Frank Kane 著

# 实用数据科学和 Python机器学习
## （影印版）

Hands-On Data Science and Python Machine Learning

Packt>

# 实用数据科学和
# Python 机器学习(影印版)
## Hands-On Data Science and
## Python Machine Learning

Frank Kane 著

# Credits

**Author**
Frank Kane

**Proofreader**
Safis Editing

**Acquisition Editor**
Ben Renow-Clarke

**Indexer**
Tejal Daruwale Soni

**Content Development Editor**
Khushali Bhangde

**Graphics**
Jason Monteiro

**Technical Editor**
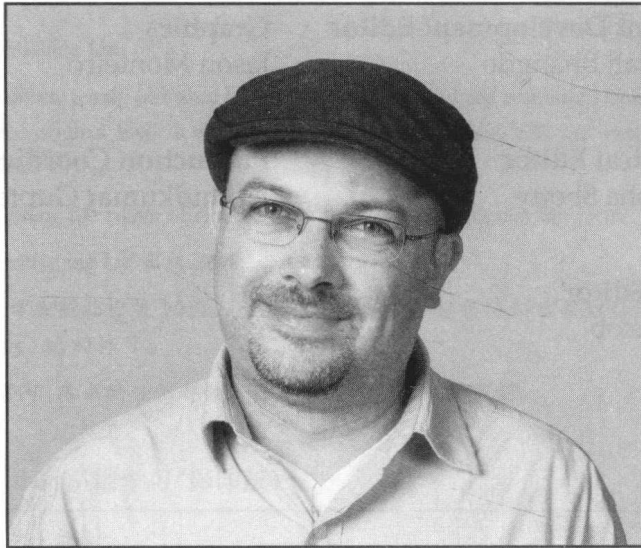Nidhisha Shetty

**Production Coordinator**
Arvindkumar Gupta

**Copy Editor**
Tom Jacob

# About the Author

My name is Frank Kane. I spent nine years at `amazon.com` and `imdb.com`, wrangling millions of customer ratings and customer transactions to produce things such as personalized recommendations for movies and products and "people who bought this also bought." I tell you, I wish we had Apache Spark back then, when I spent years trying to solve these problems there. I hold 17 issued patents in the fields of distributed computing, data mining, and machine learning. In 2012, I left to start my own successful company, Sundog Software, which focuses on virtual reality environment technology, and teaching others about big data analysis.

# www.PacktPub.com

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.comand 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.

# Mapt

https://www.packtpub.com/mapt

Get the most in-demand software skills with Mapt. Mapt gives you full access to all Packt books and video courses, as well as industry-leading tools to help you plan your personal development and advance your career.

## 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

# Customer Feedback

Thanks for purchasing this Packt book. At Packt, quality is at the heart of our editorial process. To help us improve, please leave us an honest review on this book's Amazon page at https://www.amazon.com/dp/1787280748.

If you'd like to join our team of regular reviewers, you can email us at customerreviews@packtpub.com. We award our regular reviewers with free eBooks and videos in exchange for their valuable feedback. Help us be relentless in improving our products!

# Table of Contents

# Preface

Being a data scientist in the tech industry is one of the most rewarding careers on the planet today. I went and studied actual job descriptions for data scientist roles at tech companies and I distilled those requirements down into the topics that you'll see in this course.

*Hands-On Data Science and Python Machine Learning* is really comprehensive. We'll start with a crash course on Python and do a review of some basic statistics and probability, but then we're going to dive right into over 60 topics in data mining and machine learning. That includes things such as Bayes' theorem, clustering, decision trees, regression analysis, experimental design; we'll look at them all. Some of these topics are really fun.

We're going to develop an actual movie recommendation system using actual user movie rating data. We're going to create a search engine that actually works for Wikipedia data. We're going to build a spam classifier that can correctly classify spam and nonspam emails in your email account, and we also have a whole section on scaling this work up to a cluster that runs on big data using Apache Spark.

If you're a software developer or programmer looking to transition into a career in data science, this course will teach you the hottest skills without all the mathematical notation and pretense that comes along with these topics. We're just going to explain these concepts and show you some Python code that actually works that you can dive in and mess around with to make those concepts sink home, and if you're working as a data analyst in the finance industry, this course can also teach you to make the transition into the tech industry. All you need is some prior experience in programming or scripting and you should be good to go.

The general format of this book is I'll start with each concept, explaining it in a bunch of sections and graphical examples. I will introduce you to some of the notations and fancy terminologies that data scientists like to use so you can talk the same language, but the concepts themselves are generally pretty simple. After that, I'll throw you into some actual Python code that actually works that we can run and mess around with, and that will show you how to actually apply these ideas to actual data. These are going to be presented as IPython Notebook files, and that's a format where I can intermix code and notes surrounding the code that explain what's going on in the concepts. You can take these notebook files with you after going through this book and use that as a handy-quick reference later on in your career, and at the end of each concept, I'll encourage you to actually dive into that Python code, make some modifications, mess around with it, and just gain more familiarity by getting hands-on and actually making some modifications, and seeing the effects they have.

# Who this book is for

If you are a budding data scientist or a data analyst who wants to analyze and gain actionable insights from data using Python, this book is for you. Programmers with some experience in Python who want to enter the lucrative world of Data Science will also find this book to be very useful.

# Conventions

In this book, you will find a number of text styles 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 can measure that using the `r2_score()` function from `sklearn.metrics`."

A block of code is set as follows:

```
import numpy as np
import pandas as pd
from sklearn import tree

input_file = "c:/spark/DataScience/PastHires.csv"
df = pd.read_csv(input_file, header = 0)
```

When we wish to draw your attention to a particular part of a code block, the relevant lines or items are set in bold:

```
import numpy as np
import pandas as pd
from sklearn import tree

input_file = "c:/spark/DataScience/PastHires.csv"
df = pd.read_csv(input_file, header = 0)
```

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

```
spark-submit SparkKMeans.py
```