics, signal processing, pattern recognition, genomics, and so on. In spite of such an overwhelming success and wide recognition of SVM methodology, many important VC theoretical concepts responsible for good generalization of SVMs (such as margin, VC dimension) remain rather poorly understood. For example, many recent monographs and research papers refer to SVMs as a “special case of regularization.” So in this second edition, we made a special effort to emphasize the conceptual aspects of VC theory and to contrast the VC theoretical approach to learning (i.e., *system imitation*) versus the classical statistical and function approximation approach (i.e., *system identification*). Accurate interpretation of VC theoretical concepts is important for improved understanding of inductive learning algorithms, as well as for developing emerging state-of-the-art approaches based on noninductive learning settings (as discussed in Chapter 10). In this edition, we emphasize the philosophical interpretation of predictive learning, in general, and of several VC theoretical concepts, in particular. These philosophical connections appear to be quite useful for understanding recent advanced learning methods and for motivating new noninductive types of inference. Moreover, philosophical aspects of predictive learning can be immediately related to epistemology (understanding of human knowledge), as discussed in Chapter 11.

Many people have contributed directly and indirectly to this book. First and foremost, we are greatly indebted to Vladimir Vapnik of NEC Labs for his fundamental contributions to SLT and for his patience in explaining this theory to us. We would like to acknowledge many people whose constructive feedback helped improve the quality of the second edition, including Ella Bingham, John Boik, Olivier Chapelle, David Hand, Nicol Schraudolph, Simon Haykin, David Musicant, Erinija Pranceviciene, and D. Solomatine—all of whom provided many useful comments.

This book was used in the graduate course “Predictive Learning from Data” at the University of Minnesota over the past 10 years, and we would like to thank students who took this course for their valuable feedback. In particular, we acknowledge former graduate students X. Shao, Y. Ma, T. Xiong, L. Liang, H Gao, M. Ramani, R. Singh, and Y. Kim whose research contributions are incorporated in this book in the form of several fine figures and empirical

comparisons. Finally, we would like to thank our families for their patience and support.

**Vladimir Cherkassky**  
**Filip Mulier**

*Minneapolis, Minnesota*  
*March 2007*

# NOTATION

The following uniform notation is used throughout the book. Scalars are indicated by script letters such as  $a$ . Vectors are indicated by lowercase bold letters such as  $\mathbf{w}$ . Matrices are given using uppercase bold letters  $\mathbf{V}$ . When elements of a matrix are accessed individually, we use the corresponding lowercase script letter. For example, the  $(i, j)$  element of the matrix  $\mathbf{V}$  is  $v_{ij}$ . Common notation for all chapters is as follows:

## Data

$n$	Number of samples
$d$	Number of input variables
$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$	Matrix of input samples
$\mathbf{y} = [y_1, \dots, y_n]$	Vector of output samples
$\mathbf{Z} = [\mathbf{X}, \mathbf{y}]$	Combined input–output training data or
$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]$	Representation of data points in a feature space

## Distribution

$P$	Probability
$F(\mathbf{x})$	Cumulative probability distribution function (cdf)
$p(\mathbf{x})$	Probability density function (pdf)
$p(\mathbf{x}, y)$	Joint probability density function
$p(\mathbf{x}; \omega)$	Probability density function, which is parameterized
$p(y \mathbf{x})$	Conditional density
$t(\mathbf{x})$	Target function

## Approximating Functions

$f(\mathbf{x}, \omega)$	A class of approximating functions indexed by abstract parameter $\omega$ ( $\omega$ can be a scalar, vector, or matrix). Interpretation of $f(\mathbf{x}, \omega)$ depends on the particular learning problem
-------------------------	--

$f(\mathbf{x}, \omega_0)$	The function that minimizes the expected risk (optimal solution)
$f(\mathbf{x}, \omega^*)$	Estimate of the optimal solution obtained from finite data
$f(\mathbf{x}, \mathbf{w}, \mathbf{V}) = \sum_{i=1}^m w_i g_i(\mathbf{x}, \mathbf{v}_i) + b$	Basis function expansion of approximating functions with bias term
$g_i(\mathbf{x}, \mathbf{v})$	Basis function in a basis function expansion
$w, \mathbf{w}, \mathbf{W}$	Parameters of approximating function
$\mathbf{v}, \mathbf{v}, \mathbf{V}$	Basis function parameters
$m$	Number of basis functions
$\Omega$	Set of parameters, as in $\mathbf{w} \in \Omega$
$\Delta$	Margin distance
$t(\mathbf{x})$	Target function
$\zeta$	Error between the target function and the approximating function, or error between model estimate and time output

