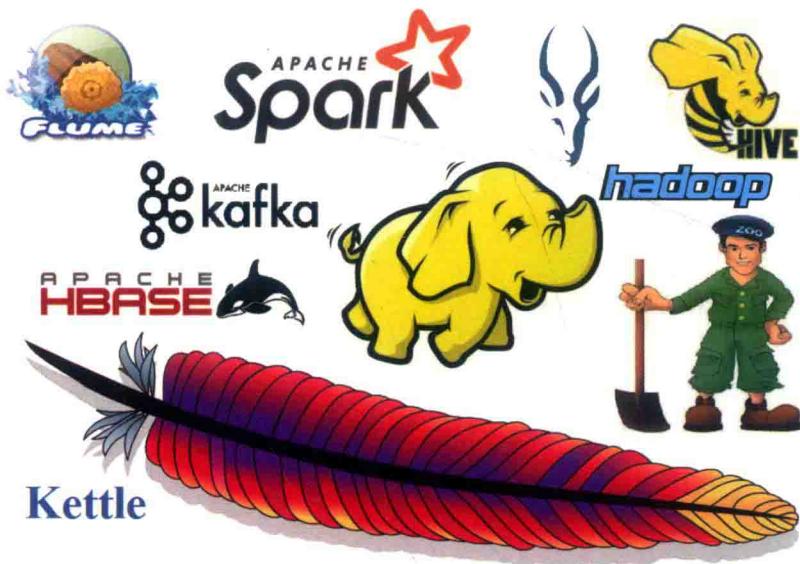


天善智能创始人 梁 勇 ·
清控道口资本管理有限公司管理合伙人 钟志新 ·
深圳纳实大数据技术有限公司CEO 吕 骏 · 联袂
斯凯奇(中国)有限公司 CIO 李宏阳 · 推荐
上海亦策软件科技有限公司总经理 邓强勇 ·

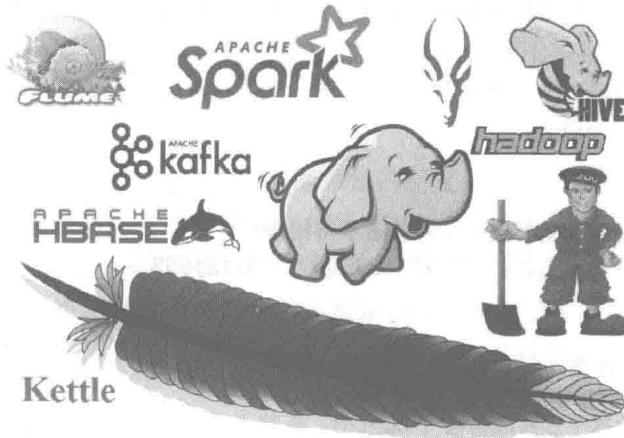


Cloudera Hadoop 大数据平台实战指南

宋立桓 陈建平 著



清华大学出版社



Cloudera Hadoop

大数据平台实战指南

宋立桓 陈建平 著

清华大学出版社
北京

内 容 简 介

对于入门和学习大数据技术的读者来说，大数据技术的生态圈和知识体系过于庞大，可能还没有开始学习就已经陷入众多的陌生名词和泛泛的概念中。本书的切入点明确而清晰，从 Hadoop 生态系统的明星 Cloudera 入手，逐步引出各类大数据基础和核心应用框架。

本书分为 18 章，系统介绍 Hadoop 生态系统大数据相关的知识，包括大数据概述、Cloudera Hadoop 平台的安装部署、HDFS 分布式文件系统、MapReduce 计算框架、资源管理调度框架 YARN、Hive 数据仓库、数据迁移工具 Sqoop、分布式数据库 HBase、ZooKeeper 分布式协调服务、准实时分析系统 Impala、日志采集工具 Flume、分布式消息系统 Kafka、ETL 工具 Kettle、Spark 计算框架等内容，最后给出两个综合实操案例，以巩固前面所学的知识点。

