

O'REILLY®

TURING

图灵程序设计丛书

第2版



# Spark

## 高级数据分析

Advanced Analytics with Spark, Second Edition

涵盖大规模数据分析中常用算法、数据集和设计模式

[美]桑迪·里扎 [美]于里·莱瑟森 著  
[英]肖恩·欧文 [美]乔希·威尔斯  
龚少成 邱鑫 译



中国工信出版集团



人民邮电出版社  
POSTS & TELECOM PRESS

# Spark高级数据分析（第2版）

Advanced Analytics with Spark, Second Edition

[美]桑迪·里扎 [美]于里·莱瑟森 [英]肖恩·欧文 [美]乔希·威尔斯 著  
龚少成 邱鑫 译



Beijing • Boston • Farnham • Sebastopol • Tokyo

O'REILLY®

O'Reilly Media, Inc.授权人民邮电出版社出版

人民邮电出版社  
北京

## 图书在版编目 (C I P ) 数据

Spark高级数据分析 : 第2版 / (美) 桑迪·里扎  
(Sandy Ryza) 等著 ; 龚少成, 邱鑫译. — 北京 : 人民  
邮电出版社, 2018. 6

(图灵程序设计丛书)

ISBN 978-7-115-48252-5

I. ①S… II. ①桑… ②龚… ③邱… III. ①数据处  
理软件 IV. ①TP274

中国版本图书馆CIP数据核字(2018)第076391号

## 内 容 提 要

本书是使用 Spark 进行大规模数据分析的实战宝典, 由知名数据科学家撰写。本书在第 1 版的基础上, 针对 Spark 近年来的发展, 对样例代码和所使用的资料进行了大量更新。新版 Spark 使用了全新的核心 API, MLlib 和 Spark SQL 两个子项目也发生了较大变化, 本书为关注 Spark 发展趋势的读者提供了与时俱进的资料, 例如 Dataset 和 DataFrame 的使用, 以及与 DataFrame API 高度集成的 Spark ML API。

本书适合从事数据分析的各类专业人员阅读。

- 
- ◆ 著 [美] 桑迪·里扎 [美] 于里·莱瑟森  
[英] 肖恩·欧文 [美] 乔希·威尔斯  
译 龚少成 邱 鑫  
责任编辑 朱 巍  
执行编辑 温 雪  
责任印制 周昇亮
- ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号  
邮编 100164 电子邮件 315@ptpress.com.cn  
网址 <http://www.ptpress.com.cn>
- 北京鑫正大印刷有限公司印刷
- ◆ 开本: 800×1000 1/16  
印张: 15.25  
字数: 360千字 2018年6月第1版  
印数: 1~4 000册 2018年6月北京第1次印刷  
著作权合同登记号 图字: 01-2018-2239号
- 

定价: 69.00元

读者服务热线: (010)51095186转600 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

## 译者介绍



### 龚少成

现任万达科技集团数据工程部总经理，清华大学自动化系研究生毕业，国内专注企业级大数据平台建设的先驱者之一，曾经在Intel和Cloudera公司担任大数据技术负责人，Cloudera公司认证大数据培训讲师。



### 邱鑫

毕业于武汉大学，目前就职于英特尔亚太研发有限公司，是Intel大数据团队高级工程师。主要研究大数据与深度学习技术，是基于Spark的深度学习框架BigDL的核心贡献者。



微信连接



回复“Spark”查看相关书单



微博连接

关注@图灵教育 每日分享IT好书



QQ连接

图灵读者官方群I: 218139230

图灵读者官方群II: 164939616

图灵社区  
iTuring.cn

在线出版，电子书，《码农》杂志，图灵访谈

# 版权声明

© 2017 by Sandy Ryza, Uri Laserson, Sean Owen, and Josh Wills.

Simplified Chinese Edition, jointly published by O'Reilly Media, Inc. and Posts & Telecom Press, 2018. Authorized translation of the English edition, 2018 O'Reilly Media, Inc., 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 O'Reilly Media, Inc. 出版，2017。

简体中文版由人民邮电出版社出版，2018。英文原版的翻译得到 O'Reilly Media, Inc. 的授权。此简体中文版的出版和销售得到出版权和销售权的所有者——O'Reilly Media, Inc. 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。

# O'Reilly Media, Inc.介绍

O'Reilly Media 通过图书、杂志、在线服务、调查研究和会议等方式传播创新知识。自 1978 年开始，O'Reilly 一直都是前沿发展的见证者和推动者。超级极客们正在开创着未来，而我们关注真正重要的技术趋势——通过放大那些“细微的信号”来刺激社会对新科技的应用。作为技术社区中活跃的参与者，O'Reilly 的发展充满了对创新的倡导、创造和发扬光大。

