field cady THE DATA SCIENCE HANDBOOK

10

15 16 17 18 19 22



The Data Science Handbook

The Data Science Handbook

Field Cady



This edition first published 2017 © 2017 John Wiley & Sons, Inc.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at http://www.wiley.com/go/permissions.

The right of Field Cady to be identified as the author(s) of this work has been asserted in accordance with law.

Registered Offices

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA

Editorial Office 111 River Street, Hoboken, NJ 07030, USA

For details of our global editorial offices, customer services, and more information about Wiley products visit us at www.wiley.com.

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

Limit of Liability/Disclaimer of Warranty

MATLAB[®] is a trademark of The MathWorks, Inc. and is used with permission. The MathWorks does not warrant the accuracy of the text or exercises in this book. This work's use or discussion of MATLAB® software or related products does not constitute endorsement or sponsorship by The MathWorks of a particular pedagogical approach or particular use of the MATLAB® software. While the publisher and authors have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and authors endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read. Neither the publisher nor authors shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

Library of Congress Cataloguing-in-Publication Data

Names: Cady, Field, 1984- author. Title: The data science handbook / Field Cady. Description: Hoboken, NJ: John Wiley & Sons, Inc., 2017. | Includes bibliographical references and index. Identifiers: LCCN 2016043329 (print) | LCCN 2016046263 (ebook) | ISBN 9781119092940 (cloth) | ISBN 9781119092933 (pdf) | ISBN 9781119092926 (epub) Subjects: LCSH: Databases--Handbooks, manuals, etc. | Statistics--Data processing--Handbooks, manuals, etc. | Big data--Handbooks, manuals, etc. | Information theory--Handbooks, manuals, etc. Classification: LCC QA76.9.D32 C33 2017 (print) | LCC QA76.9.D32 (ebook) | DDC 005.74--dc23 LC record available at https://lccn.loc.gov/2016043329 Cover image: Cepreň Хакимуллин/Gettyimages Cover design by Wiley

Printed in the United States of America

Set in 10/12pt Warnock by SPi Global, Chennai, India

 $10\,9\,8\,7\,6\,5\,4\,3\,2\,1$

To my wife, Ryna. Thank you honey, for your support and for always believing in me.

Contents

Preface xvii

1 Introduction: Becoming a Unicorn 1

- 1.1 Aren't Data Scientists Just Overpaid Statisticians? 2
- 1.2 How Is This Book Organized? 3
- 1.3 How to Use This Book? 3
- 1.4 Why Is It All in PythonTM, Anyway? 4
- 1.5 Example Code and Datasets 4
- 1.6 Parting Words 5

Part I The Stuff You'll Always Use 7

2 The Data Science Road Map 9

- 2.1 Frame the Problem 10
- 2.2 Understand the Data: Basic Questions 11
- 2.3 Understand the Data: Data Wrangling 12
- 2.4 Understand the Data: Exploratory Analysis 13
- 2.5 Extract Features 14
- 2.6 Model 15
- 2.7 Present Results 15
- 2.8 Deploy Code 16
- 2.9 Iterating 16
- 2.10 Glossary 17

3 Programming Languages 19

- 3.1 Why Use a Programming Language? What Are the Other Options? *19*
- 3.2 A Survey of Programming Languages for Data Science 20
- 3.2.1 Python 20
- 3.2.2 R *21*
- 3.2.3 MATLAB[®] and Octave 21
- 3.2.4 SAS[®] 21

