

大数据分析得天下

大数据应用赢未来

# 大数据

## 技术及应用教程

李联宁 编著

# Big Data



清华大学出版社

# 大数据

## 技术及应用教程

李联宁 编著

# Big Data



清华大学出版社  
北京

## 内 容 简 介

本书详细介绍了大数据技术的基础理论和最新主流前沿技术,全书共分为10章,分别介绍我们目前面临的数字化信息社会的大数据时代、大数据技术基本概念、云计算网络、大数据采集与预处理、大数据存储、计算模式与处理系统、查询显示与交互、大数据分析 with 数据挖掘、隐私与安全、大数据技术发展前景,同时包括行业案例研究(银行、保险、证券、金融行业),典型系统与相关大数据分析实例。

本书主要作为高等院校计算机专业、信息管理专业、经济类专业、管理类专业相关本科生和研究生专业基础课的教材,也可以作为干部培训、职业技术教育以及职业培训机构的云计算与大数据分析技术的专业训练教材。对从事云计算与大数据分析工作的财政金融、政府管理、计算机网络、软件工程的方面的管理与工程技术人员也有学习参考价值。

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

大数据技术及应用教程/李联宁编著. —北京:清华大学出版社,2016

ISBN 978-7-302-44561-6

I. ①大… II. ①李… III. ①数据处理—教材 IV. ①TP274

中国版本图书馆 CIP 数据核字(2016)第 174858 号

责任编辑:白立军 李 晔

封面设计:杨玉兰

责任校对:李建庄

责任印制:宋 林

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)

课 件 下 载: <http://www.tup.com.cn>, 010-62795954

印 装 者:北京密云胶印厂

经 销:全国新华书店

开 本:185mm×260mm 印 张:23.75 字 数:576千字

版 次:2016年10月第1版 印 次:2016年10月第1次印刷

印 数:1~2000

定 价:49.00元

---

产品编号:069155-01



本书试图在介绍大数据技术的理论基础上对大数据分析最新前沿技术做全面详细介绍,给出实际案例及行业解决方案,达到技术全面、案例教学及工程实用的目的。

本书主要分为4个部分,共10章,分别按大数据的技术架构分层次详细讲述涉及大数据分析系统的各类相关技术:

第一部分 大数据基础知识,简单介绍我们目前面临的数字化时代与信息社会的状况,大数据的定义和特点、大数据技术基础、大数据的社会价值、大数据的商业应用、大数据的基础架构、云计算网络的技术层次、典型的云计算网络平台,包括第1章“大数据技术基本概念”和第2章“基础架构——云计算网络”;

第二部分 大数据理论与技术,介绍涉及大数据分析的基本理论与技术基础,按照技术层次分别介绍大数据采集与预处理、大数据存储、大数据计算模式与处理系统、大数据查询、显示与交互、大数据分析 with 数据挖掘、大数据隐私与安全,包括第3章到第8章的内容;

第三部分 为行业案例研究,以银行、保险、证券、金融行业为例,介绍涉及大数据分析的理论与技术方法在具体行业中的应用,包括第9章“行业案例研究”;

第四部分 大数据技术发展前景,介绍大数据引发的新一代信息技术变革浪潮、大数据各个过程的最新技术与发展前景,包括第10章大数据技术发展前景。

本书主要作为高等院校计算机专业、信息管理与信息系统专业、经济类专业、管理类相关专业本科生和研究生专业基础课的教材,安排课时为48课时(3学分)。如课时缩减,可在概要叙述第一部分的基础上,主要讲解第二部分第3章到第8章的内容,并安排学生在课外自主阅读每章节后的案例及第9章“行业案例研究”。第10章“大数据技术发展前景”仅作参考性讲解。

本书的特点是紧扣实践应用需求,全面讲述云计算与大数据分析实用技术,提供了大量的实际案例、数据分析适用技术。内容新颖、用表格和结构图直观描述知识并力图反映最新主流技术。

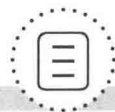
每一章在讲解相关理论外,还讲解了最新前沿技术。各章都附有案例、习题以帮助读者学习理解和实际工程应用。为方便教师教学,附有全套教学PPT课件、教学大纲、教学计划以便教师使用。

