

Python语言构建 机器学习系统 第2版(影印版)

Building Machine Learning Systems with Python Second Edition

Luis Pedro Coelho, Willi Richert 著

[PACKT]

SO 东南大学出版社 SOUTHEAST UNIVERSITY PRESS

Python 语言构建机器学习系统 第2版 (影印版)

Luis Pedro Coelho, Willi Richert 著

南京 东南大学出版社

图书在版编目(CIP)数据

Python 语言构建机器学习系统:第2版:英文/ (美)科埃略(Coelho, L.P.),(美)里克特(Richert, W.) 著.一影印本.一南京:东南大学出版社,2016.1

书名原文: Building Machine Learning Systems with Python, Second Edition

ISBN 978 - 7 - 5641 - 6062 - 3

I.①P··· □.①科··· ②里··· □.①软件工具─程序设计─英文 Ⅳ.①TP311.56

中国版本图书馆 CIP 数据核字(2015)第 241940号

© 2015 by PACKT Publishing Ltd

Reprint of the English Edition, jointly published by PACKT Publishing Ltd and Southeast University Press, 2016. Authorized reprint of the original English edition, 2015 PACKT Publishing Ltd, the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 PACKT Publishing Ltd 出版 2015。

英文影印版由东南大学出版社出版 2016。此影印版的出版和销售得到出版权和销售权的所有者——PACKT Publishing Ltd 的许可。

版权所有,未得书面许可,本书的任何部分和全部不得以任何形式重制。

Python 语言构建机器学习系统 第 2 版(影印版)

出版发行: 东南大学出版社

地 址:南京四牌楼 2号 邮编:210096

出版人: 江建中

网 址: http://www.seupress.com

电子邮件: press@seupress.com

印刷:常州市武进第三印刷有限公司

开 本: 787毫米×980毫米 16开本

印 张: 20.25

字 数: 397千字

版 次: 2016年1月第1版

印 次: 2016年1月第1次印刷

书 号: ISBN 978-7-5641-6062-3

定 价: 68.00元

Credits

Authors

Luis Pedro Coelho

Willi Richert

Reviewers

Matthieu Brucher

Maurice HT Ling

Radim Řehůřek

Commissioning Editor

Kartikey Pandey

Acquisition Editors

Greg Wild

Richard Harvey

Kartikey Pandey

Content Development Editor

Arun Nadar

Technical Editor

Pankaj Kadam

Copy Editors

Relin Hedly

Sameen Siddiqui

Laxmi Subramanian

Project Coordinator

Nikhil Nair

Proofreaders

Simran Bhogal

Lawrence A. Herman

Linda Morris

Paul Hindle

Indexer

Hemangini Bari

Graphics

Sheetal Aute

Abhinash Sahu

Production Coordinator

Arvindkumar Gupta

Cover Work

Arvindkumar Gupta

About the Authors

Luis Pedro Coelho is a computational biologist: someone who uses computers as a tool to understand biological systems. In particular, Luis analyzes DNA from microbial communities to characterize their behavior. Luis has also worked extensively in bioimage informatics—the application of machine learning techniques for the analysis of images of biological specimens. His main focus is on the processing and integration of large-scale datasets.

Luis has a PhD from Carnegie Mellon University, one of the leading universities in the world in the area of machine learning. He is the author of several scientific publications.

Luis started developing open source software in 1998 as a way to apply real code to what he was learning in his computer science courses at the Technical University of Lisbon. In 2004, he started developing in Python and has contributed to several open source libraries in this language. He is the lead developer on the popular computer vision package for Python and mahotas, as well as the contributor of several machine learning codes.

Luis currently divides his time between Luxembourg and Heidelberg.

I thank my wife, Rita, for all her love and support and my daughter, Anna, for being the best thing ever.

Willi Richert has a PhD in machine learning/robotics, where he used reinforcement learning, hidden Markov models, and Bayesian networks to let heterogeneous robots learn by imitation. Currently, he works for Microsoft in the Core Relevance Team of Bing, where he is involved in a variety of ML areas such as active learning, statistical machine translation, and growing decision trees.

This book would not have been possible without the support of my wife, Natalie, and my sons, Linus and Moritz. I am especially grateful for the many fruitful discussions with my current or previous managers, Andreas Bode, Clemens Marschner, Hongyan Zhou, and Eric Crestan, as well as my colleagues and friends, Tomasz Marciniak, Cristian Eigel, Oliver Niehoerster, and Philipp Adelt. The interesting ideas are most likely from them; the bugs belong to me.