3.2.5	Scala [®] 22
3.3	Python Crash Course 22
3.3.1	A Note on Versions 22
3.3.2	"Hello World" Script 23
3.3.3	More Complicated Script 23
3.3.4	Atomic Data Types 26
3.4	Strings 27
3.4.1	Comments and Docstrings 28
3.4.2	Complex Data Types 29
3.4.3	Lists 29
3.4.4	Strings and Lists 30
3.4.5	Tuples 31
3.4.6	Dictionaries 31
3.4.7	Sets 32
3.5	Defining Functions 32
3.5.1	For Loops and Control Structures 33
3.5.2	A Few Key Functions 34
3.5.3	Exception Handling 35
3.5.4	Libraries 35
3.5.5	Classes and Objects 35
3.5.6	GOTCHA: Hashable and Unhashable Types 36
3.6	Python's Technical Libraries 37
3.6.1	Data Frames 38
3.6.2	Series 39
3.6.3	Joining and Grouping 40
3.7	Other Python Resources 42
3.8	Further Reading 42
3.9	Glossary 43
3a	Interlude: My Personal Toolkit 45
4	Data Munging: String Manipulation, Regular Expressions, and Data Cleaning 47
4.1	The Worst Dataset in the World 48
4.2	How to Identify Pathologies 48
4.3	Problems with Data Content 49
4.3.1	Duplicate Entries 49
4.3.2	Multiple Entries for a Single Entity 49
4.3.3	Missing Entries 49
4.3.4	NULLs 50
4.3.5	Huge Outliers 50
4.3.6	Out-of-Date Data 50
4.3.7	Artificial Entries 50

viii

- 4.3.8 Irregular Spacings 51
- 4.4 Formatting Issues 51
- 4.4.1 Formatting Is Irregular between Different Tables/Columns *51*
- 4.4.2 Extra Whitespace 51
- 4.4.3 Irregular Capitalization 52
- 4.4.4 Inconsistent Delimiters 52
- 4.4.5 Irregular NULL Format 52
- 4.4.6 Invalid Characters 52
- 4.4.7 Weird or Incompatible Datetimes 52
- 4.4.8 Operating System Incompatibilities 53
- 4.4.9 Wrong Software Versions 53
- 4.5 Example Formatting Script 54
- 4.6 Regular Expressions 55
- 4.6.1 Regular Expression Syntax 56
- 4.7 Life in the Trenches 60
- 4.8 Glossary 60

5 Visualizations and Simple Metrics 61

- 5.1 A Note on Python's Visualization Tools 62
- 5.2 Example Code 62
- 5.3 Pie Charts 63
- 5.4 Bar Charts 65
- 5.5 Histograms 66
- 5.6 Means, Standard Deviations, Medians, and Quantiles 69
- 5.7 Boxplots 70
- 5.8 Scatterplots 72
- 5.9 Scatterplots with Logarithmic Axes 74
- 5.10 Scatter Matrices 76
- 5.11 Heatmaps 77
- 5.12 Correlations 78
- 5.13 Anscombe's Quartet and the Limits of Numbers 80
- 5.14 Time Series 81
- 5.15 Further Reading 85
- 5.16 Glossary 85

6 Machine Learning Overview 87

- 6.1 Historical Context 88
- 6.2 Supervised versus Unsupervised 89
- 6.3 Training Data, Testing Data, and the Great Boogeyman of Overfitting *89*
- 6.4 Further Reading 91
- 6.5 Glossary 91

- **x** Contents
 - 7 Interlude: Feature Extraction Ideas 93
 - 7.1 Standard Features 93
 - 7.2 Features That Involve Grouping 94
 - 7.3 Preview of More Sophisticated Features 95
 - 7.4 Defining the Feature You Want to Predict 95

8 Machine Learning Classification 97

- 8.1 What Is a Classifier, and What Can You Do with It? 97
- 8.2 A Few Practical Concerns 98
- 8.3 Binary versus Multiclass 99
- 8.4 Example Script 99
- 8.5 Specific Classifiers 101
- 8.5.1 Decision Trees 101
- 8.5.2 Random Forests 103
- 8.5.3 Ensemble Classifiers 104
- 8.5.4 Support Vector Machines 105
- 8.5.5 Logistic Regression 108
- 8.5.6 Lasso Regression 110
- 8.5.7 Naive Bayes 110
- 8.5.8 Neural Nets 112
- 8.6 Evaluating Classifiers 114
- 8.6.1 Confusion Matrices 114
- 8.6.2 ROC Curves 115
- 8.6.3 Area under the ROC Curve *116*
- 8.7 Selecting Classification Cutoffs 117
- 8.7.1 Other Performance Metrics 118
- 8.7.2 Lift–Reach Curves 118
- 8.8 Further Reading 119
- 8.9 Glossary 119

