



# 大数据 商业实战三部曲

---

内核解密 | 商业案例 | 性能调优

---

王家林 段智华 夏 阳◎编著

清华大学出版社



# 大数据 商业实战三部曲

内核解密 | 商业案例 | 性能调优

王家林 段智华 夏 阳◎编著

清华大学出版社  
北京

## 内 容 简 介

本书基于 Spark 2.2.X 最新版本，以 Spark 商业案例实战和 Spark 在生产环境下几乎所有类型的性能调优为核心，以 Spark 内核解密为基石，分为上篇、中篇、下篇，对企业生产环境下的 Spark 商业案例与性能调优抽丝剥茧地进行剖析。上篇基于 Spark 源码，从一个动手实战案例入手，循序渐进地全面解析了 Spark 2.2.X 新特性及 Spark 内核源码；中篇选取 Spark 开发中最具有代表的经典学习案例，深入浅出地介绍，在案例中综合应用 Spark 的大数据技术；下篇性能调优内容基本完全覆盖了 Spark 在生产环境下的所有调优技术。

本书适合所有 Spark 学习者和从业人员使用。对于有分布式计算框架应用经验的人员，本书也可以作为 Spark 高手修炼的参考书籍。同时，本书也特别适合作为高等院校的大数据教材使用。

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

### 图书在版编目（CIP）数据

Spark 大数据商业实战三部曲：内核解密|商业案例|性能调优 / 王家林，段智华，夏阳编著. —北京：清华大学出版社，2018

ISBN 978-7-302-48962-7

I. ①S… II. ①王… ②段… ③夏… III. ①数据处理 IV. ①TP274

中国版本图书馆 CIP 数据核字（2017）第 287681 号

责任编辑：袁金敏 常建丽

封面设计：刘新新

责任校对：徐俊伟

责任印制：沈 露

出版发行：清华大学出版社

网 址：<http://www.tup.com.cn>, <http://www.wqbook.com>

地 址：北京清华大学学研大厦 A 座 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈：010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 装 者：三河市金元印装有限公司

经 销：全国新华书店

开 本：185mm×260mm 印 张：72.75 字 数：1815 千字

版 次：2018 年 2 月第 1 版 印 次：2018 年 2 月第 1 次印刷

定 价：299.00 元

---

产品编号：075671-01

# 前　　言

大数据像当年的石油、人工智能（Artificial Intelligence）像当年的电力一样，正以前所未有的广度和深度影响所有的行业，现在及未来公司的核心壁垒是数据，核心竞争力来自基于大数据的人工智能的竞争。Spark 是当今大数据领域最活跃、最热门、最高效的大数据通用计算平台，2009 年诞生于美国加州大学伯克利分校 AMP 实验室，2010 年正式开源，2013 年成为 Apache 基金项目，2014 年成为 Apache 基金的顶级项目。基于 RDD，Spark 成功构建起了一体化、多元化的大数据处理体系。

在任何规模的数据计算中，Spark 在性能和扩展性上都更具优势。

(1) Hadoop 之父 Doug Cutting 指出：Use of MapReduce engine for Big Data projects will decline, replaced by Apache Spark (大数据项目的 MapReduce 引擎的使用将下降，由 Apache Spark 取代。)

(2) Hadoop 商业发行版本的市场领导者 Cloudera、HortonWorks、MapR 纷纷转投 Spark，并把 Spark 作为大数据解决方案的首选和核心计算引擎。

2014 年的 Sort Benchmark 测试中，Spark 秒杀 Hadoop，在使用十分之一计算资源的情况下，相同数据的排序上，Spark 比 MapReduce 快 3 倍！在没有官方 PB 排序对比的情况下，首次将 Spark 推到了 1PB 数据（十万亿条记录）的排序，在使用 190 个节点的情况下，工作负载在 4 小时内完成，同样远超雅虎之前使用 3800 台主机耗时 16 个小时的记录。

2015 年 6 月，Spark 最大的集群来自腾讯——8000 个节点，单个 Job 最大分别是阿里巴巴和 Databricks——1PB，震撼人心！同时，Spark 的 Contributor 比 2014 年涨了 3 倍，达到 730 人；总代码行数也比 2014 年涨了 2 倍多，达到 40 万行。IBM 于 2015 年 6 月承诺大力推进 Apache Spark 项目，并称该项目为：以数据为主导的，未来十年最重要的新的开源项目。这一承诺的核心是将 Spark 嵌入 IBM 业内领先的分析和商务平台，并将 Spark 作为一项服务，在 IBMBluemix 平台上提供给客户。IBM 还将投入超过 3500 名研究和开发人员在全球 10 余个实验室开展与 Spark 相关的项目，并将为 Spark 开源生态系统无偿提供突破性的机器学习技术——IBM SystemML。同时，IBM 还将培养超过 100 万名 Spark 数据科学家和数据工程师。