O'Reilly 为软件开发人员带来革命性的“动物书”；创建第一个商业网站（GNN）；组织了影响深远的开放源代码峰会，以至于开源软件运动以此命名；创立了 Make 杂志，从而成为 DIY 革命的主要先锋；公司一如既往地通过多种形式缔结信息与人的纽带。O'Reilly 的会议和峰会集聚了众多超级极客和高瞻远瞩的商业领袖，共同描绘出开创新产业的革命性思想。作为技术人士获取信息的选择，O'Reilly 现在还将先锋专家的知识传递给普通的计算机用户。无论是通过书籍出版、在线服务或者面授课程，每一项 O'Reilly 的产品都反映了公司不可动摇的理念——信息是激发创新的力量。

## 业界评论

“O'Reilly Radar 博客有口皆碑。”

——*Wired*

“O'Reilly 凭借一系列（真希望当初我也想到了）非凡想法建立了数百万美元的业务。”

——*Business 2.0*

“O'Reilly Conference 是聚集关键思想领袖的绝对典范。”

——*CRN*

“一本 O'Reilly 的书就代表一个有用、有前途、需要学习的主题。”

——*Irish Times*

“Tim 是位特立独行的商人，他不光放眼于最长远、最广阔的视野，并且切实地按照 Yogi Berra 的建议去做了：‘如果你在路上遇到岔路口，走小路（岔路）。’回顾过去，Tim 似乎每一次都选择了小路，而且有几次都是一闪即逝的机会，尽管大路也不错。”

——*Linux Journal*

---

# 推荐序

数据的爆炸式增长和隐藏在这些数据背后的商业价值催生了一代又一代的大数据处理技术。十余年前 Hadoop 横空出世，Doug Cutting 先生将谷歌的 MapReduce 思想用开源的方式实现出来，由此拉开了基于 MapReduce 的大数据处理框架在企业中应用的序幕。近年来，Hadoop 生态系统又发展出以 Spark 为代表的新计算框架。相比 MapReduce，Spark 速度快，开发简单，并且能同时兼顾批处理和实时数据分析。Spark 起源于加州大学伯克利分校的 AMPLab，Cloudera 公司作为大数据市场上的翘楚，很早就开始将 Spark 推广到广大企业级客户并积累了大量的经验。*Advanced Analysis with Spark* 一书正是这些经验的结晶。另一方面，企业级用户在引入 Spark 技术时碰到的最大难题之一就是能够灵活应用 Spark 技术的人才匮乏。龚少成与图灵公司将 *Advanced Analysis with Spark* 翻译成中文，让国内读者第一时间用母语感受 Spark 这一新技术在数据分析和处理方面的魔力，实在是国内技术圈的幸事。能为本书作序推荐，也算是为国内企业更好地应用 Spark 技术尽自己的一份力量！

本书开篇介绍了 Spark 的基础知识，然后详细介绍了如何将 Spark 应用到各个行业。与许多图书只着重描述最终方案不同，本书作者在介绍案例时把解决问题的整个过程也展现了出来。在介绍一个主题时，并不是一开始就给出最终方案，而是先给出一个最初并不完善的方案，然后指出方案的不足，引导读者思考并逐步改进，最终得出一个相对完善方案。这体现了工程问题的解决思路，也体现了大数据分析是一个迭代的过程。这样的论述方式更能激发读者的思考，这一点实在难能可贵。

本书英文版自第 1 版出版以来，在亚马逊网站大数据分析类图书中一直名列前茅，而且获得的多为五星级评价，可见国外读者对该书的喜爱。本书中文版译者龚少成技术扎实，在英特尔和 Cloudera 工作期间带领团队成功实施过许多国内标杆大数据平台项目，最近两年又转战万达科技集团大数据中心从零到一构建 PB 级大数据平台并支撑业务落地，而且其英语功底也相当扎实，此外我偶然得知他还是国内少数通过高级口译考试的专业人才。所以本书的中文版交给龚少成翻译实在是件让人欣慰的事情。本书中文版初稿也证实了我的判断，不仅保持了英文版的风格，而且语言也十分流畅。如果你了解 Scala 语言，还有一些统计学和机器学习基础，那么本书是你学习 Spark 时必备的图书之一！

