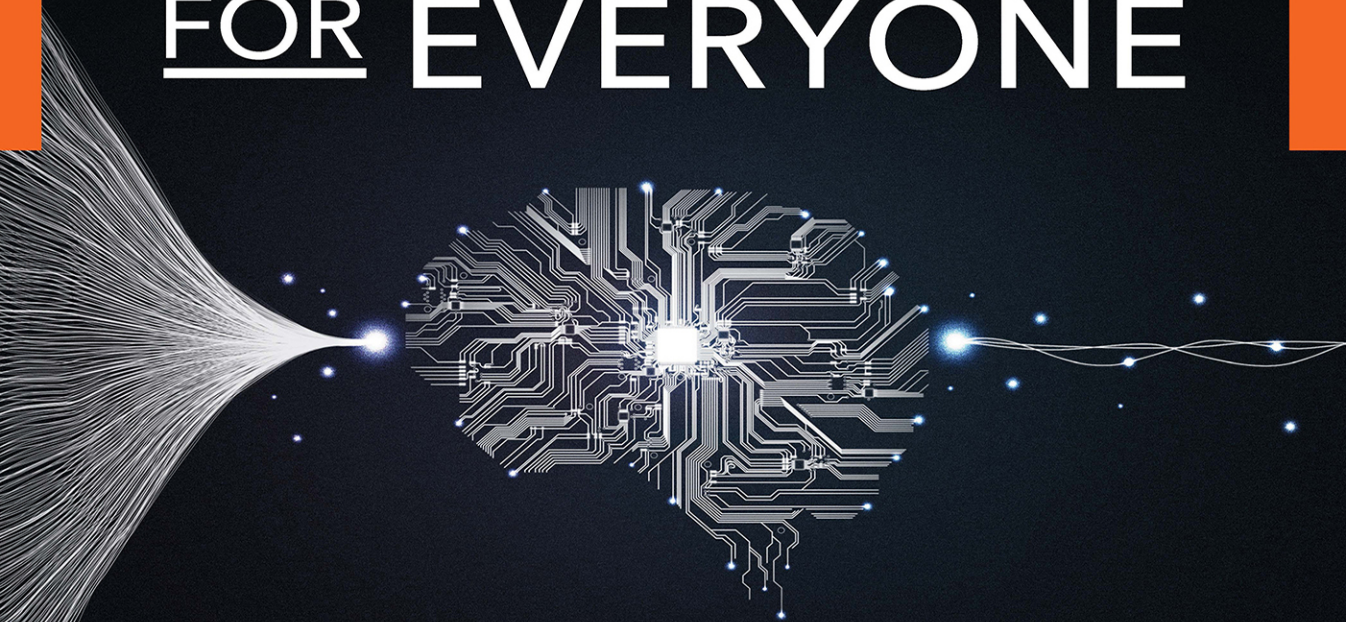


ADDISON WESLEY DATA & ANALYTICS SERIES



MACHINE LEARNING WITH PYTHON FOR EVERYONE



MARK E. FENNER

FREE SAMPLE CHAPTER

SHARE WITH OTHERS



Machine Learning with Python for Everyone

Machine Learning with Python for Everyone

Mark E. Fenner

◆◆Addison-Wesley

Boston • Columbus • New York • San Francisco • Amsterdam • Cape Town
Dubai • London • Madrid • Milan • Munich • Paris • Montreal • Toronto • Delhi • Mexico City
São Paulo • Sydney • Hong Kong • Seoul • Singapore • Taipei • Tokyo

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The author and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact intlcs@pearson.com.

Visit us on the Web: informit.com/aw

Library of Congress Control Number: 2019938761

Copyright © 2020 Pearson Education, Inc.

Cover image: [cono0430/Shutterstock](#)

Pages 58, 87: Screenshot of seaborn © 2012–2018 Michael Waskom.

Pages 167, 177, 192, 201, 278, 284, 479, 493: Screenshot of seaborn heatmap © 2012–2018 Michael Waskom.

Pages 178, 185, 196, 197, 327, 328: Screenshot of seaborn swarmplot © 2012–2018 Michael Waskom.

Page 222: Screenshot of seaborn stripplot © 2012–2018 Michael Waskom.

Pages 351, 354: Screenshot of seaborn implot © 2012–2018 Michael Waskom.

Pages 352, 353, 355: Screenshot of seaborn displot © 2012–2018 Michael Waskom.

Pages 460, 461: Screenshot of Manifold © 2007–2018, scikit-learn developers.

Page 480: Screenshot of cluster © 2007–2018, scikit-learn developers.

Pages 483, 484, 485: Image of accordion, Vereshchagin Dmitry/Shutterstock.

Page 485: Image of fighter jet, [3dgenerator/123RF](#).

Page 525: Screenshot of seaborn jointplot © 2012–2018 Michael Waskom.

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit www.pearsoned.com/permissions/.

ISBN-13: 978-0-13-484562-3

ISBN-10: 0-13-484562-5

ScoutAutomatedPrintCode

*To my son, Ethan—
with the eternal hope of a better tomorrow*

This page intentionally left blank

Contents

Foreword xxi

Preface xxiii

About the Author xxvii

I First Steps 1

1 Let's Discuss Learning 3

- 1.1 Welcome 3
- 1.2 Scope, Terminology, Prediction, and Data 4
 - 1.2.1 Features 5
 - 1.2.2 Target Values and Predictions 6
- 1.3 Putting the Machine in Machine Learning 7
- 1.4 Examples of Learning Systems 9
 - 1.4.1 Predicting Categories: Examples of Classifiers 9
 - 1.4.2 Predicting Values: Examples of Regressors 10
- 1.5 Evaluating Learning Systems 11
 - 1.5.1 Correctness 11
 - 1.5.2 Resource Consumption 12
- 1.6 A Process for Building Learning Systems 13
- 1.7 Assumptions and Reality of Learning 15
- 1.8 End-of-Chapter Material 17
 - 1.8.1 The Road Ahead 17
 - 1.8.2 Notes 17

2 Some Technical Background 19

- 2.1 About Our Setup 19
- 2.2 The Need for Mathematical Language 19

2.3	Our Software for Tackling Machine Learning	20
2.4	Probability	21
2.4.1	Primitive Events	22
2.4.2	Independence	23
2.4.3	Conditional Probability	24
2.4.4	Distributions	25
2.5	Linear Combinations, Weighted Sums, and Dot Products	28
2.5.1	Weighted Average	30
2.5.2	Sums of Squares	32
2.5.3	Sum of Squared Errors	33
2.6	A Geometric View: Points in Space	34
2.6.1	Lines	34
2.6.2	Beyond Lines	39
2.7	Notation and the Plus-One Trick	43
2.8	Getting Groovy, Breaking the Straight-Jacket, and Nonlinearity	45
2.9	NumPy versus “All the Maths”	47
2.9.1	Back to 1D versus 2D	49
2.10	Floating-Point Issues	52
2.11	EOC	53
2.11.1	Summary	53
2.11.2	Notes	54

3 Predicting Categories: Getting Started with Classification 55

3.1	Classification Tasks	55
3.2	A Simple Classification Dataset	56
3.3	Training and Testing: Don’t Teach to the Test	59
3.4	Evaluation: Grading the Exam	62
3.5	Simple Classifier #1: Nearest Neighbors, Long Distance Relationships, and Assumptions	63
3.5.1	Defining Similarity	63
3.5.2	The k in k -NN	64
3.5.3	Answer Combination	64

3.5.4	<i>k</i> -NN, Parameters, and Nonparametric Methods	65
3.5.5	Building a <i>k</i> -NN Classification Model	66
3.6	Simple Classifier #2: Naive Bayes, Probability, and Broken Promises	68
3.7	Simplistic Evaluation of Classifiers	70
3.7.1	Learning Performance	70
3.7.2	Resource Utilization in Classification	71
3.7.3	Stand-Alone Resource Evaluation	77
3.8	EOC	81
3.8.1	Sophomore Warning: Limitations and Open Issues	81
3.8.2	Summary	82
3.8.3	Notes	82
3.8.4	Exercises	83

4 Predicting Numerical Values: Getting Started with Regression 85

4.1	A Simple Regression Dataset	85
4.2	Nearest-Neighbors Regression and Summary Statistics	87
4.2.1	Measures of Center: Median and Mean	88
4.2.2	Building a <i>k</i> -NN Regression Model	90
4.3	Linear Regression and Errors	91
4.3.1	No Flat Earth: Why We Need Slope	92
4.3.2	Tilting the Field	94
4.3.3	Performing Linear Regression	97
4.4	Optimization: Picking the Best Answer	98
4.4.1	Random Guess	98
4.4.2	Random Step	99
4.4.3	Smart Step	99
4.4.4	Calculated Shortcuts	100

	4.4.5	Application to Linear Regression	101
4.5		Simple Evaluation and Comparison of Regressors	101
	4.5.1	Root Mean Squared Error	101
	4.5.2	Learning Performance	102
	4.5.3	Resource Utilization in Regression	102
4.6		EOC	104
	4.6.1	Limitations and Open Issues	104
	4.6.2	Summary	105
	4.6.3	Notes	105
	4.6.4	Exercises	105

II Evaluation 107

5 Evaluating and Comparing Learners 109

5.1	Evaluation and Why Less Is More	109
5.2	Terminology for Learning Phases	110
	5.2.1	Back to the Machines 110
	5.2.2	More Technically Speaking . . . 113
5.3	Major Tom, There's Something Wrong: Overfitting and Underfitting	116
	5.3.1	Synthetic Data and Linear Regression 117
	5.3.2	Manually Manipulating Model Complexity 118
	5.3.3	Goldilocks: Visualizing Overfitting, Underfitting, and "Just Right" 120
	5.3.4	Simplicity 124
	5.3.5	Take-Home Notes on Overfitting 124
5.4	From Errors to Costs	125
	5.4.1	Loss 125
	5.4.2	Cost 126

5.4.3	Score	127
5.5	(Re)Sampling: Making More from Less	128
5.5.1	Cross-Validation	128
5.5.2	Stratification	132
5.5.3	Repeated Train-Test Splits	133
5.5.4	A Better Way and Shuffling	137
5.5.5	Leave-One-Out Cross-Validation	140
5.6	Break-It-Down: Deconstructing Error into Bias and Variance	142
5.6.1	Variance of the Data	143
5.6.2	Variance of the Model	144
5.6.3	Bias of the Model	144
5.6.4	All Together Now	145
5.6.5	Examples of Bias-Variance Tradeoffs	145
5.7	Graphical Evaluation and Comparison	149
5.7.1	Learning Curves: How Much Data Do We Need?	150
5.7.2	Complexity Curves	152
5.8	Comparing Learners with Cross-Validation	154
5.9	EOC	155
5.9.1	Summary	155
5.9.2	Notes	155
5.9.3	Exercises	157

6 Evaluating Classifiers 159

6.1	Baseline Classifiers	159
6.2	Beyond Accuracy: Metrics for Classification	161
6.2.1	Eliminating Confusion from the Confusion Matrix	163
6.2.2	Ways of Being Wrong	164
6.2.3	Metrics from the Confusion Matrix	165
6.2.4	Coding the Confusion Matrix	166
6.2.5	Dealing with Multiple Classes: Multiclass Averaging	168

6.2.6	F_1	170
6.3	ROC Curves	170
6.3.1	Patterns in the ROC	173
6.3.2	Binary ROC	174
6.3.3	AUC: Area-Under-the-(ROC)-Curve	177
6.3.4	Multiclass Learners, One-versus-Rest, and ROC	179
6.4	Another Take on Multiclass: One-versus-One	181
6.4.1	Multiclass AUC Part Two: The Quest for a Single Value	182
6.5	Precision-Recall Curves	185
6.5.1	A Note on Precision-Recall Tradeoff	185
6.5.2	Constructing a Precision-Recall Curve	186
6.6	Cumulative Response and Lift Curves	187
6.7	More Sophisticated Evaluation of Classifiers: Take Two	190
6.7.1	Binary	190
6.7.2	A Novel Multiclass Problem	195
6.8	EOC	201
6.8.1	Summary	201
6.8.2	Notes	202
6.8.3	Exercises	203

7 Evaluating Regressors 205

7.1	Baseline Regressors	205
7.2	Additional Measures for Regression	207
7.2.1	Creating Our Own Evaluation Metric	207
7.2.2	Other Built-in Regression Metrics	208
7.2.3	R^2	209

- 7.3 Residual Plots 214
 - 7.3.1 Error Plots 215
 - 7.3.2 Residual Plots 217
- 7.4 A First Look at Standardization 221
- 7.5 Evaluating Regressors in a More Sophisticated Way: Take Two 225
 - 7.5.1 Cross-Validated Results on Multiple Metrics 226
 - 7.5.2 Summarizing Cross-Validated Results 230
 - 7.5.3 Residuals 230
- 7.6 EOC 232
 - 7.6.1 Summary 232
 - 7.6.2 Notes 232
 - 7.6.3 Exercises 234

III More Methods and Fundamentals 235

8 More Classification Methods 237

- 8.1 Revisiting Classification 237
- 8.2 Decision Trees 239
 - 8.2.1 Tree-Building Algorithms 242
 - 8.2.2 Let's Go: Decision Tree Time 245
 - 8.2.3 Bias and Variance in Decision Trees 249
- 8.3 Support Vector Classifiers 249
 - 8.3.1 Performing SVC 253
 - 8.3.2 Bias and Variance in SVCs 256
- 8.4 Logistic Regression 259
 - 8.4.1 Betting Odds 259
 - 8.4.2 Probabilities, Odds, and Log-Odds 262
 - 8.4.3 Just Do It: Logistic Regression Edition 267
 - 8.4.4 A Logistic Regression: A Space Oddity 268

- 8.5 Discriminant Analysis 269
 - 8.5.1 Covariance 270
 - 8.5.2 The Methods 282
 - 8.5.3 Performing DA 283
- 8.6 Assumptions, Biases, and Classifiers 285
- 8.7 Comparison of Classifiers: Take Three 287
 - 8.7.1 Digits 287
- 8.8 EOC 290
 - 8.8.1 Summary 290
 - 8.8.2 Notes 290
 - 8.8.3 Exercises 293

9 More Regression Methods 295

- 9.1 Linear Regression in the Penalty Box: Regularization 295
 - 9.1.1 Performing Regularized Regression 300
- 9.2 Support Vector Regression 301
 - 9.2.1 Hinge Loss 301
 - 9.2.2 From Linear Regression to Regularized Regression to Support Vector Regression 305
 - 9.2.3 Just Do It—SVR Style 307
- 9.3 Piecewise Constant Regression 308
 - 9.3.1 Implementing a Piecewise Constant Regressor 310
 - 9.3.2 General Notes on Implementing Models 311
- 9.4 Regression Trees 313
 - 9.4.1 Performing Regression with Trees 313
- 9.5 Comparison of Regressors: Take Three 314
- 9.6 EOC 318
 - 9.6.1 Summary 318
 - 9.6.2 Notes 318
 - 9.6.3 Exercises 319

10 Manual Feature Engineering: Manipulating Data for Fun and Profit 321

- 10.1 Feature Engineering Terminology and Motivation 321
 - 10.1.1 Why Engineer Features? 322
 - 10.1.2 When Does Engineering Happen? 323
 - 10.1.3 How Does Feature Engineering Occur? 324
- 10.2 Feature Selection and Data Reduction: Taking out the Trash 324
- 10.3 Feature Scaling 325
- 10.4 Discretization 329
- 10.5 Categorical Coding 332
 - 10.5.1 Another Way to Code and the Curious Case of the Missing Intercept 334
- 10.6 Relationships and Interactions 341
 - 10.6.1 Manual Feature Construction 341
 - 10.6.2 Interactions 343
 - 10.6.3 Adding Features with Transformers 348
- 10.7 Target Manipulations 350
 - 10.7.1 Manipulating the Input Space 351
 - 10.7.2 Manipulating the Target 353
- 10.8 EOC 356
 - 10.8.1 Summary 356
 - 10.8.2 Notes 356
 - 10.8.3 Exercises 357

11 Tuning Hyperparameters and Pipelines 359

- 11.1 Models, Parameters, Hyperparameters 360
- 11.2 Tuning Hyperparameters 362
 - 11.2.1 A Note on Computer Science and Learning Terminology 362
 - 11.2.2 An Example of Complete Search 362
 - 11.2.3 Using Randomness to Search for a Needle in a Haystack 368