2016 年，在有“计算界奥运会”之称的国际著名 Sort Benchmark 全球数据排序大赛中，由南京大学计算机科学与技术系 PASA 大数据实验室、阿里巴巴和 Databricks 公司组成的参赛团队 NADSort，以 144 美元的成本完成 100TB 标准数据集的排序处理，创下了每 TB 数据排序 1.44 美元成本的最新世界纪录，比 2014 年夺得冠军的加州大学圣地亚哥分校 TritonSort 团队每 TB 数据 4.51 美元的成本降低了近 70%，而这次比赛依旧使用 Apache Spark 大数据计算平台，在大规模并行排序算法以及 Spark 系统底层进行了大量的优化，以尽可能提高排序计算性能并降低存储资源开销，确保最终赢得比赛。

在 Full Stack 理想的指引下，Spark 中的 Spark SQL、SparkStreaming、MLLib、GraphX、R 五大子框架和库之间可以无缝地共享数据和操作，这不仅打造了 Spark 在当今大数据计算领域其他计算框架都无可匹敌的优势，而且使得 Spark 正在加速成为大数据处理中心首选通

用计算平台，而 Spark 商业案例和性能优化必将成为接下来的重中之重！

本书根据王家林老师亲授课程及结合众多大数据项目经验编写而成，其中王家林、段智华编写了本书近 90% 的内容，具体编写章节如下：

- 第 3 章 Spark 的灵魂：RDD 和 DataSet；
- 第 4 章 Spark Driver 启动内幕剖析；
- 第 5 章 Spark 集群启动原理和源码详解；
- 第 6 章 Spark Application 提交给集群的原理和源码详解；
- 第 7 章 Shuffle 原理和源码详解；
- 第 8 章 Job 工作原理和源码详解；
- 第 9 章 Spark 中 Cache 和 checkpoint 原理和源码详解；
- 第 10 章 Spark 中 Broadcast 和 Accumulator 原理和源码详解；
- 第 11 章 Spark 与大数据其他经典组件整合原理与实战；
- 第 12 章 Spark 商业案例之大数据电影点评系统应用案例；
- 第 13 章 Spark 2.2 实战之 Dataset 开发实战企业人员管理系统应用案例；
- 第 14 章 Spark 商业案例之电商交互式分析系统应用案例；
- 第 15 章 Spark 商业案例之 NBA 篮球运动员大数据分析系统应用案例；
- 第 16 章 电商广告点击大数据实时流处理系统案例；
- 第 17 章 Spark 在通信运营商生产环境中的应用案例；
- 第 18 章 使用 Spark GraphX 实现婚恋社交网络多维度分析案例；
- 第 23 章 Spark 集群中 Mapper 端、Reducer 端内存调优；
- 第 24 章 使用 Broadcast 实现 Mapper 端 Shuffle 聚合功能的原理和调优实战；
- 第 25 章 使用 Accumulator 高效地实现分布式集群全局计数器的原理和调优案例；
- 第 27 章 Spark 五大子框架调优最佳实践；
- 第 28 章 Spark 2.2.0 新一代鸽丝计划优化引擎；
- 第 30 章 Spark 性能调优之数据倾斜调优一站式解决方案原理与实战；
- 第 31 章 Spark 大数据性能调优实战专业之路。

其中，段智华根据自身多年的大数据工作经验对本书的案例等部分进行了扩展。

除上述章节外，剩余内容由夏阳、郑采翔、闫恒伟三位作者根据王家林老师的大数据授课内容而完成。

在阅读本书的过程中，如发现任何问题或有任何疑问，可以加入本书的阅读群（QQ：418110145）讨论，会有专人答疑。同时，该群也会提供本书所用案例源码及本书的配套学习视频。

如果读者想要了解或者学习更多大数据相关技术，可以关注 DT 大数据梦工厂微信公众号 DT\_Spark，也可以通过 YY 客户端登录 68917580 永久频道直接体验。

王家林老师的新浪微博是 <http://weibo.com/ilovepains/> 欢欢迎大家在微博上与作者进行互动。

由于时间仓促，书中难免存在不妥之处，请读者谅解，并提出宝贵意见。

王家林 2017 年中秋之夜于美国硅谷

# 目 录

## 上篇 内核解密

第 1 章 电光石火间体验 Spark 2.2 开发实战	2
1.1 通过 RDD 实战电影点评系统入门及源码阅读	2
1.1.1 Spark 核心概念图解	2
1.1.2 通过 RDD 实战电影点评系统案例	4
1.2 通过 DataFrame 和 DataSet 实战电影点评系统	7
1.2.1 通过 DataFrame 实战电影点评系统案例	7
1.2.2 通过 DataSet 实战电影点评系统案例	10
1.3 Spark 2.2 源码阅读环境搭建及源码阅读体验	11
第 2 章 Spark 2.2 技术及原理	14
2.1 Spark 2.2 综述	14
2.1.1 连续应用程序	14
2.1.2 新的 API	15
2.2 Spark 2.2 Core	16
2.2.1 第二代 Tungsten 引擎	16
2.2.2 SparkSession	16
2.2.3 累加器 API	17
2.3 Spark 2.2 SQL	19
2.3.1 Spark SQL	20
2.3.2 DataFrame 和 Dataset API	20
2.3.3 Timed Window	21
2.4 Spark 2.2 Streaming	21
2.4.1 Structured Streaming	21
2.4.2 增量输出模式	23
2.5 Spark 2.2 MLlib	27
2.5.1 基于 DataFrame 的 Machine Learning API	28
2.5.2 R 的分布式算法	28
2.6 Spark 2.2 GraphX	29
第 3 章 Spark 的灵魂：RDD 和 DataSet	30
3.1 为什么说 RDD 和 DataSet 是 Spark 的灵魂	30
3.1.1 RDD 的定义及五大特性剖析	30

