



大数据技术与应用专业规划教材
教育部-阿里云产学合作专业综合改革项目规划教材

大数据 挖掘与应用

◎ 王振武 编著

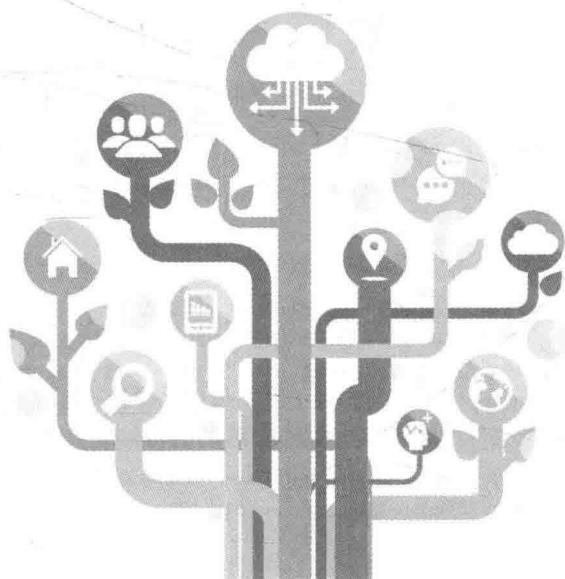


清华大学出版社

大数据技术与应用专业规划教材
教育部-阿里云产学合作专业综合改革项目规划教材

大数据 挖掘与应用

◎ 王振武 编著



清华大学出版社
北京

内 容 简 介

本书对大数据挖掘与应用的基本算法进行了系统的介绍,每种算法不仅包括对算法基本原理的介绍,而且配有大量的例题以及基于阿里云数加平台的演示,这种理论与实践相结合的方式极大地方便了读者对抽象的数据挖掘算法的理解和掌握。

本书共 17 章,内容覆盖了数据预处理、关联规则挖掘算法、分类算法和聚类算法及常见的数据挖掘应用,具体章节包括大数据简介、数据预处理技术、关联规则挖掘、逻辑回归方法、KNN 算法、朴素贝叶斯分类算法、随机森林分类算法、支持向量机、人工神经网络算法、决策树分类算法、K-means 聚类算法、K-中心点聚类算法、自组织神经网络聚类算法、DBSCAN 聚类算法以及社交网络分析方法及应用、文本分析方法及应用和推荐系统方法及应用等内容。

本书可作为高等院校数据挖掘课程的教材,也可作为从事数据挖掘工作及其他相关工程技术工作的人员的参考书。

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

大数据挖掘与应用/王振武编著. —北京: 清华大学出版社, 2017
(大数据技术与应用专业规划教材)

ISBN 978-7-302-46043-5

I. ①大… II. ①王… III. ①数据采集—高等教育—教材 IV. ①TP274

中国版本图书馆 CIP 数据核字(2016)第 315830 号

责任编辑: 刘 星 王冰飞

封面设计: 刘 键

责任校对: 梁 肖

责任印制: 王静怡

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>, 010-62795954

印 装 者: 北京鑫海金澳胶印有限公司

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 24.25 插 页: 1 字 数: 588 千字

版 次: 2017 年 6 月第 1 版 印 次: 2017 年 6 月第 1 次印刷

印 数: 1~2000

定 价: 49.50 元

产品编号: 072070-01

DT 时代的数据思维与智能思维

本套云计算大数据丛书出版正值信息科技领域进入新一轮巨变，中国经济面临转型机遇的特殊时期。全球信息科技行业伴随着云计算、大数据、物联网、人工智能的发展即将进入一个泛智能的时代，云计算成为数字经济的基础设施；数据驱动、泛在智能成为各行各业转型升级的基础，不仅传统的 IT 从业人员面临能力升级，大多数在校大学生也面临新一轮知识体系的更新，各个垂直行业面临新一轮的人才升级。新一代人才教育与培训，需要一套产学研一体的培训课程体系，这是阿里云愿意投身云计算大数据网络安全人才培养体系的时代背景。云计算、大数据、网络安全不仅关乎网络强国的大使命，也逐步成为各行各业专业人才的“元学科”，会逐步成为高等与职业教育的通识课程，一些发达国家已经在中小学立法普及编程课，已经开始指向这个趋势。“懂云计算，有数据思维，理解智能化”，未来可能是每一个工程技术人员与专业人士的必要素质。