- 11.3 Down the Recursive Rabbit Hole:
Nested Cross-Validation 370
 - 11.3.1 Cross-Validation, Redux 370
 - 11.3.2 GridSearch as a Model 371
 - 11.3.3 Cross-Validation Nested
within Cross-Validation 372
 - 11.3.4 Comments on Nested
CV 375
- 11.4 Pipelines 377
 - 11.4.1 A Simple Pipeline 378
 - 11.4.2 A More Complex
Pipeline 379
- 11.5 Pipelines and Tuning Together 380
- 11.6 EOC 382
 - 11.6.1 Summary 382
 - 11.6.2 Notes 382
 - 11.6.3 Exercises 383

IV Adding Complexity 385

12 Combining Learners 387

- 12.1 Ensembles 387
- 12.2 Voting Ensembles 389
- 12.3 Bagging and Random Forests 390
 - 12.3.1 Bootstrapping 390
 - 12.3.2 From Bootstrapping to
Bagging 394
 - 12.3.3 Through the Random
Forest 396
- 12.4 Boosting 398
 - 12.4.1 Boosting Details 399
- 12.5 Comparing the Tree-Ensemble
Methods 401
- 12.6 EOC 405
 - 12.6.1 Summary 405
 - 12.6.2 Notes 405
 - 12.6.3 Exercises 406

13 Models That Engineer Features for Us 409

- 13.1 Feature Selection 411
 - 13.1.1 Single-Step Filtering with Metric-Based Feature Selection 412
 - 13.1.2 Model-Based Feature Selection 423
 - 13.1.3 Integrating Feature Selection with a Learning Pipeline 426
- 13.2 Feature Construction with Kernels 428
 - 13.2.1 A Kernel Motivator 428
 - 13.2.2 Manual Kernel Methods 433
 - 13.2.3 Kernel Methods and Kernel Options 438
 - 13.2.4 Kernelized SVCs: SVMs 442
 - 13.2.5 Take-Home Notes on SVM and an Example 443
- 13.3 Principal Components Analysis: An Unsupervised Technique 445
 - 13.3.1 A Warm Up: Centering 445
 - 13.3.2 Finding a Different Best Line 448
 - 13.3.3 A First PCA 449
 - 13.3.4 Under the Hood of PCA 452
 - 13.3.5 A Finale: Comments on General PCA 457
 - 13.3.6 Kernel PCA and Manifold Methods 458
- 13.4 EOC 462
 - 13.4.1 Summary 462
 - 13.4.2 Notes 462
 - 13.4.3 Exercises 467

14 Feature Engineering for Domains: Domain-Specific Learning 469

- 14.1 Working with Text 470
 - 14.1.1 Encoding Text 471
 - 14.1.2 Example of Text Learning 476
- 14.2 Clustering 479
 - 14.2.1 k -Means Clustering 479

- 14.3 Working with Images 481
 - 14.3.1 Bag of Visual Words 481
 - 14.3.2 Our Image Data 482
 - 14.3.3 An End-to-End System 483
 - 14.3.4 Complete Code of BoVW Transformer 491
- 14.4 EOC 493
 - 14.4.1 Summary 493
 - 14.4.2 Notes 494
 - 14.4.3 Exercises 495

15 Connections, Extensions, and Further Directions 497

- 15.1 Optimization 497
- 15.2 Linear Regression from Raw Materials 500
 - 15.2.1 A Graphical View of Linear Regression 504
- 15.3 Building Logistic Regression from Raw Materials 504
 - 15.3.1 Logistic Regression with Zero-One Coding 506
 - 15.3.2 Logistic Regression with Plus-One Minus-One Coding 508
 - 15.3.3 A Graphical View of Logistic Regression 509
- 15.4 SVM from Raw Materials 510
- 15.5 Neural Networks 512
 - 15.5.1 A NN View of Linear Regression 512
 - 15.5.2 A NN View of Logistic Regression 515
 - 15.5.3 Beyond Basic Neural Networks 516
- 15.6 Probabilistic Graphical Models 516
 - 15.6.1 Sampling 518
 - 15.6.2 A PGM View of Linear Regression 519

15.6.3	A PGM View of Logistic Regression	523
15.7	EOC	525
15.7.1	Summary	525
15.7.2	Notes	526
15.7.3	Exercises	527

A mlwpy.py Listing 529

Index 537

This page intentionally left blank

Foreword

Whether it is called statistics, data science, machine learning, or artificial intelligence, learning patterns from data is transforming the world. Nearly every industry imaginable has been touched (or soon will be) by machine learning. The combined progress of both hardware and software improvements are driving rapid advancements in the field, though it is upon software that most people focus their attention.

While many languages are used for machine learning, including R, C/C++, Fortran, and Go, Python has proven remarkably popular. This is in large part thanks to scikit-learn, which makes it easy to not only train a host of different models but to also engineer features, evaluate the model quality, and score new data. The scikit-learn project has quickly become one of Python's most important and powerful software libraries.

While advanced mathematical concepts underpin machine learning, it is entirely possible to train complex models without a thorough background in calculus and matrix algebra. For many people, getting into machine learning through programming, rather than math, is a more attainable goal. That is precisely the goal of this book: to use Python as a hook into machine learning and then add in some math as needed. Following in the footsteps of *R for Everyone* and *Pandas for Everyone*, *Machine Learning with Python for Everyone* strives to be open and accessible to anyone looking to learn about this exciting area of math and computation.

Mark Fenner has spent years practicing the communication of science and machine learning concepts to people of varying backgrounds, honing his ability to break down complex ideas into simple components. That experience results in a form of storytelling that explains concepts while minimizing jargon and providing concrete examples. The book is easy to read, with many code samples so the reader can follow along on their computer.

With more people than ever eager to understand and implement machine learning, it is essential to have practical resources to guide them, both quickly and thoughtfully. Mark fills that need with this insightful and engaging text. *Machine Learning with Python for Everyone* lives up to its name, allowing people with all manner of previous training to quickly improve their machine learning knowledge and skills, greatly increasing access to this important field.

Jared Lander,
Series Editor

This page intentionally left blank

Preface

In 1983, the movie *WarGames* came out. I was a preteen and I was absolutely engrossed: by the possibility of a nuclear apocalypse, by the almost magical way the lead character interacted with computer systems, but mostly by the potential of machines that could *learn*. I spent years studying the strategic nuclear arsenals of the East and the West—fortunately with a naivete of a tweener—but it was almost ten years before I took my first serious steps in computer programming. Teaching a computer to do a set process was amazing. Learning the intricacies of complex systems and bending them around my curiosity was a great experience. Still, I had a large step forward to take. A few short years later, I worked with my first program that was explicitly designed to *learn*. I was blown away and I knew I found my intellectual home. I want to share the world of *computer programs that learn* with you.

Audience

Who do I think *you* are? I've written *Machine Learning with Python for Everyone* for the absolute beginner to machine learning. Even more so, you may well have very little college-level mathematics in your toolbox *and I'm not going to try to change that*. While many machine learning books are very heavy on mathematical concepts and equations, I've done my best to *minimize* the amount of mathematical luggage you'll have to carry. I do expect, given the book's title, that you'll have some basic proficiency in Python. If you can *read* Python, you'll be able to get a lot more out of our discussions. While many books on machine learning rely on mathematics, I'm relying on stories, pictures, and Python code to communicate with you. There *will* be the occasional equation. Largely, these can be skipped if you are so inclined. But, if I've done my job well, I'll have given you enough context around the equation to maybe—just *maybe*—understand what it is trying to say.

Why might you have this book in your hand? The least common denominator is that all of my readers want to *learn* about machine learning. Now, you might be coming from very different backgrounds: a student in an introductory computing class focused on machine learning, a mid-career business analyst who all of sudden has been thrust beyond the limits of spreadsheet analysis, a tech hobbyist looking to expand her interests, or a scientist needing to analyze data in a new way. Machine learning is permeating society. Depending on your background, *Machine Learning with Python for Everyone* has different things to offer you. Even a mathematically sophisticated reader who is looking to do a break-in to machine learning using Python can get a lot out of this book.

So, my goal is to take someone with an interest or need to do some machine learning and teach them the *process* and the most important *concepts* of machine learning in a concrete way using the Python scikit-learn library and some of its friends. You'll come

away with overall patterns, strategies, pitfalls, and gotchas that will be applicable in every learning system you ever study, build, or use.

Approach

Many books that try to explain mathematical topics, such as machine learning, do so by presenting equations as if they tell a story to the uninitiated. I think that leaves many of us—even those of us who like mathematics!—stuck. Personally, I build a far better mental picture of the process of machine learning by combining visual and verbal descriptions with *running code*. I’m a computer scientist at heart and by training. I love building things. Building things is how I know that I’ve reached a level where I *really* understand them. You might be familiar with the phrase, “If you really want to know something, teach it to someone.” Well, there’s a follow-on. “If you really want to know something, teach a computer to do it!” That’s my take on how I’m going to teach you machine learning. With minimal mathematics, I want to give you the concepts behind the most important and frequently used machine learning tools and techniques. Then, I want you to immediately see how to make a computer do it. One note: we won’t be programming these methods from scratch. We’ll be standing on the shoulders of giants and using some very powerful, time-saving, prebuilt software libraries (more on that shortly).

We won’t be covering all of these libraries in great detail—there is simply too much material to do that. Instead, we are going to be practical. We are going to use the best tool for the job. I’ll explain enough to orient you in the concept we’re using—and then we’ll get to using it. For our mathematically inclined colleagues, I’ll give pointers to more in-depth references they can pursue. I’ll save most of this for end-of-the-chapter notes so the rest of us can skip it easily.

If you are flipping through this introduction, deciding if you want to invest time in this book, I want to give you some insight into things that are out-of-scope for us. We aren’t going to dive into mathematical proofs or rely on mathematics to explain things. There are many books out there that follow that path and I’ll give pointers to my favorites at the ends of the chapters. Likewise, I’m going to assume that you are fluent in basic- to intermediate-level Python programming. However, for more advanced Python topics—and things that show up from third-party packages like NumPy or Pandas—I’ll explain enough of what’s going on so that you can understand each technique and its context.

Overview

In **Part I**, we establish a foundation. I’ll give you some verbal and conceptual introductions to machine learning in Chapter 1. In Chapter 2 we introduce and take a slightly different approach to some mathematical and computational topics that show up repeatedly in machine learning. Chapters 3 and 4 walk you through your first steps in building, training, and evaluating learning systems that classify examples (classifiers) and quantify examples (regressors).

Part II shifts our focus to the most important aspect of applied machine learning systems: evaluating the success of our system in a realistic way. Chapter 5 talks about general

evaluation techniques that will apply to all of our learning systems. Chapters 6 and 7 take those general techniques and add evaluation capabilities for classifiers and regressors.

Part III broadens our toolbox of learning techniques and fills out the components of a practical learning system. Chapters 8 and 9 give us additional classification and regression techniques. Chapter 10 describes *feature engineering*: how we smooth the edges of rough data into forms that we can use for learning. Chapter 11 shows how to chain multiple steps together as a single learner and how to tune a learner's inner workings for better performance.

Part IV takes us beyond the basics and discusses more recent techniques that are driving machine learning forward. We look at learners that are made up of multiple little learners in Chapter 12. Chapter 13 discusses learning techniques that incorporate automated feature engineering. Chapter 14 is a wonderful capstone because it takes the techniques we describe throughout the book and applies them to two particularly interesting types of data: images and text. Chapter 15 both reviews many of the techniques we discuss and shows how they relate to more advanced learning architectures—neural networks and graphical models.

Our main focus is on the techniques of machine learning. We will investigate a number of learning algorithms and other processing methods along the way. However, completeness is not our goal. We'll discuss the most common techniques and only glance briefly at the two large subareas of machine learning: graphical models and neural, or deep, networks. However, we will see how the techniques we focus on relate to these more advanced methods.

Another topic we won't cover is implementing specific learning algorithms. We'll build on top of the algorithms that are already available in scikit-learn and friends; we'll create larger solutions using them as components. Still, someone has to implement the gears and cogs inside the black-box we funnel data into. If you are really interested in implementation aspects, you are in good company: I love them! Have all your friends buy a copy of this book, so I can argue I need to write a follow-up that dives into these lower-level details.

Acknowledgments

I must take a few moments to thank several people that have contributed greatly to this book. My editor at Pearson, Debra Williams Cauley, has been instrumental in every phase of this book's development. From our initial meetings, to her probing for a topic that might meet both our needs, to gently shepherding me through many (many!) early drafts, to constantly giving me just enough of a push to keep going, and finally climbing the steepest parts of the mountain at its peak . . . through all of these phases, Debra has shown the highest degrees of professionalism. I can only respond with a heartfelt *thank you*.

My wife, Dr. Barbara Fenner, also deserves more praise and thanks than I can give her in this short space. In addition to the burdens that any partner of an author must bear, she *also* served as my primary draft reader *and* our intrepid illustrator. She did the hard work of drafting all of the non-computer-generated diagrams in this book. While this is not our first joint academic project, it has been turned into the longest. Her patience is, by all appearances, never ending. Barbara, *I thank you!*

My primary technical reader was Marilyn Roth. Marilyn was unfailingly positive towards even my most egregious errors. *Machine Learning with Python for Everyone* is immeasurably better for her input. *Thank you.*

I would also like to thank several members of Pearson's editorial staff: Alina Kirsanova and Dmitry Kirsanov, Julie Nahil, and many other behind-the-scenes folks that I didn't have the pleasure of meeting. This book would not exist without you and your hardworking professionalism. *Thank you.*

Publisher's Note

The text contains unavoidable references to color in figures. To assist readers of the print edition, color PDFs of figures are available for download at <http://informit.com/title/9780134845623>.

For formatting purposes, decimal values in many tables have been manually rounded to two place values. In several instances, Python code and comments have been slightly modified—all such modifications should result in valid programs.

Online resources for this book are available at <https://github.com/mfenner1>.

Register your copy of *Machine Learning with Python for Everyone* on the InformIT site for convenient access to updates and/or corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account. Enter the product ISBN (9780134845623) and click Submit. Look on the Registered Products tab for an Access Bonus Content link next to this product, and follow that link to access any available bonus materials. If you would like to be notified of exclusive offers on new editions and updates, please check the box to receive email from us.

About the Author

Mark Fenner, PhD, has been teaching computing and mathematics to adult audiences—from first-year college students to grizzled veterans of industry—since 1999. In that time, he has also done research in machine learning, bioinformatics, and computer security. His projects have addressed design, implementation, and performance of machine learning and numerical algorithms; security analysis of software repositories; learning systems for user anomaly detection; probabilistic modeling of protein function; and analysis and visualization of ecological and microscopy data. He has a deep love of computing and mathematics, history, and adventure sports. When he is not actively engaged in writing, teaching, or coding, he can be found launching himself, with abandon, through the woods on his mountain bike or sipping a post-ride beer at a swimming hole. Mark holds a *niidan* rank in judo and is a certified Wilderness First Responder. He and his wife are graduates of Allegheny College and the University of Pittsburgh. Mark holds a PhD in computer science. He lives in northeastern Pennsylvania with his family and works through his company, Fenner Training and Consulting, LLC.

This page intentionally left blank

Predicting Categories: Getting Started with Classification

In [1]:

```
# setup
from mlwpy import *
%matplotlib inline
```

3.1 Classification Tasks

Now that we've laid a bit of groundwork, let's turn our attention to the main attraction: building and evaluating learning systems. We'll start with classification and we need some data to play with. If that weren't enough, we need to establish some evaluation criteria for success. All of these are just ahead.

Let me squeeze in a few quick notes on terminology. If there are only two target classes for output, we can call a learning task *binary classification*. You can think about $\{\text{Yes}, \text{No}\}$, $\{\text{Red}, \text{Black}\}$, or $\{\text{True}, \text{False}\}$ targets. Very often, binary problems are described mathematically using $\{-1, +1\}$ or $\{0, 1\}$. Computer scientists love to encode $\{\text{False}, \text{True}\}$ into the numbers $\{0, 1\}$ as the output values. In reality, $\{-1, +1\}$ or $\{0, 1\}$ are both used for mathematical convenience, and it won't make much of a difference to us. (The two encodings often cause head-scratching if you lose focus reading two different mathematical presentations. You might see one in a blog post and the other in an article and you can't reconcile them. I'll be sure to point out any differences in *this* book.) With more than two target classes, we have a *multiclass* problem.

