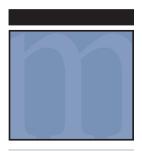
Machine Learning in Python®



Machine Learning in Python®

Essential Techniques for Predictive Analysis

Michael Bowles

WILEY

Machine Learning in Python®: Essential Techniques for Predictive Analysis

Published by John Wiley & Sons, Inc. 10475 Crosspoint Boulevard Indianapolis, IN 46256 www.wiley.com

Copyright © 2015 by John Wiley & Sons, Inc., Indianapolis, Indiana Published simultaneously in Canada

ISBN: 978-1-118-96174-2 ISBN: 978-1-118-96176-6 (ebk) ISBN: 978-1-118-96175-9 (ebk)

Manufactured in the United States of America

10987654321

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 Sections 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, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. 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/permissions.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Web site is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or website may provide or recommendations it may make. Further, readers should be aware that Internet websites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services 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 publishes in a variety of print and electronic formats and by print-on-demand. Some material included with standard print versions of this book may not be included in e-books or in print-on-demand. If this book refers to media such as a CD or DVD that is not included in the version you purchased, you may download this material at http://booksupport.wiley.com. For more information about Wiley products, visit www.wiley.com.

Library of Congress Control Number: 2015930541

Trademarks: Wiley and the Wiley logo are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. Python is a registered trademark of Python Software Foundation. All other trademarks are the property of their respective owners. John Wiley & Sons, Inc. is not associated with any product or vendor mentioned in this book.

To my children, Scott, Seth, and Cayley. Their blossoming lives and selves bring me more joy than anything else in this world.

To my close friends David and Ron for their selfless generosity and steadfast friendship.

To my friends and colleagues at Hacker Dojo in Mountain View, California, for their technical challenges and repartee.

To my climbing partners. One of them, Katherine, says climbing partners make the best friends because "they see you paralyzed with fear, offer encouragement to overcome it, and celebrate when you do."

About the Author

Dr. Michael Bowles (Mike) holds Bachelor's and Master's degrees in Mechanical Engineering, an Sc.D. in Instrumentation, and an MBA. He has worked in academia, technology, and business. Mike currently works with startup companies where machine learning is integral to success. He serves variously as part of the management team, a consultant, or advisor. He also teaches machine learning courses at Hacker Dojo, a co-working space and startup incubator in Mountain View, California.

Mike was born in Oklahoma and earned his Bachelor's and Master's degrees there. Then after a stint in Southeast Asia, Mike went to Cambridge for his Sc.D. and then held the C. Stark Draper Chair at MIT after graduation. Mike left Boston to work on communications satellites at Hughes Aircraft company in Southern California, and then after completing an MBA at UCLA moved to the San Francisco Bay Area to take roles as founder and CEO of two successful venture-backed startups.

Mike remains actively involved in technical and startup-related work. Recent projects include the use of machine learning in automated trading, predicting biological outcomes on the basis of genetic information, natural language processing for website optimization, predicting patient outcomes from demographic and lab data, and due diligence work on companies in the machine learning and big data arenas. Mike can be reached through www.mbowles.com.

About the Technical Editor

Daniel Posner holds Bachelor's and Master's degrees in Economics and is completing a Ph.D. in Biostatistics at Boston University. He has provided statistical consultation for pharmaceutical and biotech firms as well as for researchers at the Palo Alto VA hospital.

Daniel has collaborated with the author extensively on topics covered in this book. In the past, they have written grant proposals to develop web-scale gradient boosting algorithms. Most recently, they worked together on a consulting contract involving random forests and spline basis expansions to identify key variables in drug trial outcomes and to sharpen predictions in order to reduce the required trial populations.

Credits

Executive Editor Robert Elliott

Project Editor Jennifer Lynn

Technical Editor Daniel Posner

Production Editor

Dassi Zeidel

Copy Editor Keith Cline

Manager of Content Development

& Assembly

Mary Beth Wakefield

Marketing DirectorDavid Mayhew

Marketing Manager

Carrie Sherrill

Professional Technology &

Strategy DirectorBarry Pruett

Business Manager

Amy Knies

Associate Publisher

Jim Minatel

Project Coordinator, Cover

Brent Savage

Proofreader

Word One New York

Indexer

Johnna VanHoose Dinse

Cover Designer

Wiley

Acknowledgments

I'd like to acknowledge the splendid support that people at Wiley have offered during the course of writing this book. It began with Robert Elliot, the acquisitions editor, who first contacted me about writing a book; he was very easy to work with. It continued with Jennifer Lynn, who has done the editing on the book. She's been very responsive to questions and very patiently kept me on schedule during the writing. I thank you both.