3.1.2 DataSet 的定义及内部机制剖析 .....	34
3.2 RDD 弹性特性七个方面解析 .....	36
3.3 RDD 依赖关系 .....	43
3.3.1 窄依赖解析 .....	43
3.3.2 宽依赖解析 .....	45
3.4 解析 Spark 中的 DAG 逻辑视图 .....	46
3.4.1 DAG 生成的机制 .....	46
3.4.2 DAG 逻辑视图解析 .....	47
3.5 RDD 内部的计算机制 .....	49
3.5.1 Task 解析 .....	49
3.5.2 计算过程深度解析 .....	49
3.6 Spark RDD 容错原理及其四大核心要点解析 .....	57
3.6.1 Spark RDD 容错原理 .....	57
3.6.2 RDD 容错的四大核心要点 .....	57
3.7 Spark RDD 中 Runtime 流程解析 .....	59
3.7.1 Runtime 架构图 .....	59
3.7.2 生命周期 .....	60
3.8 通过 WordCount 实战解析 Spark RDD 内部机制 .....	70
3.8.1 Spark WordCount 动手实践 .....	70
3.8.2 解析 RDD 生成的内部机制 .....	72
3.9 基于 DataSet 的代码到底是如何一步步转化成为 RDD 的 .....	78
<b>第 4 章 Spark Driver 启动内幕剖析 .....</b>	<b>81</b>
4.1 Spark Driver Program 剖析 .....	81
4.1.1 Spark Driver Program .....	81
4.1.2 SparkContext 深度剖析 .....	81
4.1.3 SparkContext 源码解析 .....	82
4.2 DAGScheduler 解析 .....	96
4.2.1 DAG 的定义 .....	96
4.2.2 DAG 的实例化 .....	97
4.2.3 DAGScheduler 划分 Stage 的原理 .....	98
4.2.4 DAGScheduler 划分 Stage 的具体算法 .....	99
4.2.5 Stage 内部 Task 获取最佳位置的算法 .....	113
4.3 TaskScheduler 解析 .....	116
4.3.1 TaskScheduler 原理剖析 .....	116
4.3.2 TaskScheduler 源码解析 .....	117
4.4 SchedulerBackend 解析 .....	132
4.4.1 SchedulerBackend 原理剖析 .....	132
4.4.2 SchedulerBackend 源码解析 .....	132
4.4.3 Spark 程序的注册机制 .....	133

---

4.4.4	Spark 程序对计算资源 Executor 的管理 .....	134
4.5	打通 Spark 系统运行内幕机制循环流程 .....	135
4.6	本章总结 .....	145
<b>第 5 章</b>	<b>Spark 集群启动原理和源码详解 .....</b>	<b>146</b>
5.1	Master 启动原理和源码详解 .....	146
5.1.1	Master 启动的原理详解 .....	146
5.1.2	Master 启动的源码详解 .....	147
5.1.3	Master HA 双机切换 .....	157
5.1.4	Master 的注册机制和状态管理解密 .....	163
5.2	Worker 启动原理和源码详解 .....	170
5.2.1	Worker 启动的原理流程 .....	170
5.2.2	Worker 启动的源码详解 .....	174
5.3	ExecutorBackend 启动原理和源码详解 .....	178
5.3.1	ExecutorBackend 接口与 Executor 的关系 .....	178
5.3.2	ExecutorBackend 的不同实现 .....	179
5.3.3	ExecutorBackend 中的通信 .....	181
5.3.4	ExecutorBackend 的异常处理 .....	183
5.4	Executor 中任务的执行 .....	184
5.4.1	Executor 中任务的加载 .....	184
5.4.2	Executor 中的任务线程池 .....	185
5.4.3	任务执行失败处理 .....	186
5.4.4	揭秘 TaskRunner .....	188
5.5	Executor 执行结果的处理方式 .....	189
5.6	本章总结 .....	197
<b>第 6 章</b>	<b>Spark Application 提交给集群的原理和源码详解 .....</b>	<b>198</b>
6.1	Spark Application 到底是如何提交给集群的 .....	198
6.1.1	Application 提交参数配置详解 .....	198
6.1.2	Application 提交给集群原理详解 .....	199
6.1.3	Application 提交给集群源码详解 .....	201
6.2	Spark Application 是如何向集群申请资源的 .....	211
6.2.1	Application 申请资源的两种类型详解 .....	211
6.2.2	Application 申请资源的源码详解 .....	213
6.3	从 Application 提交的角度重新审视 Driver .....	219
6.3.1	Driver 到底是什么时候产生的 .....	220
6.3.2	Driver 和 Master 交互原理解析 .....	238
6.3.3	Driver 和 Master 交互源码详解 .....	244
6.4	从 Application 提交的角度重新审视 Executor .....	249
6.4.1	Executor 到底是什么时候启动的 .....	249
6.4.2	Executor 如何把结果交给 Application .....	254