9 Technical Communication and Documentation 121

- 9.1 Several Guiding Principles 122
- 9.1.1 Know Your Audience 122
- 9.1.2 Show Why It Matters 122
- 9.1.3 Make It Concrete 123
- 9.1.4 A Picture Is Worth a Thousand Words 123
- 9.1.5 Don't Be Arrogant about Your Tech Knowledge 124
- 9.1.6 Make It Look Decent 124
- 9.2 Slide Decks 124
- 9.2.1 C.R.A.P. Design 125
- 9.2.2 A Few Tips and Rules of Thumb 127
- 9.3 Written Reports 128
- 9.4 Speaking: What Has Worked for Me 130

- 9.5 Code Documentation 131
- 9.6 Further Reading *132*
- 9.7 Glossary 132

Part II Stuff You Still Need to Know 133

10 Unsupervised Learning: Clustering and Dimensionality Reduction 135

- 10.1 The Curse of Dimensionality 136
- 10.2 Example: Eigenfaces for Dimensionality Reduction 138
- 10.3 Principal Component Analysis and Factor Analysis 140
- 10.4 Skree Plots and Understanding Dimensionality 142
- 10.5 Factor Analysis 143
- 10.6 Limitations of PCA 143
- 10.7 Clustering 144
- 10.7.1 Real-World Assessment of Clusters 144
- 10.7.2 k-Means Clustering 145
- 10.7.3 Gaussian Mixture Models 146
- 10.7.4 Agglomerative Clustering 147
- 10.7.5 Evaluating Cluster Quality 148
- 10.7.6 Silhouette Score 148
- 10.7.7 Rand Index and Adjusted Rand Index 149
- 10.7.8 Mutual Information 150
- 10.8 Further Reading 151
- 10.9 Glossary 151

11 Regression 153

- 11.1 Example: Predicting Diabetes Progression 153
- 11.2 Least Squares 156
- 11.3 Fitting Nonlinear Curves 157
- 11.4 Goodness of Fit: R^2 and Correlation 159
- 11.5 Correlation of Residuals 160
- 11.6 Linear Regression 161
- 11.7 LASSO Regression and Feature Selection 162
- 11.8 Further Reading 164
- 11.9 Glossary 164

12 Data Encodings and File Formats 165

- 12.1 Typical File Format Categories 165
- 12.1.1 Text Files 166
- 12.1.2 Dense Numerical Arrays 166
- 12.1.3 Program-Specific Data Formats 166

xii	Contents

- 12.1.4 Compressed or Archived Data 166
- 12.2 CSV Files 167
- 12.3 JSON Files 168
- 12.4 XML Files 170
- 12.5 HTML Files 172
- 12.6 Tar Files 174
- 12.7 GZip Files 175
- 12.8 Zip Files 175
- 12.9 Image Files: Rasterized, Vectorized, and/or Compressed 176
- 12.10 It's All Bytes at the End of the Day 177
- 12.11 Integers 178
- 12.12 Floats 179
- 12.13 Text Data 180
- 12.14 Further Reading 183
- 12.15 Glossary 183
- **13 Big Data** 185
- 13.1 What Is Big Data? *185*
- 13.2 Hadoop: The File System and the Processor 187
- 13.3 Using HDFS 188
- 13.4 Example PySpark Script 189
- 13.5 Spark Overview 190
- 13.6 Spark Operations 192
- 13.7 Two Ways to Run PySpark 193
- 13.8 Configuring Spark 194
- 13.9 Under the Hood 195
- 13.10 Spark Tips and Gotchas 196
- 13.11 The MapReduce Paradigm 197
- 13.12 Performance Considerations 199
- 13.13 Further Reading 200
- 13.14 Glossary 200
- 14 Databases 203
- 14.1 Relational Databases and MySQL[®] 204
- 14.1.1 Basic Queries and Grouping 204
- 14.1.2 Joins 207
- 14.1.3 Nesting Queries 208
- 14.1.4 Running MySQL and Managing the DB 209
- 14.2 Key-Value Stores 210
- 14.3 Wide Column Stores 211
- 14.4 Document Stores 211
- 14.4.1 MongoDB[®] 212
- 14.5 Further Reading 214
- 14.6 Glossary 214