I also want to acknowledge the enormous comfort that comes from having such a sharp, thorough statistician and programmer as Daniel Posner doing the technical editing on the book. Thank you for that and thanks also for the fun and interesting discussions on machine learning, statistics, and algorithms. I don't know anyone else who'll get as deep as fast.

Contents at a Glance

Introduction	1	xxiii
Chapter 1	The Two Essential Algorithms for Making Predictions	1
Chapter 2	Understand the Problem by Understanding the Data	23
Chapter 3	Predictive Model Building: Balancing Performance, Complexity, and Big Data	75
Chapter 4	Penalized Linear Regression	121
Chapter 5	Building Predictive Models Using Penalized Linear Methods	165
Chapter 6	Ensemble Methods	211
Chapter 7	Building Ensemble Models with Python	255
Index		319

Contents

Introduction	n	xxiii
Chapter 1	The Two Essential Algorithms for Making Predictions	1
	Why Are These Two Algorithms So Useful?	2
	What Are Penalized Regression Methods?	7
	What Are Ensemble Methods?	9
	How to Decide Which Algorithm to Use	11
	The Process Steps for Building a Predictive Model	13
	Framing a Machine Learning Problem	15
	Feature Extraction and Feature Engineering	17
	Determining Performance of a Trained Model	18
	Chapter Contents and Dependencies	18
	Summary	20
Chapter 2	Understand the Problem by Understanding the Data	23
	The Anatomy of a New Problem	24
	Different Types of Attributes and Labels	
	Drive Modeling Choices	26
	Things to Notice about Your New Data Set	27
	Classification Problems: Detecting Unexploded	
	Mines Using Sonar	28
	Physical Characteristics of the Rocks Versus Mines Data Set	29
	Statistical Summaries of the Rocks versus Mines Data Set	32
	Visualization of Outliers Using Quantile-Quantile Plot	35
	Statistical Characterization of Categorical Attributes	37
	How to Use Python Pandas to Summarize the	
	Rocks Versus Mines Data Set	37