Some classifiers try to make a decision about the output in a direct fashion. The direct approach gives us great flexibility in the relationships we find, but that very flexibility means that we aren't tied down to assumptions that might lead us to better decisions. These assumptions are similar to limiting the suspects in a crime to people that were near where the crime occurred. Sure, we could start with no assumptions at all and equally consider suspects from London, Tokyo, and New York for a crime that occurred in

Nashville. But, adding an assumption that the suspect is in Tennessee should lead to a better pool of suspects.

Other classifiers break the decision into a two-step process: (1) build a model of how likely the outcomes are and (2) pick the most likely outcome. Sometimes we prefer the second approach because we care about the grades of the prediction. For example, we might want to know how likely it is that someone is sick. That is, we want to know that there is a 90% chance someone is sick, versus a more generic estimate “yes, we think they are sick.” That becomes important when the real-world cost of our predictions is high. When cost matters, we can combine the probabilities of events with the costs of those events and come up with a decision model to choose a real-world action that balances these, possibly competing, demands. We will consider one example of each type of classifier: Nearest Neighbors goes directly to an output class, while Naive Bayes makes an intermediate stop at an estimated probability.

3.2 A Simple Classification Dataset

The *iris* dataset is included with `sklearn` and it has a long, rich history in machine learning and statistics. It is sometimes called Fisher’s Iris Dataset because Sir Ronald Fisher, a mid-20th-century statistician, used it as the sample data in one of the first academic papers that dealt with what we now call classification. Curiously, Edgar Anderson was responsible for gathering the data, but his name is not as frequently associated with the data. Bummer. History aside, what is the *iris* data? Each row describes one iris—that’s a flower, by the way—in terms of the length and width of that flower’s sepals and petals (Figure 3.1). Those are the big flowery parts and little flowery parts, if you want to be highly technical. So, we have four total measurements per iris. Each of the measurements is a length of one aspect of that iris. The final column, our classification target, is the particular species—one of three—of that iris: *setosa*, *versicolor*, or *virginica*.

We’ll load the *iris* data, take a quick tabular look at a few rows, and look at some graphs of the data.

In [2]:

```
iris = datasets.load_iris()

iris_df = pd.DataFrame(iris.data,
                       columns=iris.feature_names)

iris_df['target'] = iris.target
display(pd.concat([iris_df.head(3),
                   iris_df.tail(3)]))
```

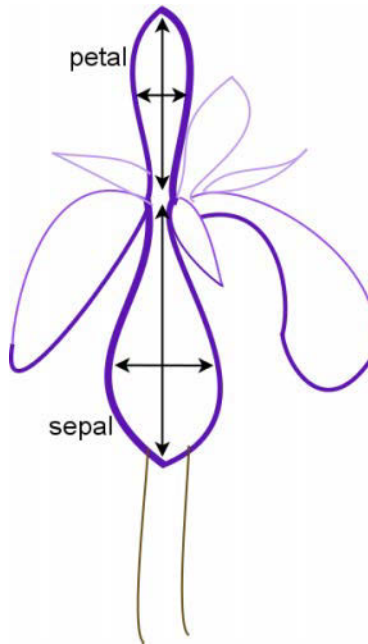
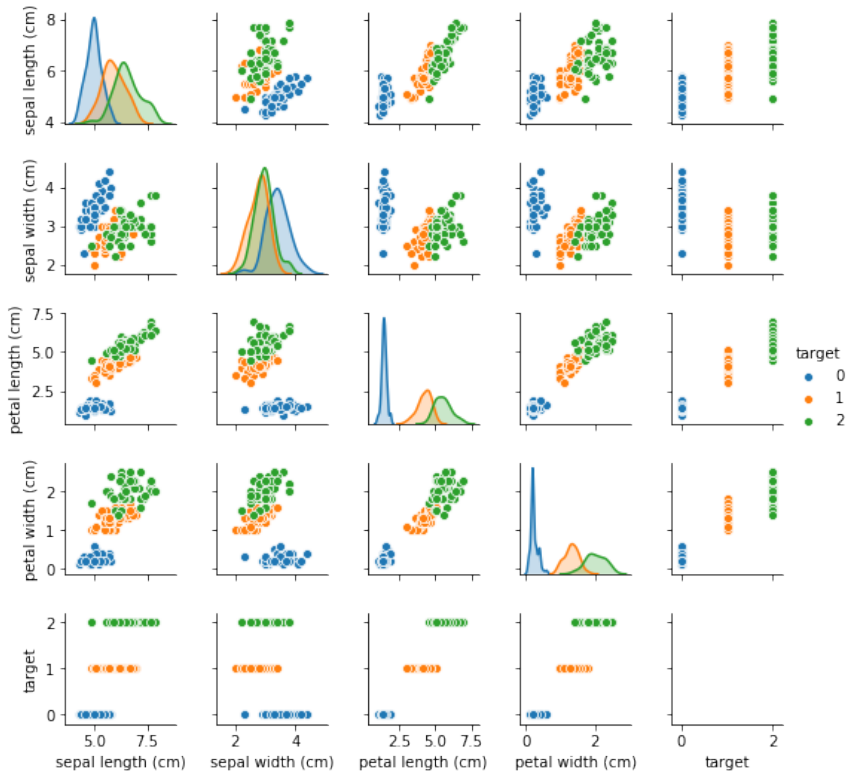


Figure 3.1 An iris and its parts.

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1000	3.5000	1.4000	0.2000	0
1	4.9000	3.0000	1.4000	0.2000	0
2	4.7000	3.2000	1.3000	0.2000	0
147	6.5000	3.0000	5.2000	2.0000	2
148	6.2000	3.4000	5.4000	2.3000	2
149	5.9000	3.0000	5.1000	1.8000	2

In [3]:

```
sns.pairplot(iris_df, hue='target', size=1.5);
```



`sns.pairplot` gives us a nice panel of graphics. Along the diagonal from the top-left to bottom-right corner, we see histograms of the frequency of the different types of iris differentiated by color. The off-diagonal entries—everything *not* on that diagonal—are scatter plots of pairs of features. You'll notice that these pairs occur twice—once above and once below the diagonal—but that each plot for a pair is flipped axis-wise on the other side of the diagonal. For example, near the bottom-right corner, we see *petal width* against *target* and then we see *target* against *petal width* (across the diagonal). When we flip the axes, we change up-down orientation to left-right orientation.

In several of the plots, the blue group (target 0) seems to stand apart from the other two groups. Which species is this?

In [4]:

```
print('targets: {}'.format(iris.target_names),
      iris.target_names[0], sep="\n")
```

```
targets: ['setosa' 'versicolor' 'virginica']
setosa
```

So, looks like *setosa* is easy to separate or partition off from the others. The *vs*, *versicolor* and *virginica*, are more intertwined.

3.3 Training and Testing: Don't Teach to the Test

Let's briefly turn our attention to how we are going to use our data. Imagine you are taking a class (Figure 3.2). Let's go wild and pretend you are studying machine learning. Besides wanting a good grade, when you take a class to learn a subject, you want to be able to use that subject in the real world. Our grade is a surrogate measure for how well we will do in the real world. Yes, I can see your grumpy faces: grades can be very bad estimates of how well we do in the real world. Well, we're in luck! We get to try to make *good* grades that really tell us how well we will do when we get out there to face reality (and, perhaps, our student loans).

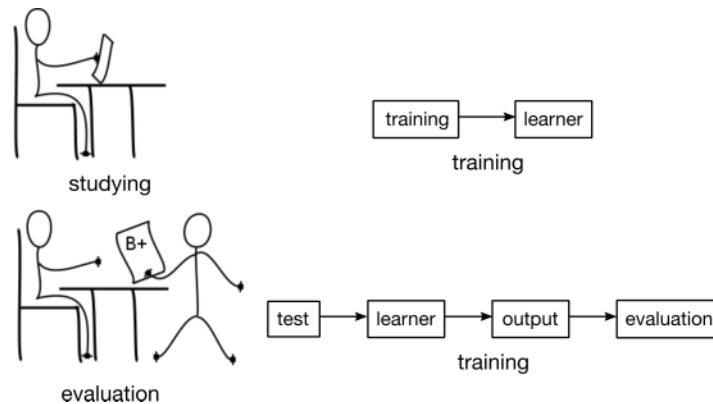


Figure 3.2 School work: training, testing, and evaluating.

So, back to our classroom setting. A common way of evaluating students is to teach them some material and then test them on it. You might be familiar with the phrase “teaching to the test.” It is usually regarded as a bad thing. Why? Because, if we teach to the test, the students will do better on the test than on other, new problems they have never seen before. They know the specific answers for the test problems, but they’ve missed out on the *general* knowledge and techniques they need to answer *novel* problems. Again, remember our goal. We want to do well in the real-world use of our subject. In a machine learning scenario, we want to do well on *unseen* examples. Our performance on unseen examples is called *generalization*. If we test ourselves on data we have already seen, we will have an overinflated estimate of our abilities on novel data.

Teachers prefer to assess students on novel problems. Why? Teachers care about how the students will do on new, never-before-seen problems. If they practice on a specific problem and figure out what’s right or wrong about their answer to it, we want that new nugget of knowledge to be something general that they can apply to other problems. If we want to estimate how well the student will do on novel problems, we have to evaluate them on novel problems. Are you starting to feel bad about studying old exams yet?

I don't want to get into too many details of too many tasks here. Still, there is one complication I feel compelled to introduce. Many presentations of learning start off using a teach-to-the-test evaluation scheme called *in-sample evaluation* or *training error*. These have their uses. However, not teaching to the test is such an important concept in learning systems that *I refuse to start you off on the wrong foot!* We just can't take an easy way out. We are going to put on our big girl and big boy pants and do this like adults with a real, *out-of-sample* or *test error* evaluation. We can use these as an estimate for our ability to generalize to unseen, future examples.

Fortunately, `sklearn` gives us some support here. We're going to use a tool from `sklearn` to avoid teaching to the test. The `train_test_split` function segments our dataset that lives in the Python variable `iris`. Remember, that dataset has two components already: the *features* and the *target*. Our new segmentation is going to split it into two buckets of examples:

1. A portion of the data that we will use to study and build up our understanding and
2. A portion of the data that we will use to test ourselves.

We will only study—that is, learn from—the *training* data. To keep ourselves honest, we will only evaluate ourselves on the *testing* data. We promise not to peek at the testing data. We started by breaking our dataset into two parts: features and target. Now, we're breaking each of those into two pieces:

1. Features → training features and testing features
2. Targets → training targets and testing targets

We'll get into more details about `train_test_split` later. Here's what a basic call looks like:

In [5]:

```
# simple train-test split
(iris_train_ftrs, iris_test_ftrs,
 iris_train_tgt, iris_test_tgt) = skms.train_test_split(iris.data,
                                                         iris.target,
                                                         test_size=.25)
print("Train features shape:", iris_train_ftrs.shape)
print("Test features shape:", iris_test_ftrs.shape)
```

Train features shape: (112, 4)

Test features shape: (38, 4)

So, our training data has 112 examples described by four features. Our testing data has 38 examples described by the same four attributes.

If you're confused about the two splits, check out Figure 3.3. Imagine we have a box drawn around a table of our total data. We identify a special column and put that special column on the right-hand side. We draw a vertical line that separates that rightmost column from the rest of the data. That vertical line is the split between our predictive

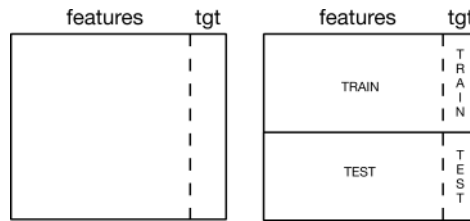


Figure 3.3 Training and testing with features and a target in a table.

features and the target feature. Now, somewhere on the box we draw a horizontal line—maybe three quarters of the way towards the bottom.

The area above the horizontal line represents the part of the data that we use for training. The area below the line is—you got it!—the testing data. And the vertical line? That single, special column is our target feature. In some learning scenarios, there might be multiple target features, but those situations don't fundamentally alter our discussion. Often, we need relatively more data to learn from and we are content with evaluating ourselves on somewhat less data, so the training part might be greater than 50 percent of the data and testing less than 50 percent. Typically, we sort data into training and testing *randomly*: imagine shuffling the examples like a deck of cards and taking the top part for training and the bottom part for testing.

Table 3.1 lists the pieces and how they relate to the *iris* dataset. Notice that I've used both some English phrases and some abbreviations for the different parts. I'll do my best to be consistent with this terminology. You'll find some differences, as you go from book A to blog B and from article C to talk D, in the use of these terms. That isn't the end of the world and there are usually close similarities. Do take a moment, however, to orient yourself when you start following a new discussion of machine learning.

Table 3.1 Relationship between Python variables and *iris* data components.

iris Python variable	Symbol	Phrase
<code>iris</code>	D_{all}	(total) dataset
<code>iris.data</code>	D_{ftrs}	train and test features
<code>iris.target</code>	D_{tgt}	train and test targets
<code>iris_train_ftrs</code>	D_{train}	training features
<code>iris_test_ftrs</code>	D_{test}	testing features
<code>iris_train_tgt</code>	D_{traintgt}	training target
<code>iris_test_tgt</code>	D_{testtgt}	testing target

One slight hiccup in the table is that `iris.data` refers to all of the input *features*. But this is the terminology that scikit-learn chose. Unfortunately, the Python variable name `data` is sort of like the mathematical x : they are both generic identifiers. `data`, as a name, can refer to just about any body of information. So, while scikit-learn is using a specific sense of the word *data* in `iris.data`, I'm going to use a more specific indicator, D_{ftrs} , for the *features* of the whole dataset.

3.4 Evaluation: Grading the Exam

We've talked a bit about how we want to design our evaluation: we don't teach to the test. So, we train on one set of questions and then evaluate on a new set of questions. How are we going to compute a grade or a score from the exam? For now—and we'll dive into this later—we are simply going to ask, "Is the answer correct?" If the answer is *true* and we predicted *true*, then we get a point! If the answer is *false* and we predicted *true*, we don't get a point. Cue :sadface:. Every correct answer will count as one point. Every missed answer will count as zero points. Every question will count equally for one or zero points. In the end, we want to know the percent we got correct, so we add up the points and divide by the number of questions. This type of evaluation is called *accuracy*, its formula being $\frac{\text{\#correct answers}}{\text{\#questions}}$. It is very much like scoring a multiple-choice exam.

So, let's write a snippet of code that captures this idea. We'll have a very short exam with four true-false questions. We'll imagine a student who finds themselves in a bind and, in a last act of desperation, answers every question with `True`. Here's the scenario:

In [6]:

```
answer_key      = np.array([True, True, False, True])
student_answers = np.array([True, True, True, True]) # desperate student!
```

We can calculate the accuracy by hand in three steps:

1. Mark each answer right or wrong.
2. Add up the correct answers.
3. Calculate the percent.

In [7]:

```
correct = answer_key == student_answers
num_correct = correct.sum() # True == 1, add them up
print("manual accuracy:", num_correct / len(answer_key))
```

manual accuracy: 0.75

Behind the scenes, sklearn's `metrics.accuracy_score` is doing an equivalent calculation:

In [8]:

```
print("sklearn accuracy:",
      metrics.accuracy_score(answer_key,
                             student_answers))
```

sklearn accuracy: 0.75

So far, we've introduced two key components in our evaluation. First, we identified which material we study from and which material we test from. Second, we decided on a method to score the exam. We are now ready to introduce our first learning method, train it, test it, and evaluate it.

3.5 Simple Classifier #1: Nearest Neighbors, Long Distance Relationships, and Assumptions

One of the simpler ideas for making predictions from a labeled dataset is:

1. Find a way to describe the similarity of two different examples.
2. When you need to make a prediction on a new, unknown example, simply take the value from the most similar known example.

This process is the nearest-neighbors algorithm in a nutshell. I have three friends *Mark*, *Barb*, *Ethan* for whom I know their favorite snacks. A new friend, Andy, is most like Mark. Mark's favorite snack is Cheetos. I predict that Andy's favorite snack is the same as Mark's: Cheetos.

There are many ways we can modify this basic template. We may consider more than *just* the single most similar example:

1. Describe similarity between pairs of examples.
2. Pick several of the most-similar examples.
3. Combine those picks to get a single answer.

3.5.1 Defining Similarity

We have complete control over what *similar* means. We could define it by calculating a *distance* between pairs of examples: `similarity = distance(example_one, example_two)`. Then, our idea of similarity becomes encoded in the way we calculate the distance. Similar things are close—a small distance apart. Dissimilar things are far away—a large distance apart.

