

内 容 简 介

本书详细介绍了大数据挖掘技术,全书分为3篇,共12章。第1篇为大数据分析基础,包括第1~4章,分别为大数据概述、大数据相关技术、数据预处理、R语言工具的使用。第2篇为大数据挖掘技术,包括第5~11章,分别为线性分类方法、分类方法、聚类分析、关联规则、预测方法与离群点诊断、时间序列分析、大数据挖掘可视化。第3篇为大数据挖掘案例,包括第12章,介绍了大数据挖掘应用案例。

本书既可作为高等学校计算机科学与技术、数据科学与大数据技术、统计学、数据分析等专业的高等教育教材,也可作为科研人员、从事大数据相关工作的技术人员的参考书。

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

大数据挖掘/赵志升主编. —北京:清华大学出版社,2019
(高等院校数据科学与大数据技术系列规划教材)
ISBN 978-7-302-51179-3

I. ①大… II. ①赵… III. ①数据采集—高等学校—教材 IV. ①TP274

中国版本图书馆CIP数据核字(2018)第210080号

责任编辑:刘翰鹏
封面设计:傅瑞学
责任校对:袁芳
责任印制:丛怀宇

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

质量反馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

课件下载: <http://www.tup.com.cn>, 010-62770175-4278

印 装 者:三河市君旺印务有限公司

经 销:全国新华书店

开 本:185mm×260mm

印 张:23.5

字 数:568千字

版 次:2019年3月第1版

印 次:2019年3月第1次印刷

定 价:59.00元

产品编号:078814-01



为什么要写这本书

大数据时代的到来,使我们的生活在政治、经济、社会、文化等各个领域都发生了很大的变化。如何从大数据中挖掘出隐含的丰富知识与价值,更好地得出结论并作出智能决策已成为相关工作者面临的机遇与挑战。

本书基于教育部“2016年产学研合作协同育人项目”——普开数据教学内容和课程体系改革项目,作为项目成果公开出版。

读者对象

本书适合作为高等教育“大数据处理”与“大数据分析”课程的教材,也可作为其他领域有数据分析需求的人员培训教材以及大数据从业人员的参考书。

如何阅读本书

本书首先介绍大数据,包括大数据的业务应用场景、云计算与大数据挖掘以及大数据挖掘过程。介绍了大数据相关技术,包括大数据获取、预处理、存储和处理、查询和分析、可视化技术以及主流大数据分析平台、R语言工具的使用。接着介绍了大数据挖掘常用的分类模型和算法,包括最基础的线性分类方法,分类器性能评价标准以及主要分类方法,内容包括K-近邻分类器、贝叶斯分类、神经网络与深度学习、支持向量机等,着重介绍了聚类分析、关联规则、时间序列分析、预测方法与离群点诊断以及大数据挖掘可视化常用技术。最后对各行各业的大数据挖掘应用案例进行了详细介绍。除了系统方法的理论讲解之外,我们在每一章给出了每种方法的R语言实现的实例。每一章的小结处按知识点提供了参考学习视频,可通过微信APP的扫一扫功能扫描观看。

作者分工与感谢

本书由赵志升撰写第1章、第2章、第12章,李静撰写第3~5章,梁俊花撰写第6章、第8章、第10章,赵志升、刘洋合写第7章、第9章、第11章。最终由赵志升、梁俊花统稿校对。感谢河北省人口健康工程技术研究中心医疗大数据研究室的人员参与本书的写作与实例算法实现,他们是靳晓松、王秀苹、吴仪、韩冰月、高雅静、李凯璇、李佳垚、樊亚宁、贾晓莹、傅轩昂、张艺璇、郭明磊、刘恬恬等。在编写的过程中也得到了刘艳霞、宋玉玺老师的帮助。本书参考了国内外学者的大量成果文献,在此一并表示诚挚的谢意。

勘误和支持

由于大数据挖掘是一个正在蓬勃发展的学科领域,涉及的内容宽泛且变化迅速,鉴于作者水平有限,在本书内容的安排、表述、推导等方面各种不当之处在所难免,敬请作者在阅读本书的过程中不吝赐教,以改进此书,读者的意见和建议请发至邮箱 zsbigdata@sina.com。