Visualizing with Parallel Coordinates Plots Visualizing Interrelationships between Attributes and Labels Visualizing Interrelationships between Attributes and Labels Visualizing Attribute and Label Correlations Using a Heat Map Summarizing the Process for Understanding Rocks versus Mines Data Set Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem For How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations Gor the Abalone Problem For Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galculate How Your Wine Tastes Galculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems		visualizing Properties of the Rocks versus Mines Data Set	40
Visualizing Attribute and Label Correlations Using a Heat Map Summarizing the Process for Understanding Rocks versus Mines Data Set Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem 60 Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary 73 Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Model and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful		Visualizing with Parallel Coordinates Plots	40
Using a Heat Map Summarizing the Process for Understanding Rocks versus Mines Data Set Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations of the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression to Control Overfitting Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful		Visualizing Interrelationships between Attributes and Labels	42
Summarizing the Process for Understanding Rocks versus Mines Data Set Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Methods Are So Useful		Visualizing Attribute and Label Correlations	
versus Mines Data Set Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of De		Using a Heat Map	49
Real-Valued Predictions with Factor Variables: How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Ochosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful		Summarizing the Process for Understanding Rocks	
How Old Is Your Abalone? Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galulate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful		versus Mines Data Set	50
Parallel Coordinates for Regression Problems—Visualize Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galulate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary 73 Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Pro		Real-Valued Predictions with Factor Variables:	
Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galculate How Your Wine Tastes Galculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Deployed Models Ochieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Methods Are So Useful		How Old Is Your Abalone?	50
Variable Relationships for Abalone Problem How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galculate How Your Wine Tastes Galculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Deployed Models Ochieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Methods Are So Useful		Parallel Coordinates for Regression Problems—Visualize	
How to Use Correlation Heat Map for Regression—Visualize Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Ochieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary Chapter 4 Penalized Linear Regression Methods Are So Useful		S S S S S S S S S S S S S S S S S S S	56
Pair-Wise Correlations for the Abalone Problem Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? 68 Summary 73 Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Methods Are So Useful			
Real-Valued Predictions Using Real-Valued Attributes: Calculate How Your Wine Tastes Galulate How Your Wine Tastes Galulate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? 68 Summary 73 Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data 75 The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data 79 Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Ochoosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful			60
Calculate How Your Wine Tastes Multiclass Classification Problem: What Type of Glass Is That? Summary Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Coefficients—Ridge Regression Methods Are So Useful Chapter 4 Penalized Linear Regression Methods Are So Useful			
Multiclass Classification Problem: What Type of Glass Is That? 68 Summary 73 Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data 75 The Basic Problem: Understanding Function Approximation 76 Working with Training Data 76 Assessing Performance of Predictive Models 78 Factors Driving Algorithm Choices and 79 Contrast Between a Simple Problem and a Complex Problem 80 Contrast Between a Simple Model and a Complex Model 82 Factors Driving Predictive Algorithm Performance 86 Choosing an Algorithm: Linear or Nonlinear? 87 Measuring the Performance of Predictive Models 88 Performance Measures for Different Types of Problems 88 Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, 80 Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting 103 Evaluating and Understanding Your Predictive Model 108 Control Overfitting by Penalizing Regression 100 Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122		~	62
Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data 75 The Basic Problem: Understanding Function Approximation 76 Working with Training Data 76 Assessing Performance of Predictive Models 78 Factors Driving Algorithm Choices and 79 Contrast Between a Simple Problem and a Complex Problem 80 Contrast Between a Simple Model and a Complex Model 82 Factors Driving Predictive Algorithm Performance 86 Choosing an Algorithm: Linear or Nonlinear? 87 Measuring the Performance of Predictive Models 88 Performance Measures for Different Types of Problems 88 Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting 103 Evaluating and Understanding Your Predictive Model 108 Control Overfitting by Penalizing Regression 100 Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122			
Chapter 3 Predictive Model Building: Balancing Performance, Complexity, and Big Data 75 The Basic Problem: Understanding Function Approximation 76 Working with Training Data 76 Assessing Performance of Predictive Models 78 Factors Driving Algorithm Choices and 79 Contrast Between a Simple Problem and a Complex Problem 80 Contrast Between a Simple Model and a Complex Model 82 Factors Driving Predictive Algorithm Performance 86 Choosing an Algorithm: Linear or Nonlinear? 87 Measuring the Performance of Predictive Models 88 Performance Measures for Different Types of Problems 88 Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting 103 Evaluating and Understanding Your Predictive Model 108 Control Overfitting by Penalizing Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122			
The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Pachieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 112 Chapter 4 Penalized Linear Regression Methods Are So Useful		Summing	, 0
The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Pachieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 112 Chapter 4 Penalized Linear Regression Methods Are So Useful	Chapter 3	Predictive Model Building: Balancing Performance,	
The Basic Problem: Understanding Function Approximation Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Hoasures for Different Types of Problems Simulating Performance of Deployed Models Ochoosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful			75
Working with Training Data Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful			76
Assessing Performance of Predictive Models Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Factors Driving Predictive Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Performance F			76
Factors Driving Algorithm Choices and Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Factors Driving Predictive Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Factors Driving Predictive Model Factors Driving Predictive Models Factors Driving Predictive Model Factors Driving Predictive Problem Problem Problem Problem Problem Problem Probl			78
Performance—Complexity and Data Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Exactors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Ochoosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary 110 Chapter 4 Penalized Linear Regression Methods Are So Useful 121			
Contrast Between a Simple Problem and a Complex Problem Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Factors Driving Predictive Algorithm Performance Reasuring an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Simulating Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary 110 Chapter 4 Penalized Linear Regression Methods Are So Useful			79
Contrast Between a Simple Model and a Complex Model Factors Driving Predictive Algorithm Performance Reference Set Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Simulating Performance of Deployed Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Simulating Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary 110 Chapter 4 Penalized Linear Regression Methods Are So Useful 121			80
Factors Driving Predictive Algorithm Performance Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 121			82
Choosing an Algorithm: Linear or Nonlinear? Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models Simulating Performance of Deployed Models Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Coefficients—Ridge Regression 110 Summary 121 Chapter 4 Penalized Linear Regression Methods Are So Useful			86
Measuring the Performance of Predictive Models Performance Measures for Different Types of Problems Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122			87
Performance Measures for Different Types of Problems Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122			88
Simulating Performance of Deployed Models 99 Achieving Harmony Between Model and Data 101 Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size 102 Using Forward Stepwise Regression to Control Overfitting 103 Evaluating and Understanding Your Predictive Model 108 Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122		~	88
Achieving Harmony Between Model and Data Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 121			99
Choosing a Model to Balance Problem Complexity, Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 121			101
Model Complexity, and Data Set Size Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary 110 Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful		·	
Using Forward Stepwise Regression to Control Overfitting Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary 110 Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful 121			102
Evaluating and Understanding Your Predictive Model Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful 108 110 110 1110 1121			103
Control Overfitting by Penalizing Regression Coefficients—Ridge Regression Summary Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful 122			108
Coefficients—Ridge Regression 110 Summary 119 Chapter 4 Penalized Linear Regression Methods Are So Useful 122			
Chapter 4 Penalized Linear Regression Why Penalized Linear Regression Methods Are So Useful 121 122			110
Why Penalized Linear Regression Methods Are So Useful 122		0 0	119
Why Penalized Linear Regression Methods Are So Useful 122			
Why Penalized Linear Regression Methods Are So Useful 122	Chapter 4	Penalized Linear Regression	121
•		<u> </u>	
Variable Importance Information 122		·	
1		Extremely Fast Evaluation When Deployed	123