Let's look at three ways of calculating the similarity of a pair of examples. The first, *Euclidean* distance, harkens back to high-school geometry or trig. We treat the two examples as points in space. Together, the two points define a line. We let that line be the hypotenuse of a right triangle and, armed with the Pythagorean theorem, use the other two sides of the triangle to calculate a distance (Figure 3.4). You might recall that $c^2 = a^2 + b^2$ or $c = \sqrt{a^2 + b^2}$. Or, you might just recall it as painful. Don't worry, we don't have to *do* the calculation. `scikit-learn` can be told, “Do that *thing* for me.” By now, you might be concerned that my next example can only get *worse*. Well, frankly, it could. The *Minkowski* distance would lead us down a path to Einstein and his theory of relativity . . . but we're going to avoid that black (rabbit) hole.

Instead, another option for calculating similarity makes sense when we have examples that consist of simple *Yes*, *No* or *True*, *False* features. With Boolean data, I can compare two examples very nicely by counting up the number of features that are *different*. This simple idea is clever enough that it has a name: the *Hamming* distance. You might recognize this as a close cousin—maybe even a sibling or evil twin—of accuracy. Accuracy is the percent *correct*—the percent of answers the *same* as the target—which is $\frac{\text{correct}}{\text{total}}$. Hamming distance is the number of *differences*. The practical implication is that when two sets of answers agree

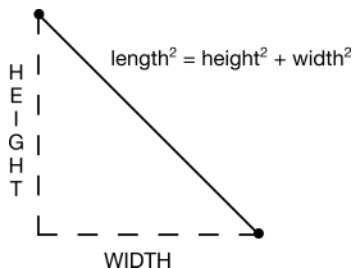


Figure 3.4 Distances from components.

completely, we want the accuracy to be high: 100%. When two sets of features are identical, we want the similarity distance between them to be low: 0.

You might have noticed that these notions of similarity have names—Euclid(-ean), Minkowski, Hamming Distance—that all fit the template of *FamousMathDude Distance*. Aside from the math dude part, the reason they share the term *distance* is because they obey the mathematical rules for what constitutes a distance. They are also called *metrics* by the mathematical wizards—that-be—as in *distance metric* or, informally, a distance measure. These mathematical terms will sometimes slip through in conversation and documentation. `sklearn`’s list of possible distance calculators is in the documentation for `neighbors.DistanceMetric`: there are about twenty metrics defined there.

3.5.2 The k in k -NN

Choices certainly make our lives complicated. After going to the trouble of choosing how to measure our local neighborhood, we have to decide how to combine the different opinions in the neighborhood. We can think about that as determining who gets to vote and how we will combine those votes.

Instead of considering only *the* nearest neighbor, we might consider some small number of nearby neighbors. Conceptually, expanding our neighborhood gives us more perspectives. From a technical viewpoint, an expanded neighborhood protects us from noise in the data (we’ll come back to this in far more detail later). Common numbers of neighbors are 1, 3, 10, or 20. Incidentally, a common name for this technique, and the abbreviation we’ll use in this book, is k -NN for “ k -Nearest Neighbors”. If we’re talking about k -NN for classification and need to clarify that, I’ll tack a C on there: k -NN- C .

3.5.3 Answer Combination

We have one last loose end to tie down. We must decide how we combine the known values (votes) from the close, or similar, neighbors. If we have an animal classification problem, four of our nearest neighbors might vote for *cat*, *cat*, *dog*, and *zebra*. How do we respond for our test example? It seems like taking the most frequent response, *cat*, would be a decent method.

In a very cool twist, we can use the exact same neighbor-based technique in *regression* problems where we try to predict a numerical value. The only thing we have to change is how we combine our neighbors' targets. If three of our nearest neighbors gave us numerical values of 3.1, 2.2, and 7.1, how do we combine them? We could use any statistic we wanted, but the mean (average) and the median (middle) are two common and useful choices. We'll come back to k -NN for regression in the next chapter.

3.5.4 k -NN, Parameters, and Nonparametric Methods

Since k -NN is the first model we're discussing, it is a bit difficult to compare it to other methods. We'll save some of those comparisons for later. There's one major difference we can dive into *right now*. I hope that grabbed your attention.

Recall the analogy of a learning model as a machine with knobs and levers on the side. Unlike many other models, k -NN outputs—the predictions—can't be computed from an input example and the values of a small, fixed set of adjustable knobs. We need *all* of the training data to figure out our output value. Really? Imagine that we throw out just one of our training examples. That example might be *the* nearest neighbor of a new test example. Surely, missing that training example will affect our output. There are other machine learning methods that have a similar requirement. Still others need some, but not *all*, of the training data when it comes to test time.

Now, you might argue that for a fixed amount of training data there could be a fixed number of knobs: say, 100 examples and 1 knob per example, giving 100 knobs. Fair enough. But then I add one example—and, poof, you now need 101 knobs, and that's a *different* machine. In this sense, the number of knobs on the k -NN machine depends on the number of examples in the training data. There is a better way to describe this dependency. Our factory machine had a side tray where we could feed additional information. We can treat the training data as this additional information. Whatever we choose, if we need either (1) a growing number of knobs or (2) the side-input tray, we say the type of machine is *nonparametric*. k -NN is a nonparametric learning method.

Nonparametric learning methods can have parameters. (Thank you for nothing, formal definitions.) What's going on here? When we call a method *nonparametric*, it means that with this method, the relationship between features and targets cannot be captured solely using a *fixed* number of parameters. For statisticians, this concept is related to the idea of parametric versus nonparametric statistics: nonparametric statistics assume less about a basket of data. However, recall that we are *not* making any assumptions about the way our black-box factory machine relates to reality. Parametric models (1) make an assumption about the form of the model and then (2) pick a specific model by setting the parameters. This corresponds to the two questions: what knobs are on the machine, and what values are they set to? We don't make assumptions like that with k -NN. However, k -NN *does* make and rely on assumptions. The most important assumption is that our similarity calculation is related to the *actual* example similarity that we want to capture.

3.5.5 Building a k -NN Classification Model

k -NN is our first example of a *model*. Remember, a supervised model is anything that captures the relationship between our features and our target. We need to discuss a few concepts that swirl around the idea of a model, so let's provide a bit of context first. Let's write down a small process we want to walk through:

1. We want to use 3-NN—three nearest neighbors—as our model.
2. We want that model to capture the relationship between the iris training features and the iris training target.
3. We want to use that model to *predict*—on previously unseen test examples—the iris target species.
4. Finally, we want to evaluate the quality of those predictions, using accuracy, by comparing predictions against reality. We didn't peek at these known answers, but we can use them as an answer key for the test.

There's a diagram of the flow of information in Figure 3.5.

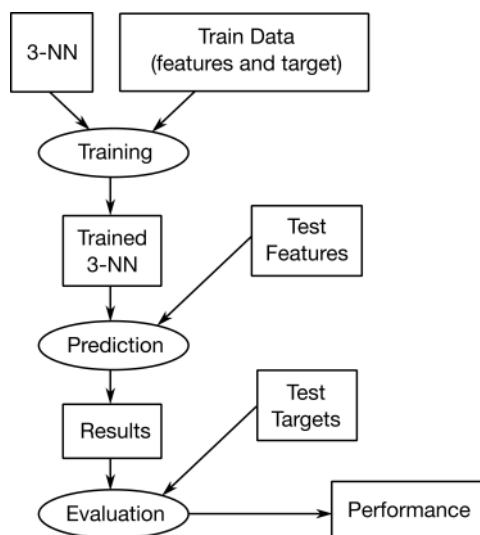


Figure 3.5 Workflow of training, testing, and evaluation for 3-NN.

As an aside on `sklearn`'s terminology, in their documentation an *estimator* is *fit* on some data and then used to *predict* on some data. If we have a training and testing split, we *fit* the *estimator* on *training data* and then use the *fit-estimator* to *predict* on the *test data*. So, let's

1. Create a 3-NN model,
2. Fit that model on the training data,
3. Use that model to predict on the test data, and
4. Evaluate those predictions using accuracy.

In [9]:

```
# default n_neighbors = 5
knn = neighbors.KNeighborsClassifier(n_neighbors=3)
fit = knn.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)

# evaluate our predictions against the held-back testing targets
print("3NN accuracy:",
      metrics.accuracy_score(iris_test_tgt, preds))
```

3NN accuracy: 1.0

Wow, 100%. We're doing great! This machine learning stuff seems pretty easy—except when it isn't. We'll come back to that shortly. We can abstract away the details of k -NN classification and write a simplified workflow template for building and assessing models in `sklearn`:

1. Build the model,
2. Fit the model using the training data,
3. Predict using the fit model on the testing data, and
4. Evaluate the quality of the predictions.

We can connect this workflow back to our conception of a model as a machine. The equivalent steps are:

1. Construct the machine, including its knobs,
2. Adjust the knobs and feed the side-inputs appropriately to capture the training data,
3. Run new examples through the machine to see what the outputs are, and
4. Evaluate the quality of the outputs.

Here's one last, quick note. The `3` in our 3-nearest-neighbors is not something that we adjust by training. It is part of the *internal* machinery of our learning machine. There is no knob on our machine for turning the `3` to a `5`. If we want a 5-NN machine, we have to build a completely different machine. The `3` is not something that is adjusted by the k -NN training process. The `3` is a *hyperparameter*. *Hyperparameters* are not trained or manipulated by the learning method they help define. An equivalent scenario is agreeing to the rules of a game and then playing the game under that *fixed* set of rules. Unless we're playing Calvinball or acting like Neo in *The Matrix*—where the flux of the rules is the point—the rules are static for the duration of the game. You can think of hyperparameters as being predetermined and fixed in place before we get a chance to do anything with them while learning. Adjusting them involves conceptually, and literally, working outside the learning box or the factory machine. We'll discuss this topic more in Chapter 11.

3.6 Simple Classifier #2: Naive Bayes, Probability, and Broken Promises

Another basic classification technique that draws directly on probability for its inspiration and operation is the Naive Bayes classifier. To give you insight into the underlying probability ideas, let me start by describing a scenario.

There's a casino that has two tables where you can sit down and play games of chance. At either table, you can play a dice game and a card game. One table is fair and the other table is rigged. Don't fall over in surprise, but we'll call these *Fair* and *Rigged*. If you sit at *Rigged*, the dice you roll have been tweaked and will only come up with six pips—the dots on the dice—one time in ten. The rest of the values are spread equally likely among 1, 2, 3, 4, and 5 pips. If you play cards, the scenario is even worse: the deck at the rigged table has no face cards—kings, queens, or jacks—in it. I've sketched this out in Figure 3.6. For those who want to nitpick, you can't tell these modifications have been made because the dice are visibly identical, the card deck is in an opaque card holder, and you make no physical contact with either the dice or the deck.

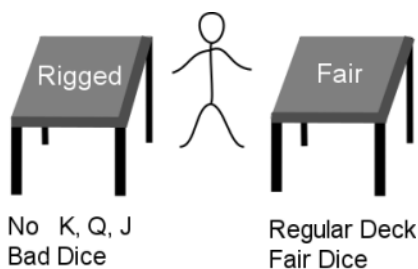


Figure 3.6 Fair and rigged tables at a casino.

Suppose I tell you—truthfully!—that you are sitting at *Rigged*. Then, when you play cards for a while and never see a face card, you aren't surprised. You also won't expect to see sixes on the die very often. Still, if you *know* you are at *Rigged*, neither of the outcomes of the dice or card events is going to *add* anything to your knowledge about the other. We *know* we are at *Rigged*, so *inferring* that we are *Rigged* doesn't add a new fact to our knowledge—although in the real world, confirmation of facts is nice.

Without knowing what table we are at, when we start seeing outcomes we receive information that indicates which table we are at. That can be turned into concrete predictions about the dice and cards. If we *know* which table we're at, that process is short-circuited and we can go directly to predictions about the dice and cards. The information about the table cuts off any gains from seeing a die or card outcome. The story is similar at *Fair*. If I tell you that you just sat down at the fair table, you would expect all the dice rolls to happen with the same probability and the face cards to come up every so often.

Now, imagine you are blindfolded and led to a table. You only know that there are two tables and you know what is happening at both—you know *Rigged* and *Fair* exist.

However, you don't know whether you are at *Rigged* or *Fair*. You sit down and the blindfold is removed. If you are dealt a face card, you immediately know you are at the *Fair* table. When we knew the table we were sitting at, knowing something about the dice didn't tell us anything additional about the cards or vice versa. Now that we don't know the table, we might get some information about the dice from the cards. If we see a face card, which doesn't exist at *Rigged*, we know we *aren't* at *Rigged*. We *must* be at *Fair*. (That's double negative logic put to good use.) As a result, we know that sixes are going to show up regularly.

Our key takeaway is that *there is no communication or causation between the dice and the cards at one of the tables*. Once we sit at *Rigged*, picking a card doesn't adjust the dice odds. The way mathematicians describe this is by saying the cards and the dice are *conditionally independent given the table*.

That scenario lets us discuss the main ideas of Naive Bayes (NB). The key component of NB is that it treats the features as if they are conditionally independent of each other given the class, just like the dice and cards at one of the tables. Knowing the table solidifies our ideas about what dice and cards we'll see. Likewise, knowing a class sets our ideas about what feature values we expect to see.

Since independence of probabilities plays out mathematically as multiplication, we get a very simple description of probabilities in a NB model. The likelihood of features for a given class can be calculated from the training data. From the training data, we store the probabilities of seeing particular features within each target class. For testing, we look up probabilities of feature values associated with a potential target class and multiply them together along with the overall class probability. We do that for each possible class. Then, we choose the class with the highest overall probability.

I constructed the casino scenario to explain what is happening with NB. However, when we use NB as our classification technique, *we assume that the conditional independence between features holds, and then we run calculations on the data*. We could be wrong. The assumptions might be broken! For example, we might not know that every time we roll a specific value on the dice, the dealers—who are *very* good card sharks—are manipulating the deck we draw from. If that were the case, there *would* be a connection between the deck and dice; our assumption that there is no connection would be *wrong*. To quote a famous statistician, George Box, “All models are wrong but some are useful.” Indeed.

Naive Bayes can be *very* useful. It turns out to be unreasonably useful in text classification. This is almost mind-blowing. It seems obvious that the words in a sentence depend on each other and on their order. We don't pick words at random; we intentionally put the right words together, in the right order, to communicate specific ideas. How can a method which *ignores* the relationship between words—which are the basis of our features in text classification—be so useful? The reasoning behind NB's success is two-fold. First, Naive Bayes is a relatively *simple* learning method that is hard to distract with irrelevant details. Second, since it is particularly simple, it benefits from having *lots* of data fed into it. I'm being slightly vague here, but you'll need to jump ahead to the discussion of *overfitting* (Section 5.3) to get more out of me.

Let's build, fit, and evaluate a simple NB model.

In [10]:

```
nb = naive_bayes.GaussianNB()
fit = nb.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)

print("NB accuracy:",
      metrics.accuracy_score(iris_test_tgt, preds))
```

NB accuracy: 1.0

Again, we are perfect. Don't be misled, though. Our success says more about the ease of the dataset than our skills at machine learning.

3.7 Simplistic Evaluation of Classifiers

We have everything lined up for the fireworks! We have data, we have methods, and we have an evaluation scheme. As the Italians say, “*Andiamo!*” Let's go!

3.7.1 Learning Performance

Shortly, we'll see a simple Python program to compare our two learners: k -NN and NB. Instead of using the names imported by our setup statement `from mlwpy import *` at the start of the chapter, it has its `imports` written out. This code is what you would write in a stand-alone script or in a notebook that *doesn't* import our convenience setup. You'll notice that we rewrote the `train_test_split` call and we also made the test set size significantly bigger. Why? Training on less data makes it a harder problem. You'll also notice that I sent an extra argument to `train_test_split`: `random_state=42` hacks the randomness of the train-test split and gives us a repeatable result. Without it, every run of the cell would result in different evaluations. Normally we want that, but here I want to be able to talk about the results *knowing* what they are.