编者

2018年11月

目 录



第 1 篇 大数据分析基础

第 1 章 大数据概述	003
1.1 大数据的业务应用场景	003
1.1.1 大数据的产生及特征	003
1.1.2 大数据现状及趋势	004
1.1.3 大数据时代面临的技术问题	007
1.2 云计算与大数据挖掘	009
1.2.1 云计算的定义与特点	009
1.2.2 云计算与大数据	010
1.2.3 大数据挖掘	011
1.3 大数据挖掘过程概述	013
1.3.1 挖掘目标的定义与数据理解	014
1.3.2 数据准备与数据理解	014
1.3.3 过程模型的建立	016
1.3.4 过程模型的评估	017
1.3.5 模型的部署与应用	017
小结	017
习题	018
第 2 章 大数据相关技术	019
2.1 大数据获取技术	019
2.1.1 分布式数据采集系统 Flume	019
2.1.2 分布式消息队列 Kafka	022
2.1.3 Sqoop 数据转移工具	024
2.1.4 网络爬虫技术	027
2.1.5 数据预处理工具 Kettle	031
2.2 大数据存储和处理技术	031
2.2.1 数据处理架构技术演进	031
2.2.2 Hadoop 分布式存储和计算平台	032
2.2.3 流式数据计算引擎 Storm	034

2.2.4	Spark 分布式内存计算引擎	035
2.2.5	大数据部署方案简介	038
2.3	大数据查询和分析技术	038
2.3.1	SQL-on-Hadoop 技术	038
2.3.2	OLAP 分析引擎 Kylin	040
2.3.3	大数据分析技术 Mahout	041
2.3.4	大数据分析技术 Spark MLlib	042
2.3.5	其他常用分析语言比较	043
2.4	大数据可视化技术	046
2.5	主流大数据分析平台简介	049
	小结	050
	习题	050
第 3 章	数据预处理	051
3.1	数据类型、数据特征与数据质量	051
3.1.1	数据类型	051
3.1.2	数据集与数据特征	052
3.1.3	探索数据结构	053
3.1.4	数据质量相关概念与数据质量分析	054
3.2	数据采集与抽样	055
3.2.1	数据采集概述	055
3.2.2	数据采集方法与应用特性	055
3.2.3	数据抽样概述	058
3.2.4	数据抽样方法与应用特性	059
3.3	数据预处理过程	062
3.3.1	数据预处理的作用与任务	062
3.3.2	数据清洗	062
3.3.3	数据集成	065
3.3.4	数据变换	067
3.3.5	数据规约	071
3.4	Hadoop 中的数据预处理应用	074
3.4.1	使用 MapReduce 进行数据预处理	074
3.4.2	使用 Kettle 和 Python 进行数据预处理	076
	小结	079
	习题	080
第 4 章	R 语言工具的使用	082
4.1	R 语言概述	082
4.1.1	下载、安装和使用	082

4.1.2 R 包的使用	084
4.2 R 语言的基本操作	085
4.2.1 数据的基本操作	085
4.2.2 R 常用函数	093
4.3 R 语言可视化绘图	097
4.3.1 R 绘图参数设置	098
4.3.2 常用图形的绘制	099
4.4 R 语言数据分析	104
4.4.1 数据处理基础函数	104
4.4.2 多元统计分析	109
4.5 RHadoop 安装与使用	117
4.5.1 环境准备	118
4.5.2 RHadoop 安装	118
4.5.3 RHadoop 程序应用	120
小结	126
习题	126

第 2 篇 大数据挖掘技术

第 5 章 线性分类方法	131
5.1 线性分类方法综述与评价准则	131
5.1.1 线性分类方法综述	131
5.1.2 分类方法评价准则	132
5.2 多元线性回归分析	132
5.2.1 回归分析原理	132
5.2.2 多元线性回归分析 R 案例	133
5.3 逻辑回归分析	139
5.3.1 逻辑回归模型	139
5.3.2 逻辑回归分析 R 案例	139
5.4 线性判别分析	142
5.4.1 线性判别分析原理	142
5.4.2 线性判别分析 R 案例	143
5.5 应用回归树和模型树进行数值预测实例	148
小结	153
习题	154
第 6 章 分类方法	155
6.1 分类方法概要	155
6.1.1 分类的基本原理	155