About the Reviewers

Matthieu Brucher holds an engineering degree from the Ecole Supérieure d'Electricité (Information, Signals, Measures), France and has a PhD in unsupervised manifold learning from the Université de Strasbourg, France. He currently holds an HPC software developer position in an oil company and is working on the next generation reservoir simulation.

Maurice HT Ling has been programming in Python since 2003. Having completed his PhD in Bioinformatics and BSc (Hons.) in Molecular and Cell Biology from The University of Melbourne, he is currently a Research Fellow at Nanyang Technological University, Singapore, and an Honorary Fellow at The University of Melbourne, Australia. Maurice is the Chief Editor for Computational and Mathematical Biology, and co-editor for The Python Papers. Recently, Maurice cofounded the first synthetic biology start-up in Singapore, AdvanceSyn Pte. Ltd., as the Director and Chief Technology Officer. His research interests lies in life — biological life, artificial life, and artificial intelligence — using computer science and statistics as tools to understand life and its numerous aspects. In his free time, Maurice likes to read, enjoy a cup of coffee, write his personal journal, or philosophize on various aspects of life. His website and LinkedIn profile are http://maurice.vodien.com and http://www.linkedin.com/in/mauriceling, respectively.

Radim Řehůřek is a tech geek and developer at heart. He founded and led the research department at Seznam.cz, a major search engine company in central Europe. After finishing his PhD, he decided to move on and spread the machine learning love, starting his own privately owned R&D company, RaRe Consulting Ltd. RaRe specializes in made-to-measure data mining solutions, delivering cutting-edge systems for clients ranging from large multinationals to nascent start-ups.

Radim is also the author of a number of popular open source projects, including gensim and smart_open.

A big fan of experiencing different cultures, Radim has lived around the globe with his wife for the past decade, with his next steps leading to South Korea. No matter where he stays, Radim and his team always try to evangelize data-driven solutions and help companies worldwide make the most of their machine learning opportunities.

www.PacktPub.com

Support files, eBooks, discount offers, and more

For support files and downloads related to your book, please visit www.PacktPub.com.

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.PacktPub.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at service@packtpub.com for more details.

At www.PacktPub.com, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



https://www2.packtpub.com/books/subscription/packtlib

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can search, access, and read Packt's entire library of books.

Why subscribe?

- · Fully searchable across every book published by Packt
- Copy and paste, print, and bookmark content
- · On demand and accessible via a web browser

Free access for Packt account holders

If you have an account with Packt at www.PacktPub.com, you can use this to access PacktLib today and view 9 entirely free books. Simply use your login credentials for immediate access.

Preface

One could argue that it is a fortunate coincidence that you are holding this book in your hands (or have it on your eBook reader). After all, there are millions of books printed every year, which are read by millions of readers. And then there is this book read by you. One could also argue that a couple of machine learning algorithms played their role in leading you to this book—or this book to you. And we, the authors, are happy that you want to understand more about the hows and whys.

Most of the book will cover the *how*. How has data to be processed so that machine learning algorithms can make the most out of it? How should one choose the right algorithm for a problem at hand?

Occasionally, we will also cover the *why*. Why is it important to measure correctly? Why does one algorithm outperform another one in a given scenario?

We know that there is much more to learn to be an expert in the field. After all, we only covered some *hows* and just a tiny fraction of the *whys*. But in the end, we hope that this mixture will help you to get up and running as quickly as possible.

What this book covers

Chapter 1, Getting Started with Python Machine Learning, introduces the basic idea of machine learning with a very simple example. Despite its simplicity, it will challenge us with the risk of overfitting.

Chapter 2, Classifying with Real-world Examples, uses real data to learn about classification, whereby we train a computer to be able to distinguish different classes of flowers.

Chapter 3, Clustering – Finding Related Posts, teaches how powerful the bag of words approach is, when we apply it to finding similar posts without really "understanding" them.

Chapter 4, Topic Modeling, moves beyond assigning each post to a single cluster and assigns them to several topics as a real text can deal with multiple topics.

Chapter 5, Classification – Detecting Poor Answers, teaches how to use the bias-variance trade-off to debug machine learning models though this chapter is mainly on using a logistic regression to find whether a user's answer to a question is good or bad.