In [11]:

```
# stand-alone code
from sklearn import (datasets, metrics,
                     model_selection as skms,
                     naive_bayes, neighbors)

# we set random_state so the results are reproducible
# otherwise, we get different training and testing sets
# more details in Chapter 5
iris = datasets.load_iris()
```

```
(iris_train_ftrs, iris_test_ftrs,
 iris_train_tgt, iris_test_tgt) = skms.train_test_split(iris.data,
                                                         iris.target,
                                                         test_size=.90,
                                                         random_state=42)

models = {'knn': neighbors.KNeighborsClassifier(n_neighbors=3),
          'NB' : naive_bayes.GaussianNB()}

for name, model in models.items():
    fit = model.fit(iris_train_ftrs, iris_train_tgt)
    predictions = fit.predict(iris_test_ftrs)

    score = metrics.accuracy_score(iris_test_tgt, predictions)
    print("{:>3s}: {:.2f}".format(name, score))
```

KNN: 0.96

NB: 0.81

With a test set size of 90% of the data, k -NN does fairly well and NB does a bit *meh* on this train-test split. If you rerun this code many times without `random_state` set and you use a more moderate amount of testing data, we get upwards of 97+% accuracy on both methods for many repeated runs. So, from a learning performance perspective, *iris* is a fairly easy problem. It is reasonably easy to distinguish the different types of flowers, based on the measurements we have, using very simple classifiers.

3.7.2 Resource Utilization in Classification

Everything we do on a computer comes with a cost in terms of processing time and memory. Often, computer scientists will talk about memory as storage space or, simply, space. Thus, we talk about the *time and space* usage of a program or an algorithm. It may seem a bit old-fashioned to worry about resource usage on a computer; today's computer are orders of magnitude faster and larger in processing and storage capabilities than their ancestors of even a few years ago—let alone the behemoth machines of the 1960s and 1970s. So why are we going down a potentially diverting rabbit hole? There are two major reasons: extrapolation and the limits of theoretical analysis.

3.7.2.1 Extrapolation

Today, much of data science and machine learning is driven by *big data*. The very nature of big data is that it pushes the limits of our computational resources. Big data is a relative term: what's big for you might not be too big for someone with the skills and budget to compute on a large cluster of machines with GPUs (graphics processing units). One possible breaking point after which I *don't* have *small* data is when the problem is so large that I can't solve it on my laptop in a “reasonable” amount of time.

If I'm doing my prototyping and development on my laptop—so I can sip a mojito under a palm tree in the Caribbean while I'm working—how can I know what sort of resources I will need when I scale up to the full-sized problem? Well, I can take measurements of smaller problems of increasing sizes and make some educated guesses about what will happen with the full dataset. To do that, I need to quantify what's happening with the smaller data in time and space. In fairness, it is only an estimate, and adding computational horsepower doesn't always get a one-to-one payback. Doubling my available memory won't always double the size of the dataset I can process.

3.7.2.2 Limits of Theory

Some of you might be aware of a subfield of computer science called *algorithm analysis* whose job is to develop equations that relate the time and memory use of a computing task to the size of that task's input. For example, we might say that the new learning method *Foo* will take $2n + 27$ steps on n input examples. (That's a drastic simplification: we almost certainly care about how many features there are in these examples.)

So, if there is a theoretical way to know the resources needed by an algorithm, why do we care about measuring them? I'm glad you asked. Algorithm analysis typically abstracts away certain mathematical details, like constant factors and terms, that can be practically relevant to real-world run times. Algorithm analysis also (1) makes certain strong or mathematically convenient assumptions, particularly regarding the average case analysis, (2) can ignore implementation details like system architecture, and (3) often uses algorithmic idealizations, devoid of real-world practicalities and necessities, to reach its conclusions.

In short, the only way to *know* how a real-world computational system is going to consume resources, short of some specialized cases that don't apply here, is to run it and measure it. Now, it is just as possible to screw this up: you could run and measure under idealized or nonrealistic conditions. We don't want to throw out algorithmic analysis altogether. My critiques are *not* failures of algorithm analysis; it's simply open-eyed understanding its limits. Algorithm analysis will always tell us some fundamental truths about how different algorithms compare and how they behave on bigger-and-bigger inputs.

I'd like to show off a few methods of comparing the resource utilization of our two classifiers. A few caveats: quantifying program behavior can be very difficult. Everything occurring on your system can potentially have a significant impact on your learning system's resource utilization. Every difference in your input can affect your system's behavior: more examples, more features, different types of features (numerical versus symbolic), and different hyperparameters can all make the same learning algorithm behave differently and consume different resources.

3.7.2.3 Units of Measure

We need to make one small digression. We're going to be measuring the resources used by computer programs. Time is measured in seconds, and space is measured in bytes. One byte is eight bits: it can hold the answers to eight yes/no questions. Eight bits can

distinguish between 256 different values—so far, so good. However, we’ll be dealing with values that are significantly larger or smaller than our normal experience. I want you to be able to connect with these values.

We need to deal with SI prefixes. SI is short for the International Standard of scientific abbreviations—but, coming from a Romance language, the adjective is *after* the noun, so the IS is swapped. The prefixes that are important for us are in Table 3.2. Remember that the exponent is the x in 10^x ; it’s also the number of “padded zeros” on the right. That is, *kilo* means $10^3 = 1000$ and 1000 has three zeros on the right. The examples are distances that would be reasonable to measure, using that prefix, applied to meters.

Table 3.2 SI prefixes and length scale examples.

Prefix	Verbal	Exponent	Example Distance
T	tera	12	orbit of Neptune around the Sun
G	giga	9	orbit of the Moon around the Earth
M	mega	6	diameter of the Moon
K	kilo	3	a nice walk
		0	1 meter \sim 1 step
m	milli	−3	mosquito
μ	micro	−6	bacteria
n	nano	−9	DNA

There is another complicating factor. Computers typically work with base-2 amounts of storage, not base-10. So, instead of 10^x we deal with 2^x . Strictly speaking—and scientists are nothing if not strict—we need to account for this difference. For memory, we have some additional prefixes (Table 3.3) that you’ll see in use soon.

Table 3.3 SI base-two prefixes and memory scale examples.

Prefix	Verbal Prefix	Number of Bytes	Example
KiB	kibi	2^{10}	a list of about 1000 numbers
MiB	mebi	2^{20}	a short song as an MP3
GiB	gibi	2^{30}	a feature-length movie
TiB	tebi	2^{40}	a family archive of photos and movies

So, 2 MiB is *two mebi-bytes* equal to 2^{20} bytes. You’ll notice that the base-2 prefixes are also pronounced differently. Ugh. You might wonder why these step up by 10s, not by 3s as in the base-10 values. Since $2^{10} = 1024 \sim 1000 = 10^3$, multiplying by ten 2s is fairly close to multiplying by three 10s. Unfortunately, these binary prefixes, defined by large standards bodies, haven’t necessarily trickled down to daily conversational use. The good news is that within one measuring system, you’ll probably only see MiB or MB, not both. When you see MiB, just know that it isn’t quite MB.

3.7.2.4 Time

In a Jupyter notebook, we have some nice tools to measure execution times. These are great for measuring the time use of small snippets of code. If we have two different ways of coding a solution to a problem and want to compare their speed, or just want to measure how long a snippet of code takes, we can use Python's `timeit` module. The Jupyter cell magic `%timeit` gives us a convenient interface to time a line of code:

In [12]:

```
%timeit -r1 datasets.load_iris()
```

1000 loops, best of 1: 1.4 ms per loop

The `-r1` tells `timeit` to measure the timing of the snippet once. If we give a higher `r`, for repeats, the code will be run multiple times and we will get statistics. Recent versions of Jupyter default to calculating the mean and standard deviation of the results. Fortunately, for a single result we just get that single value. If you are concerned about the `1000 loops`, check out my note on it at the end of the chapter.

`%%timeit`—the two-percents make it a *cell magic*—applies the same strategy to the entire block of code in a cell:

In [13]:

```
%%timeit -r1 -n1
(iris_train_ftrs, iris_test_ftrs,
 iris_train_tgt, iris_test_tgt) = skms.train_test_split(iris.data,
                                                         iris.target,
                                                         test_size=.25)
```

1 loop, best of 1: 638 μ s per loop

And now let's point our chronometer (`timeit`) at our learning workflow:

In [14]:

```
%%timeit -r1

nb = naive_bayes.GaussianNB()
fit = nb.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)

metrics.accuracy_score(iris_test_tgt, preds)
```

1000 loops, best of 1: 1.07 ms per loop

In [15]:

```
%%timeit -r1

knn = neighbors.KNeighborsClassifier(n_neighbors=3)
fit = knn.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)

metrics.accuracy_score(iris_test_tgt, preds)
```

1000 loops, best of 1: 1.3 ms per loop

If we just want to time one line in a cell—for example, we only want to see how long it takes to fit the models—we can use a single-percent version, called a *line magic*, of `timeit`:

In [16]:

```
# fitting
nb = naive_bayes.GaussianNB()
%timeit -r1 fit = nb.fit(iris_train_ftrs, iris_train_tgt)

knn = neighbors.KNeighborsClassifier(n_neighbors=3)
%timeit -r1 fit = knn.fit(iris_train_ftrs, iris_train_tgt)
```

1000 loops, best of 1: 708 μ s per loop

1000 loops, best of 1: 425 μ s per loop

In [17]:

```
# predicting
nb = naive_bayes.GaussianNB()
fit = nb.fit(iris_train_ftrs, iris_train_tgt)
%timeit -r1 preds = fit.predict(iris_test_ftrs)

knn = neighbors.KNeighborsClassifier(n_neighbors=3)
fit = knn.fit(iris_train_ftrs, iris_train_tgt)
%timeit -r1 preds = fit.predict(iris_test_ftrs)
```

1000 loops, best of 1: 244 μ s per loop

1000 loops, best of 1: 644 μ s per loop

There seems to be a bit of a tradeoff. k -NN is faster to fit, but is slower to predict. Conversely, NB takes a bit of time to fit, but is faster predicting. If you're wondering why I didn't reuse the `knn` and `nb` from the prior cell, it's because when you `%timeit`, variable assignment are trapped inside the `timeit` magic and don't leak back out to our main code. For example, trying to use `preds` as “normal” code in the prior cell will results in a `NameError`.

3.7.2.5 Memory

We can also do a very similar sequence of steps for quick-and-dirty measurements of memory use. However, two issues raise their ugly heads: (1) our tool isn't built into Jupyter, so we need to install it and (2) there are technical details—err, opportunities?—that we'll get to in a moment. As far as installation goes, install the `memory_profiler` module with `pip` or `conda` at your terminal command line:

```
pip install memory_profiler
conda install memory_profiler
```

Then, in your notebook you will be able to use `%load_ext`. This is Jupyter's command to load a Jupyter extension module—sort of like Python's `import`. For `memory_profiler`, we use it like this:

```
%load_ext memory_profiler
```

Here it goes:

In [18]:

```
%load_ext memory_profiler
```

Use it is just like `%%timeit`. Here's the cell magic version for Naive Bayes:

In [19]:

```
%%memit
nb = naive_bayes.GaussianNB()
fit = nb.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)
```

peak memory: 144.79 MiB, increment: 0.05 MiB

And for Nearest Neighbors:

In [20]:

```
%%memit
knn = neighbors.KNeighborsClassifier(n_neighbors=3)
fit = knn.fit(iris_train_ftrs, iris_train_tgt)
preds = fit.predict(iris_test_ftrs)
```

peak memory: 144.79 MiB, increment: 0.00 MiB

3.7.2.6 Complicating Factors

You may never have considered what happens with memory on your computer. In the late 2010s, you might have 4 or 8GB of system memory, RAM, on your laptop. I have 32GB

on my workhorse powerstation—or workstation powerhorse, if you prefer. Regardless, that system memory is shared by each and every running program on your computer. It is the job of the operating system—Windows, OSX, Linux are common culprits—to manage that memory and respond to applications’ requests to use it. The OS has to be a bit of a playground supervisor to enforce sharing between the different programs.

Our small Python programs, too, are playing on that playground. We have to share with others. As we request resources like memory—or time on the playground swing—the OS will respond and give us a block of memory to use. We might actually get *more* memory than we request (more on that in a second). Likewise, when we are done with a block of memory—and being the polite playground children that we are—we will return it to the playground monitor. In both our request for memory and our return of the memory, the process incurs management overhead. Two ways that OSes simplify the process and reduce the overhead are (1) by granting memory in blocks that might be more than we need and (2) by possibly letting us keep using memory, after we’ve said we’re done with it, until someone else *actively* needs it. The net result of this is that determining the actual amount of memory that we are using—versus the amount the operating system has walled off for us—can be very tricky. Measuring additional requests within a running program is even more difficult.

Another issue further complicates matters. Python is a memory-managed language: it has its own memory management facilities on top of the OS. If you were to rerun the above cells in a Jupyter notebook, you might see a memory increment of 0.00 MiB and wonder what circuits just got fried. In that case, the old memory we used was released by us—and the operating system never shuffled it off to someone else. So, when we needed more memory, we were able to reuse the old memory and didn’t need any new memory from the OS. It is almost as if the memory was released and reclaimed by us so quickly that it was never actually gone! Now, whether or not we see an increment is also dependent on (1) what the notebook cell is doing, (2) what other memory our program has claimed and is using, (3) every other program that is running on the computer, and (4) the exact details of the operating system’s memory manager. To learn more, check out a course or textbook on operating systems.

3.7.3 Stand-Alone Resource Evaluation

To minimize these concerns and to reduce confounding variables, it is extremely useful to write small, stand-alone programs when testing memory use. We can make the script general enough to be useful for stand-alone timing, as well.

In [21]:

```
!cat scripts/knn_memtest.py

import memory_profiler, sys
from mlwpy import *

@memory_profiler.profile(precision=4)
```

```
def knn_memtest(train, train_tgt, test):
    knn = neighbors.KNeighborsClassifier(n_neighbors=3)
    fit = knn.fit(train, train_tgt)
    preds = fit.predict(test)

if __name__ == "__main__":
    iris = datasets.load_iris()
    tts = skms.train_test_split(iris.data,
                                iris.target,
                                test_size=.25)
    (iris_train_ftrs, iris_test_ftrs,
     iris_train_tgt, iris_test_tgt) = tts
    tup = (iris_train_ftrs, iris_train_tgt, iris_test_ftrs)
    knn_memtest(*tup)
```

There are a few ways to use `memory_profiler`. We've seen the line and cell magics in the previous section. In `knn_memtest.py`, we use the `@memory_profiler.profile` decorator. That extra line of Python tells the memory profiler to track the memory usage of `knn_memtest` on a line-by-line basis. When we run the script, we see memory-related output for each line of `knn_memtest`:

In [22]:

```
!python scripts/knn_memtest.py
```

Filename: scripts/knn_memtest.py

output modified for formatting purposes

Line #	Mem usage	Increment	Line Contents
4	120.5430 MiB	120.5430 MiB	<code>@memory_profiler.profile(precision=4)</code>
5			<code>def knn_memtest(train, train_tgt, test):</code>
6	120.5430 MiB	0.0000 MiB	<code>knn = neighbors.</code>
			<code>KNeighborsClassifier(n_neighbors=3)</code>
7	120.7188 MiB	0.1758 MiB	<code>fit = knn.fit(train, train_tgt)</code>
8	120.8125 MiB	0.0938 MiB	<code>preds = fit.predict(test)</code>

Here's another stand-alone script to measure the memory usage of Naive Bayes:

In [23]:

```
import functools as ft
import memory_profiler
from mlwpy import *

def nb_go(train_ftrs, test_ftrs, train_tgt):
    nb = naive_bayes.GaussianNB()
```

```

fit    = nb.fit(train_ftrs, train_tgt)
preds = fit.predict(test_ftrs)

def split_data(dataset):
    split = skms.train_test_split(dataset.data,
                                   dataset.target,
                                   test_size=.25)
    return split[:-1] # don't need test tgt

def msr_mem(go, args):
    base = memory_profiler.memory_usage()[0]
    mu = memory_profiler.memory_usage((go, args),
                                       max_usage=True)[0]
    print("{:<3}: ~{:.4f} MiB".format(go.__name__, mu-base))

if __name__ == "__main__":
    msr = msr_mem
    go = nb_go

    sd = split_data(datasets.load_iris())
    msr(go, sd)

```

nb_go: ~0.0078 MiB

`nb_go` has the *model-fit-predict* pattern we saw above. `split_data` just wraps `train_test_split` in a convenient way to use with `nb_go`. The new piece is setting up the timing wrapper in `msr_mem`. Essentially, we ask what memory is used now, run `nb_go`, and then see the maximum memory used along the way. Then, we take that max, subtract what we were using before, `max-baseline`, and that's the peak memory used by `nb_go`. `nb_go` gets passed in to `msr_mem` as `go` and then finds its way to `memory_usage`.