本书既适合 Hadoop 初学者、大数据技术工程师和大数据技术爱好者自学使用，亦可作为高等院校和培训机构大数据相关课程的培训用书。

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

Cloudera Hadoop 大数据平台实战指南 / 宋立桓，陈建平著.—北京：清华大学出版社，2019
ISBN 978-7-302-51753-5

I. ①C… II. ①宋… ②陈… III. ①数据处理软件 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2019) 第 271428 号

责任编辑：夏毓彦

封面设计：王翔

责任校对：闫秀华

责任印制：刘海龙

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：<http://www.lib.tsinghua.edu.cn>

投稿与读者服务：010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈：010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者：三河市龙大印装有限公司

经 销：全国新华书店

开 本：190mm×260mm 印 张：16 字 数：410 千字

版 次：2019 年 2 月第 1 版 印 次：2019 年 2 月第 1 次印刷

定 价：59.00 元

产品编号：080579-01



推荐序一

从 2013 年起大数据的概念在国内逐步普及，经过短短几年的时间，相关的技术就在各行各业有了深入的使用和发展，并且越来越多的企业开始重视对大数据项目的规划和建设。大数据的项目建设是以 IT 信息化部门为主导、企业各部门紧密配合、企业高层驱动的一个持续的过程，其中大数据技术的相关人才尤为重要。

2013 年，我与几位志同道合的深耕于数据领域的朋友一起成立了天善智能。至今，天善智能已经成为国内最大的大数据、商业智能 BI、人工智能 AI 的垂直社区之一，来自于百度、阿里巴巴、腾讯、微软、IBM、京东等国内一众知名公司的数据专家也积极地活跃在我们的社区。这些专家广泛地参与天善智能各类线上线下有关大数据技术的布道活动，他们用自己专业的知识、精湛的技术分享极大地点燃了广大大数据技术爱好者的热情，共同推动了大数据技术在国内的普及和发展。

在天善智能的成长过程中，我也有幸结识了很多来自各行各业大数据技术圈的朋友，其中包括本书的两位作者宋立桓和陈建平。宋立桓老师以前在微软工作，由于我以前也是微软技术的程序开发者，包括之后在工作中使用到微软商业智能 BI 技术，因此让我们有了更多交流的话题。之后，宋老师去腾讯就职，从大数据到云计算，从一个很深的数据底层走向另一个更深的架构底层，这是一个很好的提升和发展。建平来自传统的行业，在传统行业的数据升级打怪过程中不断将数据运用到了一个很高的高度。

在我们和很多技术专家合作各类线上线下沙龙分享的过程中，大家都意识到了一个问题——大数据知识体系过于庞大，零零散散的知识体系终归需要有一个载体，而这个载体既可以是文字的沉淀，也可以是专业课程的沉淀。很惊喜的是这两位志同道合的朋友在精心酝酿了很长一段时间之后，终于开始行动并将过往经验一一成文。

Hadoop 走过这么多年，整个生态体系越来越庞大，作为 Hadoop 最有影响力的数据管理软件服务提供商之一的 Cloudera 无疑是一颗耀眼的明星。两位作者从这个切入点开始循序渐进地将 Hadoop 生态系统中核心的技术、框架、应用一一展开，构成一个完整的知识体系框架，不多不少入门正好。本书案例简洁清晰，不少料不拖沓，可以帮助大家快速学习掌握大数据相关的核心知识点，希望此书能够成为广大大数据技术学习爱好者的手边参考书。

最后，回归到整个大行业，我们仍然要意识到：许多传统企业在从业务信息化到数据信息化的过程并不会一帆风顺。一方面来自于传统业务与大数据结合的场景目前依旧需要实践的检验，

存在一个比较长的建设和提炼周期，需要企业在人力、物力、财力有持续的投入和保障。另一方面，每个企业的 IT 基础、数据基础、技术积累程度不同，对于选择适合自身的大数据方案也并不是那么容易。对于技术人员来说，有些问题我们可能无力解决，我们能够做到的就是不断地夯实大数据技术、用技术驱动传统业务、挖掘业务增长点，让大数据真正地为企业创造业务价值。