6.1.2	主要分类方法	156
6.1.3	分类器性能评价标准	157
6.2	K-近邻分类器	158
6.2.1	K-近邻分类算法	158
6.2.2	K-近邻算法实例	158
6.2.3	K-近邻的特点	161
6.3	贝叶斯分类	161
6.3.1	贝叶斯概述	161
6.3.2	朴素贝叶斯分类原理	163
6.3.3	朴素贝叶斯分类实例	164
6.3.4	朴素贝叶斯的特点	166
6.4	神经网络与深度学习	166
6.4.1	神经网络基本原理	166
6.4.2	深度学习	167
6.4.3	分类实例	168
6.4.4	人工神经网络及深度学习的特点	173
6.5	支持向量机	174
6.5.1	支持向量机的基本思想	174
6.5.2	支持向量机理论基础	174
6.5.3	支持向量机实例	175
6.5.4	支持向量机的特点	180
	小结	181
	习题	181
第7章 聚类分析		183
7.1	聚类分析方法概述	183
7.1.1	聚类的基本概念	183
7.1.2	类的度量方法	186
7.1.3	聚类过程与应用	187
7.2	K-Means 聚类	189
7.2.1	K-Means 聚类的原理及步骤	189
7.2.2	K-Means 特点与适用场景	190
7.2.3	K-Means 聚类的算法实例	190
7.3	层次聚类	195
7.3.1	层次聚类的原理及步骤	195
7.3.2	层次聚类算法及特点	195
7.3.3	层次聚类的算法实例	197
7.4	神经网络聚类	199
7.4.1	SOM 算法的原理及步骤	200

7.4.2	SOM 算法实例	205
7.5	模糊 FCM 算法	207
7.5.1	FCM 算法原理和步骤	207
7.5.2	FCM 应用实例	208
7.6	并行聚类分析	215
7.6.1	并行聚类的分类	215
7.6.2	并行聚类算法流程	218
7.6.3	基于 MapReduce 聚类分析	218
7.7	其他聚类分析算法	219
	小结	223
	习题	223
第 8 章	关联规则	225
8.1	关联规则概述	225
8.1.1	关联规则的基本概念	225
8.1.2	关联规则地发现步骤	226
8.1.3	关联规则挖掘算法分类	228
8.1.4	应用场景及特点	229
8.1.5	关联规则质量评价	230
8.2	Apriori 算法	231
8.2.1	Apriori 算法的基本原理	231
8.2.2	Apriori 算法步骤	231
8.2.3	Apriori 算法的频繁项集产生实例	232
8.2.4	Apriori 算法的优缺点	241
8.3	FP-Growth 算法	242
8.3.1	FP-Growth 算法的基本思想	242
8.3.2	FP-tree 表示法	242
8.3.3	FP-Growth 算法的应用实例	243
8.3.4	FP-Growth 算法的优缺点	247
8.4	关联规则的后处理与扩展	247
8.4.1	基于 RHadoop 的关联规则挖掘	247
8.4.2	基于云计算的关联规则挖掘算法	247
8.4.3	空间数据挖掘	248
	小结	249
	习题	250
第 9 章	预测方法与离群点诊断	252
9.1	预测方法概要	252
9.1.1	预测的概念及分类	253