- 15 Software Engineering Best Practices 217
- 15.1 Coding Style 217
- 15.2 Version Control and Git for Data Scientists 220
- 15.3 Testing Code 222
- 15.3.1 Unit Tests 223
- 15.3.2 Integration Tests 224
- 15.4 Test-Driven Development 225
- 15.5 AGILE Methodology 225
- 15.6 Further Reading 226
- 15.7 Glossary 226

16 Natural Language Processing 229

- 16.1 Do I Even Need NLP? 229
- 16.2 The Great Divide: Language versus Statistics 230
- 16.3 Example: Sentiment Analysis on Stock Market Articles 230
- 16.4 Software and Datasets 232
- 16.5 Tokenization 233
- 16.6 Central Concept: Bag-of-Words 233
- 16.7 Word Weighting: TF-IDF 235
- 16.8 *n*-Grams 235
- 16.9 Stop Words 236
- 16.10 Lemmatization and Stemming 236
- 16.11 Synonyms 237
- 16.12 Part of Speech Tagging 237
- 16.13 Common Problems 238
- 16.13.1 Search 238
- 16.13.2 Sentiment Analysis 239
- 16.13.3 Entity Recognition and Topic Modeling 240
- 16.14 Advanced NLP: Syntax Trees, Knowledge, and Understanding 240
- 16.15 Further Reading 241
- 16.16 Glossary 242

17 Time Series Analysis 243

- 17.1 Example: Predicting Wikipedia Page Views 244
- 17.2 A Typical Workflow 247
- 17.3 Time Series versus Time-Stamped Events 248
- 17.4 Resampling an Interpolation 249
- 17.5 Smoothing Signals 251
- 17.6 Logarithms and Other Transformations 252
- 17.7 Trends and Periodicity 252
- 17.8 Windowing 253
- 17.9 Brainstorming Simple Features 254
- 17.10 Better Features: Time Series as Vectors 255

xiv Contents

- Fourier Analysis: Sometimes a Magic Bullet 256 17.11
- 17.12 Time Series in Context: The Whole Suite of Features 259
- 17.13 Further Reading 259
- 17.14 Glossary 260

18 **Probability** 261

- 18.1 Flipping Coins: Bernoulli Random Variables 261
- 18.2 Throwing Darts: Uniform Random Variables 263
- 18.3 The Uniform Distribution and Pseudorandom Numbers 263
- 18.4 Nondiscrete, Noncontinuous Random Variables 265
- 18.5 Notation, Expectations, and Standard Deviation 267
- 18.6 Dependence, Marginal and Conditional Probability 268
- 18.7 Understanding the Tails 269
- **Binomial Distribution** 18.8 271
- 18.9 Poisson Distribution 272
- 272 18.10 Normal Distribution
- 18.11 Multivariate Gaussian 273
- 18.12 Exponential Distribution 274
- 276 18.13 Log-Normal Distribution
- 18.14 Entropy 277
- Further Reading 18.15 279
- 18.16 Glossary 279

19 Statistics 281

- Statistics in Perspective 19.1 281
- 19.2 Bayesian versus Frequentist: Practical Tradeoffs and Differing Philosophies 282
- 19.3 Hypothesis Testing: Key Idea and Example 283
- 19.4 Multiple Hypothesis Testing 285
- 19.5 Parameter Estimation 286
- 19.6 Hypothesis Testing: t-Test 287
- 19.7 Confidence Intervals 290
- 19.8 **Bayesian Statistics** 291
- 19.9 Naive Bayesian Statistics 293
- 19.10 Bayesian Networks 293
- 19.11 Choosing Priors: Maximum Entropy or Domain Knowledge 294
- 19.12 Further Reading 295
- 19.13 Glossary 295

20 Programming Language Concepts 297

- 20.1**Programming Paradigms** 297
- 20.1.1 Imperative 298
- 20.1.2 Functional 298