	Reliable Performance	123
	Sparse Solutions	123
	Problem May Require Linear Model	124
	When to Use Ensemble Methods	124
	Penalized Linear Regression: Regulating Linear	
	Regression for Optimum Performance	124
	Training Linear Models: Minimizing Errors and More	126
	Adding a Coefficient Penalty to the OLS Formulation	127
	Other Useful Coefficient Penalties—Manhattan and	
	ElasticNet	128
	Why Lasso Penalty Leads to Sparse Coefficient Vectors	129
	ElasticNet Penalty Includes Both Lasso and Ridge	131
	Solving the Penalized Linear Regression Problem	132
	Understanding Least Angle Regression and Its Relationship	
	to Forward Stepwise Regression	132
	How LARS Generates Hundreds of Models of Varying	
	Complexity	136
	Choosing the Best Model from The Hundreds	
	LARS Generates	139
	Using Glmnet: Very Fast and Very General	144
	Comparison of the Mechanics of Glmnet and	
	LARS Algorithms	145
	Initializing and Iterating the Glmnet Algorithm	146
	Extensions to Linear Regression with Numeric Input	151
	Solving Classification Problems with Penalized Regression	151
	Working with Classification Problems Having More Than	
	Two Outcomes	155
	Understanding Basis Expansion: Using Linear Methods on	
	Nonlinear Problems	156
	Incorporating Non-Numeric Attributes into Linear Methods	158
	Summary	163
		100
Chapter 5	Building Predictive Models Using Penalized	
	Linear Methods	165
	Python Packages for Penalized Linear Regression	166
	Multivariable Regression: Predicting Wine Taste	167
	Building and Testing a Model to Predict Wine Taste	168
	Training on the Whole Data Set before Deployment	172
	Basis Expansion: Improving Performance by	172
	Creating New Variables from Old Ones	178
	Binary Classification: Using Penalized Linear	
	Regression to Detect Unexploded Mines	181
	Build a Rocks versus Mines Classifier for Deployment	191
	Multiclass Classification: Classifying Crime Scene	
	Glass Samples	204
	Summary	209
	,	