2016 年开始，全球信息科技进入一个新的加速爆发周期，可能发生的大概率事件是：二十年之内，有一半的人类知识工作者会被人工智能替代，有服务能力的机器人会诞生，全世界的产业工人会少于机器人；虚拟现实和增强现实会替代今天的智能手机，变成一个新的入口；各行各业都会需要基于物联网的智能化，“中国制造”会成为广泛意义的“中国智造”。

新一轮科技带来了生活方式的变革，生产方式的变革，还有学习方式的变革，这几个趋势的背后，是云计算作为一种普惠科技的基础设施，大数据成为新能源，智能化成为一种新常识。

2016 年，全世界的短视频总量增长了 6 倍，直播业务在中国增长了 10 倍，远在偏远小镇的青年可以通过直播做电子商务，转化率可以提升十倍以上。当一个技术的使用成本趋近于零的时候，会带来广泛的社会效应。十年以前的直播只有电视台能做，需要专门的摄像机等设备，而今天的直播只需要一个手机，而且是多对多带互动的。无论是短视频，还是直播，背后都有云计算作为普惠科技的支撑作用，由此带来的，所有与知识传播有关的教育，包括整个内容行业，都会被它改变，随着大数据和人工智能的加入，人类学习的方式交互性会更强，“学习系统”会根据不同人的理解程度做个性化的推荐与辅导。

这意味着知识生产与知识传播方式的根本性转变，这个恰恰是云计算、人工智能等科技与各行各业产生化学反应的交叉点，数据是这个转变的新能源。

在 2016 年 10 月，阿里云和法院系统合作，发布了一个面向法律服务的智能应用“法小

“淘”，通过把数千万份法律判例文本化，“法小淘”智能应用可以为普通老百姓以及初级律师提供“打官司”的咨询服务，根据用户输入的案件信息给出建议，包括推荐合适的律师。貌似与科技远离的法律服务也用上了人工智能，这是垂直行业泛智能化的一个小例子。

II 中国制造进入智能时代

在工业界，阿里云跟中石化合作，协助他们做了企业的电商平台；与徐工合作，推动工厂基于工业云的智能化；与上汽合作，推出具有智能服务的互联网汽车，都收到积极的市场反馈。中国制造，面临智能化的产业机遇，借助互联网人口和产业布局两大优势成为未来的第一个智能产品制造国。

在接下来的几年，互联网+智能制造的叠加会在很多个垂直领域出现，数据智能与制造业结合，产生“跨界重混”的效果，甚至制造业就不是以制造为主，而是以服务化为主。这个巨大的重构背后依赖云和大数据。也因为这个需求，我们可预见到工业企业对云计算大数据人才的需求会越来越强烈。

“创业化生存”与共享经济的兴起

创业化，会成为一种常态，越来越多的年轻人开始告别公司，兴起中的数字经济体，都是基于云平台的网络化协作组织；云计算成为共享经济的超级容器，催生新一代创业者和“斜杠青年”。十年以后，或许一半以上的从业者都是“斜杠青年”，今天美国就有数千万人是跨工作、跨公司的“斜杠青年”。

过去十年，云计算使得创业公司的创业门槛降低了10倍，没有云计算，Airbnb、NetFlix、推特、Uber等公司不可能这么快成长壮大，新一代创业者的一个核心能力就是要懂技术，理解数据和算法的价值，缺少技术理解力的创业者将面临更大的同质化压力。一句话，无论是草根创业，还是做一个“斜杠青年”，必要的数据思维是生存本能。

创业化和共享经济的崛起，有赖于云计算作为基础设施，大数据作为新能源的全新范式，新一代创业公司需要大量的科技人才。

在未来的经济环境里，普惠云科技的基础设施化、制造的智能化、软件的泛化以及数据无处不在，是一个大趋势，并且不断向各行各业渗透。本套丛书就是希望在这个普惠科技与各行各业深度融合的时代为下一代科技人才的培养提供更多产业界的经验与实践。