### Risk Functionals

$L(y, f(\mathbf{x}, \omega))$	Discrepancy measure or loss function
$L_2$	Squared discrepancy measure
$Q(\omega)$	A set of loss functions
$R$	Risk or average loss
$R(\omega)$	Expected risk as a function of parameters
$R_{\text{emp}}(\omega)$	Empirical risk as a function of parameters

### Kernel Functions

$K(\mathbf{x}, \mathbf{x}')$	General kernel function (for kernel smothing)
$S(\mathbf{x}, \mathbf{x}')$	Equivalent kernel of a linear estimator
$H(\mathbf{x}, \mathbf{x}')$	Inner product kernel

### Miscellaneous

$(\mathbf{a} \cdot \mathbf{b})$	Inner (dot) product of two vectors
$I()$	Indicator function of a Boolean argument that takes the value 1 if its argument is true and 0 otherwise. By convention, for a real-valued argument, $I(x) = 1$ for $x > 0$ , and $I(x) = 0$ for $x \leq 0$
$\phi[f(\mathbf{x}, \omega)]$	Penalty functional
$\lambda$	Regularization parameter
$h$	VC dimension
$\gamma_k$	Learning rate for stochastic approximation at iteration step $k$
$[a]_+$	Positive argument, equals $\max(a, 0)$
$\mathcal{L}$	Lagrangian

In addition to the above notation used throughout the book, there is chapter-specific notation, which will be introduced locally in each chapter.

# CONTENTS

<b>PREFACE</b>	<b>xi</b>
<b>NOTATION</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Learning and Statistical Estimation, 2	
1.2 Statistical Dependency and Causality, 7	
1.3 Characterization of Variables, 10	
1.4 Characterization of Uncertainty, 11	
1.5 Predictive Learning versus Other Data Analytical Methodologies, 14	
<b>2 Problem Statement, Classical Approaches, and Adaptive Learning</b>	<b>19</b>
2.1 Formulation of the Learning Problem, 21	
2.1.1 Objective of Learning, 24	
2.1.2 Common Learning Tasks, 25	
2.1.3 Scope of the Learning Problem Formulation, 29	
2.2 Classical Approaches, 30	
2.2.1 Density Estimation, 30	
2.2.2 Classification, 32	
2.2.3 Regression, 34	
2.2.4 Solving Problems with Finite Data, 34	
2.2.5 Nonparametric Methods, 36	
2.2.6 Stochastic Approximation, 39	

2.3	Adaptive Learning: Concepts and Inductive Principles,	40
2.3.1	Philosophy, Major Concepts, and Issues,	40
2.3.2	A Priori Knowledge and Model Complexity,	43
2.3.3	Inductive Principles,	45
2.3.4	Alternative Learning Formulations,	55
2.4	Summary,	58
<b>3</b>	<b>Regularization Framework</b>	<b>61</b>
3.1	Curse and Complexity of Dimensionality,	62
3.2	Function Approximation and Characterization of Complexity,	66
3.3	Penalization,	70
3.3.1	Parametric Penalties,	72
3.3.2	Nonparametric Penalties,	73
3.4	Model Selection (Complexity Control),	73
3.4.1	Analytical Model Selection Criteria,	75
3.4.2	Model Selection via Resampling,	78
3.4.3	Bias–Variance Tradeoff,	80
3.4.4	Example of Model Selection,	85
3.4.5	Function Approximation versus Predictive Learning,	88
3.5	Summary,	96
<b>4</b>	<b>Statistical Learning Theory</b>	<b>99</b>
4.1	Conditions for Consistency and Convergence of ERM,	101
4.2	Growth Function and VC Dimension,	107
4.2.1	VC Dimension for Classification and Regression Problems,	110
4.2.2	Examples of Calculating VC Dimension,	111
4.3	Bounds on the Generalization,	115
4.3.1	Classification,	116
4.3.2	Regression,	118
4.3.3	Generalization Bounds and Sampling Theorem,	120
4.4	Structural Risk Minimization,	122
4.4.1	Dictionary Representation,	124
4.4.2	Feature Selection,	125
4.4.3	Penalization Formulation,	126
4.4.4	Input Preprocessing,	126
4.4.5	Initial Conditions for Training Algorithm,	127
4.5	Comparisons of Model Selection for Regression,	128
4.5.1	Model Selection for Linear Estimators,	134
4.5.2	Model Selection for $k$ -Nearest-Neighbor Regression,	137
4.5.3	Model Selection for Linear Subset Regression,	140
4.5.4	Discussion,	141
4.6	Measuring the VC Dimension,	143
4.7	VC Dimension, Occam’s Razor, and Popper’s Falsifiability,	146
4.8	Summary and Discussion,	149