本书由李联宁教授编著,在本书编写过程中,编者参考了国内外大量的云计算网络与大数据分析技术的书刊及文献资料,主要参考书籍及研究论文在书后“参考文献”中

列出。但由于大量来自网络的资料未能详尽标注作者及文献资料来源,疏漏之处在所难免,在此一并对书刊文献、科技论文的作者表示感谢。如有遗漏,恳请相应书刊文献作者及时告知,将在书籍再版时列入。如发现本书有错误或不妥之处,恳请广大读者不吝赐教。

**编 者**

2016年8月



## 第一部分 大数据基础

第1章 大数据技术基本概念	3
1.1 数据	3
1.1.1 数据的单位	4
1.1.2 数据与信息的关系	4
1.1.3 数据的分类	4
1.2 信息	6
1.2.1 信息的定义	6
1.2.2 信息资源	7
1.2.3 信息的应用意义	8
1.3 大数据	9
1.3.1 大数据发展历史	9
1.3.2 大数据的定义和特点	10
1.4 大数据技术的基本概念	15
1.4.1 传统数据处理	15
1.4.2 大数据分析的方法理论	16
1.4.3 大数据技术	17
1.5 大数据的社会价值	21
1.5.1 大数据的社会价值体现	21
1.5.2 大数据在政府管理方面的应用	22
1.5.3 大数据在公共服务领域的应用	23
1.6 大数据的商业应用	24
1.6.1 商业大数据的类型和价值挖掘方法	24
1.6.2 全球大数据市场结构	26
1.6.3 中国大数据市场	26
1.6.4 大数据给中国带来的十大商业应用场景	27
1.7 大数据与商业模式创新	32
1.7.1 商业模式的创新特点	32
1.7.2 商业模式创新可以为企业带来什么	32

1.7.3 基于大数据分析的商业模式创新 .....	33
1.8 如何成为“大数据企业” .....	35
1.8.1 驾驭企业外部大数据 .....	35
1.8.2 成为“大数据企业” .....	36
1.8.3 如何挖掘企业大数据的价值 .....	37
1.8.4 大数据实质上是一种管理思维 .....	38
1.9 大数据应用案例之：男女嘉宾《非诚勿扰》牵手数据分析 .....	39
习题与思考题 .....	42

## 第二部分 大数据技术

<b>第2章 基础架构——云计算平台</b> .....	47
2.1 大数据处理的基础架构 .....	47
2.2 云计算网络 .....	47
2.2.1 云计算简介 .....	48
2.2.2 云计算系统的体系结构 .....	50
2.2.3 云计算服务层次 .....	55
2.2.4 云计算技术层次 .....	57
2.2.5 云计算的核心技术 .....	58
2.2.6 典型云计算平台 .....	59
2.2.7 典型的云计算系统及应用 .....	64
2.2.8 大数据平台的应用 .....	67
2.3 大数据应用案例之：在“北上广”打拼是怎样一种体验 .....	69
习题与思考题 .....	72
<b>第3章 大数据采集与预处理</b> .....	74
3.1 大数据采集概念 .....	74
3.2 数据采集来源 .....	75
3.3 大数据采集方法 .....	76
3.3.1 大数据数据采集方面新方法 .....	76
3.3.2 网页数据采集方法 .....	76
3.3.3 Web 信息数据自动采集 .....	79
3.4 导入/预处理 .....	82
3.4.1 大数据导入/预处理的过程 .....	82
3.4.2 数据清洗 .....	84
3.4.3 数据采集(ETL)技术 .....	86
3.4.4 基于大数据的数据预处理 .....	88
3.4.5 数据处理的基本流程与关键技术 .....	90

