



普通高等教育“十三五”规划教材  
国家新闻出版改革发展项目库入库项目  
数据科学与大数据技术专业教材丛书

# 大数据技术基础

BIG DATA TECHNOLOGY FOUNDATION

鄂海红 宋美娜 欧中洪◎编著



北京邮电大学出版社  
www.buptpress.com



普通高等教育“十三五”规划教材  
国家新闻出版改革发展项目库入库项目  
数据科学与大数据技术专业教材丛书

# 大数据技术基础

鄂海红 宋美娜 欧中洪 编著



北京邮电大学出版社  
[www.buptpress.com](http://www.buptpress.com)

## 内 容 简 介

本书围绕大数据技术基础,重点介绍了大数据存储系统(分布式文件系统和 NoSQL 数据库)、大数据处理框架(Hadoop 的 MapReduce、Spark 及实时处理框架 Storm 和 Flink)、大数据仓库技术(Hive、Druid 等)、大数据多维分析(Kylin)、大数据可视化技术和大数据综合应用等,以及当今主流的大数据平台构建技术和开源组件实践知识,可以指导读者全面、系统地掌握大数据各层的实现方案,开展各领域的大数据实践。本书可作为计算机学科相关专业,特别是数据科学与大数据技术专业的教材。

### 图书在版编目(CIP)数据

大数据技术基础 / 鄂海红, 宋美娜, 欧中洪编著. -- 北京: 北京邮电大学出版社, 2019. 9

ISBN 978-7-5635-5878-0

I. ①大… II. ①鄂… ②宋… ③欧… III. ①数据处理—高等学校—教材 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2019) 第 204848 号

---

书 名: 大数据技术基础

作 者: 鄂海红 宋美娜 欧中洪

责任编辑: 孙宏颖

出版发行: 北京邮电大学出版社

社 址: 北京市海淀区西土城路 10 号(100876)

发 行 部: 电话: 010-62282185 传真: 010-62283578

E-mail: publish@bupt.edu.cn

经 销: 各地新华书店

印 刷: 保定市中国画美凯印刷有限公司

开 本: 787 mm×1 092 mm 1/16

印 张: 15.5

字 数: 401 千字

版 次: 2019 年 9 月第 1 版 2019 年 9 月第 1 次印刷

---

ISBN 978-7-5635-5878-0

定价: 48.00 元

· 如有印装质量问题,请与北京邮电大学出版社发行部联系 ·

# 大数据顾问委员会

宋俊德 王国胤 张云勇 郑宇  
段云峰 田世明 娄瑜 孙少隣  
王柏

## 大数据专业教材编委会

总主编：吴斌

编委：宋美娜 欧中洪 鄂海红 双锴  
于艳华 周文安 林荣恒 李静林  
袁燕妮 李劼 皮人杰

总策划：姚顺

秘书长：刘纳新

## 序 言

党的十八届五中全会明确提出实施国家大数据战略,至此大数据技术成为塑造国家竞争力的战略制高点之一。掌握和运用大数据技术的能力成为一个国家竞争力的重要体现。国内许多行业如互联网、电信、金融和交通等开始实际部署大数据平台并付诸实践,这带动了软件、硬件及服务市场的快速发展。

大数据正在成为产业发展的重要推动力,大数据相关产业的高速发展带来了大数据人才严重短缺的问题,大数据人才的培养成为当前急迫的任务。近年来大数据专业建设在全国各大高校如火如荼地开展,设立该专业的学校数量也在快速增长。截止到2019年4月,教育部累计批准486所高校设立“数据科学与大数据技术”专业,其中,2016年3所高校获批,2017年32所高校获批,2018年248所高校获批,2019年203所高校获批。如何更好地建设大数据专业和培养产业迫切需求的高水平专业人才,成为高校人才培养工作的重要挑战。

自2018年起,北京邮电大学出版社联合北京邮电大学计算机学院、网络技术研究院的多位知名教授、副教授及任课教师,共同开启“数据科学与大数据技术专业教材丛书”的出版工作。这套丛书包括《大数据技术基础》《大数据技术基础实验》《R语言编程与数据科学》《网络科学与计算》《计算机视觉》《NoSQL数据库技术》《流数据分析技术》《数据可视化》《机器学习》《分布式计算与云计算》《数据仓库与数据挖掘》《Python语言程序设计》等。这些教材的出版凝炼了众多大数据领域教学、科研专家的心得体会,为大数据创新型人才的培养奠定了基础。