- 20.1.3 Object-Oriented 301
- 20.2 Compilation and Interpretation *305*
- 20.3 Type Systems 307
- 20.3.1 Static versus Dynamic Typing 308
- 20.3.2 Strong versus Weak Typing 308
- 20.4 Further Reading 309
- 20.5 Glossary 309

21 Performance and Computer Memory 311

- 21.1 Example Script 311
- 21.2 Algorithm Performance and Big-O Notation 314
- 21.3 Some Classic Problems: Sorting a List and Binary Search 315
- 21.4 Amortized Performance and Average Performance 318
- 21.5 Two Principles: Reducing Overhead and Managing Memory 320
- 21.6 Performance Tip: Use Numerical Libraries When Applicable 322
- 21.7 Performance Tip: Delete Large Structures You Don't Need 323
- 21.8 Performance Tip: Use Built-In Functions When Possible 324
- 21.9 Performance Tip: Avoid Superfluous Function Calls 324
- 21.10 Performance Tip: Avoid Creating Large New Objects 325
- 21.11 Further Reading 325
- 21.12 Glossary 325

Part III Specialized or Advanced Topics 327

22 Computer Memory and Data Structures 329

- 22.1 Virtual Memory, the Stack, and the Heap 329
- 22.2 Example C Program 330
- 22.3 Data Types and Arrays in Memory 330
- 22.4 Structs 332
- 22.5 Pointers, the Stack, and the Heap 333
- 22.6 Key Data Structures 337
- 22.6.1 Strings 337
- 22.6.2 Adjustable-Size Arrays 338
- 22.6.3 Hash Tables 339
- 22.6.4 Linked Lists 340
- 22.6.5 Binary Search Trees 342
- 22.7 Further Reading 343
- 22.8 Glossary 343

23 Maximum Likelihood Estimation and Optimization 345

- 23.1 Maximum Likelihood Estimation 345
- 23.2 A Simple Example: Fitting a Line *346*

xvi Contents

- 23.3 Another Example: Logistic Regression 348
- 23.4 Optimization 348
- 23.5 Gradient Descent and Convex Optimization 350
- 23.6 Convex Optimization 353
- 23.7 Stochastic Gradient Descent 355
- 23.8 Further Reading 355
- 23.9 Glossary 356

24 Advanced Classifiers 357

- 24.1 A Note on Libraries 358
- 24.2 Basic Deep Learning 358
- 24.3 Convolutional Neural Networks 361
- 24.4 Different Types of Layers. What the Heck Is a Tensor? 362
- 24.5 Example: The MNIST Handwriting Dataset 363
- 24.6 Recurrent Neural Networks 366
- 24.7 Bayesian Networks 367
- 24.8 Training and Prediction 369
- 24.9 Markov Chain Monte Carlo 369
- 24.10 PyMC Example 370
- 24.11 Further Reading 373
- 24.12 Glossary 373

25 Stochastic Modeling 375

- 25.1 Markov Chains 375
- 25.2 Two Kinds of Markov Chain, Two Kinds of Questions 377
- 25.3 Markov Chain Monte Carlo 379
- 25.4 Hidden Markov Models and the Viterbi Algorithm 380
- 25.5 The Viterbi Algorithm 382
- 25.6 Random Walks 384
- 25.7 Brownian Motion 384
- 25.8 ARIMA Models 385
- 25.9 Continuous-Time Markov Processes 386
- 25.10 Poisson Processes 387
- 25.11 Further Reading 388
- 25.12 Glossary 388

25a Parting Words: Your Future as a Data Scientist 391

Index 393

Preface

This book was written to solve a problem. The people who I interview for data science jobs have sterling mathematical pedigrees, but most of them are unable to write a simple script that computes Fibonacci numbers (in case you aren't familiar with Fibonacci numbers, this takes about five lines of code). On the other side, employers tend to view data scientists as either mysterious wizards or used-car salesmen (and when data scientists can't be trusted to write a basic script, the latter impression has some merit!). These problems reflect a fundamental misunderstanding, by all parties, of what data science is (and isn't) and what skills its practitioners need.