3.5 数据集成	91
3.5.1 数据集成的概念	91
3.5.2 数据集成面临的问题	92
3.6 数据变换	92
3.6.1 异构数据交换综述	93
3.6.2 异构数据分析	94
3.6.3 异构数据交换方式	97
3.6.4 异构数据交换技术	99
3.6.5 异构数据交换与集成的研究方向	103
3.7 大数据应用案例之：互联网行业哪个职位比较有前途	103
习题与思考题	107
<b>第4章 大数据存储</b>	<b>110</b>
4.1 传统数据存储	110
4.1.1 传统数据存储介质	110
4.1.2 存储的模式	112
4.2 海量数据存储的需求	113
4.3 分布式存储系统	117
4.3.1 分布式存储系统	117
4.3.2 典型系统	118
4.4 云存储	120
4.5 数据库	123
4.5.1 数据库分类	123
4.5.2 常规 SQL 结构化关系数据库	124
4.5.3 NoSQL 非结构化数据库	124
4.5.4 NoSQL 技术	126
4.5.5 大规模并行分析数据库	129
4.6 数据仓库	131
4.6.1 数据仓库的概念	131
4.6.2 数据仓库技术发展	133
4.6.3 数据仓库原理及构成	133
4.6.4 数据仓库的基本架构	136
4.6.5 数据仓库的数据存储	136
4.6.6 数据仓库的数据应用	137
4.6.7 元数据管理	138
4.7 大数据应用案例之：一场雾霾将损失多少 GDP	138
习题与思考题	141



第 5 章 大数据计算模式与处理系统	143
5.1 数据计算	143
5.1.1 离线批处理	143
5.1.2 实时交互计算	145
5.1.3 海量数据实时计算	145
5.1.4 流计算	146
5.2 聚类算法	147
5.2.1 聚类算法的分类	147
5.2.2 数据分类与聚类	147
5.3 数据集成	148
5.3.1 数据集成概述	149
5.3.2 数据集成方案	155
5.3.3 企业数据集成应用形式	157
5.3.4 企业整体解决方案	160
5.4 机器学习	161
5.4.1 机器学习的定义和例子	162
5.4.2 机器学习的范围	164
5.4.3 机器学习的方法	165
5.4.4 机器学习的应用——大数据	170
5.4.5 机器学习的子类——深度学习	172
5.4.6 机器学习的父类——人工智能	174
5.5 数据处理语言	175
5.5.1 数据分析语言 R	175
5.5.2 大数据开发语言 Python	177
5.6 大数据应用案例之：北京的人流在哪儿？用大数据看城市	179
习题与思考题	183
第 6 章 大数据查询、显现与交互	185
6.1 数据的查询	185
6.1.1 常规数据库查询结构化数据	185
6.1.2 大数据时代的数据搜索	186
6.1.3 数据库与信息检索技术的比较	188
6.1.4 数据库技术面临的 Web 数据管理问题	189
6.2 网络数据索引与查询技术	192
6.2.1 搜索引擎技术概述	192
6.2.2 Web 搜索引擎工作原理	192
6.3 大数据索引与查询技术	200
6.3.1 大数据索引和查询	200

6.3.2	大数据处理案例：登机牌、阅卷与 MapReduce .....	201
6.4	相似性搜索工具 .....	206
6.5	数据展现与交互 .....	209
6.6	数据可视化 .....	210
6.6.1	数据可视化概念 .....	210
6.6.2	数据可视化定义与方法 .....	211
6.6.3	数据可视化分析 .....	216
6.6.4	个性化精准推荐 .....	217
6.6.5	预测和预警 .....	217
6.6.6	决策分析 .....	219
6.7	知识图谱 .....	220
6.7.1	知识图谱的概念 .....	221
6.7.2	知识图谱的表示 .....	221
6.7.3	知识图谱的存储 .....	222
6.7.4	知识图谱的应用 .....	223
6.8	大数据应用案例之：数据告诉你，上海的房子都被谁买走了 .....	229
	习题与思考题 .....	233
<b>第7章</b>	<b>大数据分析 with 数据挖掘 .....</b>	<b>235</b>
7.1	大数据的分析及应用 .....	235
7.1.1	数据处理和分析的发展 .....	235
7.1.2	大数据分析面对的数据类型 .....	236
7.1.3	大数据分析 with 处理方法 .....	237
7.1.4	数据分析的步骤 .....	237
7.1.5	大数据分析应用 .....	240
7.2	数据挖掘技术 .....	242
7.2.1	数据挖掘的定义 .....	242
7.2.2	数据挖掘的常用方法 .....	244
7.2.3	数据挖掘的功能 .....	245
7.2.4	数据挖掘技术 .....	246
7.2.5	数据挖掘的流程 .....	248
7.2.6	数据挖掘的应用 .....	250
7.2.7	“大数据自动挖掘”才是大数据的真正意义 .....	251
7.3	商业智能与数据分析 .....	252
7.3.1	商业智能技术辅助决策的发展 .....	252
7.3.2	商业智能系统架构 .....	253
7.3.3	商业智能的技术体系 .....	253

