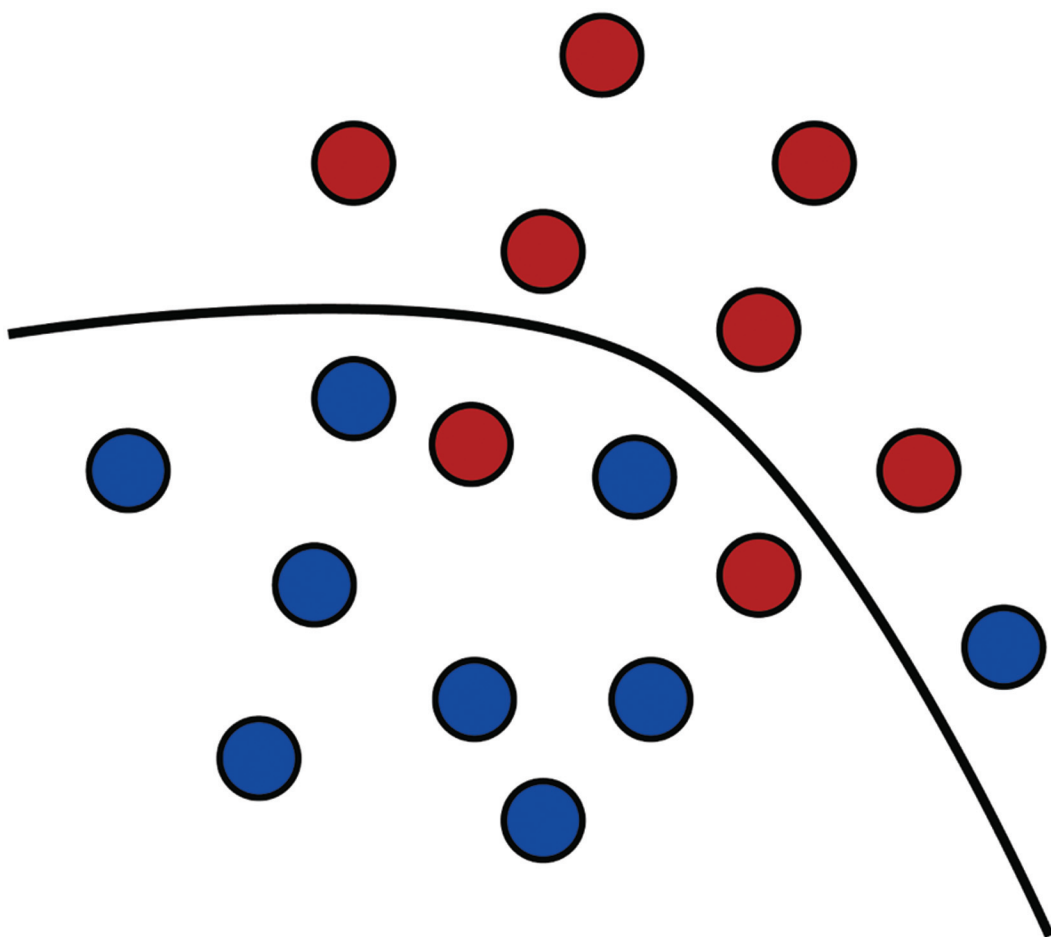


Foundations of Machine Learning



Mehryar Mohri,
Afshin Rostamizadeh,
and Ameet Talwalkar

Foundations of Machine Learning

Adaptive Computation and Machine Learning

Thomas Dietterich, Editor

Christopher Bishop, David Heckerman, Michael Jordan, and Michael Kearns,
Associate Editors

A complete list of books published in The Adaptive Computations and Machine Learning series appears at the back of this book.

Foundations of Machine Learning

Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar

The MIT Press
Cambridge, Massachusetts
London, England

© 2012 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.

MIT Press books may be purchased at special quantity discounts for business or sales promotional use. For information, please email special_sales@mitpress.mit.edu or write to Special Sales Department, The MIT Press, 55 Hayward Street, Cambridge, MA 02142.

This book was set in L^AT_EX by the authors. Printed and bound in the United States of America.