6.5	Spark 1.6 RPC 内幕解密：运行机制、源码详解、Netty 与 Akka 等.....	254
6.6	本章总结.....	267
<b>第 7 章</b>	<b>Shuffle 原理和源码详解.....</b>	<b>268</b>
7.1	概述.....	268
7.2	Shuffle 的框架.....	269
7.2.1	Shuffle 的框架演进.....	269
7.2.2	Shuffle 的框架内核.....	270
7.2.3	Shuffle 框架的源码解析.....	272
7.2.4	Shuffle 数据读写的源码解析.....	275
7.3	Hash Based Shuffle.....	281
7.3.1	概述.....	281
7.3.2	Hash Based Shuffle 内核.....	282
7.3.3	Hash Based Shuffle 数据读写的源码解析.....	285
7.4	Sorted Based Shuffle.....	290
7.4.1	概述.....	292
7.4.2	Sorted Based Shuffle 内核.....	293
7.4.3	Sorted Based Shuffle 数据读写的源码解析.....	294
7.5	Tungsten Sorted Based Shuffle.....	302
7.5.1	概述.....	302
7.5.2	Tungsten Sorted Based Shuffle 内核.....	302
7.5.3	Tungsten Sorted Based Shuffle 数据读写的源码解析.....	303
7.6	Shuffle 与 Storage 模块间的交互.....	309
7.6.1	Shuffle 注册的交互.....	310
7.6.2	Shuffle 写数据的交互.....	314
7.6.3	Shuffle 读数据的交互.....	315
7.6.4	BlockManager 架构原理、运行流程图和源码解密.....	315
7.6.5	BlockManager 解密进阶：BlockManager 初始化和注册解密、BlockManager-Master 工作解密、BlockTransferService 解密、本地数据读写解密、远程数据读写解密.....	324
7.7	本章总结.....	341
<b>第 8 章</b>	<b>Job 工作原理和源码详解.....</b>	<b>342</b>
8.1	Job 到底在什么时候产生.....	342
8.1.1	触发 Job 的原理和源码解析.....	342
8.1.2	触发 Job 的算子案例.....	344
8.2	Stage 划分内幕.....	345
8.2.1	Stage 划分原理详解.....	345
8.2.2	Stage 划分源码详解.....	346
8.3	Task 全生命周期详解.....	346
8.3.1	Task 的生命过程详解.....	347

8.3.2 Task 在 Driver 和 Executor 中交互的全生命周期原理和源码详解.....	348
8.4 ShuffleMapTask 和 ResultTask 处理结果是如何被 Driver 管理的.....	364
8.4.1 ShuffleMapTask 执行结果和 Driver 的交互原理及源码详解.....	364
8.4.2 ResultTask 执行结果与 Driver 的交互原理及源码详解.....	370
<b>第 9 章 Spark 中 Cache 和 checkpoint 原理和源码详解 .....</b>	<b>372</b>
9.1 Spark 中 Cache 原理和源码详解.....	372
9.1.1 Spark 中 Cache 原理详解.....	372
9.1.2 Spark 中 Cache 源码详解.....	372
9.2 Spark 中 checkpoint 原理和源码详解 .....	381
9.2.1 Spark 中 checkpoint 原理详解 .....	381
9.2.2 Spark 中 checkpoint 源码详解 .....	381
<b>第 10 章 Spark 中 Broadcast 和 Accumulator 原理和源码详解 .....</b>	<b>391</b>
10.1 Spark 中 Broadcast 原理和源码详解.....	391
10.1.1 Spark 中 Broadcast 原理详解.....	391
10.1.2 Spark 中 Broadcast 源码详解.....	393
10.2 Spark 中 Accumulator 原理和源码详解 .....	396
10.2.1 Spark 中 Accumulator 原理详解 .....	396
10.2.2 Spark 中 Accumulator 源码详解 .....	396
<b>第 11 章 Spark 与大数据其他经典组件整合原理与实战 .....</b>	<b>399</b>
11.1 Spark 组件综合应用.....	399
11.2 Spark 与 Alluxio 整合原理与实战.....	400
11.2.1 Spark 与 Alluxio 整合原理 .....	400
11.2.2 Spark 与 Alluxio 整合实战 .....	401
11.3 Spark 与 Job Server 整合原理与实战 .....	403
11.3.1 Spark 与 Job Server 整合原理 .....	403
11.3.2 Spark 与 Job Server 整合实战 .....	404
11.4 Spark 与 Redis 整合原理与实战 .....	406
11.4.1 Spark 与 Redis 整合原理 .....	406
11.4.2 Spark 与 Redis 整合实战 .....	407
<b>中篇 商业案例</b>	
<b>第 12 章 Spark 商业案例之大数据电影点评系统应用案例 .....</b>	<b>412</b>
12.1 通过 RDD 实现分析电影的用户行为信息.....	412
12.1.1 搭建 IDEA 开发环境 .....	412
12.1.2 大数据电影点评系统中电影数据说明 .....	425
12.1.3 电影点评系统用户行为分析统计实战 .....	428
12.2 通过 RDD 实现电影流行度分析.....	431