Contents

xix

Chapter 6	Ensemble Methods	211
-	Binary Decision Trees	212
	How a Binary Decision Tree Generates Predictions	213
	How to Train a Binary Decision Tree	214
	Tree Training Equals Split Point Selection	218
	How Split Point Selection Affects Predictions	218
	Algorithm for Selecting Split Points	219
	Multivariable Tree Training—Which Attribute to Split?	219
	Recursive Splitting for More Tree Depth	220
	Overfitting Binary Trees	221
	Measuring Overfit with Binary Trees	221
	Balancing Binary Tree Complexity for Best Performance	222
	Modifications for Classification and Categorical Features	225
	Bootstrap Aggregation: "Bagging"	226
	How Does the Bagging Algorithm Work?	226
	Bagging Performance—Bias versus Variance	229
	How Bagging Behaves on Multivariable Problem	231
	Bagging Needs Tree Depth for Performance	235
	Summary of Bagging	236
	Gradient Boosting	236
	Basic Principle of Gradient Boosting Algorithm	237
	Parameter Settings for Gradient Boosting	239
	How Gradient Boosting Iterates Toward a Predictive Model	240 240
	Getting the Best Performance from Gradient Boosting Gradient Boosting on a Multivariable Problem	
	Summary for Gradient Boosting	247
	Random Forest	247
	Random Forests: Bagging Plus Random Attribute Subsets	250
	Random Forests Performance Drivers	251
	Random Forests Summary	252
	Summary	252
Chapter 7	Building Ensemble Models with Python	255
•	Solving Regression Problems with Python	
	Ensemble Packages	255
	Building a Random Forest Model to Predict	
	Wine Taste	256
	Constructing a RandomForestRegressor Object	256
	Modeling Wine Taste with RandomForestRegressor	259
	Visualizing the Performance of a Random	
	Forests Regression Model	262
	Using Gradient Boosting to Predict Wine Taste	263
	Using the Class Constructor for GradientBoostingRegressor	263
	Using GradientBoostingRegressor to	
	Implement a Regression Model	267
	Assessing the Performance of a Gradient Boosting Model	269

	Contents	ххі
Coding Bagging to Predict Wine Taste	270	
Incorporating Non-Numeric Attributes in		
Python Ensemble Models	275	
Coding the Sex of Abalone for Input to Randon	n	
Forest Regression in Python	275	
Assessing Performance and the Importance of		
Coded Variables	278	
Coding the Sex of Abalone for Gradient Boostir	ng	
Regression in Python	278	
Assessing Performance and the Importance of	Coded	
Variables with Gradient Boosting	282	
Solving Binary Classification Problems with Py-	thon	
Ensemble Methods	284	
Detecting Unexploded Mines with Python Ran	dom Forest 285	
Constructing a Random Forests Model to Detec	et	
Unexploded Mines	287	
Determining the Performance of a Random		
Forests Classifier	291	
Detecting Unexploded Mines with Python		
Gradient Boosting	291	
Determining the Performance of a Gradient		
Boosting Classifier	298	
Solving Multiclass Classification Problems with	L	
Python Ensemble Methods	302	
Classifying Glass with Random Forests	302	
Dealing with Class Imbalances	305	
Classifying Glass Using Gradient Boosting	307	
Assessing the Advantage of Using Random For	rest	
Base Learners with Gradient Boosting	311	
Comparing Algorithms	314	
Summary	315	
Index	319	

Introduction

Extracting actionable information from data is changing the fabric of modern business in ways that directly affect programmers. One way is the demand for new programming skills. Market analysts predict demand for people with advanced statistics and machine learning skills will exceed supply by 140,000 to 190,000 by 2018. That means good salaries and a wide choice of interesting projects for those who have the requisite skills. Another development that affects programmers is progress in developing core tools for statistics and machine learning. This relieves programmers of the need to program intricate algorithms for themselves each time they want to try a new one. Among general-purpose programming languages, Python developers have been in the forefront, building state-of-the-art machine learning tools, but there is a gap between having the tools and being able to use them efficiently.

Programmers can gain general knowledge about machine learning in a number of ways: online courses, a number of well-written books, and so on. Many of these give excellent surveys of machine learning algorithms and examples of their use, but because of the availability of so many different algorithms, it's difficult to cover the details of their usage in a survey.

This leaves a gap for the practitioner. The number of algorithms available requires making choices that a programmer new to machine learning might not be equipped to make until trying several, and it leaves the programmer to fill in the details of the usage of these algorithms in the context of overall problem formulation and solution.

This book attempts to close that gap. The approach taken is to restrict the algorithms covered to two families of algorithms that have proven to give optimum performance for a wide variety of problems. This assertion is supported by their dominant usage in machine learning competitions, their early inclusion in newly developed packages of machine learning tools, and their performance in

comparative studies (as discussed in Chapter 1, "The Two Essential Algorithms for Making Predictions"). Restricting attention to two algorithm families makes it possible to provide good coverage of the principles of operation and to run through the details of a number of examples showing how these algorithms apply to problems with different structures.

The book largely relies on code examples to illustrate the principles of operation for the algorithms discussed. I've discovered in the classes I teach at Hacker Dojo in Mountain View, California, that programmers generally grasp principles more readily by seeing simple code illustrations than by looking at math.