When I first got into data science, I was part of that problem. Years of doing academic physics had trained me to solve problems in a way that was long on abstract theory but short on common sense or flexibility. Mercifully, I also knew how to code (thanks, GoogleTM internships!), and this let me limp along while I picked up the skills and mindsets that actually mattered.

Since leaving academia, I have done data science consulting for companies of every stripe. This includes web traffic analysis for tiny start-ups, manufacturing optimizations for Fortune 100 giants, and everything in between. The problems to solve are always unique, but the skills required to solve them are strikingly universal. They are an eclectic mix of computer programming, mathematics, and business savvy. They are rarely found together in one person, but in truth they can be learned by anybody.

A few interviews I have given stand out in my mind. The candidate was smart and knowledgeable, but the interview made it painfully clear that they were unprepared for the daily work of a data scientist. What do you do as an interviewer when the candidate starts apologizing for wasting your time? We ended up filling the hour with a crash course on what they were missing and how they could go out and fill the gaps in their knowledge. They went out, learned what they needed to, and are now successful data scientists.

I wrote this book in an attempt to help people like that out, by condensing data science's various skill sets into a single, coherent volume. It is hands-on

xviii Preface

and to the point: ideal for somebody who needs to come up to speed quickly or solve a problem on a tight deadline. The educational system has not yet caught up to the demands of this new and exciting field, and my hope is that this book will help you bridge the gap.

> Field Cady September 2016 Redmond, Washington

Introduction: Becoming a Unicorn

"Data science" is a very popular term these days, and it gets applied to so many things that its meaning has become very vague. So I'd like to start this book by giving you the definition that I use. I've found that this one gets right to the heart of what sets it apart from other disciplines. Here goes:

Data science means doing analytics work that, for one reason or another, requires a substantial amount of software engineering skills.

Sometimes, the final deliverable is the kind of thing a statistician or business analyst might provide, but achieving that goal demands software skills that your typical analyst simply doesn't have. For example, a dataset might be so large that you need to use distributed computing to analyze it or so convoluted in its format that many lines of code are required to parse it. In many cases, data scientists also have to write big chunks of production software that implement their analytics ideas in real time. In practice, there are usually other differences as well. For example, data scientists usually have to extract features from raw data, which means that they tackle very open-ended problems such as how to quantify the "spamminess" of an e-mail.

It's very hard to find people who can construct good statistical models, hack quality software, and relate this all in a meaningful way to business problems. It's a lot of hats to wear! These individuals are so rare that recruiters often call them "unicorns."

The message of this book is that it is not only possible but also relatively straightforward to become a "unicorn." It's just a question of acquiring the particular balance of skills required. Very few educational programs teach all of those skills, which is why unicorns are rare, but that's mostly a historical accident. It is perfectly reasonable for a single person to have the whole palette of abilities, provided they're willing to ignore the traditional boundaries between different disciplines.

This book aims to teach you everything you'll need to know to be a competent data scientist. My guess is that you're either a computer programmer

The Data Science Handbook, First Edition. Field Cady. © 2017 John Wiley & Sons, Inc. Published 2017 by John Wiley & Sons, Inc.

1

2 1 Introduction: Becoming a Unicorn

looking to learn about analytics or more of a mathematician trying to bone up on their coding. You might also be a businessperson who needs the technical skills to answer your business questions or simply an interested layman. Whoever you are though, this book will teach you the concepts you need.

This book is not comprehensive. Data science is too big an area for any person or book to cover all of it. Besides, the field is changing so fast that any "comprehensive" book would be out-of-date before it came off the presses. Instead, I have aimed for two goals. First, I want to give a solid grounding in the big picture of what data science is, how to go about doing it, and the foundational concepts that will stand the test of time. Second, I want to give a "complete" skill set, in the sense that you have the nuts-and-bolts knowledge to go out and do data science work (you can code in Python, you know the libraries to use, most of the big machine learning models, etc.), even if particular projects or companies might require that you pick up a new skill set from somewhere else.