12.3	通过 RDD 分析各种类型的最喜爱电影 TopN 及性能优化技巧.....	433
12.4	通过 RDD 分析电影点评系统仿 QQ 和微信等用户群分析及广播 背后机制解密.....	436
12.5	通过 RDD 分析电影点评系统实现 Java 和 Scala 版本的二次排序系统.....	439
12.5.1	二次排序自定义 Key 值类实现（Java）.....	440
12.5.2	电影点评系统二次排序功能实现（Java）.....	442
12.5.3	二次排序自定义 Key 值类实现（Scala）.....	445
12.5.4	电影点评系统二次排序功能实现（Scala）.....	446
12.6	通过 Spark SQL 中的 SQL 语句实现电影点评系统用户行为分析.....	447
12.7	通过 Spark SQL 下的两种不同方式实现口碑最佳电影分析.....	451
12.8	通过 Spark SQL 下的两种不同方式实现最流行电影分析.....	456
12.9	通过 DataFrame 分析最受男性和女性喜爱电影 TopN .....	457
12.10	纯粹通过 DataFrame 分析电影点评系统仿 QQ 和微信、淘宝等用户群.....	460
12.11	纯粹通过 DataSet 对电影点评系统进行流行度和不同年龄阶段兴趣分析等.....	462
12.11.1	通过 DataSet 实现某特定电影观看者中男性和女性不同年龄的人数.....	463
12.11.2	通过 DataSet 方式计算所有电影中平均得分最高 (口碑最好) 的电影 TopN.....	464
12.11.3	通过 DataSet 方式计算所有电影中粉丝或者观看人数最多 (最流行电影) 的电影 TopN.....	465
12.11.4	纯粹通过 DataSet 的方式实现所有电影中最受男性、女性喜爱的 电影 Top10.....	466
12.11.5	纯粹通过 DataSet 的方式实现所有电影中 QQ 或者微信核心目标 用户最喜爱电影 TopN 分析 .....	467
12.11.6	纯粹通过 DataSet 的方式实现所有电影中淘宝核心目标用户最喜爱电影 TopN 分析 .....	469
12.12	大数据电影点评系统应用案例涉及的核心知识点原理、源码及案例代码 .....	470
12.12.1	知识点：广播变量 Broadcast 内幕机制 .....	470
12.12.2	知识点：SQL 全局临时视图及临时视图 .....	473
12.12.3	大数据电影点评系统应用案例完整代码 .....	474
12.13	本章总结 .....	496
<b>第 13 章</b>	<b>Spark 2.2 实战之 Dataset 开发实战企业人员管理系统应用案例 .....</b>	<b>498</b>
13.1	企业人员管理系统应用案例业务需求分析 .....	498
13.2	企业人员管理系统应用案例数据建模 .....	499
13.3	通过 SparkSession 创建案例开发实战上下文环境 .....	500
13.3.1	Spark 1.6.0 版本 SparkContext .....	500
13.3.2	Spark 2.0.0 版本 SparkSession .....	501
13.3.3	DataFrame、DataSet 剖析与实战 .....	507
13.4	通过 map、flatMap、mapPartitions 等分析企业人员管理系统 .....	510
13.5	通过 dropDuplicate、coalesce、repartition 等分析企业人员管理系统 .....	512

13.6 通过 sort、join、joinWith 等分析企业人员管理系统.....	514
13.7 通过 randomSplit、sample、select 等分析企业人员管理系统.....	515
13.8 通过 groupBy、agg、col 等分析企业人员管理系统.....	517
13.9 通过 collect_list、collect_set 等分析企业人员管理系统.....	518
13.10 通过 avg、sum、countDistinct 等分析企业人员管理系统.....	519
13.11 Dataset 开发实战企业人员管理系统应用案例代码 .....	519
13.12 本章总结.....	522
<b>第 14 章 Spark 商业案例之电商交互式分析系统应用案例.....</b>	<b>523</b>
14.1 纯粹通过 DataSet 进行电商交互式分析系统中特定时段访问次数 TopN.....	523
14.1.1 电商交互式分析系统数据说明 .....	523
14.1.2 特定时段内用户访问电商网站排名 TopN.....	525
14.2 纯粹通过 DataSet 分析特定时段购买金额 Top10 和访问次数增长 Top10 .....	527
14.3 纯粹通过 DataSet 进行电商交互式分析系统中各种类型 TopN 分析实战详解.....	530
14.3.1 统计特定时段购买金额最多的 Top5 用户.....	530
14.3.2 统计特定时段访问次数增长最多的 Top5 用户.....	530
14.3.3 统计特定时段购买金额增长最多的 Top 5 用户.....	531
14.3.4 统计特定时段注册之后前两周内访问次数最多的 Top 10 用户 .....	533
14.3.5 统计特定时段注册之后前两周内购买总额最多的 Top 10 用户 .....	534
14.4 电商交互式分析系统应用案例涉及的核心知识点原理、源码及案例代码 .....	535
14.4.1 知识点：Functions.scala .....	535
14.4.2 电商交互式分析系统应用案例完整代码 .....	548
14.5 本章总结.....	555
<b>第 15 章 Spark 商业案例之 NBA 篮球运动员大数据分析系统应用案例.....</b>	<b>556</b>
15.1 NBA 篮球运动员大数据分析系统架构和实现思路.....	556
15.2 NBA 篮球运动员大数据分析系统代码实战：数据清洗和初步处理.....	561
15.3 NBA 篮球运动员大数据分析代码实战之核心基础数据项编写 .....	565
15.3.1 NBA 球员数据每年基础数据项记录 .....	565
15.3.2 NBA 球员数据每年标准分 Z-Score 计算 .....	567
15.3.3 NBA 球员数据每年归一化计算 .....	568
15.3.4 NBA 历年比赛数据按球员分组统计分析 .....	572
15.3.5 NBA 球员年龄值及经验值列表获取 .....	575
15.3.6 NBA 球员年龄值及经验值统计分析 .....	576
15.3.7 NBA 球员系统内部定义的函数、辅助工具类 .....	578
15.4 NBA 篮球运动员大数据分析完整代码测试和实战 .....	582
15.5 NBA 篮球运动员大数据分析系统应用案例涉及的核心知识点、原理、源码 .....	594
15.5.1 知识点：StatCounter 源码分析 .....	594
15.5.2 知识点：StatCounter 应用案例 .....	598
15.6 本章总结.....	601