7.3.4	商务智能=数据+分析+决策+利益	254
7.4	电商大数据分析技术	257
7.4.1	移动互联网应用数据分析基础	257
7.4.2	用户规模和质量	258
7.4.3	参与度分析	259
7.4.4	渠道分析	260
7.4.5	功能分析	261
7.4.6	用户属性分析	262
7.5	大数据营销业务模型	263
7.5.1	大数据对业务模式的影响	263
7.5.2	大数据时代的网络化精确营销	264
7.5.3	移动互联和大数据时代的电子商务	265
7.5.4	大数据营销的定义与特点	266
7.5.5	网络营销大数据实际操作	268
7.5.6	数据营销方法论	270
7.6	基于社会媒体的分析预测技术	273
7.6.1	基于空间大数据的社会感知	273
7.6.2	基于社会媒体的预测技术	278
7.6.3	基于消费意图挖掘的预测	279
7.6.4	基于事件抽取的预测	282
7.6.5	基于因果分析的预测	282
7.7	大数据应用案例之：如何用大数据看风水？以星巴克和海底捞的 选址为例	286
	习题与思考题	287
<b>第8章</b>	<b>大数据隐私与安全</b>	<b>290</b>
8.1	大数据面临的问题	290
8.1.1	大数据面临的安全问题	290
8.1.2	使用大数据分析安全与隐私的问题	295
8.2	大数据安全与隐私保护关键技术	296
8.2.1	基于大数据的威胁发现技术	296
8.2.2	基于大数据的认证技术	297
8.2.3	基于大数据的数据真实性分析	298
8.2.4	大数据与“安全即服务”	298
8.3	大数据安全的防护策略	298
8.4	大数据应用案例之：电影《爸爸去哪儿》大卖有前兆么？	300
	习题与思考题	305

## 第三部分 大数据分析案例

第9章 行业案例研究——银行、保险、证券、金融行业 .....	309
9.1 银行业应用 .....	309
9.1.1 大数据时代：银行如何玩转数据挖掘 .....	309
9.1.2 工商银行客户关系管理案例 .....	311
9.1.3 银行风险管理 .....	314
9.2 保险业应用 .....	318
9.2.1 保险产业拥抱“大数据时代”或带来颠覆性变革 .....	318
9.2.2 保险欺诈识别 .....	320
9.3 证券期货应用 .....	322
9.3.1 安徽使用大数据监管证券期货 .....	322
9.3.2 “大数据”分析挖出基金“老鼠仓”的启示 .....	323
9.4 金融行业应用 .....	324
9.4.1 汽车金融公司怎么实现大数据管理 .....	324
9.4.2 大数据决定互联网金融未来 .....	326
9.4.3 移动大数据在互联网金融反欺诈领域的应用 .....	329
9.5 大数据应用案例之：大吃一惊！大数据下的中国原来是这样的 .....	331

## 第四部分 大数据技术现状及发展展望

第10章 大数据技术发展前景 .....	339
10.1 大数据引发新一代信息技术变革浪潮 .....	339
10.2 大数据采集与预处理技术发展前景 .....	341
10.3 大数据存储与管理技术发展前景 .....	342
10.4 大数据计算模式与系统技术发展前景 .....	347
10.5 大数据分析 with 挖掘技术发展前景 .....	351
10.6 大数据可视化分析技术发展前景 .....	353
10.7 大数据隐私与安全技术发展前景 .....	357
10.8 大数据应用案例之：数据解读城市：北京本地人 VS 外地人 .....	360
参考文献 .....	366



