

Spark快速大数据分析 (图灵程序设计丛书)

作者: [美] 卡劳 (Holden Karau) [美] 肯维尼斯科 (Andy Konwinski) [美] 温德尔 (Patrick Wendell) [加] 扎哈里亚 (Matei Zaharia)

版权信息

书名: Spark快速大数据分析

作者: [美] Holden Karau Andy Konwinski,Patrick Wendell [加] Matei Zaharia

译者: 王道远

ISBN: 978-7-115-40309-4

本书由北京图灵文化发展有限公司发行数字版。版权所有, 侵权必究。

您购买的图灵电子书仅供您个人使用, 未经授权, 不得以任何方式复制和传播本书内容。

我们愿意相信读者具有这样的良知和觉悟, 与我们共同保护知识产权。

如果购买者有侵权行为, 我们可能对该用户实施包括但不限于关闭该帐号等维权措施, 并可能追究法律责任。

图灵社区会员 张海川 (zhanghaichuan@ptpress.com.cn) 专享 尊重版权

[版权声明](#)

[O'Reilly Media, Inc. 介绍](#)

[业界评论](#)

[推荐序](#)

[译者序](#)

[序](#)

[前言](#)

[读者对象](#)

[本书结构](#)

[相关书籍](#)

[排版约定](#)

[使用代码示例](#)

[Safari® Books Online](#)

[联系我们](#)

[致谢](#)

[第1章 Spark 数据分析导论](#)

[1.1 Spark是什么](#)

[1.2 一个大一统的软件栈](#)

[1.2.1 Spark Core](#)

[1.2.2 Spark SQL](#)

[1.2.3 Spark Streaming](#)

[1.2.4 MLlib](#)

[1.2.5 GraphX](#)

[1.2.6 集群管理器](#)

[1.3 Spark的用户和用途](#)

[1.3.1 数据科学任务](#)

[1.3.2 数据处理应用](#)

[1.4 Spark简史](#)

[1.5 Spark的版本和发布](#)

[1.6 Spark的存储层次](#)

[第2章 Spark 下载与入门](#)

[2.1 下载Spark](#)

[2.2 Spark中Python和Scala的shell](#)

[2.3 Spark核心概念简介](#)

[2.4 独立应用](#)

[2.4.1 初始化SparkContext](#)

[2.4.2 构建独立应用](#)

[2.5 总结](#)

[第3章 RDD 编程](#)

[3.1 RDD基础](#)

[3.2 创建RDD](#)

[3.3 RDD操作](#)

[3.3.1 转化操作](#)

[3.3.2 行动操作](#)

[3.3.3 惰性求值](#)

[3.4 向Spark传递函数](#)

[3.4.1 Python](#)

- [3.4.2 Scald](#)
- [3.4.3 Java](#)
- [3.5 常见的转化操作和行动操作](#)
 - [3.5.1 基本RDD](#)
 - [3.5.2 在不同RDD类型间转换](#)
 - [3.6 持久化\(缓存\)](#)
 - [3.7 总结](#)
- [第4章 键值对操作](#)
 - [4.1 动机](#)
 - [4.2 创建Pair RDD](#)
 - [4.3 Pair RDD的转化操作](#)
 - [4.3.1 聚合操作](#)
 - [4.3.2 数据分组](#)
 - [4.3.3 连接](#)
 - [4.3.4 数据排序](#)
 - [4.4 Pair RDD的行动操作](#)
 - [4.5 数据分区\(进阶\)](#)
 - [4.5.1 获取RDD的分区方式](#)
 - [4.5.2 从分区中获益的操作](#)
 - [4.5.3 影响分区方式的操作](#)
 - [4.5.4 示例: PageRank](#)
 - [4.5.5 自定义分区方式](#)
 - [4.6 总结](#)
- [第5章 数据读取与保存](#)
 - [5.1 动机](#)
 - [5.2 文件格式](#)
 - [5.2.1 文本文件](#)
 - [5.2.2 JSON](#)
 - [5.2.3 逗号分隔值与制表符分隔值](#)
 - [5.2.4 SequenceFile](#)
 - [5.2.5 对象文件](#)
 - [5.2.6 Hadoop输入输出格式](#)
 - [5.2.7 文件压缩](#)
 - [5.3 文件系统](#)
 - [5.3.1 本地“常规”文件系统](#)
 - [5.3.2 Amazon S3](#)
 - [5.3.3 HDFS](#)
 - [5.4 Spark SQL中的结构化数据](#)
 - [5.4.1 Apache Hive](#)
 - [5.4.2 JSON](#)
 - [5.5 数据库](#)
 - [5.5.1 Java数据库连接](#)
 - [5.5.2 Cassandra](#)
 - [5.5.3 HBase](#)
 - [5.5.4 Elasticsearch](#)
 - [5.6 总结](#)
- [第6章 Spark 编程进阶](#)