《大数据技术基础》是“数据科学与大数据技术”专业重要的基础教材之一,主要讲授大数据知识体系中理论与工程实践结合的技术基础。该书涵盖大数据采集、存储、处理、分析、可视化及应用等一整套全流程所需的基础理论知识。为了使读者能够快速掌握大数据工程实践的知识,书中还介绍了多种开源大数据实践工具组件的技术架构和使用方法。可以说,该书所设计的内容一方面体现了对学生理论知识培养的重视,另一方面强调了计算机专业背景下数据科学的系统观,注重学生实际应用能力的培养。

该书的作者一直在大数据领域从事一线的教学和科研工作,这些工作基础为大数据专业人才的培养和大数据专业教材的出版提供着有力的支撑。该书作为北京邮电大学计算机学院“数据科学与大数据技术”专业的第一批正式出版教材,我很期待在以后的教学和科研实践中该书能够得到不断升华,也恳请全国同行在使用该书的同时予以批评指正,让我们一起为中国的大数据事业添砖加瓦。

国家“万人计划”领军人才

北京邮电大学计算机学院执行院长

博士生导师、教授

苏 森

本书一共分为 9 章。

第 1 章为大数据概述。本章首先介绍了大数据的发展历程、大数据的定义与特征、大数据与传统数据的区别；然后介绍了大数据平台应具备的能力和大数据平台架构；最后介绍了 Hadoop 生态开源组件和大数据技术的应用领域。

第 2 章为大数据存储技术。本章主要介绍主流的分布式存储系统，包括相关概念、体系结构、存储机制和操作方法，主要涵盖了分布式文件系统 HDFS 以及 4 种 NoSQL 数据库。

第 3 章、第 4 章、第 5 章为大数据处理技术。第 3 章介绍了 Hadoop 的 MapReduce 并行计算框架，第 4 章介绍了 Spark 内存计算框架，第 5 章介绍了实时计算框架。

第 6 章为大数据仓库技术。本章介绍了分布式数据仓库和数据查询技术，主要包括 3 个组件：Hive 分布式数据仓库、Druid 时序数据仓储和 Drill 分布式实时查询。

第 7 章为大数据多维分析技术。本章的主要内容包括大数据多维分析技术演进的需求和背景、开源 Kylin 的基本概念与原理、技术架构和实战操作方法。

第 8 章为大数据可视化技术。本章详细介绍了数据可视化的定义及其分类、可视化流程，以及时空数据可视化、层次和网络数据可视化、文本和文档可视化的概念，并对商业智能中的数据可视化及其应用进行了介绍；同时讲解了常见的数据可视化的实现技术和方法。

第 9 章为大数据应用案例。本章选择了某电影大数据平台案例，结合某电影大数据平台的技术体系架构，对大数据应用的构建流程进行了介绍，可以帮助读者整体性地理解和掌握本书知识内容的实践方法。

本书可以作为数据科学与大数据技术专业的本科高年级专业课教材，也可以作为研究生相关课程的参考材料。同时本书还配套了《大数据技术基础实验》，用于指导读者学习具体的实践课程知识，以使读者掌握实际大数据平台和大数据应用系统的研发能力。

本书的编写得到了北京邮电大学 PCN&CAD 中心、教育部信息网络工程研究中心和北

京邮电大学计算机学院数据科学与服务中心教师与研究生的支持,他们分别是宋美娜、欧中洪、宋俊德、毕秋波、韩鹏昊、田川、孔慧慧、赵淑晨、吴金盛、温宇飞、万仁山、谭泽华、陈小康、韦帅丽、朱永波,在此一并表示感谢。

感谢国家重点研发计划项目“大数据征信及智能评估技术”和“基于大数据的科技咨询技术与服务平台研发”、国家科技条件平台计划项目“国家人类遗传资源共享服务平台北京创新中心建设”的支持。