Library of Congress Cataloging-in-Publication Data

Mohri, Mehryar.

Foundations of machine learning / Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar.

p. cm. - (Adaptive computation and machine learning series)

Includes bibliographical references and index.

ISBN 978-0-262-01825-8 (hardcover : alk. paper) 1. Machine learning. 2. Computer algorithms. I. Rostamizadeh, Afshin. II. Talwalkar, Ameet. III. Title.

Q325.5.M64 2012

006.3'1-dc23

2012007249

10 9 8 7 6 5 4 3 2 1

Contents

Preface	xi
1 Introduction	1
1.1 Applications and problems	1
1.2 Definitions and terminology	3
1.3 Cross-validation	5
1.4 Learning scenarios	7
1.5 Outline	8
2 The PAC Learning Framework	11
2.1 The PAC learning model	11
2.2 Guarantees for finite hypothesis sets — consistent case	17
2.3 Guarantees for finite hypothesis sets — inconsistent case	21
2.4 Generalities	24
2.4.1 Deterministic versus stochastic scenarios	24
2.4.2 Bayes error and noise	25
2.4.3 Estimation and approximation errors	26
2.4.4 Model selection	27
2.5 Chapter notes	28
2.6 Exercises	29
3 Rademacher Complexity and VC-Dimension	33
3.1 Rademacher complexity	34
3.2 Growth function	38
3.3 VC-dimension	41
3.4 Lower bounds	48
3.5 Chapter notes	54
3.6 Exercises	55
4 Support Vector Machines	63
4.1 Linear classification	63
4.2 SVMs — separable case	64

4.2.1	Primal optimization problem	64
4.2.2	Support vectors	66
4.2.3	Dual optimization problem	67
4.2.4	Leave-one-out analysis	69
4.3	SVMs — non-separable case	71
4.3.1	Primal optimization problem	72
4.3.2	Support vectors	73
4.3.3	Dual optimization problem	74
4.4	Margin theory	75
4.5	Chapter notes	83
4.6	Exercises	84
5	Kernel Methods	89
5.1	Introduction	89
5.2	Positive definite symmetric kernels	92
5.2.1	Definitions	92
5.2.2	Reproducing kernel Hilbert space	94
5.2.3	Properties	96
5.3	Kernel-based algorithms	100
5.3.1	SVMs with PDS kernels	100
5.3.2	Representer theorem	101
5.3.3	Learning guarantees	102
5.4	Negative definite symmetric kernels	103
5.5	Sequence kernels	106
5.5.1	Weighted transducers	106
5.5.2	Rational kernels	111
5.6	Chapter notes	115
5.7	Exercises	116
6	Boosting	121
6.1	Introduction	121
6.2	AdaBoost	122
6.2.1	Bound on the empirical error	124
6.2.2	Relationship with coordinate descent	126
6.2.3	Relationship with logistic regression	129
6.2.4	Standard use in practice	129
6.3	Theoretical results	130
6.3.1	VC-dimension-based analysis	131
6.3.2	Margin-based analysis	131
6.3.3	Margin maximization	136
6.3.4	Game-theoretic interpretation	137

6.4	Discussion	140
6.5	Chapter notes	141
6.6	Exercises	142
7	On-Line Learning	147
7.1	Introduction	147
7.2	Prediction with expert advice	148
7.2.1	Mistake bounds and Halving algorithm	148
7.2.2	Weighted majority algorithm	150
7.2.3	Randomized weighted majority algorithm	152
7.2.4	Exponential weighted average algorithm	156
7.3	Linear classification	159
7.3.1	Perceptron algorithm	160
7.3.2	Winnnow algorithm	168
7.4	On-line to batch conversion	171
7.5	Game-theoretic connection	174
7.6	Chapter notes	175
7.7	Exercises	176
8	Multi-Class Classification	183
8.1	Multi-class classification problem	183
8.2	Generalization bounds	185
8.3	Uncombined multi-class algorithms	191
8.3.1	Multi-class SVMs	191
8.3.2	Multi-class boosting algorithms	192
8.3.3	Decision trees	194
8.4	Aggregated multi-class algorithms	198
8.4.1	One-versus-all	198
8.4.2	One-versus-one	199
8.4.3	Error-correction codes	201
8.5	Structured prediction algorithms	203
8.6	Chapter notes	206
8.7	Exercises	207
9	Ranking	209
9.1	The problem of ranking	209
9.2	Generalization bound	211
9.3	Ranking with SVMs	213
9.4	RankBoost	214
9.4.1	Bound on the empirical error	216
9.4.2	Relationship with coordinate descent	218