Chapter 6, Classification II – Sentiment Analysis, explains how Naïve Bayes works, and how to use it to classify tweets to see whether they are positive or negative.

Chapter 7, Regression, explains how to use the classical topic, regression, in handling data, which is still relevant today. You will also learn about advanced regression techniques such as the Lasso and ElasticNets.

Chapter 8, Recommendations, builds recommendation systems based on costumer product ratings. We will also see how to build recommendations just from shopping data without the need for ratings data (which users do not always provide).

Chapter 9, Classification – Music Genre Classification, makes us pretend that someone has scrambled our huge music collection, and our only hope to create order is to let a machine learner classify our songs. It will turn out that it is sometimes better to trust someone else's expertise than creating features ourselves.

Chapter 10, Computer Vision, teaches how to apply classification in the specific context of handling images by extracting features from data. We will also see how these methods can be adapted to find similar images in a collection.

Chapter 11, Dimensionality Reduction, teaches us what other methods exist that can help us in downsizing data so that it is chewable by our machine learning algorithms.

Chapter 12, Bigger Data, explores some approaches to deal with larger data by taking advantage of multiple cores or computing clusters. We also have an introduction to using cloud computing (using Amazon Web Services as our cloud provider).

Appendix, Where to Learn More Machine Learning, lists many wonderful resources available to learn more about machine learning.

What you need for this book

This book assumes you know Python and how to install a library using easy_install or pip. We do not rely on any advanced mathematics such as calculus or matrix algebra.

We are using the following versions throughout the book, but you should be fine with any more recent ones:

- Python 2.7 (all the code is compatible with version 3.3 and 3.4 as well)
- NumPy 1.8.1
- SciPy 0.13
- scikit-learn 0.14.0

Who this book is for

This book is for Python programmers who want to learn how to perform machine learning using open source libraries. We will walk through the basic modes of machine learning based on realistic examples.

This book is also for machine learners who want to start using Python to build their systems. Python is a flexible language for rapid prototyping, while the underlying algorithms are all written in optimized C or C++. Thus the resulting code is fast and robust enough to be used in production as well.

Conventions

In this book, you will find a number of styles of text that distinguish between different kinds of information. Here are some examples of these styles, and an explanation of their meaning.

Code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles are shown as follows: "We then use poly1d() to create a model function from the model parameters."

A block of code is set as follows:

```
[aws info]

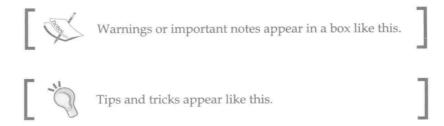
AWS_ACCESS_KEY_ID = AAKIIT7HHF6IUSN3OCAA

AWS_SECRET_ACCESS_KEY = <your secret key>
```

Any command-line input or output is written as follows:

```
>>> import numpy
>>> numpy.version.full_version
1.8.1
```

New terms and important words are shown in bold. Words that you see on the screen, in menus or dialog boxes for example, appear in the text like this: "Once the machine is stopped, the Change instance type option becomes available."



Reader feedback

Feedback from our readers is always welcome. Let us know what you think about this book—what you liked or may have disliked. Reader feedback is important for us to develop titles that you really get the most out of.

To send us general feedback, simply send an e-mail to feedback@packtpub.com, and mention the book title via the subject of your message. If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, see our author guide on www.packtpub.com/authors.

Customer support

Now that you are the proud owner of a Packt book, we have a number of things to help you to get the most from your purchase.

Downloading the example code

You can download the example code files from your account at http://www.packtpub.com for all the Packt Publishing books you have purchased. If you purchased this book elsewhere, you can visit http://www.packtpub.com/support and register to have the files e-mailed directly to you.

The code for this book is also available on GitHub at https://github.com/luispedro/BuildingMachineLearningSystemsWithPython. This repository is kept up-to-date so that it will incorporate both errata and any necessary updates for newer versions of Python or of the packages we use in the book.

Errata

Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you find a mistake in one of our books—maybe a mistake in the text or the code—we would be grateful if you could report this to us. By doing so, you can save other readers from frustration and help us improve subsequent versions of this book. If you find any errata, please report them by visiting http://www.packtpub.com/submit-errata, selecting your book, clicking on the Errata Submission Form link, and entering the details of your errata. Once your errata are verified, your submission will be accepted and the errata will be uploaded to our website or added to any list of existing errata under the Errata section of that title.