This book focuses on Python because it offers a good blend of functionality and specialized packages containing machine learning algorithms. Python is an often-used language that is well known for producing compact, readable code. That fact has led a number of leading companies to adopt Python for prototyping and deployment. Python developers are supported by a large community of fellow developers, development tools, extensions, and so forth. Python is widely used in industrial applications and in scientific programming, as well. It has a number of packages that support computationally-intensive applications like machine learning, and it is a good collection of the leading machine learning algorithms (so you don't have to code them yourself). Python is a better general-purpose programming language than specialized statistical languages such as R or SAS (Statistical Analysis System). Its collection of machine learning algorithms incorporates a number of top-flight algorithms and continues to expand.

Who This Book Is For

This book is intended for Python programmers who want to add machine learning to their repertoire, either for a specific project or as part of keeping their toolkit relevant. Perhaps a new problem has come up at work that requires machine learning. With machine learning being covered so much in the news these days, it's a useful skill to claim on a resume.

This book provides the following for Python programmers:

- A description of the basic problems that machine learning attacks
- Several state-of-the-art algorithms
- The principles of operation for these algorithms
- Process steps for specifying, designing, and qualifying a machine learning system
- Examples of the processes and algorithms
- Hackable code

To get through this book easily, your primary background requirements include an understanding of programming or computer science and the ability to read and write code. The code examples, libraries, and packages are all Python, so the book will prove most useful to Python programmers. In some cases, the book runs through code for the core of an algorithm to demonstrate the operating principles, but then uses a Python package incorporating the algorithm to apply the algorithm to problems. Seeing code often gives programmers an intuitive grasp of an algorithm in the way that seeing the math does for others. Once the understanding is in place, examples will use developed Python packages with the bells and whistles that are important for efficient use (error checking, handling input and output, developed data structures for the models, defined predictor methods incorporating the trained model, and so on).

In addition to having a programming background, some knowledge of math and statistics will help get you through the material easily. Math requirements include some undergraduate-level differential calculus (knowing how to take a derivative and a little bit of linear algebra), matrix notation, matrix multiplication, and matrix inverse. The main use of these will be to follow the derivations of some of the algorithms covered. Many times, that will be as simple as taking a derivative of a simple function or doing some basic matrix manipulations. Being able to follow the calculations at a conceptual level may aid your understanding of the algorithm. Understanding the steps in the derivation can help you to understand the strengths and weaknesses of an algorithm and can help you to decide which algorithm is likely to be the best choice for a particular problem.

This book also uses some general probability and statistics. The requirements for these include some familiarity with undergraduate-level probability and concepts such as the mean value of a list of real numbers, variance, and correlation. You can always look through the code if some of the concepts are rusty for you.

This book covers two broad classes of machine learning algorithms: penalized linear regression (for example, Ridge and Lasso) and ensemble methods (for example, Random Forests and Gradient Boosting). Each of these families contains variants that will solve regression and classification problems. (You learn the distinction between classification and regression early in the book.)

Readers who are already familiar with machine learning and are only interested in picking up one or the other of these can skip to the two chapters covering that family. Each method gets two chapters—one covering principles of operation and the other running through usage on different types of problems. Penalized linear regression is covered in Chapter 4, "Penalized Linear Regression," and Chapter 5, "Building Predictive Models Using Penalized Linear Methods." Ensemble methods are covered in Chapter 6, "Ensemble Methods," and Chapter 7, "Building Predictive Models with Python." To

familiarize yourself with the problems addressed in the chapters on usage of the algorithms, you might find it helpful to skim Chapter 2, "Understand the Problem by Understanding the Data," which deals with data exploration. Readers who are just starting out with machine learning and want to go through from start to finish might want to save Chapter 2 until they start looking at the solutions to problems in later chapters.

What This Book Covers

As mentioned earlier, this book covers two algorithm families that are relatively recent developments and that are still being actively researched. They both depend on, and have somewhat eclipsed, earlier technologies.

Penalized linear regression represents a relatively recent development in ongoing research to improve on ordinary least squares regression. Penalized linear regression has several features that make it a top choice for predictive analytics. Penalized linear regression introduces a tunable parameter that makes it possible to balance the resulting model between overfitting and underfitting. It also yields information on the relative importance of the various inputs to the predictions it makes. Both of these features are vitally important to the process of developing predictive models. In addition, penalized linear regression yields best prediction performance in some classes of problems, particularly underdetermined problems and problems with very many input parameters such as genetics and text mining. Furthermore, there's been a great deal of recent development of coordinate descent methods, making training penalized linear regression models extremely fast.