9.1.2	预测性能评价	254
9.1.3	常用的预测方法	255
9.2	灰色预测	256
9.2.1	灰色预测原理及应用场景	257
9.2.2	灰色预测实例	260
9.3	马尔科夫预测	262
9.3.1	马尔科夫预测原理	262
9.3.2	马尔科夫预测实例	265
9.4	离群点诊断	267
9.4.1	离群点的定义、来源及分类	267
9.4.2	各种离群点诊断技术	268
9.4.3	基于聚类的离群点技术	271
9.4.4	其他的离群点检测方法	273
	小结	276
	习题	276
第 10 章	时间序列分析	279
10.1	时间序列的基本概念	279
10.2	时间序列的组成因素及分类	280
10.3	时间序列分析方法	282
10.3.1	平稳时间序列分析方法	282
10.3.2	季节指数预测法	283
10.4	时间序列模型	283
10.4.1	ARMA 模型	283
10.4.2	ARIMA 模型	284
10.4.3	ARCH 模型	284
10.4.4	GARCH 模型	285
10.5	偏差检测	286
	小结	293
	习题	294
第 11 章	大数据挖掘可视化	296
11.1	大数据挖掘可视化概述	296
11.1.1	常规数据可视化方法	297
11.1.2	大数据可视化趋势与应用	298
11.2	数据可视化技术	300
11.3	可视化工具	302
11.3.1	常用可视化工具简介	302
11.3.2	大数据可视化面临的挑战	306

小结	307
习题	308

第 3 篇 大数据挖掘案例

第 12 章 大数据挖掘应用案例	311
12.1 社交网络分析	311
12.1.1 社交网络分析应用概述	311
12.1.2 社交网络应用案例	312
12.2 推荐系统	313
12.2.1 推荐系统概述	313
12.2.2 推荐系统应用案例	314
12.3 零售行业大数据解决方案	315
12.3.1 大数据在零售行业的创新性应用	315
12.3.2 零售行业大数据应用案例	316
12.4 金融：大数据理财时代	317
12.4.1 大数据时代下金融业的机遇和面临的挑战	317
12.4.2 金融行业大数据应用案例	317
12.4.3 信用卡反欺诈预测模型构建案例	319
12.5 临床医学大数据分析	333
12.5.1 医疗行业大数据应用	333
12.5.2 医疗行业大数据应用案例	334
12.5.3 威斯康星乳腺癌数据分析实例	335
12.6 交通行业大数据应用	347
12.6.1 大数据在智能交通行业的挑战	348
12.6.2 交通行业大数据应用案例	348
12.7 生产制造业大数据应用	349
12.7.1 大数据对生产制造业的影响及前景	349
12.7.2 生产制造业大数据应用案例	350
12.8 信息通信大数据解决方案	351
12.8.1 信息通信大数据应用	351
12.8.2 信息通信大数据应用案例	351
12.9 精准营销的大数据企业管理	352
12.9.1 大数据精准营销	352
12.9.2 精准营销大数据应用案例	353
12.9.3 基于大数据的中文舆情分析案例	354
12.10 教育领域大数据应用案例	356
12.10.1 教育领域大数据应用	356
12.10.2 教育大数据应用案例	356

12.11	互联网大数据应用	358
12.11.1	大数据使生活更智能	358
12.11.2	互联网大数据应用案例	358
12.12	其他行业大数据应用	360
12.12.1	能源业大数据应用	360
12.12.2	公共事业管理大数据应用	360
	小结	361
	习题	361
	参考文献	362

大数据概述

【内容摘要】 本章对大数据的产生及特征、现状及趋势及其面临的技术问题业务应用场景进行了简述,并对云计算与大数据挖掘进行了比较,对大数据挖掘过程进行了概述。

【学习目标】 理解大数据、云计算与大数据挖掘的基本概念与特征,掌握大数据挖掘过程与方法,了解大数据现状及趋势以及面临的技术问题。

1.1 大数据的业务应用场景

1.1.1 大数据的产生及特征

1. 什么是大数据

大数据(Big Data)或称巨量资料,是指需要用新处理模式才能具有更强的决策力、洞察力和流程优化能力的海量、高增长率和多样化的信息资产。大数据是指无法用现有的软件工具提取、存储、搜索、共享、分析和处理的海量的复杂的数据集合。

