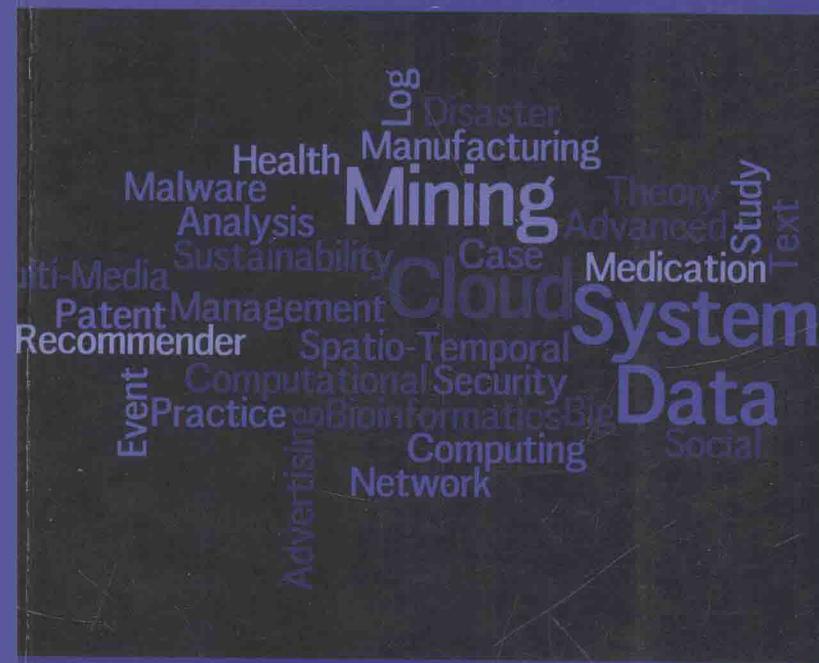


DATA MINING Where Theory Meets Practice

数据挖掘的应用与实践 ——大数据时代的案例分析



国际数据挖掘领域知名专家

李涛 等 著



厦门大学出版社 国家一级出版社
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位

DATA MINING
Where Theory Meets Practice

数据挖掘的应用与实践

——大数据时代的案例分析

李涛 等 著



厦门大学出版社 国家一级出版社
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位

图书在版编目(CIP)数据

数据挖掘的应用与实践:大数据时代的案例分析/李涛等著. —厦门:厦门大学出版社,

2013.10

ISBN 978-7-5615-4294-1

I. ①数… II. ①李… III. ①数据采集 IV. ①TP274

中国版本图书馆 CIP 数据核字(2012)第 101326 号

厦门大学出版社出版发行

(地址:厦门市软件园二期望海路 39 号 邮编:361008)

<http://www.xmupress.com>

xmup @ xmupress.com

厦门市明亮彩印有限公司印刷

2013 年 10 月第 1 版 2013 年 10 月第 1 次印刷

开本:787×1092 1/16 印张:27.25 插页:2

字数:736 千字 印数:1~2 000 册

定价:49.00 元

如有印装质量问题请与承印厂调换

内容简介

本书以笔者所带领团队的数据挖掘工作为基石，架设起研究和应用的桥梁，帮助读者们从应用实例中学习数据挖掘。

本书的宗旨是以各个领域的实际应用为导向，始终以实际案例来讲解应用之下的技术与理论。具体而言，本书对每个案例都有详细的解析，全面介绍了如何将一个实际问题抽象和转化为数据挖掘的问题，然后利用数据挖掘的理论和方法加以解决，让读者明白来龙去脉。目的是切实指导数据挖掘的应用实践，建立起研究和应用的桥梁。

本书注重原理和思想，不过多纠缠技术细节，尽量简化数学公式和模型，强调其背后的基本思想和出发点。本书不按理论和技术来划分章节，而是以实际的应用案例来贯穿始终，通过数据挖掘应用的实例来介绍如何应用和学习数据挖掘技术。