感谢清华大学出版社出版本套云计算与大数据方面的系列教材。感谢各位高校老师的辛苦努力和用心付出，使得本系列教材能够付梓出版。

——阿里云业务总经理 刘松

序言 2

近年来,随着移动互联网、物联网以及云计算的迅猛发展,大数据成为了世界各地学术界、产业界以及政府部门的关注热点。早在 2009 年,联合国就启动了“全球脉动计划”,拟通过大数据推动落后地区的发展,2012 年,美国政府提出“大数据研究和发展倡议”,发起了全球开放政府数据的运动,2015 年,中国政府也通过了“促进大数据发展行动纲要”;在学术界,美国麻省理工大学计算机科学与人工智能实验室建立了大数据科学技术中心(ISTC),英国牛津大学成立了首个综合运用大数据的医药卫生科研中心;在产业界,IBM、Microsoft、阿里巴巴等都提出了各自的大数据解决方案或应用。

在对大数据的研究和应用过程中,数据挖掘始终占据着核心位置。与传统的数据相比,大数据的特征可以总结为 5 个 V,即体量大(Volume)、速度快(Velocity)、模态多(Variety)、难辨识(Veracity)和价值大(Value),这些特征给数据挖掘工作带来了新的挑战。

本教材以数据挖掘的研究任务为主线,系统介绍了常见的各类算法及应用。在介绍过程中,本书不仅通过深入浅出的例题讲解、完整规范的源程序实现展示了众多算法的原理,而且基于阿里云数加平台,给出了算法平台级的实现及应用,以便读者更深入地理解大数据挖掘的原理及应用。

相信此书可以帮助读者学习和掌握大数据挖掘的常用算法,更好地理解大数据挖掘的应用场景,进而推动大数据在更加广泛范围内的应用。

——中国矿业大学(北京)机电与信息工程学院院长 钱旭

大数据泛指大规模、超大规模数据集,因可从中挖掘出有价值的信息而备受关注。数据挖掘是一个涉及数据库技术、人工智能、统计学、机器学习等多个学科的领域,并且已经在各行各业有着非常广泛的应用。为适应我国数据挖掘的教学工作,笔者在数据挖掘教学实践的基础上,参阅了多种国内外最新版本的教材,编写了本书。本书可以作为高等院校研究生的教材,也可以为相关行业的工程技术人员提供有益的参考。

本书是教育部阿里云产学合作项目,在内容安排上循序渐进,对大数据挖掘的基本算法进行详细的讲解。本书的最大特点是理论与实践相结合,算法理论与产业一线实践相结合,全书几乎所有的算法都配有实例和基于阿里云数加平台的演示。这种理论与实践相结合的方法克服了重理论、轻实践的内容组织方式,极大地方便了读者的理解。具体而言,本书17章内容之间的关系如下图所示。



本书提供教学课件,读者可从 www.tup.com.cn 网站自行下载。由于编者水平有限,本书必定存在不妥和不足之处,恳请专家和读者批评指正。

目 录

第一篇 基 础 篇

第1章 大数据简介	3
1.1 大数据	3
1.1.1 大数据的定义	3
1.1.2 大数据的特点	3
1.1.3 大数据处理的挑战	4
1.2 大数据挖掘	5
1.2.1 大数据挖掘的定义	6
1.2.2 大数据挖掘的特点	6
1.3 大数据挖掘的相关方法	7
1.3.1 数据预处理技术	7
1.3.2 关联规则挖掘	7
1.3.3 分类	7
1.3.4 聚类	8
1.3.5 孤立点挖掘	8
1.3.6 演变分析	8
1.3.7 特异群组分析	8
1.4 大数据挖掘类型	9
1.4.1 Web 数据挖掘	9
1.4.2 空间数据挖掘	10
1.4.3 流数据挖掘	11
1.5 大数据挖掘的常见应用	12
1.5.1 社交网络分析	12
1.5.2 文本分析	13
1.5.3 推荐系统	13
1.6 常用的大数据统计分析方法	14
1.6.1 百分位	14
1.6.2 皮尔森相关系数	15
1.6.3 直方图	16