——苗凯翔，思科中国研发公司首席技术官，前 Cloudera 公司副总裁

# 译者序

大数据是这几年科技和应用领域炙手可热的话题，而 Spark 又是大数据领域里最活跃的技术。随着人工智能的崛起，业内对大数据的需求不再局限于一般意义上的大数据存储、加工和分析，如何挖掘大数据的潜在价值成为新的热点。本书四位作者均在 Cloudera 公司担任过数据科学家，长期为客户提供专业的数据分析和挖掘服务。可以说，本书的出版将为 Spark 在数据分析和挖掘领域起到巨大的推动作用。

同时我们也注意到，国内介绍 Spark 数据分析方面的图书还比较匮乏，而且许多图书都停留在源代码研究的层面上。当然，这些书中也不乏非常优秀的作品，但我们认为 Spark 真正的力量在于其开发的大数据应用。所以早在本书还处于初期编写过程中时，我们就自告奋勇和作者联系中文版事宜，希望以此为中国的数据分析事业略尽绵力。

本书在翻译过程中得到了许多人的帮助。首先要感谢我在 Cloudera 公司的前同事，也就是本书的 4 位作者。在本书的翻译过程中，由于不同语言的习惯问题，4 位作者桑迪·里扎、于里·莱瑟森、肖恩·欧文和乔希·威尔斯花了许多时间和我交流。本人之所以有幸负责本书的中文版翻译，也是承蒙肖恩·欧文的引荐。其次要感谢星环信息科技有限公司创始人孙元浩先生将我带入到大数据这个领域，让我的人生轨迹发生变化；感谢思科中国研发公司首席技术官苗凯翔博士在英特尔和 Cloudera 工作期间曾经给我的指导，让我有了端正的工作态度和价值观；感谢我的前同事田占凤博士和陈建忠的鼓励，中文版的翻译工作才得以开始。同时本书在翻译过程中还得到了 Cloudera 公司中国区前同事刘贺峰、糜君、陈飚、陈新江、李大超和张莉萍的鼎力帮助。感谢图灵公司的李松峰、岳新欣、温雪编辑在翻译过程中的指导和仔细审阅。由于本书的翻译都是在周末完成的，所以特别感谢我的妻子周幼琼在每个周末对我的照顾。

龚少成

首先非常感谢龚少成给我这次机会，使我有幸成为本书第 2 版的译者之一。

其次要感谢英特尔大数据团队的同事们，是你们带领我走进了 Spark 的时代。

最后要感谢我的妻子和孩子对我工作的理解和支持，让我腾出业余时间完成此次翻译工作。

由于译者水平有限，同时本书涉及许多课题，所以现有译文中难免存在纰漏之处。希望读者能够不吝赐教，发现问题时麻烦和译者联系。邮件请发送至 [gongshaocheng@gmail.com](mailto:gongshaocheng@gmail.com) 或 [qiuxin2012cs@gmail.com](mailto:qiuxin2012cs@gmail.com)。

邱鑫

---

# 序

自从在加州大学伯克利分校创立 Spark 项目起，我就时常心潮澎湃。不仅因为 Spark 可以帮助人们快速构建并行系统，更因为 Spark 帮助了越来越多的人使用大规模计算。因此看到这本介绍 Spark 高级分析的书，我非常欣慰！该书由数据科学领域 4 位专家桑迪、千里、肖恩和乔希携手打造。4 位作者研习 Spark 已久，他们在本书中跟读者分享了关于 Spark 的大量精彩内容，同时本书的案例部分同样出众！

对于这本书，我最钟爱的是它强调案例，而且这些案例都源于现实数据和实际应用。找到 1 个像样的、能在笔记本电脑上运行的大数据案例已经很难，更遑论 10 个了。但本书作者做到了！作者为大家准备好了一切，只等你在 Spark 中运行它们。更难能可贵的是，作者不仅讨论了核心算法，更倾心于数据准备和模型调优，没有这些工作，实际项目中就无法得到好的结果。认真研读此书，你应该可以吸收这些案例中的概念并直接将其运用在自己的项目中！

大数据处理无疑是当今计算领域最激动人心的方向之一，发展非常迅猛，新思想层出不穷。愿本书能帮助你在这个崭新的领域中扬帆启航！

Matei Zaharia  
Databricks 公司 CTO 兼 Apache Spark 项目副总裁

---

# 前言

作者：桑迪·里扎

我不希望我的人生有很多遗憾。2011年，某个慵懒的时刻，我正在绞尽脑汁地想如何把高难度的离散优化问题最优地分配给计算机集群处理，真是很难想到有什么好方法。我的导师跟我讲，他听说有个叫Apache Spark的新技术，可我基本上没当回事。Spark的想法太好了，让人觉得有点儿不靠谱。就这样，我很快又回去接着写MapReduce的本科毕业论文了。时光荏苒，Spark和我都渐渐成熟，不过令我望尘莫及的是，Spark已然成为冉冉之星，这让人不禁感叹“Spark”（星星之火）这个双关语是多么贴切。若干年后，Spark的价值举世皆知！