## 第一部分

# 大数据基础

## 第1章 大数据技术基本概念



# 第 1 章 大数据技术基本概念

当今,信息技术为人类步入智能社会开启了大门,带动了互联网、物联网、电子商务、现代物流、网络金融等现代服务业发展,催生了车联网、智能电网、新能源、智能交通、智能城市、高端装备制造等新兴产业发展。现代信息技术正成为各行各业运营和发展的引擎。但这个引擎正面临着大数据这个巨大的考验。各种业务数据正以几何级数的形式爆发,其格式、收集、储存、检索、分析、应用等诸多问题,不再能以传统的信息处理技术加以解决,对人类实现数字社会、网络社会和智能社会带来了极大的障碍。

大数据的出现将影响各行各业以及每个人生活。以下十个事实会让你相信,每个人都必须注意大数据:

(1) 全球数据的 90% 产生于过去 2 年内。

(2) 当前数据产生的速度非常快,以今天的数据生产速度,我们可以在 2 天内生产出 2003 年以前的所有数据。

(3) 行业内获取并且存储的数据量每 1.2 年就会翻一番。

(4) 到 2020 年,全球数据量将由现在的 3.2ZB 变为 40ZB(1ZB=1024EB,1EB=1024PB,1PB=1024TB)。

(5) 仅 Google 一家搜索引擎,每秒就处理 4 万次搜索查询,一天之内更是超过 35 亿次。

(6) 最近的统计报告显示,我们每分钟在 Facebook 上贡献 180 万次赞,上传 20 万张照片。与此同时,我们每分钟还发送 2.04 亿封邮件,发送 27.8 万个推文。

(7) 每分钟大约有 100 小时的视频被传上类似 YouTube 这样的视频网站。更有趣的是,要花费 15 年才能看完一天之内被传到 YouTube 上的全部视频。

(8) AT&T 被认为是能够用单一数据库存储最多数据量的数据中心。

(9) 在美国,很多新的 IT 工作将被创造出来以处理即将到来的大数据工程潮,而每个这样的职位都将需要 3 个额外职位的支持,这将会带来总计 600 万个新增工作岗位。

(10) 全球每分钟会新增 570 个网站。这一统计数字至关重要,也具有颠覆性。

预测是:数据以及数据分析能力正与日俱增,未来五年,无论何等规模的企业都将使用某种形式的数据分析来影响其商业运作。

## 1.1 数据

数据(data)是对客观事物的逻辑归纳,用符号、字母等方式对客观事物进行直观描述。数据是进行各种统计、计算、科学研究或技术设计等所依据的数值,是表达知识的字符的集合。数据是信息的表现形式。数据可以是连续的值,例如声音,称为模拟数据;也

可以是不连续(离散)的值,例如成绩,称为数字数据。

### 1.1.1 数据的单位

数据最小的基本单位是 bit,按顺序给出所有单位: bit、Byte、KB、MB、GB、TB、PB、EB、ZB、YB、BB、NB、DB。

它们按照进率  $1024(2$  的十次方)来计算:

$$\begin{aligned}1\text{Byte} &= 8\text{bit} \\1\text{KB} &= 1024\text{Bytes} = 8192\text{bit} \\1\text{MB} &= 1024\text{KB} = 1\,048\,576\text{Bytes} \\1\text{GB} &= 1024\text{MB} = 1\,048\,576\text{KB} \\1\text{TB} &= 1024\text{GB} = 1\,048\,576\text{MB} \\1\text{PB} &= 1024\text{TB} = 1\,048\,576\text{GB} \\1\text{EB} &= 1024\text{PB} = 1\,048\,576\text{TB} \\1\text{ZB} &= 1024\text{EB} = 1\,048\,576\text{PB} \\1\text{YB} &= 1024\text{ZB} = 1\,048\,576\text{EB} \\1\text{BB} &= 1024\text{YB} = 1\,048\,576\text{ZB} \\1\text{NB} &= 1024\text{BB} = 1\,048\,576\text{YB} \\1\text{DB} &= 1024\text{NB} = 1\,048\,576\text{BB}\end{aligned}$$