大数据是一个宽泛的概念,麦肯锡咨询公司研究大数据的先驱,在其报告 *Big data: The next frontier for innovation, competition, and productivity* 中给出的大数据定义:大数据指的是大小超出常规的数据库工具获取、存储、管理和分析能力的数据集。它同时强调,并不是说一定要超过特定 TB 值的数据集才能算是大数据。国际数据公司(IDC)从大数据的四个特征来定义,即海量的数据规模(Volume)、快速的数据流转和动态的数据体系(Velocity)、多样的数据类型(Variety)和巨大的数据价值(Value)。亚马逊公司的大数据科学家 John Rauser 给出了一个简单的定义:大数据是任何超过了一台计算机处理能力的的数据量。在维克托·迈尔-舍恩伯格及肯尼斯·库克耶编写的《大数据时代》中,大数据是指不用随机分析法(抽样调查)这样的捷径,而采用所有数据进行分析处理。

前面几个定义都是从大数据本身出发,我们的定义更关心大数据的功用,即大数据是在多样的或者大量数据中迅速获取信息的能力。在这个定义中,重心是能力。大数据的核心能力,是发现规律和预测未来。

大数据技术是指从各种各样类型的数据中快速获得有价值信息的能力。适用于大数据的技术,包括大规模并行处理(MPP)数据库、数据挖掘技术、分布式文件系统、分布式数据库、云计算平台、互联网和可扩展的存储系统。

2. 大数据的特征

作为《大数据时代》一书的作者,牛津大学网络学院互联网治理与监管专业教授、大数据权威咨询顾问维克托·迈尔-舍恩伯格博士认为,大数据有三个主要特点,分别是全体、混杂

和相关关系。

首先是全体,即收集和分析更多的数据。这个数据都是有关研究问题的数据,其中数据点绝对的数字并不重要,重要的是有多少数据点和研究的现象相关。如果想要研究的现象只有 6000 个数据点,抓住 6000 个数据点就是大数据,因为这就抓住了所有数据。通过这种方式可以看到很多细节,这些细节在之前随机抽样是得不到的。

其次是混杂,即接受混杂。在小数据时代人们总试图收集一些非常干净的数据、高质量的数据,花费很多金钱和精力来确定这些数据是好数据,是高质量的数据。可是在大数据时代,不用去追求那种特别的精确性,当宏观上失去了精确性,微观上却能获得准确性。

最后是相关关系。因为更加混杂,因果关系转向相关关系。人们不要认为可以真正地、容易地找到因果关系,其实那只是发现相关关系。我们应该关注是什么,而不是关注为什么。

业界通常用 4 个 V (Volume、Variety、Value 及 Velocity) 来概括大数据的特征。大数据呈现出“4V+1C”的特点。

(1) 数据体量巨大 (Volume)。通过各种设备产生的海量数据,其数据规模极为庞大,远大于目前互联网上的信息流量,PB 级别将是常态。截至目前,人类生产的所有印刷材料的数据量是 200PB (1PB=210TB),而历史上全人类说过的所有的话的数据量大约是 5EB (1EB=210PB)。当前,典型个人计算机硬盘的容量为 TB 量级,而一些大企业的数据量已经接近 EB 量级。

(2) 数据类型繁多 (Variety)。大数据种类繁多,在编码方式、数据格式、应用特征等多个方面存在差异性,多信息源并发形成大量的异构数据,这种类型的多样性让数据被分为结构化数据和非结构化数据。相对于以往便于存储的以文本为主的结构化数据,非结构化数据越来越多,包括网络日志、音频、视频、图片、地理位置信息等,这些多类型的数据对数据的处理能力提出了更高要求。

(3) 价值密度低 (Value)。价值密度的高低与数据总量的大小成反比,大数据量反而价值密度低。以视频为例,一部时长 1 小时的视频,在连续不间断的监控中,有用数据可能仅有一两秒。如何通过强大的机器算法更迅速地完成数据的价值“提纯”成为目前大数据背景下亟待解决的难题。

(4) 处理速度快 (Velocity)。涉及感知、传输、决策、控制开放式循环的大数据,对数据实时处理有着极高的要求,这是大数据区别于传统数据挖掘的最显著特征。根据 IDC 的“数字宇宙”的报告,预计到 2020 年,全球数据使用量将达到 35.2ZB。