<b>5</b>	<b>Nonlinear Optimization Strategies</b>	<b>151</b>
5.1	Stochastic Approximation Methods, 154	
5.1.1	Linear Parameter Estimation, 155	
5.1.2	Backpropagation Training of MLP Networks, 156	
5.2	Iterative Methods, 161	
5.2.1	EM Methods for Density Estimation, 161	
5.2.2	Generalized Inverse Training of MLP Networks, 164	
5.3	Greedy Optimization, 169	
5.3.1	Neural Network Construction Algorithms, 169	
5.3.2	Classification and Regression Trees, 170	
5.4	Feature Selection, Optimization, and Statistical Learning Theory, 173	
5.5	Summary, 175	
<b>6</b>	<b>Methods for Data Reduction and Dimensionality Reduction</b>	<b>177</b>
6.1	Vector Quantization and Clustering, 183	
6.1.1	Optimal Source Coding in Vector Quantization, 184	
6.1.2	Generalized Lloyd Algorithm, 187	
6.1.3	Clustering, 191	
6.1.4	EM Algorithm for VQ and Clustering, 192	
6.1.5	Fuzzy Clustering, 195	
6.2	Dimensionality Reduction: Statistical Methods, 201	
6.2.1	Linear Principal Components, 202	
6.2.2	Principal Curves and Surfaces, 205	
6.2.3	Multidimensional Scaling, 209	
6.3	Dimensionality Reduction: Neural Network Methods, 214	
6.3.1	Discrete Principal Curves and Self-Organizing Map Algorithm, 215	
6.3.2	Statistical Interpretation of the SOM Method, 218	
6.3.3	Flow-Through Version of the SOM and Learning Rate Schedules, 222	
6.3.4	SOM Applications and Modifications, 224	
6.3.5	Self-Supervised MLP, 230	
6.4	Methods for Multivariate Data Analysis, 232	
6.4.1	Factor Analysis, 233	
6.4.2	Independent Component Analysis, 242	
6.5	Summary, 247	
<b>7</b>	<b>Methods for Regression</b>	<b>249</b>
7.1	Taxonomy: Dictionary versus Kernel Representation, 252	
7.2	Linear Estimators, 256	
7.2.1	Estimation of Linear Models and Equivalence of Representations, 258	
7.2.2	Analytic Form of Cross-Validation, 262	



7.2.3	Estimating Complexity of Penalized Linear Models,	263
7.2.4	Nonadaptive Methods,	269
7.3	Adaptive Dictionary Methods,	277
7.3.1	Additive Methods and Projection Pursuit Regression,	279
7.3.2	Multilayer Perceptrons and Backpropagation,	284
7.3.3	Multivariate Adaptive Regression Splines,	293
7.3.4	Orthogonal Basis Functions and Wavelet Signal Denoising,	298
7.4	Adaptive Kernel Methods and Local Risk Minimization,	309
7.4.1	Generalized Memory-Based Learning,	313
7.4.2	Constrained Topological Mapping,	314
7.5	Empirical Studies,	319
7.5.1	Predicting Net Asset Value (NAV) of Mutual Funds,	320
7.5.2	Comparison of Adaptive Methods for Regression,	326
7.6	Combining Predictive Models,	332
7.7	Summary,	337
<b>8</b>	<b>Classification</b>	<b>340</b>
8.1	Statistical Learning Theory Formulation,	343
8.2	Classical Formulation,	348
8.2.1	Statistical Decision Theory,	348
8.2.2	Fisher's Linear Discriminant Analysis,	362
8.3	Methods for Classification,	366
8.3.1	Regression-Based Methods,	368
8.3.2	Tree-Based Methods,	378
8.3.3	Nearest-Neighbor and Prototype Methods,	382
8.3.4	Empirical Comparisons,	385
8.4	Combining Methods and Boosting,	390
8.4.1	Boosting as an Additive Model,	395
8.4.2	Boosting for Regression Problems,	400
8.5	Summary,	401
<b>9</b>	<b>Support Vector Machines</b>	<b>404</b>
9.1	Motivation for Margin-Based Loss,	408
9.2	Margin-Based Loss, Robustness, and Complexity Control,	414
9.3	Optimal Separating Hyperplane,	418
9.4	High-Dimensional Mapping and Inner Product Kernels,	426
9.5	Support Vector Machine for Classification,	430
9.6	Support Vector Implementations,	438
9.7	Support Vector Regression,	439
9.8	SVM Model Selection,	445
9.9	Support Vector Machines and Regularization,	453