We can write a similar `msr_time` driver to evaluate time, and we can write a similar `knn_go` to kick off a *k*-NN classifier for measuring time and memory. Here are all four pieces in a single script:

In [24]:

```

!cat scripts/perf_01.py

import timeit, sys
import functools as ft
import memory_profiler
from mlwpy import *

def knn_go(train_ftrs, test_ftrs, train_tgt):
    knn    = neighbors.KNeighborsClassifier(n_neighbors=3)
    fit    = knn.fit(train_ftrs, train_tgt)

```

```

preds = fit.predict(test_ftrs)

def nb_go(train_ftrs, test_ftrs, train_tgt):
    nb = naive_bayes.GaussianNB()
    fit = nb.fit(train_ftrs, train_tgt)
    preds = fit.predict(test_ftrs)

def split_data(dataset):
    split = skms.train_test_split(dataset.data,
                                   dataset.target,
                                   test_size=.25)
    return split[:-1] # don't need test tgt

def msr_time(go, args):
    call = ft.partial(go, *args)
    tu = min(timeit.Timer(call).repeat(repeat=3, number=100))
    print("{:<6}: ~{:.4f} sec".format(go.__name__, tu))

def msr_mem(go, args):
    base = memory_profiler.memory_usage()[0]
    mu = memory_profiler.memory_usage((go, args),
                                       max_usage=True)[0]
    print("{:<3}: ~{:.4f} MiB".format(go.__name__, mu-base))

if __name__ == "__main__":
    which_msr = sys.argv[1]
    which_go = sys.argv[2]

    msr = {'time': msr_time, 'mem': msr_mem}[which_msr]
    go = {'nb': nb_go, 'knn': knn_go}[which_go]

    sd = split_data(datasets.load_iris())
    msr(go, sd)

```

With all this excitement, let's see where we end up using Naive Bayes:

In [25]:

```

!python scripts/perf_01.py mem nb
!python scripts/perf_01.py time nb

```

```

nb_go: ~0.1445 MiB
nb_go : ~0.1004 sec

```

And with k -NN:

In [26]:

```
!python scripts/perf_01.py mem knn
!python scripts/perf_01.py time knn
```

knn_go: ~0.3906 MiB

knn_go: ~0.1035 sec

In summary, our learning and resource performance metrics look like this (the numbers may vary a bit):

Method	Accuracy	~Time(s)	~Memory (MiB)
k -NN	0.96	0.10	.40
NB	0.80	0.10	.14

Don't read too much into the accuracy scores! I'll tell you why in a minute.

3.8 EOC

3.8.1 Sophomore Warning: Limitations and Open Issues

There are several caveats to what we've done in this chapter:

- We compared these learners on a single dataset.
- We used a very simple dataset.
- We did *no* preprocessing on the dataset.
- We used a single train-test split.
- We used accuracy to evaluate the performance.
- We didn't try different numbers of neighbors.
- We only compared two simple models.

Each one of these caveats is great! It means we have more to talk about in the forthcoming chapters. In fact, discussing *why* these are concerns and figuring out *how* to address them is the point of this book. Some of these issues have no fixed answer. For example, no one learner is best on *all* datasets. So, to find a good learner *for a particular problem*, we often try several different learners and pick the one that does the best *on that particular problem*. If that sounds like teaching-to-the-test, you're right! We have to be very careful in how we select the model we use from many potential models. Some of these issues, like our use of accuracy, will spawn a long discussion of how we quantify and visualize the performance of classifiers.

3.8.2 Summary

Wrapping up our discussion, we've seen several things in this chapter:

1. *iris*, a simple real-world dataset
2. Nearest-neighbors and Naive Bayes classifiers
3. The concept of training and testing data
4. Measuring learning performance with accuracy
5. Measuring time and space usage within a Jupyter notebook and via stand-alone scripts

3.8.3 Notes

If you happen to be a botanist or are otherwise curious, you can read Anderson's original paper on irises: www.jstor.org/stable/2394164. The version of the *iris* data with `sklearn` comes from the UCI Data repository: <https://archive.ics.uci.edu/ml/datasets/iris>.

The Minkowski distance isn't really as scary as it seems. There's another distance called the Manhattan distance. It is the distance it would take to walk as directly as possible from one point to the other, if we were on a fixed grid of streets like in Manhattan. It simply adds up the absolute values of the feature differences without squares or square roots. All Minkowski does is extend the formulas so we can pick Manhattan, Euclidean, or other distances by varying a value p . The weirdness comes in when we make p very, very big: $p \rightarrow \infty$. Of course, that has its own name: the Chebyshev distance.

If you've seen theoretical resource analysis of algorithms before, you might remember the terms *complexity analysis* or *Big-O* notation. The Big-O analysis simplifies statements on the upper bounds of resource use, as input size grows, with mathematical statements like $\mathcal{O}(n^2)$ —hence the name Big-O.

I briefly mentioned graphics processing units (GPUs). When you look at the mathematics of computer graphics, like the visuals in modern video games, it is all about describing points in space. And when we play with data, we often talk about examples as points in space. The “natural” mathematical language to describe this is *matrix algebra*. GPUs are designed to perform matrix algebra at warp speed. So, it turns out that machine learning algorithms can be run very, very efficiently on GPUs. Modern projects like Theano, TensorFlow, and Keras are designed to take advantage of GPUs for learning tasks, often using a type of learning model called a *neural network*. We'll briefly introduce these in Chapter 15.

In this chapter, we used Naive Bayes on discrete data. Therefore, learning involved making a table of how often values occurred for the different target classes. When we have continuous numerical values, the game is a bit different. In that case, learning means figuring out the center and spread of a distribution of values. Often, we assume that a *normal* distribution works well with the data; the process is then called *Gaussian Naive Bayes*—Gaussian and normal are essentially synonyms. Note that we are making an *assumption*—it might work well but we might also be *wrong*. We'll talk more about GNB in Section 8.5.

In any chapter that discusses performance, I would be remiss if I didn't tell you that “premature optimization is the root of all evil . . . in programming.” This quote is from an essay form of Donald Knuth's 1974 Turing Award—the Nobel Prize of Computer Science—acceptance speech. Knuth is, needless to say, a giant in the discipline. There are two points that underlie his quote. Point one: in a computer system, the majority of the execution time is usually tied up in a small part of the code. This observation is a form of the Pareto principle or the 80–20 rule. Point two: optimizing code is hard, error-prone, and makes the code more difficult to understand, maintain, and adapt. Putting these two points together tells us that we can waste an awful lot of programmer time optimizing code that isn't contributing to the overall performance of our system. So, what's the better way? (1) Write a good, solid, *working* system and then measure its performance. (2) Find the bottlenecks—the slow and/or calculation-intensive portions of the program. (3) Optimize those bottlenecks. We only do the work that we *know* needs to be done and has a chance at meeting our goals. We also do as little of this intense work as possible. One note: *inner loops*—the innermost nestings of repetition—are often the most fruitful targets for optimization because they are, by definition, code that is repeated the most times.

Recent versions of Jupyter now report a mean and standard deviation for `%timeit` results. However, the Python core developers and documenters prefer a different strategy for analyzing `timeit` results: they prefer either (1) taking the minimum of several repeated runs to give an idea of best-case performance, which will be more consistent for comparison sake, or (2) looking at all of the results as a whole, without summary. I think that (2) is *always* a good idea in data analysis. The mean and standard deviation are not *robust*; they respond poorly to outliers. Also, while the mean and standard deviation completely characterize normally distributed data, other distributions will be characterized in very different ways; see Chebyshev's inequality for details. I would be far happier if Jupyter reported medians and inter-quartile ranges (those are the 50th percentile and the 75th–25th percentiles). These are robust to outliers and are not based on distributional assumptions about the data.

What was up with the `1000 loops` in the `timeit` results? Essentially, we are stacking multiple runs of the same, potentially short-lived, task one after the other so we get a longer-running pseudo-task. This longer-running task plays more nicely with the level of detail that the timing functions of the operating system support. Imagine measuring a 100-yard dash using a sundial. It's going to be very hard because there's a mismatch between the time scales. As we repeat the task multiple times—our poor sprinters might get worn out but, fortunately, Python keeps chugging along—we may get more meaningful measurements. Without specifying a `number`, `timeit` will attempt to find a good number for you. In turn, this may take a while because it will try increasing values for `number`. There's also a `repeat` value you can use with `timeit`; `repeat` is an *outer loop* around the whole process. That's what we discussed computing statistics on in the prior paragraph.

3.8.4 Exercises

You might be interested in trying some classification problems on your own. You can follow the model of the sample code in this chapter with some other classification datasets

from `sklearn`: `datasets.load_wine` and `datasets.load_breast_cancer` will get you started. You can also download numerous datasets from online resources like:

- The UCI Machine Learning Repository,
<https://archive.ics.uci.edu/ml/datasets.html>
- Kaggle, www.kaggle.com/datasets

Index

Symbols

+1 trick, 38, 43–45, 336, 521
1-NN model, 145, 154–156
 for the circle problem, 464
3-NN model, 66, 193
3D datasets, 460–461
80–20 rule, 83
 Σ , in math, 30

A

accuracy, 15, 27, 163
 calculating, 62
 fundamental limits of, 163
`accuracy_score`, 62
AdaBoost, 400, 405
`AdaBoostClassifier`, 400, 403–406
additive model, 318
aggregation, 390
algorithms
 analysis of, 72
 genetic, 101
 less important than data, 15
amoeba (StackExchange user), 465
analytic learning, 18
Anderson, Edgar, 56
ANOVAs test, 463
area under the curve (AUC), 177–178,
 182–193, 202
arguments, of a function, 362
arithmetic mean, 170, *see also* average
`array` (NumPy), 276, 494
assessment, 113–115

assumptions, 55, 270, 282, 286–287, 439
attributes, 4–5
average
 computing from confusion matrix, 170
 simple, 30
 weighted, 31–32, 34, 89
average centered dot product, *see*
 covariance

B

background knowledge, 322, 331, 439
bag of global visual words (BoGVW), 483,
 488–490
bag of visual words (BoVW), 481–483
 transformer for, 491–493
bag of words (BOW), 471–473
 normalizing, 474–476
bagged classifiers
 creating, 394
 implementing, 407
bagging, 390, 394
 basic algorithm for, 395
 bias-variance in, 396
`BaggingRegressor`, 407
base models
 overfitting, 396
 well-calibrated, 407
`BaseEstimator`, 311
baseline methods, 159–161, 189, 191
baseline regressors, 205–207
baseline values, 356
basketball players, 397
Bayes optimal classifier, 464

betting odds, 259–262
 bias, 110, 144–145, 292
 addressing, 350–351
 in combined models, 390
 in SVCs, 256–259
 number of, 148
 reducing, 396, 400, 406
 bias-variance tradeoffs, 145–149, 154, 396
 in decision trees, 249
 in performance estimating, 382
 big data, 71
 Big-O analysis, 82
 bigrams, 471
 binary classification, 55, 174, 267
 confusion matrix for, 164
 binomials, 524
 bivariate correlation, 415
 black holes, models of, 467
 body mass index (BMI), 322, 410–411
 boosting, 398–401, 406
 bootstrap aggregation, *see* bagging
 bootstrap mean, 391–393
 bootstrapping, 157, 390–394
 Box, George, 69

C

C4.5, C5.0, CART algorithms, 244
 calculated shortcut strategy, 100–101, 104
Caltech101 dataset, 482–483
 Calvinball game, 67
 card games, 21
 rigged, 68–69
 case-based reasoning, 18
 categorical coding, 332–341
 categorical features, 5–7, 18, 346
 numerical values for, 85–86
 categories, 332
 predicting, 9–10
 Cauchy-Schwarz inequality, 463
 causality, 233
 Celsius, converting to Fahrenheit, 325–326
 center, *see* mean
 classification, 7, 55–58
 binary, 55, 164, 174, 267
 nonlinear, 418–419
 classification_report, 169–170
 ClassifierMixin, 202
 classifiers
 baseline, 159–161, 189, 191
 comparing, 287–290, 368
 evaluating, 70–71, 159–203, 238–239
 making decisions, 55–56
 simple, 63–81
 smart, 189
 closures, 382, 394
 clustering, 18, 479–481
 on subsets of features, 494
 coefficient of determination, 130
 coin flipping, 21
 and binomials, 524
 increasing number of, 25–27
 collections, 20
 collinearity, 340, 356
 combinations, 41
 combinatorics, 423
 complexity, 124–125
 cost of increasing, 12
 evaluating, 152–154, 363–365
 manipulating, 119–123
 penalizing, 300, 306, 502
 trading off for errors, 125–126, 295–301
 complexity analysis, 82
 compound events, 22–23
 compression, 13
 computational learning theory, 15
 computer graphics, 82
 computer memory, *see* memory
 computer science, 362
 confounding factors, 233

confusion matrix, 164, 171–178
 computing averages from, 168, 170
constant, 160–161
 constant linear model, 146
 constants, 35–38
 contrast coding, 356
 conveyor belt, 377
 convolutional neural network, 516
 corpus, 472
corrcoef (NumPy), 416
 correctness, 11–12
 correlation, 415–417, 423, 464
 squared, 415–417
 Cortez, Paulo, 195, 203
 cosine similarity, 462–463
 cost, 126–127
 comparing, for different models, 127
 lowering, 299, 497–500
 of predictions, 56
CountVectorizer, 473
 covariance, 270–292, 415–417
 between all pairs of features, 278
 exploring graphically, 292
 length-normalized, 463
 not affected by data shifting, 451
 visualizing, 275–281
 covariance matrix (CM), 279–283, 451, 456
 computing, 455
 diagonal, 281
 eigendecomposition of, 452, 456, 459
 for multiple classes, 281–282
 CRISP-DM process, 18
 cross-validation (CV), 128–131
 2-fold, 132–133
 3-fold, 128–130
 5-fold, 129–130, 132
 as a single learner, 230
 comparing learners with, 154–155
 extracting scores from, 192
 feature engineering during, 323
 flat, 370–371, 376

 leave-one-out, 140–142
 minimum number of examples for, 152
 nested, 157, 370–377
 on multiple metrics, 226–229
 with boosting, 403
 wrapping methods inside, 370–372
cross_val_predict, 192, 230
cross_val_score, 130, 132, 137, 196, 207, 379
 Cumulative Response curve, 189
 curves, 45–47
 using kernels with, 461
cut, 328, 330–331

D

data
 accuracy of, 15
 big, 71
 centering, 221, 322, 325, 445–447, 451, 457
 cleaning, 323
 collecting, 14
 converting to tabular, 470
 fuzzy towards the tails, 88
 geometric view of, 410
 incomplete, 16
 making assumptions about, 270
 modeling, 14
 more important than algorithms, 15
 multimodal, 327, 357
 noisiness of, 15
 nonlinear, 285
 preparing, 14
 preprocessing, 341
 reducing, 250–252, 324–325, 461
 redundant, 324, 340, 411
 scaling, 85, 221, 445, 447
 sparse, 333, 356, 471, 473
 standardized, 105, 221–225, 231, 315–316, 447
 synthetic, 117

data (*continued*)

- total amount of variation in, 451
- transforming, *see* feature engineering
- variance of, 143, 145, 445
- weighted, 399–400

DataFrame, 323, 363–364

datasets

- 3D, 460–461
- applying learners to, 394
- examples in, 5
- features in, 5
- finding relationships in, 445
- missing values in, 322
- multiple, 128, 156
- poorly represented classes in, 133
- reducing, 449
- single, distribution from, 390
- testing, *see* testing datasets
- training, *see* training datasets

datasets.load_boston, 105, 234

datasets.load_breast_cancer, 84, 203

datasets.load_digits, 319

datasets.load_wine, 84, 203

decision stumps, 399, 401–403

decision trees (DT), 239–249, 290–291, 464

- bagged, 395
- bias-variance tradeoffs in, 249
- building, 244, 291
- depth of, 241, 249
- flexibility of, 313
- for nonlinear data, 285–286
- performance of, 429–430
- prone to overfitting, 241
- selecting features in, 325, 412
- unique identifiers in, 241, 322
- viewed as ensembles, 405
- vs. random forests, 396