第 16 章 电商广告点击大数据实时流处理系统案例 .....	602
16.1 电商广告点击综合案例需求分析和技术架构 .....	602
16.1.1 电商广告点击综合案例需求分析 .....	602
16.1.2 电商广告点击综合案例技术架构 .....	603
16.1.3 电商广告点击综合案例整体部署 .....	606
16.1.4 生产数据业务流程及消费数据业务流程 .....	607
16.1.5 Spark JavaStreamingContext 初始化及启动 .....	607
16.1.6 Spark Streaming 使用 No Receivers 方式读取 Kafka 数据及监控 .....	609
16.2 电商广告点击综合案例在线点击统计实战 .....	612
16.3 电商广告点击综合案例黑名单过滤实现 .....	615
16.3.1 基于用户广告点击数据表，动态过滤黑名单用户 .....	616
16.3.2 黑名单的整个 RDD 进行去重操作 .....	617
16.3.3 将黑名单写入到黑名单数据表 .....	618
16.4 电商广告点击综合案例底层数据层的建模和编码实现（基于 MySQL） .....	618
16.4.1 电商广告点击综合案例数据库链接单例模式实现 .....	619
16.4.2 电商广告点击综合案例数据库操作实现 .....	622
16.5 电商广告点击综合案例动态黑名单过滤真正的实现代码 .....	624
16.5.1 从数据库中获取黑名单封装成 RDD .....	624
16.5.2 黑名单 RDD 和批处理 RDD 进行左关联，过滤掉黑名单 .....	625
16.6 动态黑名单基于数据库 MySQL 的真正操作代码实战 .....	627
16.6.1 MySQL 数据库操作的架构分析 .....	627
16.6.2 MySQL 数据库操作的代码实战 .....	628
16.7 通过 updateStateByKey 等实现广告点击流量的在线更新统计 .....	634
16.8 实现每个省份点击排名 Top5 广告 .....	639
16.9 实现广告点击 Trend 趋势计算实战 .....	643
16.10 实战模拟点击数据的生成和数据表 SQL 的建立 .....	648
16.10.1 电商广告点击综合案例模拟数据的生成 .....	648
16.10.2 电商广告点击综合案例数据表 SQL 的建立 .....	651
16.11 电商广告点击综合案例运行结果 .....	654
16.11.1 电商广告点击综合案例 Hadoop 集群启动 .....	654
16.11.2 电商广告点击综合案例 Spark 集群启动 .....	655
16.11.3 电商广告点击综合案例 Zookeeper 集群启动 .....	656
16.11.4 电商广告点击综合案例 Kafka 集群启动 .....	658
16.11.5 电商广告点击综合案例 Hive metastore 集群启动 .....	660
16.11.6 电商广告点击综合案例程序运行 .....	660
16.11.7 电商广告点击综合案例运行结果 .....	661
16.12 电商广告点击综合案例 Scala 版本关注点 .....	663
16.13 电商广告点击综合案例课程的 Java 源码 .....	666
16.14 电商广告点击综合案例课程的 Scala 源码 .....	694
16.15 本章总结 .....	711