To help you understand penalized linear regression, this book recapitulates ordinary linear regression and other extensions to it, such as stepwise regression. The hope is that these will help cultivate intuition.

Ensemble methods are one of the most powerful predictive analytics tools available. They can model extremely complicated behavior, especially for problems that are vastly overdetermined, as is often the case for many web-based prediction problems (such as returning search results or predicting ad click-through rates). Many seasoned data scientists use ensemble methods as their first try because of their performance. They are also relatively simple to use, and they also rank variables in terms of predictive performance.

Ensemble methods have followed a development path parallel to penalized linear regression. Whereas penalized linear regression evolved from overcoming the limitations of ordinary regression, ensemble methods evolved to overcome the limitations of binary decision trees. Correspondingly, this book's coverage of ensemble methods covers some background on binary decision trees because ensemble methods inherit some of their properties from binary

decision trees. Understanding them helps cultivate intuition about ensemble methods.

How This Book Is Structured

This book follows the basic order in which you would approach a new prediction problem. The beginning involves developing an understanding of the data and determining how to formulate the problem, and then proceeds to try an algorithm and measure the performance. In the midst of this sequence, the book outlines the methods and reasons for the steps as they come up. Chapter 1 gives a more thorough description of the types of problems that this book covers and the methods that are used. The book uses several data sets from the UC Irvine data repository as examples, and Chapter 2 exhibits some of the methods and tools that you can use for developing insight into a new data set. Chapter 3, "Predictive Model Building: Balancing Performance, Complexity, and Big Data," talks about the difficulties of predictive analytics and techniques for addressing them. It outlines the relationships between problem complexity, model complexity, data set size, and predictive performance. It discusses overfitting and how to reliably sense overfitting. It talks about performance metrics for different types of problems. Chapters 4 and 5, respectively, cover the background on penalized linear regression and its application to problems explored in Chapter 2. Chapters 6 and 7 cover background and application for ensemble methods.

What You Need to Use This Book

To run the code examples in the book, you need to have Python 2.x, SciPy, NumPy, Pandas, and scikit-learn. These can be difficult to install due to crossdependencies and version issues. To make the installation easy, I've used a free distribution of these packages that's available from Continuum Analytics (http://continuum.io/). Their Anaconda product is a free download and includes Python 2.x and all the packages you need to run the code in this book (and more). I've run the examples on Ubuntu 14.04 Linux but haven't tried them on other operating systems.

Conventions

To help you get the most from the text and keep track of what's happening, we've used a number of conventions throughout the book.

WARNING Boxes like this one hold important, not-to-be forgotten information that is directly relevant to the surrounding text.

NOTE Notes, tips, hints, tricks, and asides to the current discussion are offset and appear like this.

As for styles in the text:

- We *highlight* new terms and important words when we introduce them.
- We show keyboard strokes like this: Ctrl+A.
- We show filenames, URLs, and code within the text like so: persistence.properties.
- We present code in two different ways:

We use a monofont type with no highlighting for most code examples. We use bold to emphasize code that's particularly important in the present context.

Source Code

As you work through the examples in this book, you may choose either to type in all the code manually or to use the source code files that accompany the book. All the source code used in this book is available for download from http://www.wiley.com/go/pythonmachinelearning. You will find the code snippets from the source code are accompanied by a download icon and note indicating the name of the program so that you know it's available for download and can easily locate it in the download file. Once at the site, simply locate the book's title (either by using the Search box or by using one of the title lists) and click the Download Code link on the book's detail page to obtain all the source code for the book.

NOTE Because many books have similar titles, you may find it easiest to search by ISBN; this book's ISBN is 978-1-118-96174-2.

After you download the code, just decompress it with your favorite compression tool.

Errata

We make every effort to ensure that no errors appear in the text or in the code. However, no one is perfect, and mistakes do occur. If you find an error in one

of our books, like a spelling mistake or faulty piece of code, we would be very grateful for your feedback. By sending in errata, you might save another reader hours of frustration, and at the same time you will be helping us provide even higher-quality information.

To find the errata page for this book, go to http://www.wiley.com and locate the title using the Search box or one of the title lists. Then, on the book details page, click the Book Errata link. On this page, you can view all errata that has been submitted for this book and posted by Wiley editors.

Machine Learning in Python®