9.4.3	Margin bound for ensemble methods in ranking	220
9.5	Bipartite ranking	221
9.5.1	Boosting in bipartite ranking	222
9.5.2	Area under the ROC curve	224
9.6	Preference-based setting	226
9.6.1	Second-stage ranking problem	227
9.6.2	Deterministic algorithm	229
9.6.3	Randomized algorithm	230
9.6.4	Extension to other loss functions	231
9.7	Discussion	232
9.8	Chapter notes	233
9.9	Exercises	234
10	Regression	237
10.1	The problem of regression	237
10.2	Generalization bounds	238
10.2.1	Finite hypothesis sets	238
10.2.2	Rademacher complexity bounds	239
10.2.3	Pseudo-dimension bounds	241
10.3	Regression algorithms	245
10.3.1	Linear regression	245
10.3.2	Kernel ridge regression	247
10.3.3	Support vector regression	252
10.3.4	Lasso	257
10.3.5	Group norm regression algorithms	260
10.3.6	On-line regression algorithms	261
10.4	Chapter notes	262
10.5	Exercises	263
11	Algorithmic Stability	267
11.1	Definitions	267
11.2	Stability-based generalization guarantee	268
11.3	Stability of kernel-based regularization algorithms	270
11.3.1	Application to regression algorithms: SVR and KRR	274
11.3.2	Application to classification algorithms: SVMs	276
11.3.3	Discussion	276
11.4	Chapter notes	277
11.5	Exercises	277
12	Dimensionality Reduction	281
12.1	Principal Component Analysis	282

12.2	Kernel Principal Component Analysis (KPCA)	283
12.3	KPCA and manifold learning	285
12.3.1	Isomap	285
12.3.2	Laplacian eigenmaps	286
12.3.3	Locally linear embedding (LLE)	287
12.4	Johnson-Lindenstrauss lemma	288
12.5	Chapter notes	290
12.6	Exercises	290
13	Learning Automata and Languages	293
13.1	Introduction	293
13.2	Finite automata	294
13.3	Efficient exact learning	295
13.3.1	Passive learning	296
13.3.2	Learning with queries	297
13.3.3	Learning automata with queries	298
13.4	Identification in the limit	303
13.4.1	Learning reversible automata	304
13.5	Chapter notes	309
13.6	Exercises	310
14	Reinforcement Learning	313
14.1	Learning scenario	313
14.2	Markov decision process model	314
14.3	Policy	315
14.3.1	Definition	315
14.3.2	Policy value	316
14.3.3	Policy evaluation	316
14.3.4	Optimal policy	318
14.4	Planning algorithms	319
14.4.1	Value iteration	319
14.4.2	Policy iteration	322
14.4.3	Linear programming	324
14.5	Learning algorithms	325
14.5.1	Stochastic approximation	326
14.5.2	TD(0) algorithm	330
14.5.3	Q-learning algorithm	331
14.5.4	SARSA	334
14.5.5	TD(λ) algorithm	335
14.5.6	Large state space	336
14.6	Chapter notes	337

Conclusion	339
A Linear Algebra Review	341
A.1 Vectors and norms	341
A.1.1 Norms	341
A.1.2 Dual norms	342
A.2 Matrices	344
A.2.1 Matrix norms	344
A.2.2 Singular value decomposition	345
A.2.3 Symmetric positive semidefinite (SPSD) matrices	346
B Convex Optimization	349
B.1 Differentiation and unconstrained optimization	349
B.2 Convexity	350
B.3 Constrained optimization	353
B.4 Chapter notes	357
C Probability Review	359
C.1 Probability	359
C.2 Random variables	359
C.3 Conditional probability and independence	361
C.4 Expectation, Markov's inequality, and moment-generating function	363
C.5 Variance and Chebyshev's inequality	365
D Concentration inequalities	369
D.1 Hoeffding's inequality	369
D.2 McDiarmid's inequality	371
D.3 Other inequalities	373
D.3.1 Binomial distribution: Slud's inequality	374
D.3.2 Normal distribution: tail bound	374
D.3.3 Khintchine-Kahane inequality	374
D.4 Chapter notes	376
D.5 Exercises	377
E Notation	379
References	381
Index	397