### 1.1.2 数据与信息的关系

数据是一种未经加工的原始资料。数字、文字、符号、图像都是数据。数据是客观对象的表示,而信息则是数据内涵的意义,是数据的内容和解释。综上所述,数据就是指能够客观反映事实的数字和资料。

信息与数据的关系是:信息与数据是不可分离的,数据是信息的表达,信息是数据的内涵。数据本身并没有意义数据只有对实体行为产生影响时才成为信息。

### 1.1.3 数据的分类

在信息社会,信息可以划分为两大类:一类信息能够用数据或统一的结构加以表示,我们称之为结构化数据,如数字、符号;另一类信息无法用数字或统一的结构表示,如文本、图像、声音、网页等,我们称之为非结构化数据。结构化数据属于非结构化数据的一部分,是非结构化数据的特例。

#### 1. 结构化数据

结构化信息是指信息经过分析后可分解成多个互相关联的组成部分,各组成部分间有明确的层次结构,其使用和维护通过数据库进行管理,并有一定的操作规范。我们通常接触的,包括生产、业务、交易、客户信息等方面的记录都属于结构化信息。

结构化数据简单来说就是存储在结构化数据库里的数据,可以用二维表结构来逻辑表达实现的数据。结合到典型场景中更容易理解,比如企业 ERP、财务系统;医疗 HIS 数



数据库;教育一卡通;政府行政审批;其他核心数据库等。这些应用需要包括高速存储应用需求、数据备份需求、数据共享需求以及数据容灾需求。

## 2. 非结构化数据

不方便用数据库二维逻辑表来表现的数据即称为非结构化数据,包括所有格式的办公文档、文本、图片、标准通用标记语言下的子集 XML、HTML、各类报表、图像和音频/视频信息等等。

所谓非结构化数据库,是指数据库的变长记录由若干不可重复和可重复的字段组成,而每个字段又可由若干不可重复和可重复的子字段组成。用它不仅可以处理结构化数据(如数字、符号等信息)而且更适合处理非结构化数据(全文文本、图像、声音、影视、超媒体等信息)。简单地说,非结构化数据库就是字段可变的数据库。

非结构化 Web 数据库主要是针对非结构化数据而产生的,与以往流行的关系数据库相比,其最大区别在于它突破了关系数据库结构定义不易改变和数据定长的限制,支持重复字段、子字段以及变长字段并实现了对变长数据和重复字段进行处理和数据项的变长存储管理,在处理连续信息(包括全文信息)和非结构化信息(包括各种多媒体信息)中有着传统关系型数据库所无法比拟的优势。

## 3. 半结构化数据

所谓半结构化数据,就是介于完全结构化数据(如关系型数据库、面向对象数据库中的数据)和完全无结构的数据(如声音、图像文件等)之间的数据,HTML 文档就属于半结构化数据。它一般是自描述的,数据的结构和内容混在一起,没有明显的区分。

## 4. 各类数据的区别

结构化数据:行数据,存储在数据库里,可以用二维表结构来逻辑表达实现的数据。

非结构化数据:包括所有格式的办公文档、文本、图片、XML、HTML、各类报表、图像和音频/视频信息等等。

半结构化数据:介于完全结构化数据和完全无结构的数据之间的数据,它一般是自描述的,数据的结构和内容混在一起。

### 1) 数据模型

各类数据的数据模型和基本特征如下:

结构化数据:二维表(关系型)。

半结构化数据:树、图。

非结构化数据:无。

### 2) 关系型数据库系统 RMDBS 的数据模型

RMDBS 的数据模型包括网状数据模型、层次数据模型、关系型。

### 3) 不同类型数据的形成过程

结构化数据:先有结构,再有数据。

半结构化数据:先有数据,再有结构。

## 5. 互联网信息分类

互联网上出现的海量信息,同样分为结构化、半结构化和非结构化三种。