作者作为在计算机领域从事科研和教学的教师,由于在专业知识的深度和广度上的局限性使得本书存在不足之处,欢迎广大读者反馈对本书的意见和建议,我们将随着“大数据技术基础”专业课程的建设,不断地改进本书的质量。

鄂海红  
于北京

第 1 章 大数据概述	1
本章思维导图	1
1.1 大数据简介	2
1.1.1 大数据的发展历程	2
1.1.2 大数据的定义与特征	2
1.1.3 大数据与传统数据的区别	3
1.2 大数据平台应具备的能力	3
1.3 大数据平台架构	5
1.4 Hadoop 生态系统	8
1.5 大数据应用	10
1.5.1 互联网大数据应用	10
1.5.2 金融行业大数据应用	10
1.5.3 医疗行业大数据应用	11
1.5.4 智慧交通大数据应用	11
本章课后习题	12
本章参考文献	12
第 2 章 大数据存储——分布式文件系统及 NoSQL 数据库	14
本章思维导图	14
2.1 分布式文件系统	15
2.1.1 HDFS 相关概念	15
2.1.2 HDFS 体系结构	16
2.1.3 HDFS 存储机制	18
2.1.4 HDFS 读/写操作	20
2.1.5 HDFS 数据导入	21

2.2	NoSQL 数据库 .....	22
2.2.1	Key-Value 模型 .....	22
2.2.2	Key-Document 模型 .....	23
2.2.3	Key-Column 模型 .....	24
2.2.4	图模型 .....	25
2.3	列族数据库 .....	25
2.3.1	列族数据库简介 .....	25
2.3.2	HBase 的基本原理 .....	26
2.3.3	HBase 的数据模型 .....	30
2.4	键值数据库 .....	33
2.4.1	键值数据库简介 .....	33
2.4.2	选择键值数据库的原因 .....	33
2.4.3	Redis 的数据结构简介 .....	34
2.4.4	Redis 的数据持久化 .....	36
2.4.5	Redis 的数据复制 .....	37
2.5	文档数据库 .....	38
2.5.1	文档数据库简介 .....	38
2.5.2	MongoDB 的数据类型 .....	39
2.5.3	MongoDB 的数据复制 .....	40
2.6	图数据库 .....	42
2.6.1	图数据库简介 .....	42
2.6.2	图数据库的优势 .....	43
2.6.3	Neo4j 的基本元素与概念 .....	44
2.6.4	Cypher 简介 .....	46
	本章课后习题 .....	47
	本章参考文献 .....	47

### 第 3 章 大数据处理——MapReduce 处理框架 .....

	本章思维导图 .....	48
3.1	MapReduce 的发展背景 .....	49
3.2	MapReduce 框架 .....	50
3.3	MapReduce 的编程模型 .....	52
3.3.1	MapReduce 初析 .....	52
3.3.2	MapReduce 的运行机制 .....	57
3.3.3	MapReduce 的相关问题 .....	59

3.4 MapReduce 的集群调度 .....	60
3.4.1 Hadoop1.x 的传统集群调度框架 .....	60
3.4.2 Hadoop2.x 的集群调度框架 YARN .....	61
3.4.3 Hadoop 作业调度器 .....	64
本章课后习题 .....	67
本章参考文献 .....	67
<b>第 4 章 大数据处理——分布式内存处理框架 Spark .....</b>	<b>68</b>
本章思维导图 .....	68
4.1 Spark 简介 .....	69
4.1.1 Spark 介绍 .....	69
4.1.2 提出 Spark 的原因 .....	70
4.1.3 Spark 中的关键术语 .....	70
4.1.4 Spark 的优点 .....	71
4.2 Spark 框架 .....	72
4.2.1 Spark 框架图 .....	72
4.2.2 Spark 运行图 .....	73
4.2.3 Spark 任务调度方法 .....	73
4.3 RDD 概念理解 .....	74
4.3.1 RDD 介绍 .....	74
4.3.2 RDD 的操作 .....	75
4.3.3 RDD 的存储 .....	75
4.3.4 RDD 分区 .....	76
4.3.5 RDD 优先位置 .....	76
4.3.6 RDD 依赖关系 .....	76
4.4 RDD 操作 .....	78
4.4.1 RDD 创建 .....	78
4.4.2 转换操作 .....	78
4.4.3 行动操作 .....	80
4.5 Scala 语言 .....	81
4.5.1 Scala 介绍 .....	81
4.5.2 Scala 基本语法 .....	82
4.5.3 Scala 编写 Spark 示例 .....	86
4.6 Spark SQL 简介 .....	86
4.6.1 Spark SQL 与 Shark 的对比 .....	86