Preface

This book is a general introduction to machine learning that can serve as a textbook for students and researchers in the field. It covers fundamental modern topics in machine learning while providing the theoretical basis and conceptual tools needed for the discussion and justification of algorithms. It also describes several key aspects of the application of these algorithms.

We have aimed to present the most novel theoretical tools and concepts while giving concise proofs, even for relatively advanced results. In general, whenever possible, we have chosen to favor succinctness. Nevertheless, we discuss some crucial complex topics arising in machine learning and highlight several open research questions. Certain topics often merged with others or treated with insufficient attention are discussed separately here and with more emphasis: for example, a different chapter is reserved for multi-class classification, ranking, and regression.

Although we cover a very wide variety of important topics in machine learning, we have chosen to omit a few important ones, including graphical models and neural networks, both for the sake of brevity and because of the current lack of solid theoretical guarantees for some methods.

The book is intended for students and researchers in machine learning, statistics and other related areas. It can be used as a textbook for both graduate and advanced undergraduate classes in machine learning or as a reference text for a research seminar. The first three chapters of the book lay the theoretical foundation for the subsequent material. Other chapters are mostly self-contained, with the exception of chapter 5 which introduces some concepts that are extensively used in later ones. Each chapter concludes with a series of exercises, with full solutions presented separately.

The reader is assumed to be familiar with basic concepts in linear algebra, probability, and analysis of algorithms. However, to further help him, we present in the appendix a concise linear algebra and a probability review, and a short introduction to convex optimization. We have also collected in the appendix a number of useful tools for concentration bounds used in this book.

To our knowledge, there is no single textbook covering all of the material presented here. The need for a unified presentation has been pointed out to us

every year by our machine learning students. There are several good books for various specialized areas, but these books do not include a discussion of other fundamental topics in a general manner. For example, books about kernel methods do not include a discussion of other fundamental topics such as boosting, ranking, reinforcement learning, learning automata or online learning. There also exist more general machine learning books, but the theoretical foundation of our book and our emphasis on proofs make our presentation quite distinct.

Most of the material presented here takes its origins in a machine learning graduate course (*Foundations of Machine Learning*) taught by the first author at the Courant Institute of Mathematical Sciences in New York University over the last seven years. This book has considerably benefited from the comments and suggestions from students in these classes, along with those of many friends, colleagues and researchers to whom we are deeply indebted.

We are particularly grateful to Corinna Cortes and Yishay Mansour who have both made a number of key suggestions for the design and organization of the material presented with detailed comments that we have fully taken into account and that have greatly improved the presentation. We are also grateful to Yishay Mansour for using a preliminary version of the book for teaching and for reporting his feedback to us.

We also thank for discussions, suggested improvement, and contributions of many kinds the following colleagues and friends from academic and corporate research laboratories: Cyril Allauzen, Stephen Boyd, Spencer Greenberg, Lisa Hellerstein, Sanjiv Kumar, Ryan McDonald, Andres Muñoz Medina, Tyler Neylon, Peter Norvig, Fernando Pereira, Maria Pershina, Ashish Rastogi, Michael Riley, Umar Syed, Csaba Szepesvári, Eugene Weinstein, and Jason Weston.

Finally, we thank the MIT Press publication team for their help and support in the development of this text.

1 Introduction

Machine learning can be broadly defined as computational methods using experience to improve performance or to make accurate predictions. Here, *experience* refers to the past information available to the learner, which typically takes the form of electronic data collected and made available for analysis. This data could be in the form of digitized human-labeled training sets, or other types of information obtained via interaction with the environment. In all cases, its quality and size are crucial to the success of the predictions made by the learner.