- [6.1 简介](#)
- [6.2 累加器](#)
 - [6.2.1 累加器与容错性](#)
 - [6.2.2 自定义累加器](#)
- [6.3 广播变量](#)
- [广播的优化](#)
- [6.4 基于分区进行操作](#)
- [6.5 与外部程序间的管道](#)
- [6.6 数值RDD的操作](#)
- [6.7 总结](#)
- [第7章 在集群上运行 Spark](#)
 - [7.1 简介](#)
 - [7.2 Spark运行时架构](#)
 - [7.2.1 驱动器节点](#)
 - [7.2.2 执行器节点](#)
 - [7.2.3 集群管理器](#)
 - [7.2.4 启动一个程序](#)
 - [7.2.5 小结](#)
 - [7.3 使用spark-submit部署应用](#)
 - [7.4 打包代码与依赖](#)
 - [7.4.1 使用Maven构建的用Java编写的Spark应用](#)
 - [7.4.2 使用sbt构建的用Scala编写的Spark应用](#)
 - [7.4.3 依赖冲突](#)
 - [7.5 Spark应用内与应用间调度](#)
 - [7.6 集群管理器](#)
 - [7.6.1 独立集群管理器](#)
 - [7.6.2 Hadoop YARN](#)
 - [7.6.3 Apache Mesos](#)
 - [7.6.4 Amazon EC2](#)
 - [7.7 选择合适的集群管理器](#)
 - [7.8 总结](#)
- [第8章 Spark 调优与调试](#)
 - [8.1 使用SparkConf配置Spark](#)
 - [8.2 Spark执行的组成部分：作业、任务和步骤](#)
 - [8.3 查找信息](#)
 - [8.3.1 Spark网页用户界面](#)
 - [8.3.2 驱动器进程和执行器进程的日志](#)
 - [8.4 关键性能考量](#)
 - [8.4.1 并行度](#)
 - [8.4.2 序列化格式](#)
 - [8.4.3 内存管理](#)
 - [8.4.4 硬件供给](#)
 - [8.5 总结](#)
- [第9章 Spark SQL](#)
 - [9.1 连接Spark SQL](#)
 - [9.2 在应用中使用Spark SQL](#)
 - [9.2.1 初始化Spark SQL](#)

[9.2.2 基本查询示例](#)

[9.2.3 SchemaRDD](#)

[9.2.4 缓存](#)

[9.3 读取和存储数据](#)

[9.3.1 Apache Hive](#)

[9.3.2 Parquet](#)

[9.3.3 JSON](#)

[9.3.4 基于RDD](#)

[9.4 JDBC/ODBC服务器](#)

[9.4.1 使用Beeline](#)

[9.4.2 长生命周期的表与查询](#)

[9.5 用户自定义函数](#)

[9.5.1 Spark SQL UDF](#)

[9.5.2 Hive UDF](#)

[9.6 Spark SQL性能](#)

[性能调优选项](#)

[9.7 总结](#)

[第 10 章 Spark Streaming](#)

[10.1 一个简单的例子](#)

[10.2 架构与抽象](#)

[10.3 转化操作](#)

[10.3.1 无状态转化操作](#)

[10.3.2 有状态转化操作](#)

[10.4 输出操作](#)

[10.5 输入源](#)

[10.5.1 核心数据源](#)

[10.5.2 附加数据源](#)

[10.5.3 多数据源与集群规模](#)

[10.6 24/7不间断运行](#)

[10.6.1 检查点机制](#)

[10.6.2 驱动程序程序容错](#)

[10.6.3 工作节点容错](#)

[10.6.4 接收器容错](#)

[10.6.5 处理保证](#)

[10.7 Streaming用户界面](#)

[10.8 性能考量](#)

[10.8.1 批次和窗口大小](#)

[10.8.2 并行度](#)

[10.8.3 垃圾回收和内存使用](#)

[10.9 总结](#)

[第 11 章基于 MLlib 的机器学习](#)

[11.1 概述](#)

[11.2 系统要求](#)

[11.3 机器学习基础](#)

[示例：垃圾邮件分类](#)

[11.4 数据类型](#)

[操作向量](#)

[11.5 算法](#)

[11.5.1 特征提取](#)

[11.5.2 统计](#)

[11.5.3 分类与回归](#)

[11.5.4 聚类](#)

[11.5.5 协同过滤与推荐](#)

[11.5.6 降维](#)

[11.5.7 模型评估](#)

[11.6 一些提示与性能考量](#)

[11.6.1 准备特征](#)

[11.6.2 配置算法](#)

[11.6.3 缓存RDD以重复使用](#)

[11.6.4 识别稀疏程度](#)

[11.6.5 并行度](#)

[11.7 流水线API](#)

[11.8 总结](#)

[作者简介](#)

[封面介绍](#)

版权声明

© 2015 by O'Reilly Media, Inc.