Spark的前辈有很多，从MPI到MapReduce。利用这些计算框架，我们写的程序可以充分利用大量资源，但不需要关心分布式系统的实现细节。数据处理的需求促进了这些技术框架的发展。同样，大数据领域也和这些框架关系密切，这些框架界定了大数据的范围。Spark有望更进一步，让写分布式程序就像写普通程序一样。

Spark能大大提升ETL流水作业的性能，并把MapReduce程序员从每天问天天不灵、问地地不应的绝望痛苦中解救出来。对我而言，Spark的激动人心之处在于，它真正打开了复杂数据分析的大门。Spark带来了支持迭代式计算和交互式探索的模式。利用这一开源计算框架，数据科学家终于可以在大数据集上高效地工作了。

我认为数据科学教学最有效的方法是利用实例。为此，我和同事一起编写了这本关于实际应用的书，希望它能涵盖大规模数据分析中最常用的算法、数据集和设计模式。阅读本书时不必一页一页地看，可以根据工作需要或按兴趣直接翻到相关章节。

# 本书内容

第1章结合数据科学和大数据分析的广阔背景来讨论Spark。随后各章在介绍Spark数据分析时都自成一体。第2章通过数据清洗这一使用场景来介绍用Spark和Scala进行数据处理的基础知识。接下来几章深入讨论如何将Spark用于机器学习，介绍了常见应用中几个最常用的算法。其余几章则收集了一些更新颖的应用，比如通过文本隐含语义关系来查询Wikipedia或分析基因数据。

## 第2版说明

自本书第1版出版以来，Spark进行了一次重大的版本更新：使用了一个全新的核心API；MLlib和Spark SQL两个子项目也发生了翻天覆地的变化。第2版根据新版Spark的最佳实践，对样例代码和所使用的资料进行了大量更新。

## 使用代码示例

补充材料（代码示例、练习、勘误表等）可以从<https://github.com/sryza/aas>下载<sup>1</sup>。

本书是要帮你完成工作的。一般来说，如果本书提供了示例代码，你可以把它用在你的程序或文档中。除非你使用了很大一部分代码，否则无须联系我们获得许可。比如，用本书的几个代码片段写一个程序就无须获得许可，销售或分发O'Reilly图书的示例光盘则需要获得许可；引用本书中的示例代码回答问题无须获得许可，将书中大量的代码放到你的产品文档中则需要获得许可。

我们很希望但并不强制要求你在引用本书内容时加上引用说明。引用说明一般包括书名、作者、出版社和ISBN。比如：“Advanced Analytics with Spark by Sandy Ryza, Uri Laserson, Sean Owen, and Josh Wills (O'Reilly). Copyright 2015 Sandy Ryza, Uri Laserson, Sean Owen, and Josh Wills, 978-1-491-91276-8.”

如果你觉得自己对示例代码的用法超出了上述许可的范围，欢迎你通过permissions@oreilly.com与我们联系。

## O'Reilly Safari



Safari（前身为Safari Books Online）是为企业、政府、教育机构和个人提供的会员制培训和参考平台。

注1：本书中文版勘误提交及资料下载，请访问本书图灵社区页面：<http://www.ituring.com.cn/book/2039>。——编者注

会员可以访问来自 250 多家出版商的上千种图书、培训视频、学习路径、互动教程和精选播放列表。这些出版商包括 O'Reilly Media、Harvard Business Review、Prentice Hall Professional、Addison-Wesley Professional、Microsoft Press、Sams、Que、Peachpit Press、Adobe、Focal Press、Cisco Press、John Wiley & Sons、Syngress、Morgan Kaufmann、IBM Redbooks、Packt、Adobe Press、FT Press、Apress、Manning、New Riders、McGraw-Hill、Jones & Bartlett、Course Technology，等等。