天善智能创始人 梁勇

推荐序二

最近几年，Hadoop 作为最基础和最流行的大数据技术平台，被广泛地应用。作为一个开源技术平台，Hadoop 平台发展非常迅猛，在此基础上，也发展出很多的商业版本和分支。Cloudera 作为最早开源平台的贡献者，迅速成为这个技术的领导者，无论是社区版还是商业版都被大量的客户广泛采用。

亦策软件是一家国内领先的专注于大数据整体解决方案的高科技企业，为客户提供大数据分析平台端到端的解决方案。在和客户交流的过程中，我们发现很多用户和技术人员在学习、使用和运维 Cloudera Hadoop 平台的时候都面临一个严峻的挑战，就是中文的教材和资料非常少，相关的课程和培训又非常昂贵，所以只能靠自己摸索或者学习英文的资料。这给入门和掌握 Cloudera Hadoop 平台造成很大的不便，降低了效率。听闻宋立桓老师要出版《Cloudera Hadoop 大数据平台实战指南》这本书，我很高兴受邀为这本书写序。这本书梳理了大数据基础和核心技术框架、注重实操、案例简洁，很适合广大大数据技术爱好者参考学习。我想本书的出版一定能帮助大家解决学习资料的燃眉之急。

感谢宋立桓老师为大数据技术的推广做出贡献，并真诚希望广大读者通过本书能缩短相关技术从入门到精通的过程，最后祝本书大卖！

上海亦策软件科技有限公司总经理 邓强勇

推荐序三

2011 年秋，我在西雅图参加微软的技术大会，在之前的一年，EMC 刚刚提出了 Big Data 和 Data Lake 的理念，我和同事们疯狂地讨论着这个也许会颠覆软件业的新理念。我们意识到，大量的数据和非结构化的社会网络信息会成为这个全新时代的重要资源高地。我们讨论着 2005 年诞生的 Hadoop，2007 年 Linus Torvalds 于 10 天之内开发 GitHub 的传奇，2008 年中本聪发表比特币论文，2010 年张小龙开发了微信，那是伟大的事件频繁发生的 5 年。

我与本书作者宋立桓老师是在那时候认识的，我们谈了很多关于未来数据应用的场景，我们一致认为，数据的价值将会在未来 10 年内超过那些陈旧的系统，会有更新一代的应用基于显而易见、唾手可得的数据而诞生。我们将不再依赖复杂的流程和权限，会让每个社会实体都有自知之明和洞见之能。

从那以后，宋立桓老师一直致力于探索数据的价值，以及如何实现数据的价值，从数据分析到数据业务探索，从数据整合到数据共享，他研究得非常系统和完整。《Cloudera Hadoop 大数据平台实战指南》是他继《人人都是数据分析师：微软 Power BI 实践指南》之后的又一力作，通俗易懂，概念清晰，是对大数据架构和相关大数据系统普及的好教程。在本书中，从大数据概念到原理，从理论到实战，从部署到操作，无一不凝聚着他严谨的学习态度和实践精神，是很好的打开大数据宝藏的一把钥匙，也是近年来在该领域不可多得的学习材料。

希望你们和我一起加入本书的阅览与实战的旅程中。

感谢宋立桓老师的作序邀请，预祝有更多推动大数据行业的真知灼见早日发表。

深圳纳实大数据技术有限公司 CEO 吕骏

前 言

大数据这个词也许几年前你听着还有点陌生，但我相信你现在听到 Hadoop 这个词时会觉得“熟悉”！你会发现身边从事 Hadoop 开发或者正在学习 Hadoop 的人越来越多。

最早提出“大数据”时代到来的是全球知名咨询公司麦肯锡，麦肯锡称：“数据，已经渗透到当今每一个行业和业务职能领域，成为重要的生产因素。互联网技术发展到现今阶段，大量日常、工作等事务产生的数据比以前有了爆炸式的增长，以前的传统数据处理技术已经无法胜任，需求催生技术——一套用来处理海量数据的软件技术框架 Hadoop 应运而生！”

我本人一直从事云计算、大数据方面的咨询和培训工作。大数据产业高速发展促使 Hadoop 人才的需求井喷式增长，但 Hadoop 大数据工程师培养数量远远无法满足市场的需求。为了不被淹没在大数据技术的浪潮中，我们只有坚持学习，通过增加知识来实现对自我价值的挖掘和体现。