4.6.2	Spark SQL 的优势 .....	87
4.6.3	Spark SQL 生态 .....	87
4.7	MLlib 简介 .....	88
4.7.1	MLlib 介绍 .....	88
4.7.2	MLlib 支持机器学习算法 .....	88
	本章课后习题 .....	89
	本章参考文献 .....	89
<b>第 5 章</b>	<b>大数据处理——实时处理框架</b> .....	<b>90</b>
	本章思维导图 .....	90
5.1	实时处理架构 .....	91
5.1.1	基本概念 .....	91
5.1.2	批量和流式计算 .....	92
5.1.3	系统生态简介 .....	92
5.2	Storm 框架 .....	93
5.2.1	Storm 的基本术语和概念 .....	93
5.2.2	Storm 特性及运行原理 .....	94
5.2.3	消息的生命周期 .....	95
5.2.4	消息的可靠性保障 .....	96
5.3	Flume 分布式日志收集 .....	98
5.3.1	Flume 的基本术语和概念 .....	98
5.3.2	源 .....	99
5.3.3	通道 .....	100
5.3.4	接收器 .....	100
5.4	Kafka 分布式消息队列 .....	101
5.4.1	Kafka 的基本术语和概念 .....	102
5.4.2	生产者 .....	103
5.4.3	消费者 .....	104
5.4.4	数据传递的可靠性保障 .....	105
5.5	Spark Streaming 框架 .....	107
5.5.1	Spark Streaming 架构 .....	107
5.5.2	输入数据源 .....	108
5.5.3	DStream 的转换操作 .....	108
5.5.4	输出存储 .....	110
5.5.5	容错机制 .....	110

5.6 Flink 框架 .....	112
5.6.1 Flink 架构 .....	112
5.6.2 Client .....	112
5.6.3 JobManager .....	113
5.6.4 TaskManager .....	114
本章课后习题 .....	115
本章参考文献 .....	115
<b>第 6 章 大数据查询——分布式数据查询</b> .....	<b>116</b>
本章思维导图 .....	116
6.1 分布式数据查询简介 .....	117
6.2 Hive 分布式数据仓库 .....	118
6.2.1 Hive 概述 .....	118
6.2.2 Hive 内部介绍 .....	118
6.2.3 Hive 架构介绍 .....	119
6.2.4 HiveQL:数据定义 .....	119
6.2.5 HiveQL:数据导入 .....	121
6.2.6 HiveQL:查询 .....	123
6.3 Druid 时序数据仓储 .....	129
6.3.1 Druid 概述 .....	129
6.3.2 架构详解 .....	132
6.3.3 数据摄入 .....	135
6.3.4 数据查询 .....	141
6.4 Drill 分布式实时查询 .....	156
6.4.1 使用 Apache Drill 的原因 .....	156
6.4.2 Drill 架构与原理 .....	157
6.4.3 Drill 核心模块 .....	160
6.4.4 使用 Drill 实现查询 .....	161
本章课后习题 .....	168
本章参考文献 .....	168
<b>第 7 章 大数据分析——Kylin 分布式多维数据分析</b> .....	<b>170</b>
本章思维导图 .....	170
7.1 使用 Apache Kylin 的原因 .....	171
7.2 Kylin 学习的前奏 .....	172