欲知更多信息，请访问 <https://www.safaribooksonline.com/>。

## 联系我们

请把对本书的评价和问题发给出版社。

美国：

O'Reilly Media, Inc.  
1005 Gravenstein Highway North  
Sebastopol, CA 95472

中国：

北京市西城区西直门南大街 2 号成铭大厦 C 座 807 室（100035）  
奥莱利技术咨询（北京）有限公司

O'Reilly 的每一本书都有专属网页，你可以在那儿找到本书的相关信息，包括勘误表、示例代码以及其他信息。本书的网站地址是：

<http://shop.oreilly.com/product/0636920056591.do>

对于本书的评论和技术性问题，请发送电子邮件到：

[bookquestions@oreilly.com](mailto:bookquestions@oreilly.com)

要了解更多 O'Reilly 图书、培训课程、会议和新闻的信息，请访问以下网站：

<http://www.oreilly.com>

我们在 Facebook 的地址如下：

<https://facebook.com/oreilly>

请关注我们的 Twitter 动态：

<https://twitter.com/oreillymedia>

我们的 YouTube 视频地址如下：

<https://www.youtube.com/oreillymedia>

# 致谢

如果没有 Apache Spark 和 MLlib，就没有本书。所以我们应该感谢开发了 Spark 和 MLlib 并将其开源的团体，也要感谢那些添砖加瓦的数以百计的代码贡献者。

我们还要感谢本书的每一位审阅者，感谢他们花费了大量的时间来审阅本书的内容，感谢他们的专业视角，他们是 Michael Bernico、Adam Breindel、Ian Buss、Parvis Deyhim、Jeremy Freeman、Chris Fregly、Debashish Ghosh、Juliet Hougland、Jonathan Keebler、Nisha Muktewar、Frank Nothaft、Nick Pentreath、Kostas Sakellis、Tom White、Marcelo Vanzin 和另一位 Juliet Hougland。谢谢你们所有人！我们欠你们一个大人情！你们的努力大大改进了本书的结构和质量。

我（桑迪）还要感谢 Jordan Pinkus 和 Richard Wang，你们帮助我完成了风险分析章节的原理部分。

感谢 Marie Beaugureau 和 O'Reilly 出版社在本书出版和发行过程中贡献的宝贵经验和大力支持！

## 电子版

扫描如下二维码，即可购买本书电子版。



# 目录

推荐序 .....	ix
译者序 .....	xi
序 .....	xiii
前言 .....	xv
<b>第 1 章 大数据分析 .....</b>	<b>1</b>
1.1 数据科学面临的挑战 .....	2
1.2 认识 Apache Spark .....	4
1.3 关于本书 .....	5
1.4 第 2 版说明 .....	6
<b>第 2 章 用 Scala 和 Spark 进行数据分析 .....</b>	<b>8</b>
2.1 数据科学家的 Scala .....	9
2.2 Spark 编程模型 .....	10
2.3 记录关联问题 .....	10
2.4 小试牛刀：Spark shell 和 SparkContext .....	11
2.5 把数据从集群上获取到客户端 .....	16
2.6 把代码从客户端发送到集群 .....	19
2.7 从 RDD 到 DataFrame .....	20
2.8 用 DataFrame API 来分析数据 .....	23
2.9 DataFrame 的统计信息 .....	27
2.10 DataFrame 的转置和重塑 .....	29
2.11 DataFrame 的连接和特征选择 .....	32
2.12 为生产环境准备模型 .....	33
2.13 评估模型 .....	35
2.14 小结 .....	36

第3章 音乐推荐和Audioscrobbler数据集	37
3.1 数据集	38
3.2 交替最小二乘推荐算法	39
3.3 准备数据	41
3.4 构建第一个模型	44
3.5 逐个检查推荐结果	47
3.6 评价推荐质量	50
3.7 计算AUC	51
3.8 选择超参数	53
3.9 产生推荐	55
3.10 小结	56
第4章 用决策树算法预测森林植被	58
4.1 回归简介	59
4.2 向量和特征	59
4.3 样本训练	60
4.4 决策树和决策森林	61
4.5 Covtype数据集	63
4.6 准备数据	64
4.7 第一棵决策树	66
4.8 决策树的超参数	72
4.9 决策树调优	73
4.10 重谈类别型特征	77
4.11 随机决策森林	79
4.12 进行预测	81
4.13 小结	82
第5章 基于K均值聚类的网络流量异常检测	84
5.1 异常检测	85
5.2 K均值聚类	85
5.3 网络入侵	86
5.4 KDD Cup 1999数据集	86
5.5 初步尝试聚类	87
5.6 $k$ 的选择	90
5.7 基于SparkR的可视化	92
5.8 特征的规范化	96
5.9 类别型变量	98
5.10 利用标号的熵信息	99
5.11 聚类实战	100
5.12 小结	102