(5) 数据复杂 (Complexity)。通过数据库处理持久存储的数据不再适用于大数据处理,需要有新的方法来满足异构数据统一接入和实时数据处理的需求。

1.1.2 大数据现状及趋势

数据价值的凸显和数据获取手段、数据处理技术的改进是大数据应用爆发的根源。随着数据生产要素化,数据科学、数据科技的不断发展和数据价值的深度挖掘及应用,一场大数据革命正在进行,它将带动国家战略及区域经济发展,智慧城市建设,企业转型升级,社会管理及个人工作、生活等各个领域的创新和变革。如何真正应用好大数据,发挥大数据的威力,是当前所有人都在共同研究和探索的问题。

大数据在数据科学理论的指导下,改变创新模式和理念,各个国家都积极推进大数据的战略性产业,利用大数据来提高国家的经济决策和社会服务能力,保障国家安全。互联网、物联网每天都在产生大量的数据,据调查,2015年有近200亿个设备连接到互联网上,这些设备不仅是计算机、智能终端设备,更有汽车、工厂设备、数字标牌等。从产业拓展的角度看,大数据是继云计算、物联网之后的一个新产业领域,其蕴含的机会和挑战将大大多于云计算和物联网。大数据产业(数据产业)具有很强的蜂箱效应,除了产业自身的经济蕴藏量之外,还将大大撬动其他产业的跨越升级。

2009年,联合国启动“全球脉动计划”,借大数据推动落后地区发展。美国从开放政府数据、开展关键技术研究 and 推动大数据应用三方面布局大数据产业。美国在开放政府上非常积极,通过Data.gov开放37万个数据集,并开放网站的API和源代码,提供上千个数据应用。除了推动本国政府数据开放,美国倡导发起全球开放政府数据运动,已有41个国家响应。美国联邦政府下属的国防部、能源部、卫生总署等7部委联合推动,于2012年3月底发布了大数据研发专项研究计划(Big Data Initiative),投入2亿美元用于研究开发科学探索、环境和生物医学、教育和国家安全等重大领域和行业所急需的大数据处理技术和工具,把大数据研究上升为国家发展战略。

在我国,2011年以来,中国计算机学会、中国通信学会先后成立了大数据委员会,研究大数据中的科学与工程问题。2015年9月国务院出台了《促进大数据发展行动纲要》,通过开放、产业和安全三位一体建设数据强国。2016年3月我国发布的“十三五”规划纲要又对实施网络强国战略、“互联网+”行动计划、大数据战略等作了部署。实施国家大数据战略,把大数据作为基础性战略资源,全面实施促进大数据发展行动,加快推动数据资源共享开放和开发应用,助力产业转型升级和社会治理创新。全面推进重点领域大数据高效采集、有效整合,深化政府数据和社会数据关联分析、融合利用,提高宏观调控、市场监管、社会治理和公共服务精准性和有效性。

1. “大数据资源”成为重要战略资源,将成为最有价值的资产

互联网时代,“资源”的含义正在发生极大的变化,它不仅是指煤、石油、矿产等一些看得见、摸得着的实体,大数据正在演变成不可或缺的战略资源,数据成为新的战略制高点,成为一种新的资产类别,就像货币或黄金一样。大数据已经被视为一种资产、一种财富、一种可以被衡量和计算的价值。一个国家拥有数据的规模和运用数据的能力将成为综合国力的重要组成部分,对数据的占有和控制也将成为国家间和企业间新的争夺焦点。

2. “大数据决策”成为一种新决策方式

依据大数据进行决策,从数据中获取价值,让数据主导决策,是一种前所未有的决策方式,正在推动着人类信息管理准则的重新定位。随着大数据分析和预测分析对管理决策影响力的逐渐加大,依靠直觉做决定的状况将会被彻底改变。

3. “大数据应用”促进信息技术与各行业深度融合

有专家指出,大数据及其分析会在未来10年改变几乎每一个行业的业务功能,在制造业、医疗与健康、交通、能源、材料、商业和服务等行业领域甚至在新闻传媒领域,也都在以大数据为发展契机,加速这些行业与信息技术的深度融合。大数据和传统商业智能融合产生大数据商业智能,从而形成一个全面、完整的数据价值发展平台。大数据服务提供商将会以更加定制化的适用于各行业的商业智能解决方案提供大数据服务,在业务运营智能监控、精