关于本书

Hadoop 的发行版除了社区的 Apache Hadoop 外，Cloudera、Hortonworks、华为等公司都提供了自己的商业版本。因为企业通常使用的是 Hadoop 商业版本，所以本书实操的运行环境采用 Cloudera 的 CDH。本书定位是大数据从入门到应用的简明系统教程，特色是理论联系实践、实战实用为主、内容全面系统、讲解深入浅出，是大数据技术爱好者入门的最佳图书。

本书分为 18 章（宋立桓老师撰写第 1~12 章、陈建平撰写第 13~18 章），分别从大数据概述、Cloudera Hadoop 平台的安装部署、大数据 Hadoop 组件三方面进行介绍，内容包括 HDFS 分布式文件系统、MapReduce 计算框架、资源管理框架 YARN、Hive 数据仓库、数据迁移工具 Sqoop、分布式数据库 HBase、ZooKeeper 分布式协调服务、准实时分析系统 Impala、日志采集工具 Flume、分布式消息系统 Kafka、ETL 工具 Kettle、Spark 计算框架等知识，最后用两个综合实操案例把所有知识点串起来。

本书使用的操作环境是 Hadoop 商业发行版的 Cloudera Express（Express 是免费版本，企业版需付费）。全书秉承“实践为主、理论够用”的原则，将演示实验融入各个知识点讲解中。

本书另提供丰富的案例源文件和大数据工具软件下载，供读者亲自操作练习，在作者博客 <http://blog.51cto.com/lihuansong> 中有下载地址。

学习本书之前，希望大家应该具有如下基础：有一定计算机网络基础知识，熟悉常用 Linux 操作命令，对 Java 语言和数据库理论有基本的了解。

资源下载与技术支持

本书提供详细的案例资源文件，在作者博客置顶文章中提供下载地址，便于读者动手实践：

<http://blog.51cto.com/lihuansong/2317021>

欢迎读者来信互动，宋立桓的邮箱是 songlihuan@hotmail.com，陈建平的邮箱是 daxia1520@163.com。

致谢

感谢我的妻子，她是我完成此书的坚强后盾。

感谢我的朋友和同事，他们让我学会知识的增值和变现。

感谢清华大学出版社的编辑夏毓彦和其他工作人员帮助我出版了这本有意义的著作。

阿基米德有一句名言：“给我一个支点，我就能撬起地球。”谨以此书献给那些为大数据与商业智能分析铺路的人，让更多的人享受到大数据时代到来的红利。

宋立桓

云计算架构师、大数据咨询顾问

2018年11月

目 录

第1章 大数据概述	1
1.1 大数据时代的数据特点	1
1.2 大数据时代的发展趋势——数据将成为资产	2
1.3 大数据时代处理数据理念的改变	3
1.3.1 要全体不要抽样	3
1.3.2 要效率不要绝对精确	3
1.3.3 要相关不要因果	4
1.4 大数据时代的关键技术	5
1.5 大数据时代的典型应用案例	5
1.5.1 塔吉特超市精准营销案例	5
1.5.2 谷歌流感趋势案例	6
1.5.3 证券行业案例	6
1.5.4 某运营商大数据平台案例	7
1.6 Hadoop 概述和介绍	7
1.6.1 Hadoop 发展历史和应用现状	7
1.6.2 Hadoop 的特点	8
1.6.3 Hadoop 的生态系统	8
第2章 Cloudera 大数据平台介绍	10
2.1 Cloudera 简介	10
2.2 Cloudera 的 Hadoop 发行版 CDH 简介	11
2.2.1 CDH 概述	11
2.2.2 CDH 和 Apache Hadoop 对比	12
2.3 Cloudera Manager 大数据管理平台介绍	12
2.3.1 Cloudera Manager 概述和整体架构	12
2.3.2 Cloudera Manager 的基本核心功能	14
2.3.3 Cloudera Manager 的高级功能	18
2.4 Cloudera 平台参考部署架构	19
2.4.1 Cloudera 的软件体系结构	19
2.4.2 群集硬件规划配置	19
2.4.3 Hadoop 集群角色分配	21
2.4.4 网络拓扑	23