<b>第 17 章 Spark 在通信运营商生产环境中的应用案例</b>	712
17.1 Spark 在通信运营商融合支付系统日志统计分析中的综合应用案例	712
17.1.1 融合支付系统日志统计分析综合案例需求分析	712
17.1.2 融合支付系统日志统计分析数据说明	714
17.1.3 融合支付系统日志清洗中 Scala 正则表达式与模式匹配结合的代码实战	718
17.1.4 融合支付系统日志在大数据 Splunk 中的可视化展示	722
17.1.5 融合支付系统日志统计分析案例涉及的正则表达式知识点及案例代码	733
17.2 Spark 在光宽用户流量热力分布 GIS 系统中的综合应用案例	742
17.2.1 光宽用户流量热力分布 GIS 系统案例需求分析	742
17.2.2 光宽用户流量热力分布 GIS 应用的数据说明	742
17.2.3 光宽用户流量热力分布 GIS 应用 Spark 实战	744
17.2.4 光宽用户流量热力分布 GIS 应用 Spark 实战成果	748
17.2.5 光宽用户流量热力分布 GIS 应用 Spark 案例代码	749
17.3 本章总结	752
<b>第 18 章 使用 Spark GraphX 实现婚恋社交网络多维度分析案例</b>	753
18.1 Spark GraphX 发展演变历史和在业界的使用案例	753
18.2 Spark GraphX 设计实现的核心原理	757
18.3 Table operator 和 Graph Operator	760
18.4 Vertices、edges、triplets	762
18.5 以最原始的方式构建 Graph	765
18.6 第一个 Graph 代码实例并进行 Vertices、edges、triplets 操作实战	765
18.7 数据加载成为 Graph 并进行操作实战	775
18.8 图操作之 Property Operators 实战	782
18.9 图操作之 Structural Operators 实战	784
18.10 图操作之 Computing Degree 实战	788
18.11 图操作之 Collecting Neighbors 实战	791
18.12 图操作之 Join Operators 实战	793
18.13 图操作之 aggregateMessages 实战	796
18.14 图算法之 Pregel API 原理解析与实战	799
18.15 图算法之 ShortestPaths 原理解析与实战	804
18.16 图算法之 PageRank 原理解析与实战	805
18.17 图算法之 TriangleCount 原理解析与实战	807
18.18 使用 Spark GraphX 实现婚恋社交网络多维度分析实战	809
18.18.1 婚恋社交网络多维度分析实战图的属性演示	811
18.18.2 婚恋社交网络多维度分析实战图的转换操作	814
18.18.3 婚恋社交网络多维度分析实战图的结构操作	815
18.18.4 婚恋社交网络多维度分析实战图的连接操作	816

18.18.5	婚恋社交网络多维度分析实战图的聚合操作	818
18.18.6	婚恋社交网络多维度分析实战图的实用操作	822
18.19	婚恋社交网络多维度分析案例代码	823
18.20	本章总结	832

## 下篇 性能调优

<b>第 19 章</b>	<b>对运行在 YARN 上的 Spark 进行性能调优</b>	834
19.1	运行环境 Jar 包管理及数据本地性原理调优实践	834
19.1.1	运行环境 Jar 包管理及数据本地性原理	834
19.1.2	运行环境 Jar 包管理及数据本地性调优实践	835
19.2	Spark on YARN 两种不同的调度模型及其调优	836
19.2.1	Spark on YARN 的两种不同类型模型优劣分析	836
19.2.2	Spark on YARN 的两种不同类型调优实践	837
19.3	YARN 队列资源不足引起的 Spark 应用程序失败的原因及调优方案	838
19.3.1	失败的原因剖析	838
19.3.2	调优方案	838
19.4	Spark on YARN 模式下 Executor 经常被杀死的原因及调优方案	838
19.4.1	原因剖析	838
19.4.2	调优方案	839
19.5	YARN-Client 模式下网卡流量激增的原因及调优方案	839
19.5.1	原因剖析	839
19.5.2	调优方案	840
19.6	YARN-Cluster 模式下 JVM 栈内存溢出的原因及调优方案	840
19.6.1	原因剖析	841
19.6.2	调优方案	841
<b>第 20 章</b>	<b>Spark 算子调优最佳实践</b>	842
20.1	使用 mapPartitions 或者 mapPartitionWithIndex 取代 map 操作	842
20.1.1	mapPartitions 内部工作机制和源码解析	842
20.1.2	mapPartitionWithIndex 内部工作机制和源码解析	842
20.1.3	使用 mapPartitions 取代 map 案例和性能测试	843
20.2	使用 foreachPartition 把 Spark 数据持久化到外部存储介质	844
20.2.1	foreachPartition 内部工作机制和源码解析	844
20.2.2	使用 foreachPartition 写数据到 MySQL 中案例和性能测试	845
20.3	使用 coalesce 取代 repartition 操作	845
20.3.1	coalesce 和 repartition 工作机制和源码剖析	845
20.3.2	通过测试对比 coalesce 和 repartition 的性能	847
20.4	使用 repartitionAndSortWithinPartitions 取代 repartition 和 sort 的联合操作	847
20.4.1	repartitionAndSortWithinPartitions 工作机制和源码解析	847