To view the previously submitted errata, go to https://www.packtpub.com/books/content/support and enter the name of the book in the search field. The required information will appear under the Errata section.

Another excellent way would be to visit www.TwoToReal.com where the authors try to provide support and answer all your questions.

Piracy

Piracy of copyright material on the Internet is an ongoing problem across all media. At Packt, we take the protection of our copyright and licenses very seriously. If you come across any illegal copies of our works, in any form, on the Internet, please provide us with the location address or website name immediately so that we can pursue a remedy.

Please contact us at copyright@packtpub.com with a link to the suspected pirated material.

We appreciate your help in protecting our authors, and our ability to bring you valuable content.

Questions

You can contact us at questions@packtpub.com if you are having a problem with any aspect of the book, and we will do our best to address it.

Table of Contents

Preface	vii
Chapter 1: Getting Started with Python Machine Learning	1
Machine learning and Python – a dream team	2
What the book will teach you (and what it will not)	3
What to do when you are stuck	4
Getting started	5
Introduction to NumPy, SciPy, and matplotlib	6
Installing Python	6
Chewing data efficiently with NumPy and intelligently with SciPy	6
Learning NumPy	7
Indexing	9
Handling nonexisting values	10
Comparing the runtime	11
Learning SciPy	12
Our first (tiny) application of machine learning	13
Reading in the data	14
Preprocessing and cleaning the data	15
Choosing the right model and learning algorithm	17
Before building our first model	18
Starting with a simple straight line	18
Towards some advanced stuff Stepping back to go forward – another look at our data	20 22
Training and testing	26
Answering our initial question	27
Summary	28
Chapter 2: Classifying with Real-world Examples	29
The Iris dataset	30
Visualization is a good first step	30
Building our first classification model	32
Evaluation – holding out data and cross-validation	36
r:1	

FT 7 7	10	~	
Table	nt	(on	tents
THULL	U	-012	POTTED

Building more complex classifiers	39
A more complex dataset and a more complex classifier	41
Learning about the Seeds dataset	41
Features and feature engineering	42
Nearest neighbor classification	43
Classifying with scikit-learn	43
Looking at the decision boundaries	45
Binary and multiclass classification	47
Summary	49
Chapter 3: Clustering – Finding Related Posts	51
Measuring the relatedness of posts	52
How not to do it	52
How to do it	53
Preprocessing – similarity measured as a similar	
number of common words	54
Converting raw text into a bag of words	54
Counting words	55
Normalizing word count vectors	58
Removing less important words	59
Stemming Stop words on steroids	60 63
Our achievements and goals	65
Clustering	66
K-means	66
Getting test data to evaluate our ideas on	70
Clustering posts	72
Solving our initial challenge	73
Another look at noise	75
Tweaking the parameters	76
Summary	77
Chapter 4: Topic Modeling	79
Latent Dirichlet allocation	80
Building a topic model	81
Comparing documents by topics	86
Modeling the whole of Wikipedia	89
Choosing the number of topics	92
Summary	94
Chapter 5: Classification – Detecting Poor Answers	95
Sketching our roadmap	96
Learning to classify classy answers	96
Tuning the instance	96

Tuning the classifier	96
Fetching the data	97
Slimming the data down to chewable chunks	98
Preselection and processing of attributes	98
Defining what is a good answer	100
Creating our first classifier	100
Starting with kNN	100
Engineering the features	101
Training the classifier	103
Measuring the classifier's performance	103
Designing more features	104
Deciding how to improve	107
Bias-variance and their tradeoff	108
Fixing high bias	108
Fixing high variance	109
High bias or low bias	109
Using logistic regression	112
A bit of math with a small example	112
Applying logistic regression to our post classification problem	114
Looking behind accuracy – precision and recall	116
Slimming the classifier	120
Ship it!	121
Summary	121
Chapter 6: Classification II – Sentiment Analysis	123
Sketching our roadmap	123
Fetching the Twitter data	124
Introducing the Naïve Bayes classifier	124
Getting to know the Bayes' theorem	125
Being naïve	126
Using Naïve Bayes to classify	127
Accounting for unseen words and other oddities	131
Accounting for arithmetic underflows	132
Creating our first classifier and tuning it	134
Solving an easy problem first	135
Using all classes	138
Tuning the classifier's parameters	141
Cleaning tweets	146
Taking the word types into account	148
Determining the word types	148
Successfully cheating using SentiWordNet	150