第 3 章 Cloudera Manager 及 CDH 离线安装部署	25
3.1 安装前的准备工作	25
3.2 Cloudera Manager 及 CDH 安装	30
3.3 添加其他大数据组件	35
第 4 章 分布式文件系统 HDFS	37
4.1 HDFS 简介	37
4.2 HDFS 体系结构	38
4.2.1 HDFS 架构概述	38
4.2.2 HDFS 命名空间管理	38
4.2.3 NameNode	39
4.2.4 SecondaryNameNode	39
4.3 HDFS 2.0 新特性	41
4.3.1 HDFS HA	41
4.3.2 HDFS Federation	42
4.4 HDFS 操作常用 shell 命令	43
4.4.1 HDFS 目录操作和文件处理命令	43
4.4.2 HDFS 的 Web 管理界面	44
4.4.3 dfsadmin 管理维护命令	45
4.4.4 namenode 命令	47
4.5 Java 编程操作 HDFS 实践	47
4.6 HDFS 的参数配置和规划	49
4.7 使用 Cloudera Manager 启用 HDFS HA	51
4.7.1 HDFS HA 高可用配置	51
4.7.2 HDFS HA 高可用功能测试	54
第 5 章 分布式计算框架 MapReduce	57
5.1 MapReduce 概述	57
5.2 MapReduce 原理介绍	58
5.2.1 工作流程概述	58
5.2.2 MapReduce 框架的优势	58
5.2.3 MapReduce 执行过程	59
5.3 MapReduce 编程——单词示例解析	59
5.4 MapReduce 应用开发	60
5.4.1 配置 MapReduce 开发环境	60
5.4.2 编写和运行 MapReduce 程序	61
第 6 章 资源管理调度框架 YARN	65
6.1 YARN 产生背景	65
6.2 YARN 框架介绍	66

6.3 YARN 工作原理	67
6.4 YARN 框架和 MapReduce1.0 框架对比	69
6.5 CDH 集群的 YARN 参数调整	69
第 7 章 数据仓库 Hive	72
7.1 Hive 简介	72
7.2 Hive 体系架构和应用场景	73
7.2.1 Hive 体系架构	73
7.2.2 Hive 应用场景	74
7.3 Hive 的数据模型	75
7.3.1 内部表	75
7.3.2 外部表	75
7.3.3 分区表	75
7.3.4 桶	75
7.4 Hive 实战操作	76
7.4.1 Hive 内部表操作	77
7.4.2 Hive 外部表操作	77
7.4.3 Hive 分区表操作	79
7.4.4 桶表	80
7.4.5 Hive 应用实例 WordCount	82
7.4.6 UDF	84
7.5 基于 Hive 的应用案例	86
第 8 章 数据迁移工具 Sqoop	88
8.1 Sqoop 概述	88
8.2 Sqoop 工作原理	89
8.3 Sqoop 版本和架构	91
8.4 Sqoop 实战操作	93
第 9 章 分布式数据库 HBase	100
9.1 HBase 概述	100
9.2 HBase 数据模型	101
9.3 HBase 生态地位和系统架构	101
9.3.1 HBase 的生态地位解析	101
9.3.2 HBase 系统架构	102
9.4 HBase 运行机制	103
9.4.1 Region	103
9.4.2 Region Server 工作原理	103
9.4.3 Store 工作原理	104
9.5 HBase 操作实战	104