20.4.2 repartitionAndSortWithinPartitions 性能测试	848
20.5 使用 treeReduce 取代 reduce 的原理和源码	848
20.5.1 treeReduce 进行 reduce 的工作原理和源码	848
20.5.2 使用 treeReduce 进行性能测试	849
20.6 使用 treeAggregate 取代 Aggregate 的原理和源码	851
20.6.1 treeAggregate 进行 Aggregate 的工作原理和源码	851
20.6.2 使用 treeAggregate 进行性能测试	852
20.7 reduceByKey 高效运行的原理和源码解密	853
20.8 使用 AggregateByKey 取代 groupByKey 的原理和源码	857
20.8.1 使用 AggregateByKey 取代 groupByKey 的原理	857
20.8.2 源码剖析	858
20.8.3 使用 AggregateByKey 取代 groupByKey 进行性能测试	859
20.9 Join 不产生 Shuffle 的情况及案例实战	860
20.9.1 Join 在什么情况下不产生 Shuffle 及其运行原理	860
20.9.2 Join 不产生 Shuffle 的情况案例实战	860
20.10 RDD 复用性能调优最佳实践	861
20.10.1 什么时候需要复用 RDD	861
20.10.2 如何复用 RDD 算子	862
<b>第 21 章 Spark 频繁遇到的性能问题及调优技巧</b>	<b>864</b>
21.1 使用 BroadCast 广播大变量和业务配置信息原理 和案例实战	864
21.1.1 使用 BroadCast 广播大变量和业务配置信息原理	864
21.1.2 使用 BroadCast 广播大变量和业务配置信息案例实战	865
21.2 使用 Kryo 取代 Scala 默认的序列器原理和案例实战	865
21.2.1 使用 Kryo 取代 Scala 默认的序列器原理	865
21.2.2 使用 Kryo 取代 Scala 默认的序列器案例实战	866
21.3 使用 FastUtil 优化 JVM 数据格式解析和案例实战	866
21.3.1 使用 FastUtil 优化 JVM 数据格式解析	866
21.3.2 使用 FastUtil 优化 JVM 数据格式案例实战	867
21.4 Persist 及 checkpoint 使用时的正误方式	868
21.5 序列化导致的报错原因解析和调优实战	870
21.5.1 报错原因解析	870
21.5.2 调优实战	870
21.6 算子返回 NULL 产生的问题及解决办法	874
<b>第 22 章 Spark 集群资源分配及并行度调优最佳实践</b>	<b>875</b>
22.1 实际生产环境下每个 Executor 内存及 CPU 的具体配置 及原因	875
22.1.1 内存的具体配置及原因	875
22.1.2 实际生产环境下一般每个 Executor 的 CPU 的具体配置及原因	877
22.2 Spark 并行度设置最佳实践	878
22.2.1 并行度设置的原理和影响因素	878

22.2.2 并行度设置最佳实践	878
<b>第 23 章 Spark 集群中 Mapper 端、Reducer 端内存调优</b>	<b>880</b>
23.1 Spark 集群中 Mapper 端内存调优实战	880
23.1.1 内存使用详解	880
23.1.2 内存性能调优实战	881
23.2 Spark 集群中 Reducer 端内存调优实战	881
23.2.1 内存使用详解	881
23.2.2 内存性能调优实战	883
<b>第 24 章 使用 Broadcast 实现 Mapper 端 Shuffle 聚合功能的原理和调优实战</b>	<b>885</b>
24.1 使用 Broadcast 实现 Mapper 端 Shuffle 聚合功能的原理	885
24.2 使用 Broadcast 实现 Mapper 端 Shuffle 聚合功能调优实战	885
<b>第 25 章 使用 Accumulator 高效地实现分布式集群全局计数器的原理和调优案例</b>	<b>887</b>
25.1 Accumulator 内部工作原理	887
25.2 Accumulator 自定义实现原理和源码解析	887
25.3 Accumulator 作全局计数器案例实战	888
<b>第 26 章 Spark 下 JVM 性能调优最佳实践</b>	<b>889</b>
26.1 JVM 内存架构详解及调优	889
26.1.1 JVM 的堆区、栈区、方法区等详解	889
26.1.2 JVM 线程引擎及内存共享区域详解	890
26.1.3 JVM 中年轻代和老年代及元空间原理详解	891
26.1.4 JVM 进行 GC 的具体工作流程详解	895
26.1.5 JVM 常见调优参数详解	895
26.2 Spark 中对 JVM 使用的内存原理图详解及调优	896
26.2.1 Spark 中对 JVM 使用的内存原理图说明	896
26.2.2 Spark 中对 JVM 使用的内存原理图内幕详解	897
26.2.3 Spark 下的常见的 JVM 内存调优参数最佳实践	899
26.3 Spark 下 JVM 的 On-Heap 和 Off-Heap 解密	900
26.3.1 JVM 的 On-Heap 和 Off-Heap 详解	901
26.3.2 Spark 是如何管理 JVM 的 On-Heap 和 Off-Heap 的	902
26.3.3 Spark 下 JVM 的 On-Heap 和 Off-Heap 调优最佳实践	903
26.4 Spark 下的 JVM GC 导致的 Shuffle 拉取文件失败及调优方案	905
26.4.1 Spark 下的 JVM GC 导致的 Shuffle 拉取文件失败原因解密	905
26.4.2 Spark 下的 JVM GC 导致的 Shuffle 拉取文件失败时调优	906
26.5 Spark 下的 Executor 对 JVM 堆外内存连接等待时长调优	906
26.5.1 Executor 对堆外内存等待工作过程	906
26.5.2 Executor 对堆外内存等待时长调优	907
26.6 Spark 下的 JVM 内存降低 Cache 内存占比的调优	907
26.6.1 什么时候需要降低 Cache 的内存占用	907