1.6.4 T 检验	17
1.6.5 卡方检验	20
1.7 常用的大数据挖掘评估方法	24
1.8 大数据平台相关技术	25
1.8.1 分布式存储技术	25
1.8.2 分布式任务调度技术	28
1.8.3 并行计算技术	29
1.8.4 其他技术	32
1.9 大数据平台实例——阿里云数加平台	33
1.9.1 数加平台简介	33
1.9.2 数加平台产品简介	34
1.9.3 数加平台优势特色	37
1.9.4 机器学习平台简介	37
1.9.5 机器学习平台功能	38
1.9.6 机器学习平台操作流程	39
1.10 小结	48
思考题	49

第二篇 技术篇

第2章 数据预处理技术	53
2.1 数据预处理的目的	53
2.2 数据采样	54
2.2.1 加权采样	54
2.2.2 随机采样	56
2.2.3 分层采样	56
2.3 数据清理	57
2.3.1 填充缺失值	57
2.3.2 光滑噪声数据	57
2.3.3 数据清理过程	58
2.4 数据集成	59
2.4.1 数据集成简介	59
2.4.2 常用数据集成方法	60
2.5 数据变换	61
2.5.1 数据变换简介	61
2.5.2 数据规范化	62
2.6 数据归约	63
2.6.1 数据立方体聚集	63
2.6.2 维归约	63

2.6.3 数据压缩	64
2.6.4 数值归约	65
2.6.5 数据离散化与概念分层	68
2.7 特征选择.....	70
2.7.1 特征选择简介	70
2.7.2 Relief 算法	72
2.7.3 Fisher 判别法	76
2.7.4 基于 GBDT 的过滤式特征选择	82
2.8 特征提取.....	84
2.8.1 特征提取简介	84
2.8.2 DKLST 特征提取方法	84
2.8.3 主成分分析法	86
2.9 基于阿里云数加平台的数据采样与特征选择实例	93
2.10 小结	98
思考题	98
第3章 关联规则挖掘.....	100
3.1 基本概念	100
3.2 关联规则挖掘算法——Apriori 算法原理	101
3.2.1 Apriori 算法原理解析	101
3.2.2 Apriori 算法应用举例	103
3.3 Apriori 算法源代码结果分析	105
3.4 Apriori 算法的特点及应用	111
3.4.1 Apriori 算法的特点	111
3.4.2 Apriori 算法的应用	112
3.5 小结	112
思考题	113
第4章 逻辑回归方法.....	114
4.1 基本概念	114
4.1.1 回归概述	114
4.1.2 线性回归简介	114
4.2 逻辑回归	116
4.2.1 二分类逻辑回归	116
4.2.2 多分类逻辑回归	117
4.2.3 逻辑回归应用举例	117
4.2.4 逻辑回归方法的特点	119
4.2.5 逻辑回归方法的应用	119
4.3 逻辑回归源代码结果分析	120

4.3.1 线性回归.....	120
4.3.2 多分类逻辑回归.....	123
4.4 基于阿里云数加平台的逻辑回归实例	129
4.4.1 二分类逻辑回归应用实例.....	129
4.4.2 多分类逻辑回归应用实例.....	132
4.5 小结	134
思考题.....	135
第 5 章 KNN 算法	136
5.1 KNN 算法简介	136
5.1.1 KNN 算法原理	136
5.1.2 KNN 算法应用举例	138
5.2 KNN 算法的特点及改进	141
5.2.1 KNN 算法的特点	141
5.2.2 KNN 算法的改进策略	141
5.3 KNN 源代码结果分析	142
5.4 基于阿里云数加平台的 KNN 算法应用实例	147
5.5 小结	148
思考题.....	149
第 6 章 朴素贝叶斯分类算法	150
6.1 基本概念	150
6.1.1 主观概率.....	150
6.1.2 贝叶斯定理.....	151
6.1.3 朴素贝叶斯分类模型.....	152
6.1.4 朴素贝叶斯分类器实例分析.....	154
6.2 朴素贝叶斯算法的特点及应用	156
6.2.1 朴素贝叶斯算法的特点	156
6.2.2 朴素贝叶斯算法的应用场景	157
6.3 朴素贝叶斯源代码结果分析	157
6.4 基于阿里云数加平台的朴素贝叶斯实例	162
6.5 小结	164
思考题.....	164
第 7 章 随机森林分类算法	165
7.1 随机森林算法简介	165
7.1.1 随机森林算法原理.....	165
7.1.2 随机森林算法应用举例	166
7.2 随机森林算法的特点及应用	171