9.5.1 HBase 常用 shell 命令	104
9.5.2 HBase 编程实践	107
9.5.3 HBase 参数调优的案例分享	109
第 10 章 分布式协调服务 ZooKeeper	111
10.1 ZooKeeper 的特点	111
10.2 ZooKeeper 的工作原理	112
10.2.1 基本架构	112
10.2.2 ZooKeeper 实现分布式 Leader 节点选举	112
10.2.3 ZooKeeper 配置文件重点参数详解	112
10.3 ZooKeeper 典型应用场景	115
10.3.1 ZooKeeper 实现 HDFS 的 NameNode 高可用 HA	115
10.3.2 ZooKeeper 实现 HBase 的 HMaster 高可用	116
10.3.3 ZooKeeper 在 Storm 集群中的协调者作用	116
第 11 章 准实时分析系统 Impala	118
11.1 Impala 概述	118
11.2 Impala 组件构成	119
11.3 Impala 系统架构	119
11.4 Impala 的查询处理流程	120
11.5 Impala 和 Hive 的关系和对比	121
11.6 Impala 安装	122
11.7 Impala 入门实战操作	124
第 12 章 日志采集工具 Flume	128
12.1 Flume 概述	128
12.2 Flume 体系结构	129
12.2.1 Flume 外部结构	129
12.2.2 Flume 的 Event 事件概念	130
12.2.3 Flume 的 Agent	130
12.3 Flume 安装和集成	131
12.3.1 搭建 Flume 环境	131
12.3.2 Kafka 与 Flume 集成	132
12.4 Flume 操作实例介绍	132
12.4.1 例子概述	132
12.4.2 第一步：配置数据流向	132
12.4.3 第二步：启动服务	133
12.4.4 第三步：新建空数据文件	133
12.4.5 第四步：运行 flume-ng 命令	133
12.4.6 第五步：运行命令脚本	134

12.4.7 最后一步：测试结果	134
第 13 章 分布式消息系统 Kafka	135
13.1 Kafka 架构设计	135
13.1.1 基本架构	135
13.1.2 基本概念	136
13.1.3 Kafka 主要特点	136
13.2 Kafka 原理解析	137
13.2.1 主要的设计理念	137
13.2.2 ZooKeeper 在 Kafka 的作用	137
13.2.3 Kafka 在 ZooKeeper 的执行流程	137
13.3 Kafka 安装和部署	138
13.3.1 CDH5 完美集成 Kafka	138
13.3.2 Kafka 部署模式和配置	139
13.4 Java 操作 Kafka 消息处理实例	141
13.4.1 例子概述	141
13.4.2 第一步：新建工程	141
13.4.3 第二步：编写代码	141
13.4.4 第三步：运行发送数据程序	142
13.4.5 最后一步：运行接收数据程序	143
13.5 Kafka 与 HDFS 的集成	143
13.5.1 与 HDFS 集成介绍	143
13.5.2 与 HDFS 集成实例	144
13.5.3 第一步：编写代码——发送数据	144
13.5.4 第二步：编写代码——接收数据	145
13.5.5 第三步：导出文件	146
13.5.6 第四步：上传文件	146
13.5.7 第五步：运行程序——发送数据	146
13.5.8 第六步：运行程序——接收数据	147
13.5.9 最后一步：查看执行结果	147
第 14 章 大数据 ETL 工具 Kettle	148
14.1 ETL 原理	148
14.1.1 ETL 简介	148
14.1.2 ETL 在数据仓库中的作用	149
14.2 Kettle 简介	149
14.3 Kettle 完整案例实战	150
14.3.1 案例介绍	150
14.3.2 最终效果	150
14.3.3 表说明	150

14.3.4 第一步：准备数据库数据	151
14.3.5 第二步：新建转换	152
14.3.6 第三步：新建数据库连接	153
14.3.7 第四步：拖动表输入组件	153
14.3.8 第五步：设置属性——order 表	154
14.3.9 第六步：设置属性——user 表	155
14.3.10 第七步：拖动流查询并设置属性——流查询	155
14.3.11 第八步：设置属性——product 表	156
14.3.12 第九步：连接组件	156
14.3.13 第十步：设置属性——文本输出	156
14.3.14 最后一步：运行程序并查看结果	157
14.4 Kettle 调度和命令	158
14.4.1 通过页面调度	158
14.4.2 通过脚本调度	159
14.5 Kettle 使用原则	161
第 15 章 大规模数据处理计算引擎 Spark	162
15.1 Spark 简介	162
15.1.1 使用背景	162
15.1.2 Spark 特点	163
15.2 Spark 架构设计	163
15.2.1 Spark 整体架构	163
15.2.2 关键运算组件	164
15.2.3 RDD 介绍	164
15.2.4 RDD 操作	165
15.2.5 RDD 依赖关系	166
15.2.6 RDD 源码詳解	167
15.2.7 Scheduler	168
15.2.8 Storage	168
15.2.9 Shuffle	169
15.3 Spark 编程实例	170
15.3.1 实例概述	170
15.3.2 第一步：编辑数据文件	170
15.3.3 第二步：编写程序	171
15.3.4 第三步：上传 JAR 文件	171
15.3.5 第四步：远程执行程序	172
15.3.6 最后一步：查看结果	172
15.4 Spark SQL 实战	173
15.4.1 例子概述	173
15.4.2 第一步：编辑数据文件	173