Machine learning consists of designing efficient and accurate prediction *algorithms*. As in other areas of computer science, some critical measures of the quality of these algorithms are their time and space complexity. But, in machine learning, we will need additionally a notion of *sample complexity* to evaluate the sample size required for the algorithm to learn a family of concepts. More generally, theoretical learning guarantees for an algorithm depend on the complexity of the concept classes considered and the size of the training sample.

Since the success of a learning algorithm depends on the data used, machine learning is inherently related to data analysis and statistics. More generally, learning techniques are data-driven methods combining fundamental concepts in computer science with ideas from statistics, probability and optimization.

1.1 Applications and problems

Learning algorithms have been successfully deployed in a variety of applications, including

- Text or document classification, e.g., spam detection;
- Natural language processing, e.g., morphological analysis, part-of-speech tagging, statistical parsing, named-entity recognition;
- Speech recognition, speech synthesis, speaker verification;
- Optical character recognition (OCR);
- Computational biology applications, e.g., protein function or structured predic-

tion;

- Computer vision tasks, e.g., image recognition, face detection;
- Fraud detection (credit card, telephone) and network intrusion;
- Games, e.g., chess, backgammon;
- Unassisted vehicle control (robots, navigation);
- Medical diagnosis;
- Recommendation systems, search engines, information extraction systems.

This list is by no means comprehensive, and learning algorithms are applied to new applications every day. Moreover, such applications correspond to a wide variety of learning problems. Some major classes of learning problems are:

- *Classification*: Assign a category to each item. For example, document classification may assign items with categories such as *politics*, *business*, *sports*, or *weather* while image classification may assign items with categories such as *landscape*, *portrait*, or *animal*. The number of categories in such tasks is often relatively small, but can be large in some difficult tasks and even unbounded as in OCR, text classification, or speech recognition.
- *Regression*: Predict a real value for each item. Examples of regression include prediction of stock values or variations of economic variables. In this problem, the penalty for an incorrect prediction depends on the magnitude of the difference between the true and predicted values, in contrast with the classification problem, where there is typically no notion of closeness between various categories.
- *Ranking*: Order items according to some criterion. Web search, e.g., returning web pages relevant to a search query, is the canonical ranking example. Many other similar ranking problems arise in the context of the design of information extraction or natural language processing systems.
- *Clustering*: Partition items into homogeneous regions. Clustering is often performed to analyze very large data sets. For example, in the context of social network analysis, clustering algorithms attempt to identify “communities” within large groups of people.
- *Dimensionality reduction* or *manifold learning*: Transform an initial representation of items into a lower-dimensional representation of these items while preserving some properties of the initial representation. A common example involves preprocessing digital images in computer vision tasks.

The main practical objectives of machine learning consist of generating accurate predictions for unseen items and of designing efficient and robust algorithms to produce these predictions, even for large-scale problems. To do so, a number of algorithmic and theoretical questions arise. Some fundamental questions include:

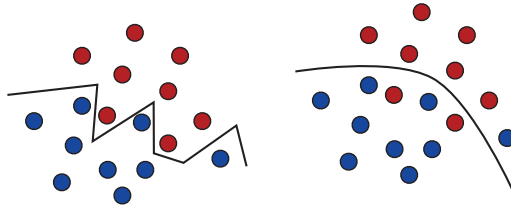


Figure 1.1 The zig-zag line on the left panel is consistent over the blue and red training sample, but it is a complex separation surface that is not likely to generalize well to unseen data. In contrast, the decision surface on the right panel is simpler and might generalize better in spite of its misclassification of a few points of the training sample.

Which concept families can actually be learned, and under what conditions? How well can these concepts be learned computationally?

1.2 Definitions and terminology

We will use the canonical problem of spam detection as a running example to illustrate some basic definitions and to describe the use and evaluation of machine learning algorithms in practice. Spam detection is the problem of learning to automatically classify email messages as either SPAM or non-SPAM.