7.2.1	数据仓库的概念与产生需求 .....	172
7.2.2	数据仓库与数据分析型系统 .....	174
7.2.3	多维数据分析 .....	175
7.2.4	OLAP 与数据立方体 .....	176
7.3	Kylin 工作原理 .....	178
7.3.1	Cube 与 Cuboid .....	178
7.3.2	工作流程 .....	178
7.4	Kylin 架构 .....	179
7.5	Kylin 快速入门 .....	181
7.5.1	在 Hive 中准备数据 .....	181
7.5.2	设计数据模型 .....	181
7.5.3	创建 Cube .....	183
7.5.4	构建 Cube .....	186
7.5.5	查询 Cube .....	188
7.6	增量构建 .....	188
7.6.1	设计增量 Cube .....	189
7.6.2	触发增量构建 .....	190
7.6.3	管理 Cube 碎片 .....	190
7.7	查询和可视化 .....	192
7.7.1	Web GUI .....	192
7.7.2	Rest API .....	194
7.7.3	ODBC .....	197
7.7.4	通过 Tableau 访问 Kylin .....	197
7.8	Cube 优化 .....	201
	本章课后习题 .....	204
	本章参考文献 .....	204

**第 8 章 数据可视化** .....

	本章思维导图 .....	205
8.1	数据可视化定义及分类 .....	206
8.1.1	数据可视化定义 .....	206
8.1.2	数据可视化分类 .....	206
8.2	数据可视化基础 .....	208
8.2.1	数据可视化流程 .....	208
8.2.2	可视化中的数据 .....	209

8.2.3 可视化的基本图表 .....	210
8.2.4 视图的交互 .....	211
8.3 信息可视化分类 .....	212
8.3.1 时空数据可视化 .....	212
8.3.2 层次和网络数据可视化 .....	213
8.3.3 文本和文档可视化 .....	214
8.4 在商业智能中的数据可视化应用 .....	214
8.4.1 商业智能可视化的基本元素 .....	215
8.4.2 仪表盘的设计准则 .....	215
8.5 数据可视化的实现 .....	216
8.5.1 数据可视化工具 .....	216
8.5.2 ECharts .....	217
8.5.3 Plotly .....	218
本章课后习题 .....	220
本章参考文献 .....	221
<b>第9章 大数据应用系统案例——互联网应用大数据系统构建</b> .....	<b>222</b>
本章思维导图 .....	222
9.1 互联网业务背景介绍 .....	223
9.2 案例的大数据平台技术体系架构 .....	223
9.2.1 数据采集 .....	224
9.2.2 数据存储 .....	226
9.2.3 数据计算 .....	227
9.2.4 数据应用 .....	229
本章课后习题 .....	230
本章参考文献 .....	230

# 第 1 章

## 大数据概述

### 本章思维导图

随着数据的爆炸式增长和计算机技术的迅速发展,大数据技术迎来了前所未有的发展,它使人们的生活发生新变化的同时,也给人们带来了许多挑战,包括如何存储、查询、计算这些海量数据等,因此构建一个统一的大数据平台显得尤为重要。目前业界普遍认为大数据平台应具有数据源、数据采集、存储、处理、分析、可视化及其应用这 6 个层次。Hadoop 作为一个开源的大数据平台,目前已成为大数据领域的技术标准,它具有高可靠性、高扩展性、高效性和高容错性等优点,这些优点使它应对大数据领域的大部分问题。

本章首先介绍了大数据的发展历程、大数据的定义与特征、大数据与传统数据的区别,使读者对大数据概念有个整体的了解;然后介绍了大数据平台应具备的能力和大数据平台架构,使读者对大数据平台的架构有大体的轮廓;接着介绍了 Hadoop 生态系统,使读者能够认识其基础的组件;最后介绍了大数据应用,使读者能了解目前现实生活中大数据应用的例子。本章思维导图如图 1-0 所示。