本书包括16个章节，第一章对数据挖掘做一个简单的介绍，其余的15章分为四个部分：

1. 第一部分：数据挖掘在计算机系统方面的应用，包括系统日志和事件的挖掘、数据挖掘在云计算中的应用、恶意软件智能检测；
2. 第二部分：数据挖掘在社会服务方面的应用，包括社交媒体挖掘、推荐系统、智能广告、灾难信息管理；
3. 第三部分：不同的数据类型下的典型应用，包括文本挖掘、多媒体数据挖掘、空间数据挖掘；
4. 第四部分：数据挖掘的一些综合应用，包括生物信息学和健康医疗、数据挖掘在建筑业中的应用、数据挖掘在高端制造业的应用、数据挖掘在可持续发展的应用和在专利领域中的应用。

本书各章相对独立，用户可以直接阅读跟自己具体应用领域相关的部分，而不用按照顺序进行阅读。值得一提的是，本书中的绝大部分内容是基于笔者团队的科研项目和研究积累。我们尽量提供书中涉及到原始数据和开源的软件平台，方便读者学习和使用。对笔者来说，本书既是一个阶段性的研究和应用工作总结，也是对未来大数据应用研究的铺垫。我可以无愧地说，写这本书是很严肃和很有诚意的。

前言

笔者长期从事数据挖掘研究和教学工作，经历了从最初数据挖掘基础研究的兴起，到如今数据挖掘应用百花齐放这样一个时代的变迁，深刻体会到研究和应用两者间不可分割的联系：数据挖掘研究源于实践中的实际应用需求，用具体的应用数据作为驱动，以方法、工具和系统作为支撑，最终将发现的知识和信息运用到实践中去，从而提供量化的、合理的、可行的，并且能够产生巨大价值的信息。数据挖掘是理论技术和实际应用的完美结合，所以，数据挖掘践行者们都要时刻坚定——应用是检验研究的最高标准这样的理念。

数据挖掘是大数据中最关键和最有价值的工作

国际知名的行业战略咨询公司麦肯锡（McKinsey & Company）在其2011的大数据研究报告中明确指出，在诸如卫生保健、公共部门管理、个人位置信息服务等领域，有效利用大数据能够带来每年超过千亿美元的经济价值。同时，在零售、制造等行业应用大数据解决方案，能够给企业带来相当巨大的资金效率和生产效率提升。IBM、谷歌、微软、阿里巴巴等IT巨头也将大数据描述成一种颠覆性的技术，其力量在将来足以影响和改变我们每一个人，甚至一个行业和一个国家。充分发挥大数据的巨大潜力，数据的产生和收集是基本，数据挖掘（知识发现）是工具和手段，是大数据中最关键和最有价值的工作。

在大数据时代，利用数据挖掘提升竞争力已成为各行各业都在追逐和挑战的目标，精彩的故事也层出不穷。在2012年美国总统大选中，奥巴马最终连任，其大数据挖掘与分析团队居功至伟。该团队利用两年时间收集、处理与整合了海量数据，将以往选举数据、居民基本信息、社交网络等数据整合在一个数据仓库中，利用数据挖掘算法与统计模型，预测有效选民、进行精准广告投放、优化资源配置，最终帮助奥巴马团队募集到10亿美金资金且最终赢得选举。

实践出真知

了解和学习数据挖掘，可以让我们为迎接大数据时代的激烈竞争作好准备。在长期的数据挖掘研究和教学工作中，笔者发现学习数据挖掘主要有两大难点：其一是数据挖掘是一个交叉学科，融合了统计分析、模式识别、机器学习、信息检索、数据库、信息论和最优化算法等领域的学术思想，所以其技术理论基础比较多并且分散。初学者往往很难把握数据挖掘的整个脉络，将技术理论的众多知识点系统有机地联系起来。其二是技术理论和应用实践容易脱节，初学者往往不能很好地将两者相结合。数据挖掘的应用性很强，包括诸如关联规则挖掘、时间

序列模式挖掘、分类预测、聚类分析、链接分析和异常检测等多种功能。与此同时，不同的功能通常有不同的理论和技术基础，而每一个具体的应用案例往往涉及多个不同的功能。对于有兴趣进行数据挖掘应用实践的读者们来说，他们常常有这样的困惑，如何将实际问题和已经学到的方法原理联系起来，如何将数据挖掘技术有效地运用在实际应用中，给使用者带来价值。