- *Examples*: Items or instances of data used for learning or evaluation. In our spam problem, these examples correspond to the collection of email messages we will use for learning and testing.
- *Features*: The set of attributes, often represented as a vector, associated to an example. In the case of email messages, some relevant features may include the length of the message, the name of the sender, various characteristics of the header, the presence of certain keywords in the body of the message, and so on.
- *Labels*: Values or categories assigned to examples. In classification problems, examples are assigned specific categories, for instance, the SPAM and non-SPAM categories in our binary classification problem. In regression, items are assigned real-valued labels.
- *Training sample*: Examples used to train a learning algorithm. In our spam problem, the training sample consists of a set of email examples along with their associated labels. The training sample varies for different learning scenarios, as described in section 1.4.
- *Validation sample*: Examples used to tune the parameters of a learning algorithm

when working with labeled data. Learning algorithms typically have one or more free parameters, and the validation sample is used to select appropriate values for these model parameters.

- *Test sample*: Examples used to evaluate the performance of a learning algorithm. The test sample is separate from the training and validation data and is not made available in the learning stage. In the spam problem, the test sample consists of a collection of email examples for which the learning algorithm must predict labels based on features. These predictions are then compared with the labels of the test sample to measure the performance of the algorithm.
- *Loss function*: A function that measures the difference, or loss, between a predicted label and a true label. Denoting the set of all labels as \mathcal{Y} and the set of possible predictions as \mathcal{Y}' , a loss function L is a mapping $L: \mathcal{Y} \times \mathcal{Y}' \rightarrow \mathbb{R}_+$. In most cases, $\mathcal{Y}' = \mathcal{Y}$ and the loss function is bounded, but these conditions do not always hold. Common examples of loss functions include the zero-one (or misclassification) loss defined over $\{-1, +1\} \times \{-1, +1\}$ by $L(y, y') = 1_{y' \neq y}$ and the squared loss defined over $I \times I$ by $L(y, y') = (y' - y)^2$, where $I \subseteq \mathbb{R}$ is typically a bounded interval.
- *Hypothesis set*: A set of functions mapping features (feature vectors) to the set of labels \mathcal{Y} . In our example, these may be a set of functions mapping email features to $\mathcal{Y} = \{\text{SPAM}, \text{non-SPAM}\}$. More generally, hypotheses may be functions mapping features to a different set \mathcal{Y}' . They could be linear functions mapping email feature vectors to real numbers interpreted as *scores* ($\mathcal{Y}' = \mathbb{R}$), with higher score values more indicative of SPAM than lower ones.

We now define the learning stages of our spam problem. We start with a given collection of labeled examples. We first randomly partition the data into a training sample, a validation sample, and a test sample. The size of each of these samples depends on a number of different considerations. For example, the amount of data reserved for validation depends on the number of free parameters of the algorithm. Also, when the labeled sample is relatively small, the amount of training data is often chosen to be larger than that of test data since the learning performance directly depends on the training sample.

Next, we associate relevant features to the examples. This is a critical step in the design of machine learning solutions. Useful features can effectively guide the learning algorithm, while poor or uninformative ones can be misleading. Although it is critical, to a large extent, the choice of the features is left to the user. This choice reflects the user's *prior knowledge* about the learning task which in practice can have a dramatic effect on the performance results.

Now, we use the features selected to train our learning algorithm by fixing different values of its free parameters. For each value of these parameters, the algorithm

selects a different hypothesis out of the hypothesis set. We choose among them the hypothesis resulting in the best performance on the validation sample. Finally, using that hypothesis, we predict the labels of the examples in the test sample. The performance of the algorithm is evaluated by using the loss function associated to the task, e.g., the zero-one loss in our spam detection task, to compare the predicted and true labels.

Thus, the performance of an algorithm is of course evaluated based on its test error and not its error on the training sample. A learning algorithm may be *consistent*, that is it may commit no error on the examples of the training data, and yet have a poor performance on the test data. This occurs for consistent learners defined by very complex decision surfaces, as illustrated in figure 1.1, which tend to memorize a relatively small training sample instead of seeking to generalize well. This highlights the key distinction between memorization and generalization, which is the fundamental property sought for an accurate learning algorithm. Theoretical guarantees for consistent learners will be discussed with great detail in chapter 2.