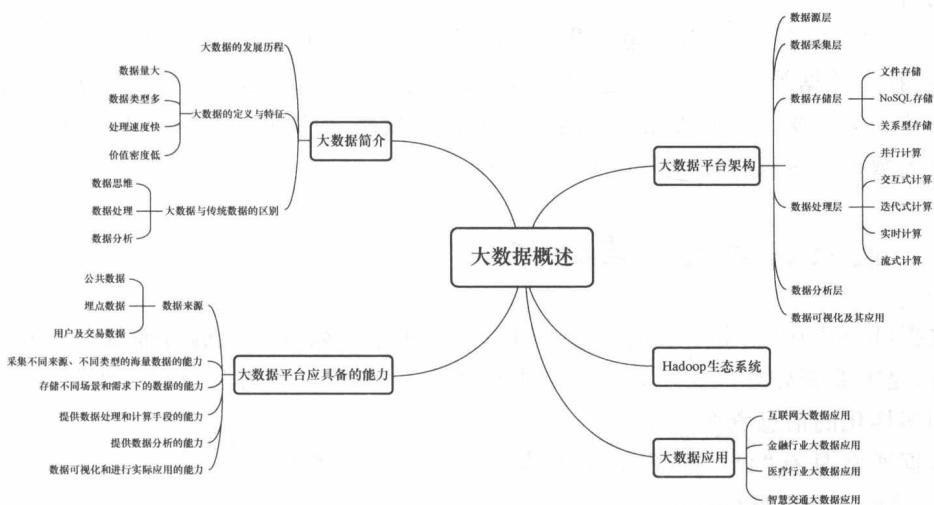


图 1-0 本章思维导图

## 1.1 大数据简介

21 世纪以来,随着计算机技术,尤其是互联网和移动技术的发展,使得数据规模呈爆炸性增长,因此“大数据”概念应运而生。大数据是继云计算、物联网之后信息技术产业领域的又一重大技术革新,它使人们的生活发生了新的变化。本节首先帮助读者更好地认识和了解大数据的发展历程、大数据的定义与特征以及大数据与传统数据的区别等<sup>[1]</sup>。

### 1.1.1 大数据的发展历程

2005 年 Hadoop 项目诞生。Hadoop 是由多个软件产品组成的一个生态系统,这些软件产品共同实现全面功能和灵活的大数据分析<sup>[2]</sup>。

2008 年 9 月, *Nature* 推出 *Big Data* 专刊<sup>[3]</sup>, 并邀请一些研究人员和企业家预测大数据所带来的革新。同年, 计算社区联盟发表了报告“Big-data computing: creating revolutionary breakthroughs in commerce, science, and society”<sup>[4]</sup>, 阐述了在数据驱动的研究背景下, 解决大数据问题所需的技术以及大数据在商业、科研和社会领域所面临的一些挑战。

2011 年 2 月, *Science* 推出 *Dealing with Data* 专刊<sup>[5]</sup>, 该专刊围绕着科学研究中大数据的问题展开讨论。麦肯锡公司在同年 5 月份发布了“Big data: the next frontier for innovation, competition, and productivity”<sup>[6]</sup>, 对大数据的影响、关键技术和应用领域等进行了详细的介绍。

2012 年 3 月, 美国政府在白宫网站发布了“Big data research and development initiative”, 这一举动标志着大数据已经成为重要的时代特征<sup>[7]</sup>。同年 7 月, 联合国在纽约发布了一本关于大数据政务的白皮书 *Big Data for Development: Opportunities & Challenges*<sup>[8]</sup>, 标志着全球大数据的研究和发展进入了前所未有的高潮阶段。

2014 年, “大数据”一词首次写入我国《政府工作报告》, 报告中指出, 要设立新兴产业创业创新平台, 在大数据等方面赶超先进, 引领未来产业发展<sup>[2]</sup>。2017 年 12 月, 习近平主席在中共中央政治局第二次会议时提出“实施国家大数据战略 加快建设数字中国”的目标, 这代表着我国对大数据的重视程度上升到了一个新的高度<sup>[9]</sup>。

拓展阅读



工业和信息化部——  
《大数据产业发展规划  
(2016—2020 年)》

### 1.1.2 大数据的定义与特征

大数据(big data)是指无法在一定时间范围内用常规软件工具进行捕捉、管理和处理的数据集合, 是需要新处理模式才能具有更强的决策力、洞察发现力和流程优化能力的海量、高增长率和多样化的信息资产<sup>[10]</sup>。

大数据通常具有“4V”特征, 即数据量大(volume)、数据类型多(variety)、处理速度快(velocity)和价值密度低(value)<sup>[7]</sup>。

① 数据体量庞大。采集、存储和计算的量都非常大。数据时代刚刚来临的时候, 一般的数