DecisionTreeClassifier, 247

decomposition, 452, 455

deep neural networks, 481

democratic legislature, 388

dependent variables, *see* targets

deployment, 14

Descartes, René, 170

design matrix, 336, 347

diabetes dataset, 85, 105, 322, 416

diagonal covariance matrix, 281

Diagonal Linear Discriminant Analysis (DLDA), 282–285, 292

diagrams, drawing, 245

dice rolling, 21–24

- expected value of, 31–32
- rigged, 68–69

Dietterich, Tom, 375

digits dataset, 287–290, 401

Dijkstra, Edsger, 54

directions, 441, 445

- finding the best, 449, 459
- with PCA, 451

discontinuous target, 308

discretization, 329–332

discriminant analysis (DA), 269–287, 290–292

- performing, 283–285
- variations of, 270, 282–285

distances, 63–64

- as weights, 90
- sum product of, 275
- total, 94

distractions, 109–110, 117

distributions, 25–27

- binomial, 524
- from a single dataset, 390
- normal, 27, 520–524
- of the mean, 390–391
- random, 369

domain knowledge, *see* background knowledge

dot, 29–30, 38, 47–52, 245, 455

dot products, 29–30, 38, 47–52

- advantages of, 43
- and kernels, 438–441, 458–459, 461
- average centered, *see* covariance

length-normalized, 462–463
 double cross strategy, 375
 dual problem, solving, 459
 dummy coding, *see* one-hot coding
 dummy methods, *see* baseline methods

E

edit distance, 439, 464
 educated guesses, 71
 eigendecomposition (EIGD), 452, 456, 458, 465–466
 eigenvalues and eigenvectors, 456–457
 Einstein, Albert, 124
ElasticNet, 318
 empirical loss, 125
 ensembles, 387–390
 enterprises, competitive advantages of, 16
 entropy, 464
enumerate, 494
enumerate_outer, 491–492, 494
 error plots, 215–217
 errors
 between predictions and reality, 350
 ignoring, 302–305
 in data collection process, 322
 in measurements, 15, 142–143, 241
 margin, 254
 measuring, 33
 minimizing, 448–449, 451
 negating, 207
 positive, 33
 sources of, 145
 trading off for complexity, 125–126, 295–301
 vs. residuals, 218
 vs. score, 207
 weighted, 399
 estimated values, *see* predicted values
 estimators, 66
 Euclidean distance, 63, 367
 Euclidean space, 466

evaluation, 14, 62, 109–157
 deterministic, 142
 events
 compound vs. primitive, 22–23
 probability distribution of, 25–27
 random, 21–22
 examples, 5
 dependent vs. independent, 391
 distance between, 63–64, 438–439
 duplicating by weight, 399
 focusing on hard, 252, 398
 grouping together, 479
 learning from, 4
 quantity of, 15
 relationships between, 434
 supporting, 252
 tricky vs. bad, 144
 execution time, *see* time
 expected value, 31–32
 extract-transform-load (ETL), 323
 extrapolation, 71
 extreme gradient boosting, 406
 extreme random forest, 397–398

F

F_1 calculation, 170
f_classif, 422
f_regression, 416–417
 Facebook, 109, 388
 factor analysis (FA), 466
 factorization, 452, 455
 factory machines, 7–9, 114
 choosing knob values for, 115, 144, 156, 337
 stringing together, 377
 testing, 110–113
 with a side tray, 65
 Fahrenheit, converting to Celsius, 325–326
 failures, in a legal system, 12
 fair bet, 259

- false negative rate (FNR), 164–166
- false positive rate (FPR), 164–166, 173–181
- Fawcett, Tom, 18
- feature construction, 322, 341–350, 410–411
 - manual, 341–343
 - with kernels, 428–445
- feature engineering, 321–356
 - how to perform, 324
 - limitations of, 377
 - when to perform, 323–324
- feature extraction, 322, 470
- feature selection, 322, 324–325, 410–428, 449
 - by importance, 425
 - formal statistics for, 463
 - greedy, 423–424
 - integrating with a pipeline, 426–428
 - model-based, 423–426
 - modelless, 464
 - random, 396–397, 423, 425
 - recursive, 425–426
- feature-and-split
 - finding the best, 244, 397
 - random, 397–398
- feature-pairwise Gram matrix, 464
- `feature_names`, 413–414
- features, 5
 - categorical, 7, 346
 - causing targets, 233
 - conditionally independent, 69
 - correlation between, 415–417
 - counterproductive, 322
 - covariant, 270
 - different, 63
 - evaluating, 462–463
 - interactions between, 343–348
 - irrelevant, 15, 241, 324, 411
 - number of, 146–148
 - numerical, 6–7, 18, 225, 343–344, 346
 - relationships between, 417
 - scaling, 322, 325–329
 - scoring, 412–415
 - sets of, 423
 - standardizing, 85
 - training vs. testing, 60–61
 - transforming, 348–353
 - useful, 15, 412
 - variance of, 412–415
- Fenner, Ethan, 237–238
- Fisher's Iris Dataset, *see iris* dataset
- Fisher, Sir Ronald, 56
- `fit`, 224–225, 337, 363, 367–368, 371–372, 379, 381
- fit-estimators, 66
- `fit_intercept`, 340
- `fit_transform`, 326, 413
- flash cards, 398
- flashlights, messaging with, 417–418
- flat surface, *see* planes
- flipping coins, 21
 - and binomials, 524
 - increasing number of, 25–27
- `float`, 52–53
- floating-point numbers, 52–53
- `fmin`, 500
- folds, 128
- forward stepwise selection, 463
- `fromiter` (NumPy), 494
- full joint distribution, 148
- functions
 - parameters of
 - vs. arguments, 362
 - vs. values, 360–361
 - wrapping, 361, 502
- `FunctionTransformer`, 348–349
- `functools`, 20
- fundraising campaign, 189
- future, predicting, 7
- fuzzy specialist scenario, 405

G

gain curve, *see* Lift Versus Random curve

games

- expected value of, 32
- fair, 259
- sets of rules for, 67

Gaussian Naive Bayes (GNB), 82, 282–287

generalization, 59, 126

genetic algorithms, 101

geometric mean, 170

`get_support`, 413

Ghostbusters, 218

Gini index, 202, 245, 464

Glaton regression, 7

global visual words, 483, 487–490

good old-fashioned (GOF) linear regression, 300–301, 519–521

- and complex problems, 307

gradient descent (GD), 101, 292

`GradientBoostingClassifier`, 400, 403–406

Gram matrix, 464

graphics processing units (GPUs), 71, 82

greediness, for feature selection, 423–424

`GridSearch`, 363, 368, 377, 382, 405, 427–428

- wrapped inside CV, 370–372

`GridSearchCV`, 368, 371–377

H

Hamming distance, 63

Hand and Till *M* method, 183–185, 197, 200, 202

handwritten digits, 287–290

harmonic mean, 170

Hettinger, Raymond, 54

hinge loss, 301–305, 465

`hist`, 22

`histogram`, 21

hold-out test set (HOT), 114–115

hyperparameters, 67, 115

- adjusting, 116
- choosing, 359
- cross-validation for, 371–377, 380–382
- evaluating, 363–368
- for tradeoffs between complexity and errors, 126
- overfitting, 370
- random combinations of, 368–370
- tuning, 362–369, 380–382

hyperplanes, 39

I

IBM, 3

ID3 algorithm, 244

identification variables, 241, 322, 324

identity matrix, 456, 465

illusory correlations, 233

images, 481–493

- BoVW transformer for, 491–493
- classification of, 9
- describing, 488–490
- predicting, 490–491
- processing, 485–487

`import`, 19

in-sample evaluation, 60

independence, 23

independence assumptions, 148

independent component analysis (ICA), 466

independent variables, *see* features

indicator function, 243

inductive logic programming, 18

infinity-norm, 367

information gain, 325

information theory, 417

input features, 7

inputs, *see* features

intercept, 336–341

- avoiding, 356

International Standard of scientific abbreviations (SI), 73
iris dataset, 56–58, 60–61, 82, 133, 166–168, 174, 190–195, 242, 245, 329–332, 336, 480, 495
 IsoMap, 462
 iteratively reweighted least squares (IRLS), 291
itertools, 20, 41

J

jackknife resampling, 157
jointplot, 524–525
 Jupyter notebooks, 19

K

k-Cross-Validation (CV), 129–131
 with repeated train-test splits, 137
k-Means Clustering (*k*-MC), 479–481
k-Nearest Neighbors (*k*-NN), 64–67
 1-NN model, 145, 154–156, 464
 3-NN model, 66, 193
 algorithm of, 63
 bias-variance for, 145
 building models, 66–67, 91
 combining values from, 64
 evaluating, 70–71
 metrics for, 162–163
 for nonlinear data, 285
 performance of, 74–76, 78–81, 429–430
 picking the best *k*, 113, 116, 154, 363–365
k-Nearest Neighbors classification (*k*-NN-C), 64
k-Nearest Neighbors regression (*k*-NN-R), 87–91
 comparing to linear regression, 102–104, 147–229
 evaluating, 221
 vs. piecewise constant regression, 310

Kaggle website, 406
Karate Kid, The, 182, 250
 Keras, 82
 kernel matrix, 438
 kernel methods, 458
 automated, 437–438
 learners used with, 438
 manual, 433–437
 mock-up, 437
 kernels, 438–445
 and dot products, 438–441, 458–459, 461
 approximate vs. full, 436
 feature construction with, 428–445
 linear, 253, 438
 polynomial, 253, 437
KFold, 139–140, 368
KNeighborsClassifier, 66, 362–363
KNeighborsRegressor, 91
knn_statistic, 394–395
 Knuth, Donald, 83
 kurtosis, 466

L

L_1 regularization, *see* lasso regression
 L_2 regularization, *see* ridge regression
label_binarize, 179, 183
Lasso, 300
 lasso regression (L_1), 300, 307
 blending with ridge regression, 318
 selecting features in, 325, 411, 424
 learning algorithms, 8
 learning curves, 131, 150–152
 in *sklearn*, 157
 learning methods
 incremental/decremental, 130
 nonparametric, 65
 parameters of, 115
 requiring normalization, 221
 learning models, *see* models

- learning systems, 9–10
 - building, 13–15, 366
 - choosing, 81
 - combining multiple, *see* ensembles
 - evaluating, 11–13, 109–157
 - from examples, 4, 9–11
 - performance of, 102
 - overestimating, 109
 - tolerating mistakes in data, 16
 - used with kernel methods, 438
- `learning_curve`, 150–152
- least-squares fitting, 101
- leave-one-out cross-validation (LOOCV), 140–142
- length-normalized covariance, 463
- length-normalized dot product, 462–463
- Levenshtein distance, 464
- `liblinear`, 291–292
- `libsvm`, 291, 443, 465
- Lift Versus Random curve, 189, 193
- limited capacity, 109–110, 117
- limited resources, 187
- `linalg.svd` (NumPy), 455
- line magic, 75
- linear algebra, 452, 457, 465
- linear combination, 28
- Linear Discriminant Analysis (LDA), 282–285, 495
- linear kernel, 253, 438
- linear regression (LR), 91–97, 305
 - bias of, 350
 - bias-variance for, 146–147
 - calculating predicted values with, 97, 265
 - comparing to k -NN-R, 102–104, 229
 - complexity of, 119–123
 - default metric for, 209
 - example of, 118
 - for nonlinear data, 285
 - from raw materials, 500–504
 - good old-fashioned (GOF), 300–301, 307, 519–521
 - graphical presentation of, 504
 - performing, 97
 - piecewise, 309–313
 - regularized, 296–301
 - relating to k -NN, 147–148
 - selecting features in, 425
 - using standardized data for, 105
 - viewed as ensembles, 405
- linear relationships, 415, 417
- linearity, 285
- `LinearRegression`, 371
- `LinearSVC`, 253, 291, 465
- lines, 34–39
 - between classes, 250
 - drawing through points, 92, 237–238
 - finding the best, 98–101, 253, 268–269, 350, 410, 448–449, 457, 465
 - piecewise, 313
 - sloped, 37, 94–97
 - straight, 91
 - limited capacity of, 122
- local visual words, 483–488
 - extracting, 485–487
 - finding synonyms for, 487–488
- log-odds, 259, 262–266
 - predicting, 505–508
- logistic regression (LogReg), 259–269, 287, 290–292
 - and loss, 526
 - calculating predicted values with, 265
 - for nonlinear data, 285
 - from raw materials, 504–509
 - kernelized, 436
 - performance of, 429
 - PGM view of, 523–525
 - solving perfectly separable classification problems with, 268–269
- `LogisticRegression`, 267, 292
- `logreg_loss_01`, 507
- lookup tables, 13

loss, 125–126, 295
 defining, 501
 hinge, 301–305, 465
 minimizing, 526
 vs. score, 127, 207

M

M method, 183–185, 197, 200, 202

machine learning

 and math, 19–20
 definition of, 4
 limits of, 15
 running on GPUs, 82

`macro`, 168

macro precision, 168

`magical_minimum_finder`, 500–511

`make_cost`, 502–503

`make_scorer`, 185, 196, 208

Manhattan distance, 82, 367

manifolds, 459–462

 differentiable, 466–467

Mann-Whitney *U* statistic, 202

margin errors, 254

mathematics

 1-based indexing in, 54

\sum notation, 30

 derivatives, 526

 eigenvalues and eigenvectors, 456–457

 linear algebra, 452, 457, 465

 matrix algebra, 82, 465–466

 optimization, 500

 parameters, 318

`matplotlib`, 20, 22, 222–223

matrices, 456

 breaking down, 457

 decomposition (factorization), 452, 455

 identity, 465

 multiplication of, 82, 465

 orthogonal, 465–466

 squaring, 466

 transposing, 465

Matrix, The, 67

`matshow`, 275–277

`max_depth`, 242

maximum margin separator, 252

mean, 54, 85, 271, 446

 arithmetic, 170, *see also* average

 bootstrap, 391–393

 computing, 390–391, 395

 definition of, 88

 distribution of, 390–391

 empirical, 457

 for two variables, multiplying, 271

 geometric, 170

 harmonic, 170

 multiple, for each train-test split, 231

 predicting, 147, 205

 weighted, 89–90

mean absolute error (MAE), 209

mean squared error (MSE), 91, 101, 130, 209

`mean_squared_error`, 91, 126

measurements

 accuracy of, 27

 critical, 16

 errors in, 15, 142–143, 241

 levels of, 18

 overlapping, 410

 rescaling, 328, 414

 scales of, 412–414

median, 206, 446

 computing on training data, 349

 definition of, 88

 predicting, 205

median absolute error, 209

medical diagnosis, 10

 assessing correctness of, 11–12

 confusion matrix for, 165–166

 example of, 6–7

 for rare diseases, 160, 163, 178

memory

 constraints of, 325

- cost of, 71
- measuring, 12, 76
- relating to input size, 72
- shared between programs, 76–77
- testing usage of, 77–81, 102–104
- `memory_profiler`, 78
- `merge`, 334
- meta level, 4, 17
- methods
 - baseline, 159–161
 - chaining, 166
- `metrics.accuracy_score`, 62
- `metrics.mean_squared_error`, 91
- `metrics.roc_curve`, 174, 179
- `metrics.SCORERS.keys()`, 161–162, 208
- micro, 168
- Minkowski distance, 63, 82, 367
- `MinMaxScaler`, 327
- mistakes, *see* errors
- Mitchell, Tom, 18
- Moby Dick*, 13
- mode value, 446
- models, 8, 66
 - additive, 318
 - bias of, 144–145
 - building, 14
 - combining, 390–398
 - comparing, 14
 - concrete, 371
 - evaluating, 14, 110
 - features working well with, 423–426, 464
 - fitting, 359–361, 363, 367, 370
 - fully defined, 371
 - keeping simple, 126, 295
 - not modifying the internal state of, 8, 361
 - performance of, 423
 - selecting, 113–114, 361–362
 - variability of, 144–145
 - workflow template for, 67, 90
- Monte Carlo, *see* randomness

- Monte Carlo cross-validation, *see* repeated train-test splitting (RTTS)
- Morse code, 417
- `most_frequent`, 160–161
- multiclass learners, 179–185, 195–201
 - averaging, 168–169
- mutual information, 418–423, 464
 - minimizing, 466
- `mutual_info_classif`, 419, 421–422
- `mutual_info_regression`, 420–421