1.1 Aren't Data Scientists Just Overpaid Statisticians?

Nate Silver, a statistician famous for accurate forecasting of US elections, once famously said: "I think data scientist is a sexed-up term for statistician." He has a point, but what he said is only partly true. The discipline of statistics deals mostly with rigorous mathematical methods for solving well-defined problems. Data scientists spend most of their time getting data into a form where statistical methods could even be applied. This involves making sure that the analytics problem is a good match to business objectives, extracting meaningful features from the raw data and coping with any pathologies of the data or weird edge cases. Once that heavy lifting is done, you can apply statistical tools to get the final results, although, in practice, you often don't even need them. Professional statisticians need to do a certain amount of preprocessing themselves, but there is a massive difference in degree.

Historically, data science emerged as a field independently from statistics. Most of the first data scientists were computer programmers or machine learning experts who were working on Big Data problems. They were analyzing datasets of the kind that statisticians don't touch: HTML pages, image files, e-mails, raw output logs of web servers, and so on. These datasets don't fit the mold of relational databases or statistical tools, so for decades, they were just piling up without being analyzed. Data science came into being as a way to finally milk them for insights.

In 20 years, I suspect that statistics, data science, and machine learning will blur into a single discipline. The differences between them are, after all, really just a matter of degree and/or historical accident. But in practical terms, for the time being, solving data science problems requires skills that a normal statistician does not have. In fact, these skills, which include extensive software engineering and domain-specific feature extraction, constitute the overwhelming majority of the work that needs to be done. In the daily work of a data scientist, statistics plays second fiddle.

1.2 How Is This Book Organized?

This book is organized into three sections. The first, The Stuff You'll Always Use, covers topics that, in my experience, you will end up using in almost any data science project. They are core skills, which are absolutely indispensable for data science at any level.

The first section was also written with an eye toward people who need data science to answer a specific question but do not aspire to become full-fledged data scientists. If you are in this camp, then there is a good chance that Part I of the book will give you everything you need.

The second section, Stuff You Still Need to Know, covers additional core skills for a data scientist. Some of these, such as clustering, are so common that they almost made it into the first section, and they could easily play a role in any project. Others, such as natural language processing, are somewhat specialized subjects that are critical in certain domains but superfluous in others. In my judgment, a data scientist should be conversant in all of these subjects, even if they don't always use them all.

The final section, Stuff That's Good to Know, covers a variety of topics that are optional. Some of these chapters are just expansions on topics from the first two sections, but they give more theoretical background and discuss some additional topics. Others are entirely new material, which does come up in data science, but which you could go through a career without ever running into.

1.3 How to Use This Book?

This book was written with three use cases in mind:

- 1) You can read it cover-to-cover. If you do that, it should give you a self-contained course in data science that will leave you ready to tackle real problems. If you have a strong background in computer programming, or in mathematics, then some of it will be review.
- 2) You can use it to come quickly up to speed on a specific subject. I have tried to make the different chapters pretty self-contained, especially the chapters after the first section.

1 Introduction: Becoming a Unicorn

3) The book contains a lot of sample codes, in pieces that are large enough to use as a starting point for your own projects.

1.4 Why Is It All in Python[™], Anyway?

The example code in this book is all in Python, except for a few domain-specific languages such as SQL. My goal isn't to push you to use Python; there are lots of good tools out there, and you can use whichever ones you want.

However, I wanted to use one language for all of my examples. This keeps the book readable, and it also lets readers follow the whole book while only knowing one language. Of the various languages available, there are two reasons why I chose Python:

- 1) Python is the most popular language for data scientists. R is its only major competitor, at least when it comes to free tools. I have used both extensively, and I think that Python is flat-out better (except for some obscure statistics packages that have been written in R and that are rarely needed anyway).
- 2) I like to say that for any task, Python is the second-best language. It's a jackof-all-trades. If you only need to worry about statistics, or numerical computation, or web parsing, then there are better options out there. But if you need to do all of these things within a single project, then Python is your best option. Since data science is so inherently multidisciplinary, this makes it a perfect fit.

As a note of advice, it is much better to be proficient in one language, to the point where you can reliably churn out code that is of high quality, than to be mediocre at several.