现今市面上已经有书籍全面地介绍数据挖掘的技术理论基础，详细解析各种挖掘算法的原理和细节。同时还有书籍专门介绍各种数据挖掘算法的实现和相关工具的使用。但这些书籍侧重于介绍单个数据挖掘功能及其相关算法原理，并没有涉及如何将数据挖掘应用到具体实践。目前，关于数据挖掘技术的应用案例分析都零星分散在一些会议论文、期刊和报告之中，并没有专门的书籍来详述。

鉴于此，笔者希望本书能填补这个空白。本书以笔者所带领团队的数据挖掘工作为基石，架设起研究和应用的桥梁，帮助读者们从应用实例中学习数据挖掘。具体而言，本书以各个领域的实际应用为导向，始终以实际案例来讲解应用之下的技术和理论。本书对每个案例都有详细的解析，全面介绍了如何将一个实际问题抽象和转化为数据挖掘的问题，然后利用数据挖掘的理论和方法加以解决，让读者明白来龙去脉。目的是切实指导数据挖掘的应用实践，建立起研究和应用的桥梁。

在编写本书时，笔者制订了两个原则。首先是内容要尽量全面，即覆盖当前数据挖掘的主要应用。在介绍每个应用案例时，详细阐述应用的背景，该领域中数据的来源和特点，数据采集与预处理方式，应用领域中数据挖掘的任务和实施数据挖掘技术的难点。同时提供相应的数据挖掘算法分析、工具设计以及系统实现。其次是要条理清晰、便于理解。一方面本书主要面向的是热爱和关心数据挖掘技术的学术界和工业界读者，帮助他们更好地理解研究的目的和应用的基础；另一方面，笔者也争取让没有太多相关技术背景的读者可以通过阅读本书了解数据挖掘的意义和价值，可以看出数据挖掘是如何被广泛地应用于实际案例并成为解决各种问题的核心工具。

笔者希望这是一本既通俗易懂，适合不同背景的读者阅读，同时又比较全面，且融入最新前沿技术和应用的数据挖掘书籍，也欢迎各大高校的师生把此书作为数据挖掘和机器学习课堂的实践教材和参考书籍。

致谢

本书中涉及的数据挖掘研究项目得到了美国佛罗里达国际大学计算机学院(School of Computer Science, Florida International University)、美国国家自然科学基金(National Science Foundation, NSF)、美国国土安全部(Department of Homeland Security, DHS)、美国军方研究实验室(Army Research Office, ARO)、中国国家自然科学研究基金、中国福建省自然科学基金、美国国际商业机器公司研究中心(IBM Research)、日本电气股份有限公司研究中心(NEC Research)、美国施乐公司研究中心(Xerox Research)、中国四川长虹电子集团公司、中国福建省厦门人才市场、中国福建省厦门市信息技术中心的资助。厦门大学信息科学与技术学院和厦门理工学院计算机与信息工程学院对本书的编写和出版给予了大力

支持，在此一并致谢！

关于作者

笔者2004年7月毕业于美国罗彻斯特大学（University of Rochester），获计算机科学博士学位。现为美国佛罗里达国际大学（Florida International University, FIU）计算机学院终身教授，数据挖掘实验室主任，是国内多家高校的客座教授。笔者长期从事关于大数据分析、数据挖掘和信息检索等方面的研究，在基于矩阵方法的数据挖掘和学习，智能推荐系统，音乐信息检索，系统日志数据挖掘，数据挖掘的各种应用等方面做出了有影响力的工作，在国际著名会议及期刊上已发表超过两百篇文章。同时，笔者是数据挖掘和知识发现的国际权威期刊ACM Transactions on Knowledge Discovery from Data (ACM TKDD), IEEE Transactions on Knowledge and Data Engineering (IEEE TKDE), 和Knowledge and Information Systems (KAIS) 的副主编。