7.2.1 随机森林算法的特点	171
7.2.2 随机森林算法的应用	172
7.3 随机森林算法源程序结果分析	172
7.4 基于阿里云数加平台的随机森林分类实例	184
7.5 小结	185
思考题	185
第8章 支持向量机	186
8.1 基本概念	186
8.1.1 支持向量机理论基础	186
8.1.2 统计学习核心理论	186
8.1.3 学习过程的一致性条件	186
8.1.4 函数集的 VC 维	187
8.1.5 泛化误差界	188
8.1.6 结构风险最小化归纳原理	188
8.2 支持向量机原理	189
8.2.1 支持向量机核心理论	189
8.2.2 最大间隔分类超平面	189
8.2.3 支持向量机	190
8.2.4 核函数分类	193
8.3 支持向量机的特点及应用	194
8.3.1 支持向量机的特点	194
8.3.2 支持向量机的应用	194
8.4 支持向量机分类实例分析	195
8.5 基于阿里云数加平台的支持向量机分类实例	197
8.6 小结	199
思考题	199
第9章 人工神经网络算法	200
9.1 基本概念	200
9.1.1 生物神经元模型	200
9.1.2 人工神经元模型	201
9.1.3 主要的神经网络模型	202
9.2 BP 算法的原理	204
9.2.1 Delta 学习规则的基本原理	204
9.2.2 BP 神经网络的结构	204
9.2.3 BP 神经网络的算法描述	205
9.2.4 标准 BP 神经网络的工作过程	206
9.3 BP 神经网络实例分析	207

9.4 BP 神经网络的特点及应用	210
9.4.1 BP 神经网络的特点	210
9.4.2 BP 神经网络的应用	212
9.5 BP 神经网络算法源代码结果分析	212
9.6 小结	215
思考题	215
第 10 章 决策树分类算法	216
10.1 基本概念	216
10.1.1 决策树分类算法简介	216
10.1.2 决策树基本算法概述	216
10.2 决策树分类算法——ID3 算法原理	218
10.2.1 ID3 算法原理	218
10.2.2 熵和信息增益	219
10.2.3 ID3 算法	221
10.3 ID3 算法实例分析	221
10.4 ID3 算法的特点及应用	225
10.4.1 ID3 算法的特点	225
10.4.2 ID3 算法的应用	225
10.5 ID3 算法源程序结果分析	226
10.6 决策树分类算法——C4.5 算法原理	230
10.6.1 C4.5 算法	230
10.6.2 C4.5 算法的伪代码	232
10.7 C4.5 算法实例分析	233
10.8 C4.5 算法的特点及应用	234
10.8.1 C4.5 算法的特点	234
10.8.2 C4.5 算法的应用	235
10.9 C4.5 源程序结果分析	235
10.10 小结	244
思考题	244
第 11 章 K-means 聚类算法	245
11.1 K-means 聚类算法原理	245
11.1.1 K-means 聚类算法原理解析	245
11.1.2 K-means 聚类算法应用举例	246
11.2 K-means 聚类算法的特点及应用	250
11.2.1 K-means 聚类算法的特点	250
11.2.2 K-means 聚类算法的应用	250
11.3 K 均值聚类算法源程序结果分析	250