1.3 Cross-validation

In practice, the amount of labeled data available is often too small to set aside a validation sample since that would leave an insufficient amount of training data. Instead, a widely adopted method known as *n-fold cross-validation* is used to exploit the labeled data both for *model selection* (selection of the free parameters of the algorithm) and for training.

Let θ denote the vector of free parameters of the algorithm. For a fixed value of θ , the method consists of first randomly partitioning a given sample S of m labeled examples into n subsamples, or folds. The i th fold is thus a labeled sample $((x_{i1}, y_{i1}), \dots, (x_{im_i}, y_{im_i}))$ of size m_i . Then, for any $i \in [1, n]$, the learning algorithm is trained on all but the i th fold to generate a hypothesis h_i , and the performance of h_i is tested on the i th fold, as illustrated in figure 1.2a. The parameter value θ is evaluated based on the average error of the hypotheses h_i , which is called the *cross-validation error*. This quantity is denoted by $\widehat{R}_{CV}(\theta)$ and defined by

$$\widehat{R}_{CV}(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{m_i} \sum_{j=1}^{m_i} L(h_i(x_{ij}), y_{ij})}_{\text{error of } h_i \text{ on the } i\text{th fold}} .$$

The folds are generally chosen to have equal size, that is $m_i = m/n$ for all $i \in [1, n]$. How should n be chosen? The appropriate choice is subject to a trade-off and the topic of much learning theory research that we cannot address in this introductory

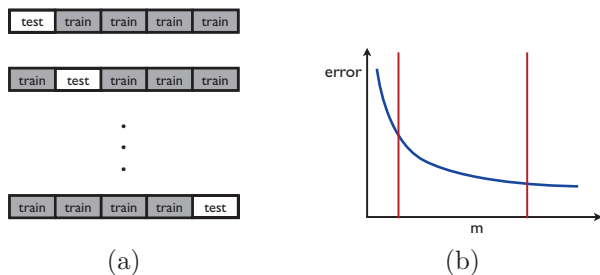


Figure 1.2 n -fold cross validation. (a) Illustration of the partitioning of the training data into 5 folds. (b) Typical plot of a classifier’s prediction error as a function of the size of the training sample: the error decreases as a function of the number of training points.

chapter. For a large n , each training sample used in n -fold cross-validation has size $m - m/n = m(1 - 1/n)$ (illustrated by the right vertical red line in figure 1.2b), which is close to m , the size of the full sample, but the training samples are quite similar. Thus, the method tends to have a small bias but a large variance. In contrast, smaller values of n lead to more diverse training samples but their size (shown by the left vertical red line in figure 1.2b) is significantly less than m , thus the method tends to have a smaller variance but a larger bias.

In machine learning applications, n is typically chosen to be 5 or 10. n -fold cross validation is used as follows in model selection. The full labeled data is first split into a training and a test sample. The training sample of size m is then used to compute the n -fold cross-validation error $\hat{R}_{CV}(\theta)$ for a small number of possible values of θ . θ is next set to the value θ_0 for which $\hat{R}_{CV}(\theta)$ is smallest and the algorithm is trained with the parameter setting θ_0 over the full training sample of size m . Its performance is evaluated on the test sample as already described in the previous section.

The special case of n -fold cross validation where $n = m$ is called *leave-one-out cross-validation*, since at each iteration exactly one instance is left out of the training sample. As shown in chapter 4, the average leave-one-out error is an approximately unbiased estimate of the average error of an algorithm and can be used to derive simple guarantees for some algorithms. In general, the leave-one-out error is very costly to compute, since it requires training n times on samples of size $m - 1$, but for some algorithms it admits a very efficient computation (see exercise 10.9).

In addition to model selection, n -fold cross validation is also commonly used for performance evaluation. In that case, for a fixed parameter setting θ , the full labeled sample is divided into n random folds with no distinction between training and test samples. The performance reported is the n -fold cross-validation on the full sample as well as the standard deviation of the errors measured on each fold.

1.4 Learning scenarios

We next briefly describe common machine learning scenarios. These scenarios differ in the types of training data available to the learner, the order and method by which training data is received and the test data used to evaluate the learning algorithm.