15.4.3 第二步：编写代码	174
15.4.4 第三步：上传文件到服务器	174
15.4.5 第四步：远程执行程序	174
15.4.6 最后一步：查看结果	175
15.5 Spark Streaming 实战	175
15.5.1 例子概述	175
15.5.2 第一步：编写代码	175
15.5.3 第二步：上传文件到服务器	176
15.5.4 第三步：远程执行程序	177
15.5.5 第四步：上传数据	177
15.5.6 最后一步：查看结果	177
15.6 Spark MLlib 实战	178
15.6.1 例子步骤	178
15.6.2 第一步：编写代码	178
15.6.3 第二步：上传文件到服务器	179
15.6.4 第三步：远程执行程序	179
15.6.5 第四步：上传数据	180
15.6.6 最后一步：查看结果	180
第 16 章 大数据全栈式开发语言 Python	182
16.1 Python 简介	182
16.2 Python 安装和配置	183
16.2.1 Anaconda 介绍	183
16.2.2 Anaconda 下载	183
16.2.3 Anaconda 安装	184
16.2.4 Anaconda 包管理	185
16.2.5 PyCharm 下载	185
16.2.6 PyCharm 安装	185
16.2.7 PyCharm 使用	187
16.3 Python 入门	190
16.3.1 例子概述	190
16.3.2 第一步：新建 Python 文件	190
16.3.3 第二步：设置字体大小	191
16.3.4 第三步：编写代码	191
16.3.5 第四步：执行程序	192
16.3.6 最后一步：改变输入	192
16.4 Python 数据科学库 pandas 入门	193
16.4.1 例子概述	193
16.4.2 pandas 包介绍	194
16.4.3 第一步：打开 Jupyter Notebook	194

16.4.4 第二步：导入包	194
16.4.5 第三步：定义数据集	195
16.4.6 第四步：过滤数据	195
16.4.7 最后一步：获取数据	196
16.5 Python 绘图库 matplotlib 入门	197
16.5.1 例子概述	197
16.5.2 第一步：新建一个 Python 文件	197
16.5.3 第二步：引入画图包	197
16.5.4 第三步：组织数据	198
16.5.5 第四步：画图	198
16.5.6 最后一步：查看结果	199
第 17 章 大数据实战案例：实时数据流处理项目	200
17.1 项目背景介绍	200
17.2 业务需求分析	200
17.3 项目技术架构	201
17.4 项目技术组成	202
17.5 项目实施步骤	202
17.5.1 第一步：运用 Kafka 产生数据	202
17.5.2 第二步：运用 Spark 接收数据	208
17.5.3 第三步：安装 Redis 软件	211
17.5.4 第四步：准备程序运行环境	214
17.5.5 第五步：远程执行 Spark 程序	216
17.5.6 第六步：编写 Python 实现可视化	218
17.5.7 最后一步：执行 Python 程序	221
17.6 项目总结	222
第 18 章 大数据实战案例：用户日志综合分析项目	223
18.1 项目背景介绍	223
18.2 项目设计目的	223
18.3 项目技术架构和组成	224
18.4 项目实施步骤	225
18.4.1 第一步：本地数据 FTP 到 Linux 环境	225
18.4.2 第二步：Linux 数据上传到 HDFS	225
18.4.3 第三步：使用 Hive 访问 HDFS 数据	226
18.4.4 第四步：使用 Kettle 把数据导入 HBase	228
18.4.5 第五步：使用 Sqoop 把数据导入 MySQL	234
18.4.6 第六步：编写 Python 程序实现可视化	236
18.4.7 最后一步：执行 Python 程序	238