N

- Naive Bayes (NB), 68–70, 292
 - bias-variance for, 148
 - evaluating, 70–71
 - in text classification, 69
 - performance of, 74–76, 78–81, 191
- natural language processing (NLP), 9
- nearest neighbors, *see* *k*-Nearest Neighbors
- Nearest Shrunken Centroids (NSC), 292
- `NearestCentroids`, 292
- negative outcome, 163–164
- nested cross-validation, 157, 370–377
- Netflix, 117
- neural networks, 512–516, 526
- newsgroups, 476
- Newton’s Method, 292
- No Free Lunch Theorem, 290
- noise, 15, 117
 - addressing, 350, 353–356
 - capturing, 122, 124, 126
 - distracting, 109–110, 296
 - eliminating, 144
 - manipulating, 117
- non-normality, 350
- nonic, 120
- nonlinearity, 285
- nonparametric learning methods, 65
- nonprimitive events, *see* compound events
- normal distribution, 27, 520–524

normal equations, 101
 normalization, 221, 322, 356, 474–476
Normalizer, 475
np.array.fromiter, 491–492, 494–495
np.cartesian_product, 41
 numbers
 binary vs. decimal, 53
 floating-point, 52–53
 numerical features, 6–7, 18, 225, 343–344, 346
 predicting, 10–11
NumPy, 20
 np.corrcoef, 416
 floating-point numbers in, 52–53
 np.array, 276, 494
 np.dot, 29–30, 38, 47–52
 np.fromiter, 494
 np.histogram, 21
 np.linalg.svd, 455
 np.polyfit, 119
 np.random.randint, 21
 np.searchsorted, 310
NuSVC, 253–257, 291
 Nystroem kernel, 436

O

Occam's razor, 124, 284
 odds
 betting, 259–262
 probability of, 262–266
 one-hot coding, 333–341, 347, 356, 526
 one-versus-all (OvA), 169
 one-versus-one (OvO), 181–182, 253
 one-versus-rest (OvR), 168, 179–182, 253, 267
OneHotEncoder, 333
 OpenCV library, 485
 optimization, 156, 497–500, 526
 premature, 83
 ordinal regression, 18
 outcome, outputs, *see* targets

overconfidence, 109–110
 and resampling, 128
 overfitting, 117, 122–126, 290, 296
 of base models, 396

P

pairplot, 86
pandas, 20
 pd.cut, 328, 330–331
 DataFrame, 323
 one-hot coding in, 333–334
 vs. **sklearn**, 323, 332
 parabolas, 45
 finding the best fit, 119–123
 piecewise, 313
 parameters, 115
 adjusting, 116
 choosing, 359
 in computer science vs. math, 318
 shuffling, 368
 tuning, 362
 vs. arguments, 362
 vs. explicit values, 360–361
 Pareto principle, 83
 partitions, 242
patsy, 334–340, 344–347
 connecting **sklearn** and, 347–348
 documentation for, 356
 PayPal, 189
PCA, 449–452
 peeking, 225
 penalization, *see* complexity
 penalties, 300, 306, 502
 percentile, 206
 performance, 102
 estimating, 382
 evaluating, 131, 150–152, 382
 measuring, 74–76, 78–81, 173, 178
 overestimating, 109
 physical laws, 17

piecewise constant regression, 309–313, 318

implementing, 310

preprocessing inputs in, 341

vs. k -NN-R, 310

PiecewiseConstantRegression, 313

Pipeline, 378–379

pipelines, 224–225, 377–382

integrating feature selection with, 426–428

plain linear model, 146, 147

planes, 39–41

finding the best, 410, 457

playing cards, 21

plots, 40, 41

plus-one trick, 38, 43–45, 336, 521

points in space, 34–43, 82

polyfit, 119

polynomial kernel, 253

polynomials

degree of, 119, 124

quadratic, 45

positive outcome, 163–164

precision, 165

macro, 168

tradeoffs between recall and, 168, 170–173, 185–187, 202

precision-recall curve (PRC), 185–187, 202

predict, 224–225, 379, 490–491

predict_proba, 174–175

predicted values, 10–11, 33

calculating, 97, 265

prediction bar, 170–177, 186

predictions, 165

combining, 389, 395, 405

evaluating, 215–217

flipping, 202

probability of, 170

real-world cost of, 56

predictive features, 7

predictive residuals, 219

predictors, *see* features

premature optimization, 83

presumption of innocence, 12

prime factorization, 452

primitive events, 22–23

principal components analysis (PCA), 445–462, 465–466

feature engineering in, 324

using dot products, 458–459, 461

prior, 160–161

probabilistic graphical models (PGMs), 516–525

and linear regression, 519–523

and logistic regression, 523–525

probabilistic principal components analysis (PPCA), 466

probabilities, 21–27

conditional, 24, 25

distribution of, 25–27, 290

expected value of, 31–32

of independent events, 23, 69

of primitive events, 22

of winning, 259–266

processing time, *see* time

programs

bottlenecks in, 83

memory usage of, 76–77

Provost, Foster, 18

purchasing behavior, predicting, 11

pydotplus, 245

pymc3, 519–521

Pythagorean theorem, 63

Python

indexing semantics in, 21, 54

list comprehension in, 136

memory management in, 77

using modules in the book, 20

Q

Quadratic Discriminant Analysis (QDA), 282–285

quadratic polynomials, *see* parabolas

quantile, 206

Quinlan, Ross, 239, 244

R

R^2 metric, 209–214

for mean model, 229

limitations of, 214, 233–234

misusing, 130

`randint`, 369

random events, 21–22

random forests (RFs), 396–398

comparing, 403

extreme, 397–398

selecting features in, 425

random guess strategy, 98–99, 101

random sampling, 325

random step strategy, 99, 101

`random.randint`, 21

`random_state`, 139–140

`RandomForestClassifier`, 425

`RandomizedSearchCV`, 369

randomness, 16

affecting data, 143

for feature selection, 423

for hyperparameters, 368–370

inherent in decisions, 241

pseudo-random, 139

to generate train-test splits, 133, 138–139

rare diseases, 160, 163, 178

`rbf`, 467

reality, 165

comparing to predictions, 215–217

recall, 165

tradeoffs between precision and, 168, 170–173, 185–187, 202

Receiver Operating Characteristic (ROC) curves, 172–181, 192, 202

and multiclass problem, 179–181

area under, 177–178, 182–193, 202

binary, 174–177

patterns in, 173–174

recentering, *see* data, centering

rectangles

areas of, 275

drawing, 275–278

overlapping, 243

recursive feature elimination, 425–426

redundancy, 324, 340

regression, 7, 64, 85–105

comparing methods of, 306–307

definition of, 85

examples of, 10–11

metrics for, 208–214

ordinal, 18

regression trees, 313–314

`RegressorMixin`, 311

regressors

baseline, 205–207

comparing, 314–317

default metric for, 209

evaluating, 205–234

implementing, 311–313

performance of, 317

scoring for, 130

regularization, 296–301

performing, 300–301

regularized linear regression, 296–301, 305

reinforcement learning, 18

repeated train-test splitting (RTTS), 133–139, 156

resampling, 128, 156, 390

with replacement, 157, 391–392

without replacement, 391

rescaling, *see* scaling, standardizing

`reshape`, 333

residual plots, 217–221, 232

residuals, 218, 230–232, 350

predictive, 219

Studentized, 232

resources

consumption of, 12–13, 71

- limited, 187
- measuring, 71–77
- needed by an algorithm, 72
- utilization in regression, 102–104
- RFE**, 425
- Ridge**, 300
- ridge regression (L_2), 300, 307
 - blending with lasso regression, 318
- rolling dice, 21–24
 - expected value of, 31–32
 - rigged, 68–69
- root mean squared error (RMSE), 101
 - calculating, 119
 - comparing regressors on, 315
 - high, 142
 - size of values in, 136
- rvs**, 369

S

- sampling, *see* resampling
- Samuel, Arthur, 3–4, 17
- scaling, 322, 325–329
 - statistical, 326
- scipy.stats**, 369
- scores, 127, 130
 - extracting from CV classifiers, 192
 - for each class, 181
 - vs. loss, 207
- scoring function, 184
- Seaborn, 20
 - pairplot**, 86
 - tspplot**, 151
- searchsorted**, 310
- SelectFromModel**, 424–425
- selection, 113–114
- SelectPercentile**, 422
- sensitivity, 173, 185
- SGDClassifier**, 267, 292
- shrinkage, *see* complexity
- shuffle**, 368
- ShuffleSplit**, 137–139
- shuffling, 137–140, 382
- SIFT_create**, 485
- signed area, 275
- Silva, Alice, 195, 203
- similarity, 63–64
- simple average, 30
- simplicity, 124
- singular value decomposition (SVD), 452, 465–466
- sklearn**, 19–20
 - 3D datasets in, 460–461
 - baseline models in, 205
 - boosters in, 400
 - classification metrics in, 161–163, 208–209
 - classifiers in, 202
 - common interface of, 379
 - confusion matrix in, 173
 - connecting **patsy** and, 347–348
 - consistency of, 225
 - cross-validation in, 129–130, 132, 184
 - custom models in, 311
 - distance calculators in, 64
 - documentation of, 368
 - feature correlation in, 416–417
 - feature evaluation in, 463
 - feature selection in, 425
 - kernels in, 435–437, 481
 - learners in, 318
 - linear regression in, 300, 310
 - logistic regression in, 267
 - naming conventions in, 207, 362
 - normalization in, 356
 - PCA in, 449–452
 - pipelines in, 224–225
 - plotting learning curves in, 157
 - R^2 in, 210–214, 233–234
 - random forests in, 396, 407
 - sparse-aware methods in, 356
 - storing data in, 333
 - SVC in, 253

sklearn (*continued*)

- SVR in, 307
- terminology of, 61, 66, 127, 160
- text representation in, 471–479, 494
- thresholds in, 176
- using alternative systems instead, 119
- using OvR, 253
- vs. **pandas**, 323, 332
- workflow in, 67, 90
- skms.cross_validate**, 226–227
- skpre.Normalizer**, 495
- Skynet, 389
- smart step strategy, 99–101, 267
- smoothness, 308, 406, *see also* complexity, regularization
- sns.pairplot**, 58
- softmax function, 526
- sorted lists, 465
- sparsity, 333, 356
- specificity, 165, 173, 185
- splines, 318
- spread, *see* standard deviation
- square root of the sum of squared errors, 93
- squared error loss, 301
- squared error points, 209
- ss.geom**, 369
- ss.normal**, 369
- ss.uniform**, 369
- StackExchange, 465
- stacking, 390
- StackOverflow, 292
- standard deviation, 54, 85, 221, 327
- standardization, 85, 105, 221–225, 231, 327
- StandardScaler**, 223–225, 326–327
- stationary learning tasks, 16
- statistics, 87
 - coefficient of determination, 130, 209
 - distribution of the mean, 391
 - dummy coding, 334
 - for feature selection, 463

- Studentized residuals, 232

- variation in data, 451

statsmodels, 292, 338–341

- documentation for, 356

- Stochastic Gradient Descent (SGD), 267

stocks

- choosing action for, 9

- predicting pricing for, 11

- stop words, 472–473, 494

storage space

- cost of, 12–13, 71

- measuring, 72

- stratification, 132–133

stratified, 160–161

- StratifiedKFold**, 130, 403

- strings, comparing, 438–439

- stripplots, 135, 155

- student performance, 195–201, 203, 225–226

- comparing regressors on, 314–317

- predicting, 10

- Studentized residuals, 232

- studying for a test, 109, 116–117

- sum, weighted, 28, 31

- sum of probabilities of events

- all primitive, 22

- independent, 23

- sum of squared errors (SSE), 33–34, 93–94, 210–212, 271, 301

- smallest, 100

- sum of squares, 32–33

- sum product, 30

- summary statistic, 87

- supervised learning from examples, 4, 9–11

- Support Vector Classifiers (SVCs), 252–259, 290–291, 301, 442

- bias-variance in, 256–259

- boundary in, 252

- computing, 291

- for nonlinear data, 285–287

- maximum margin separator in, 305

- parameters for, 254–256
- performance of, 429
- Support Vector Machines (SVMs), 252, 291, 442, 465
 - feature engineering in, 324
 - from raw materials, 510–511
 - vs. the polynomial kernel, 437
- Support Vector Regression (SVR), 301–307
 - main options for, 307
- support vectors, 252, 254
- supporting examples, 252
- SVC, 253–259, 291, 438
- synonyms, 482–483, 487–488

T

- T-distributed Stochastic Neighbor Embedding (TSNE), 462
- t-test, 463
- tabular data, 470
- targets, 6–7
 - cooperative values of, 296
 - discontinuous, 308
 - predicting, 397
 - training vs. testing, 60–61
 - transforming, 350, 353–356
- task understanding, 14
- tax brackets, 322, 331
- teaching to the test, 59–60, 114
 - in picking a learner, 112–113
 - protecting against, 110–111, 372, 377
- TensorFlow, 82
- term frequency-inverse document frequency (TF-IDF), 475–477, 495
- testing datasets, 60–61, 110, 114
 - predicting on, 66
 - resampling, 128
 - size of, 115, 130
- testing phase, *see* assessment, selection tests
 - positive vs. negative, 163–166

- specificity of, 165
- text, 470–479
 - classification of, 69
 - encoding, 471–476
 - representing as table rows, 470–471
- TfidfVectorizer, 475, 478, 495
- Theano, 82
- time
 - constraints of, 325
 - cost of, 13, 71
 - measuring, 12, 72, 74–75
 - relating to input size, 72
- time series, plotting, 151
- timeit, 74–75, 83
- todense, 333, 473
- Tolkien, J. R. R., 290
- total distance, 94
- tradeoffs, 13
 - between bias and variance, *see* bias-variance tradeoffs
 - between complexity and errors, 126
 - between false positives and negatives, 172
 - between precision and recall, 168, 170–173
- train-test splits, 60, 110, 115
 - evaluating, 70–71, 152
 - for cross-validation, 132
 - multiple, 128
 - randomly selected, 370
 - repeated, 133–139, 156
- train_test_split, 60, 70–71, 79, 349
- training datasets, 60–61, 110, 114
 - duplicating examples by weight in, 399
 - fitting estimators on, 66
 - randomly selected, 370
 - resampling, 128
 - size of, 115, 130–131, 150
 - unique identifiers in, 241, 322
- training error, 60
- training loss, 125–126, 296
- training phase, 113

transform, 224–225
Transformer, 435–436
TransformerMixin, 348, 379
 transformers, 348–350
 for images, 491–493
 treatment coding, *see* one-hot coding
 tree-building algorithms, 244
 trigrams, 471
 true negative rate (TNR), 164–166
 true positive rate (TPR), 164–166,
 173–181
 Trust Region Newton’s Method, 292
tsplot, 151
Twenty Newsgroups dataset, 476
 two-humped camel, *see* data, multimodal

U

unaccounted-for differences, 350
 underfitting, 117, 122–125, 296
uniform, 160–161
 unigrams, 471
 unique identifiers, 241, 322, 324
 univariate feature selection, 415
 unsupervised activities, 445

V

validation, 110, 156, *see also*
 cross-validation
 validation sets (ValS), 114
 randomly selected, 370
 size of, 115
 values
 accuracy of, 15
 actual, 33
 baseline, 356
 definition of, 5
 discrete, 5–6
 explicit, vs. function parameters,
 360–361
 finding the best, 98–101, 267

 missing, 18, 322
 numerical, 6–7, 18, 86, 225
 predicting, 64, 85, 87, 91
 predicted, 10–11, 33, 97, 265
 target, 6–7
 cooperative, 296
 transforming, 350
 under- vs. overestimating, 33
 variance, 110, 271, 292
 always positive, 272
 in feature values, 412–415
 in SVCs, 256–259
 maximizing, 448–449, 451
 not affected by data shifting, 451
 of data, 143, 145, 445
 of model, 144–145
 reducing, 396, 400, 406
VarianceThreshold, 413
 vectorizers, 495
 verification, 156
 vocabularies, 482
 global, 487
 votes, weighted, 390
VotingClassifier, 407

W

warp functions, 440
 weighted
 average, 31–32, 34, 89
 data, 399–400
 errors, 399
 mean, 89–90
 sum, 28, 31
 votes, 390
 weights
 adjusting, 497–500
 distributions of, 524
 pairs of, 524
 restricting, 105, 146
 total size of, 297

whuber (StackOverflow user), 292
wine dataset, 412–414, 426–428, 449
winning, odds of, 259–262
Wittgenstein, Ludwig, 18
words
 adjacent, 471
 counts of, 471, 473
 frequency of, 474–476
 in a document, 471
 stop, 472–473, 494
 visual, 491
 global, 483, 487–490
 local, 483–488
World War II, 172
wrapping functions, 361, 502

X

`xgboost`, 406
xor function, 341–343

Y

YouTube, 54, 109

Z

z-scoring, *see* standardizing
zip, 30