Simplified Chinese Edition, jointly published by O'Reilly Media, Inc. and Posts & Telecom Press, 2015. Authorized translation of the English edition, 2015 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，2015。简体中文版由人民邮电出版社出版，2015。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

近年来大数据逐渐升温，经常有人问起大数据为何重要。我们处在一个数据爆炸的时代，大量涌现的智能手机、平板、可穿戴设备及物联网设备每时每刻都在产生新的数据。当今世界，有 90% 的数据是在过去短短两年内产生的。到 2020 年，将有 500 多亿台的互联设备产生 Zeta 字节级的数据。带来革命性改变的并非海量数据本身，而是我们如何利用这些数据。大数据解决方案的强大在于它们可以快速处理大规模、复杂的数据集，可以比传统方法更快、更好地生成洞见。

一套大数据解决方案通常包含多个重要组件，从存储、计算和网络等硬件层，到数据处理引擎，再到利用改良的统计和计算算法、数据可视化来获得商业洞见的分析层。这中间，数据处理引擎起到了十分重要的作用。毫不夸张地说，数据处理引擎之于大数据就像 CPU 之于计算机，或大脑之于人类。

早在 2009 年，Matei Zaharia 在加州大学伯克利分校的 AMPLab 进行博士研究时创立了 Spark 大数据处理和计算框架。不同于传统的数据处理框架，Spark 基于内存的基本类型 (primitive) 为一些应用程序带来了 100 倍的性能提升。Spark 允许用户程序将数据加载到集群内存中用于反复查询，非常适用于大数据和机器学习，日益成为最广泛采用的大数据模块之一。包括 Cloudera 和 MapR 在内的大数据发行版也在发布时添加了 Spark。

目前，Spark 正在促使 Hadoop 和大数据生态系统发生演变，以更好地支持端到端的大数据分析需求，例如：Spark 已经超越 Spark 核心，发展到了 Spark streaming、SQL、MLlib、GraphX、SparkR 等模块。学习 Spark 和它的各个内部构件不仅有助于改善大数据处理速度，还能帮助开发者和数据科学家更轻松创建分析应用。从企业、医疗、交通到零售业，Spark 这样的大数据解决方案正以前所未见的力量推进着商业洞见的形成，带来更多更好的洞见以加速决策制定。

在过去几年中，我的部门有机会与本书的作者合作，向 Apache Spark 社区贡献成果，并在英特尔架构上优化各种大数据和 Spark 应用。《Spark 快速大数据分析》的出版为开发者和数据科学家提供了丰富的 Spark 知识。更重要的是，这本书不是简单地教开发者如何使用 Spark，而是更深入介绍了 Spark 的内部构成，并通过各种实例展示了如何优化大数据应用。我向大家推荐这本书，或更具体点，推荐这本书里提倡的优化方法和思路，相信它们能帮助你创建出更好的大数据应用。

英特尔软件服务事业部全球大数据技术中心总经理 马子雅

2015 年 7 月于加州圣克拉拉

Big data is getting hot in recent years. Quite often, folks ask why big data is a big deal. We are in the era of data explosion, with the emergence of smart phones, tablets, wearables, IoT devices, etc. Ninety percent of the data in the world today was generated in just the past two years. By 2020, we will see >50B devices connected and Zeta byte data created. It is not the quantity of the data that is revolutionary. It is that we can now do something with it that's revolutionary. The power of big data solutions is they can process large and complex data sets very fast, generate better and faster insights than conventional methods.

A big data solution suite can consist of several critical components, from the hardware layer like storage, compute and network, to data processing engine, to analytics layer where business insights are generated using improved statistical & computational algorithms and data visualization. Among all, the data processing engine is one most critical player. It is not overstating that the data processing engine for big data is like CPU for a computer or brain for a human being.

Spark was initially started for the purpose of creating a big data processing and computing framework, when Matei Zaharia was doing his Ph.D. research at UC Berkeley AMPLab in 2009. Different from the traditional data processing framework, Spark's in-memory primitives provide performance up to 100 times faster for certain applications. By allowing user programs to load data into a cluster's memory and query it repeatedly, Spark is well-suited for big data and machine learning use cases. Spark is becoming one best adopted among all big data modules. Big Data Distributions like Cloudera, MapR now all include Spark into their distributions.

Spark is now evolving the Hadoop and big data ecosystem to better support the end-to-end big data analytics needs, e.g. Spark grew beyond Spark core to Spark streaming, SQL, MLlib, GraphX, SparkR, etc. Learning Spark and its internals will not just help improve the processing speed for big data, but also help developers and data scientists create analytics applications with more ease. With big data solutions like Spark, we expect to see significant improvement with business insights which will help expedite the decision making—like we've never seen before, from enterprise, healthcare, transportation, and retail.

Over the years, my organization had the opportunities to work with authors of this book, contribute to Apache Spark, and optimize various Big Data and Spark application on Intel Architecture. The publication of Learning Spark offers developers and data scientists extensive knowledge on Spark. Moreover, Learning Spark does not simply try to tell the developers how to use Spark, it also addresses the internals and shows various examples of how to improve your big data applications. I recommend Learning Spark—that this book, and, more specifically, the method it espouses, will change your big data application for the better.

Ziya Ma, General Manager of the global Big Data Technologies organization,

SSG STO, Intel Corp.

Santa Clara, California, July 2015

欢迎访问：电子书学习和下载网站 (<https://www.shgis.cn>)

文档名称：《Spark快速大数据分析》 [美] 卡劳 (Holden Karau) 等. epub

请登录 <https://shgis.cn/post/1851.html> 下载完整文档。

手机端请扫码查看：