1.5 Example Code and Datasets

This book is rich in example code, in fairly long chunks. This was done for two reasons:

- 1) As a data scientist, you need to be able to read longish pieces of code. This is a nonoptional skill, and if you aren't used to it, then this will give you a chance to practice.
- 2) I wanted to make it easier for you to poach the code from this book, if you feel so inclined.

You can do whatever you want with the code, with or without attribution. I release it into the public domain in the hope that it can give some people a small leg up. You can find it on my GitHub page at www.github.com/field-cady.

The sample data that I used comes in two forms:

- 1) Test datasets that are built into Python's scientific libraries
- 2) Data that is pulled off the Internet, from sources such as Yahoo and Wikipedia. When I do this, the example scripts will include code that pulls the data.

1.6 Parting Words

It is my hope that this book not only teaches you how to do nut-and-bolts data science but also gives you a feel of how exciting this deeply interdisciplinary subject is. Please feel free to reach out to me at www.fieldcady.com or field. cady@gmail.com with comments, errata, or any other feedback.

Part 1

The Stuff You'll Always Use

The first section of this book covers core topics that everybody doing data science should know. This includes people who are not interested in being professional data scientists, but need to know just enough to solve some specific problem. These are the subjects that will likely arise in every data science project you do.

The Data Science Road Map

In this chapter, I will give you a high-level overview of the process of data science. I will focus on the different stages of data science work, including common pain points, key things to get right, and where data science parts ways from other disciplines.

The process of solving a data science problem is summarized in the following figure, which I called the Data Science Road Map.



The first step is always to frame the problem: understand the business use case and craft a well-defined analytics problem (or problems) out of it. This is followed by an extensive stage of grappling with the data and the real-world things that it describes, so that we can extract meaningful features. Finally, these features are plugged into analytical tools that give us hard numerical results.

Before I go into more detail about the different stages of the roadmap, I want to point out two things.

The first is that "Model and Analyze" loops back to framing the problem. This is one of the key features of data science that differentiate it from traditional software engineering. Data scientists write code, and they use many of the same tools as software engineers. However, there is a tight feedback loop between data science work and the real world. Questions are always being reframed as new insights become available, and, as a result, data scientists

The Data Science Handbook, First Edition. Field Cady. © 2017 John Wiley & Sons, Inc. Published 2017 by John Wiley & Sons, Inc.

10 *2* The Data Science Road Map

must keep their code base extremely flexible and always have an eye toward the real-world problem they are solving. Ideally, you will follow the loop back many times, constantly refining your methods and producing new insights.

The second point is that there are two different (although not mutually exclusive) ways to exit the road map: presenting results and deploying code. My friend Michael Li, a data scientist who founded The Data Incubator, likened this to having two different types of clients: humans and machines. They require distinct skill sets and modifications to every stage of the data science road map.

If your clients are humans, then usually you are trying to use available data sources to answer some kind of business problem. Examples would be the following:

- Identifying leading indicators of spikes in the price of a stock, so that people can understand what causes price spikes
- Determining whether customers break down into natural subtypes and what characteristics each type has
- Assessing whether traffic to one website can be used to predict traffic to another site.

Typically, the final deliverable for work such as this will be a PowerPoint slide deck or a written report. The goal is to give business insights, and often these insights will be used for making key decisions. This kind of data science also functions as a way to test the waters and see whether some analytics approach is worth a larger follow-up project that may result in production software.

If your clients are machines, then you are doing something that blends into software engineering, where the deliverable is a piece of software that performs some analytics work. Examples would be the following:

- Implementing the algorithm that chooses which ad to show to a customer and training it on real data
- Writing a batch process that generates daily reports based on company records generated that day, using some kind of analytics to point out salient patterns

In these cases, your main deliverable is a piece of software. In addition to performing a useful task, it had better work well in terms of performance, robustness to bad inputs, and so on.

Once you understand who your clients are, the next step is to determine what you'll be doing for them. In the next section, I will show you how to do this all-important step.

2.1 Frame the Problem

The difference between great and mediocre data science is not about math or engineering: it is about asking the right question(s). Alternately, if you're trying