- *Supervised learning*: The learner receives a set of labeled examples as training data and makes predictions for all unseen points. This is the most common scenario associated with classification, regression, and ranking problems. The spam detection problem discussed in the previous section is an instance of supervised learning.
- *Unsupervised learning*: The learner exclusively receives unlabeled training data, and makes predictions for all unseen points. Since in general no labeled example is available in that setting, it can be difficult to quantitatively evaluate the performance of a learner. Clustering and dimensionality reduction are example of unsupervised learning problems.
- *Semi-supervised learning*: The learner receives a training sample consisting of both labeled and unlabeled data, and makes predictions for all unseen points. Semi-supervised learning is common in settings where unlabeled data is easily accessible but labels are expensive to obtain. Various types of problems arising in applications, including classification, regression, or ranking tasks, can be framed as instances of semi-supervised learning. The hope is that the distribution of unlabeled data accessible to the learner can help him achieve a better performance than in the supervised setting. The analysis of the conditions under which this can indeed be realized is the topic of much modern theoretical and applied machine learning research.
- *Transductive inference*: As in the semi-supervised scenario, the learner receives a labeled training sample along with a set of unlabeled test points. However, the objective of transductive inference is to predict labels only for these particular test points. Transductive inference appears to be an easier task and matches the scenario encountered in a variety of modern applications. However, as in the semi-supervised setting, the assumptions under which a better performance can be achieved in this setting are research questions that have not been fully resolved.
- *On-line learning*: In contrast with the previous scenarios, the online scenario involves multiple rounds and training and testing phases are intermixed. At each round, the learner receives an unlabeled training point, makes a prediction, receives the true label, and incurs a loss. The objective in the on-line setting is to minimize the cumulative loss over all rounds. Unlike the previous settings just discussed, no distributional assumption is made in on-line learning. In fact, instances and their labels may be chosen adversarially within this scenario.

- *Reinforcement learning*: The training and testing phases are also intermixed in reinforcement learning. To collect information, the learner actively interacts with the environment and in some cases affects the environment, and receives an immediate reward for each action. The object of the learner is to maximize his reward over a course of actions and iterations with the environment. However, no long-term reward feedback is provided by the environment, and the learner is faced with the *exploration versus exploitation* dilemma, since he must choose between exploring unknown actions to gain more information versus exploiting the information already collected.
- *Active learning*: The learner adaptively or interactively collects training examples, typically by querying an oracle to request labels for new points. The goal in active learning is to achieve a performance comparable to the standard supervised learning scenario, but with fewer labeled examples. Active learning is often used in applications where labels are expensive to obtain, for example computational biology applications.

In practice, many other intermediate and somewhat more complex learning scenarios may be encountered.

1.5 Outline

This book presents several fundamental and mathematically well-studied algorithms. It discusses in depth their theoretical foundations as well as their practical applications. The topics covered include:

- Probably approximately correct (PAC) learning framework; learning guarantees for finite hypothesis sets;
- Learning guarantees for infinite hypothesis sets, Rademacher complexity, VC-dimension;
- Support vector machines (SVMs), margin theory;
- Kernel methods, positive definite symmetric kernels, representer theorem, rational kernels;
- Boosting, analysis of empirical error, generalization error, margin bounds;
- Online learning, mistake bounds, the weighted majority algorithm, the exponential weighted average algorithm, the Perceptron and Winnow algorithms;
- Multi-class classification, multi-class SVMs, multi-class boosting, one-versus-all, one-versus-one, error-correction methods;
- Ranking, ranking with SVMs, RankBoost, bipartite ranking, preference-based

ranking;

- Regression, linear regression, kernel ridge regression, support vector regression, Lasso;
- Stability-based analysis, applications to classification and regression;
- Dimensionality reduction, principal component analysis (PCA), kernel PCA, Johnson-Lindenstrauss lemma;
- Learning automata and languages;
- Reinforcement learning, Markov decision processes, planning and learning problems.

The analyses in this book are self-contained, with relevant mathematical concepts related to linear algebra, convex optimization, probability and statistics included in the appendix.