# Spark机器学习（影印版）

## Machine Learning with Spark

ick Pentreath　著

PACKT
PUBLISHING

# Spark 机器学习(影印版)

Nick Pentreath 著

# Credits

**Author**
Nick Pentreath

**Reviewers**
Andrea Mostosi
Hao Ren
Krishna Sankar

**Commissioning Editor**
Rebecca Youé

**Acquisition Editor**
Rebecca Youé

**Content Development Editor**
Susmita Sabat

**Technical Editors**
Vivek Arora
Pankaj Kadam

**Copy Editor**
Karuna Narayanan

**Project Coordinator**
Milton Dsouza

**Proofreaders**
Simran Bhogal
Maria Gould
Ameesha Green
Paul Hindle

**Indexer**
Priya Sane

**Graphics**
Sheetal Aute
Abhinash Sahu

**Production Coordinator**
Nitesh Thakur

**Cover Work**
Nitesh Thakur

# About the Author

**Nick Pentreath** has a background in financial markets, machine learning, and software development. He has worked at Goldman Sachs Group, Inc.; as a research scientist at the online ad targeting start-up Cognitive Match Limited, London; and led the Data Science and Analytics team at Mxit, Africa's largest social network.

He is a cofounder of Graphflow, a big data and machine learning company focused on user-centric recommendations and customer intelligence. He is passionate about combining commercial focus with machine learning and cutting-edge technology to build intelligent systems that learn from data to add value to the bottom line.

Nick is a member of the Apache Spark Project Management Committee.

# Acknowledgments

# About the Reviewers

**Andrea Mostosi** is a technology enthusiast. An innovation lover since he was a child, he started a professional job in 2003 and worked on several projects, playing almost every role in the computer science environment. He is currently the CTO at The Fool, a company that tries to make sense of web and social data. During his free time, he likes traveling, running, cooking, biking, and coding.

> I would like to thank my geek friends: Simone M, Daniele V, Luca T, Luigi P, Michele N, Luca O, Luca B, Diego C, and Fabio B. They are the smartest people I know, and comparing myself with them has always pushed me to be better.

**Hao Ren** is a software developer who is passionate about Scala, distributed systems, machine learning, and Apache Spark. He was an exchange student at EPFL when he learned about Scala in 2012. He is currently working in Paris as a backend and data engineer for ClaraVista — a company that focuses on high-performance marketing. His work responsibility is to build a Spark-based platform for purchase prediction and a new recommender system.

Besides programming, he enjoys running, swimming, and playing basketball and badminton. You can learn more at his blog http://www.invkrh.me.

**Krishna Sankar** is a chief data scientist at BlackArrow, where he is focusing on enhancing user experience via inference, intelligence, and interfaces. Earlier stints include working as a principal architect and data scientist at Tata America International Corporation, director of data science at a bioinformatics start-up company, and as a distinguished engineer at Cisco Systems, Inc. He has spoken at various conferences about data science (http://goo.gl/9pyJMH), machine learning (http://goo.gl/sSem2Y), and social media analysis (http://goo.gl/D9YpVQ). He has also been a guest lecturer at the Naval Postgraduate School. He has written a few books on Java, wireless LAN security, Web 2.0, and now on Spark. His other passion is LEGO robotics. Earlier in April, he was at the St. Louis FLL World Competition as a robots design judge.

# 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.

# PACKTLiB™

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.

# Table of Contents

# Preface

In recent years, the volume of data being collected, stored, and analyzed has exploded, in particular in relation to the activity on the Web and mobile devices, as well as data from the physical world collected via sensor networks. While previously large-scale data storage, processing, analysis, and modeling was the domain of the largest institutions such as Google, Yahoo!, Facebook, and Twitter, increasingly, many organizations are being faced with the challenge of how to handle a massive amount of data.

When faced with this quantity of data and the common requirement to utilize it in real time, human-powered systems quickly become infeasible. This has led to a rise in the so-called big data and machine learning systems that learn from this data to make automated decisions.

In answer to the challenge of dealing with ever larger-scale data without any prohibitive cost, new open source technologies emerged at companies such as Google, Yahoo!, Amazon, and Facebook, which aimed at making it easier to handle massive data volumes by distributing data storage and computation across a cluster of computers.

The most widespread of these is Apache Hadoop, which made it significantly easier and cheaper to both store large amounts of data (via the Hadoop Distributed File System, or HDFS) and run computations on this data (via Hadoop MapReduce, a framework to perform computation tasks in parallel across many nodes in a computer cluster).

However, MapReduce has some important shortcomings, including high overheads to launch each job and reliance on storing intermediate data and results of the computation to disk, both of which make Hadoop relatively ill-suited for use cases of an iterative or low-latency nature. Apache Spark is a new framework for distributed computing that is designed from the ground up to be optimized for low-latency tasks and to store intermediate data and results in memory, thus addressing some of the major drawbacks of the Hadoop framework. Spark provides a clean, functional, and easy-to-understand API to write applications and is fully compatible with the Hadoop ecosystem.

Furthermore, Spark provides native APIs in Scala, Java, and Python. The Scala and Python APIs allow all the benefits of the Scala or Python language, respectively, to be used directly in Spark applications, including using the relevant interpreter for real-time, interactive exploration. Spark itself now provides a toolkit (called MLlib) of distributed machine learning and data mining models that is under heavy development and already contains high-quality, scalable, and efficient algorithms for many common machine learning tasks, some of which we will delve into in this book.

Applying machine learning techniques to massive datasets is challenging, primarily because most well-known machine learning algorithms are not designed for parallel architectures. In many cases, designing such algorithms is not an easy task. The nature of machine learning models is generally iterative, hence the strong appeal of Spark for this use case. While there are many competing frameworks for parallel computing, Spark is one of the few that combines speed, scalability, in-memory processing, and fault tolerance with ease of programming and a flexible, expressive, and powerful API design.

Throughout this book, we will focus on real-world applications of machine learning technology. While we may briefly delve into some theoretical aspects of machine learning algorithms, the book will generally take a practical, applied approach with a focus on using examples and code to illustrate how to effectively use the features of Spark and MLlib, as well as other well-known and freely available packages for machine learning and data analysis, to create a useful machine learning system.

# What this book covers

*Chapter 1*, *Getting Up and Running with Spark*, shows how to install and set up a local development environment for the Spark framework as well as how to create a Spark cluster in the cloud using Amazon EC2. The Spark programming model and API will be introduced, and a simple Spark application will be created using each of Scala, Java, and Python.