11.4	基于阿里云数加平台的 K 均值聚类算法实例	257
11.5	基于 MaxCompute Graph 模型的 K-means 算法源程序分析.....	259
11.6	小结	264
	思考题	264
	第 12 章 K-中心点聚类算法	265
12.1	K-中心点聚类算法原理	265
12.1.1	K-中心点聚类算法原理解析	265
12.1.2	K-中心点聚类算法实例分析	266
12.2	K-中心点聚类算法的特点及应用	267
12.2.1	K-中心点聚类算法的特点	267
12.2.2	K-中心点聚类算法的应用	268
12.3	K-中心点算法源程序结果分析	268
12.4	小结	275
	思考题	275
	第 13 章 自组织神经网络聚类算法	276
13.1	SOM 网络简介	276
13.2	竞争学习算法基础	276
13.2.1	SOM 网络结构	276
13.2.2	SOM 网络概述	277
13.3	SOM 网络原理	279
13.3.1	SOM 网络的拓扑结构	279
13.3.2	SOM 权值调整域	279
13.3.3	SOM 网络运行原理	280
13.3.4	SOM 网络学习方法	281
13.4	SOM 网络应用举例	281
13.4.1	问题描述	281
13.4.2	网络设计及学习结果	282
13.4.3	输出结果分析	282
13.5	SOM 网络的特点及应用	283
13.5.1	SOM 网络的特点	283
13.5.2	SOM 网络的应用	283
13.6	SOM 神经网络源程序结果分析	284
13.7	小结	299
	思考题	300
	第 14 章 DBSCAN 聚类算法	301
14.1	DBSCAN 算法的原理	301

14.1.1 DBSCAN 算法原理解析	301
14.1.2 DBSCAN 算法应用举例	302
14.2 DBSCAN 算法的特点与应用	303
14.2.1 DBSCAN 算法的特点	303
14.2.2 DBSCAN 算法的应用	304
14.3 DBSCAN 源程序结果分析	304
14.4 小结	309
思考题	309

第三篇 综合应用篇

第 15 章 社交网络分析方法及应用	313
15.1 社交网络简介	313
15.2 K-核方法	313
15.2.1 K-核方法原理	313
15.2.2 基于阿里云数加平台的 K-核方法实例	314
15.3 单源最短路径方法	315
15.3.1 单源最短路径方法原理	315
15.3.2 基于阿里云数加平台的单源最短路径方法实例	317
15.4 PageRank 算法	318
15.4.1 PageRank 算法原理	318
15.4.2 PageRank 算法的特点及应用	319
15.4.3 基于阿里云数加平台的 PageRank 算法实例	320
15.5 标签传播算法	321
15.5.1 标签传播算法原理	321
15.5.2 基于阿里云数加平台的标签传播聚类应用实例	325
15.6 最大联通子图算法	326
15.7 聚类系数算法	328
15.7.1 聚类系数算法原理	328
15.7.2 基于阿里云数加平台的聚类系数算法应用实例	329
15.8 基于阿里云数加平台的社交网络分析实例	331
15.9 小结	335
思考题	336
第 16 章 文本分析方法及应用	337
16.1 文本分析简介	337
16.2 TF-IDF 方法	337
16.3 中文分词方法	338
16.3.1 基于字典或词库匹配的分词方法	338

16.3.2 基于词的频度统计的分词方法	339
16.3.3 其他中文分词方法	340
16.4 PLDA 方法	341
16.4.1 主题模型	341
16.4.2 PLDA 方法原理	342
16.5 Word2Vec 基本原理	344
16.5.1 词向量的表示方式	344
16.5.2 统计语言模型	344
16.5.3 霍夫曼编码	348
16.5.4 Word2Vec 原理简介	349
16.6 基于阿里云数加平台的文本分析实例	350
16.7 小结	354
思考题	354
第 17 章 推荐系统方法及应用	355
17.1 推荐系统简介	355
17.2 基于内容的推荐算法	355
17.2.1 基于内容的推荐算法原理	355
17.2.2 基于内容的推荐算法的特点	359
17.3 协同过滤推荐算法	359
17.3.1 协同过滤推荐算法简介	359
17.3.2 协同过滤推荐算法的特点	362
17.4 混合推荐算法	362
17.5 基于阿里云数加平台的推荐算法实例	364
17.6 小结	365
思考题	366
参考文献	367