从2005年至今，笔者作为独立的项目承担人，申请到超过230万美金的科研项目（包括五项美国自然科学基金项目以及两项美国军方实验室科研项目），作为重要的项目承担人和合作者，申请到超过1000万美金的科研项目。于2006年获美国自然科学基金委颁发的杰出青年教授奖（NSF CAREER Award, 2006-2010），并多次获得IBM学院研究奖(IBM Faculty Research Awards) 和2010 IBM大规模数据分析创新奖（2010 IBM Scalable Data Analytics Innovation Award）。

本书由笔者的数据挖掘团队成员执笔编写，欢迎读者积极反馈。各个章节的作者如下：

- 系统日志和事件的挖掘（唐良，李涛）
- 数据挖掘在云计算中的应用（姜页希，李涛）
- 数据挖掘在恶意软件检测中的研究与应用（叶艳芳，李涛）
- 社交媒体挖掘（李晶轩，李涛）
- 推荐系统（李磊，郑思婷，姜姗，洪文兴，李涛）
- 智能广告（李磊，李涛）
- 灾难信息管理（郑理，李涛）
- 文本挖掘（沈超，李涛）
- 多媒体数据挖掘（陆文婷，李晶轩，李涛）
- 空间数据挖掘（周武柏，李泓泰，朱顺痣，李涛）
- 数据挖掘在生物信息和健康医疗中的应用（曾而良，李涛）
- 数据挖掘在建筑业中的应用（周绮凤，杨帆，李涛）
- 数据挖掘在高端制造业的应用（曾春秋，郑理，李磊，李晶轩，李涛）

- 数据挖掘在可持续发展的应用（薛维，李涛）
- 数据挖掘在专利领域中的应用（张龙晖，李涛）

网站和联系方式

与本书配套的网站地址是：<http://bigdata-node01.cs.fiu.edu/dm-book>（美国网站）和<http://bigdata.xmu.edu.cn/dm-book>（中国网站）。该网站不仅收录了多个相关资源的链接，还提供了一些相关数据、程序和工具给大家下载使用。我们也欢迎大家将更多的反馈意见和修改建议发邮件到towerlee@xmu.edu.cn。

李涛

2013年09月16日

（薛维）美国佛罗里达国际大学计算机系副教授，博士生导师，IEEE Transactions on Knowledge and Data Engineering (IEEE TKDE) 和 ACM Transactions on Knowledge Discovery from Data (ACM TKDD) 编委，以及 ACM SIGKDD 和 KDD 等会议主席。主要研究兴趣包括数据挖掘、数据仓库、数据可视化、数据融合等。主持和参与了多项国家自然科学基金项目，包括面上项目“面向大规模数据流的数据挖掘方法”（2008-2010），青年基金项目“基于语义的多源异构数据集成”（2006-2008），以及 IBM 大数据奖（IBM Big Data Research Award, 2010）。

（张龙晖）厦门大学信息科学与工程学院教授，博士生导师，厦门市学术带头人，厦门市青年拔尖人才，福建省杰出青年科学基金获得者，福建省科技进步奖二等奖获得者。

（李涛）厦门大学信息科学与工程学院教授，博士生导师，厦门市学术带头人，厦门市青年拔尖人才，福建省杰出青年科学基金获得者，福建省科技进步奖二等奖获得者。

（薛维，李涛）感谢出版社对本书的支持。

（薛维，张龙晖）感谢审稿人中莫仁云和孙群慧对本书的审阅。

（薛维，张龙晖）感谢吕海燕对本书的审阅。

（薛维，孙晶华）感谢林敏文对本书的审阅。

（李涛，张龙晖）感谢陈文光对本书的审阅。

（李涛，薛维）感谢王静对本书的审阅。

（薛维，李涛）感谢高大伟对本书的审阅。

（薛维，张龙晖）感谢本文审阅人对本书的审阅。

（薛维，孙晶华，张龙晖）感谢吴晓东对本书的审阅。

（薛维，李涛）感谢朱晓东对本书的审阅。

（薛维，孙晶华，张龙晖）感谢孙晓东对本书的审阅。

（薛维，孙晶华，张龙晖）感谢孙晓东对本书的审阅。

（薛维，孙晶华，张龙晖，孙晓东）感谢高玉琪对本书的审阅。

目 录

第一章 数据挖掘简介	1
1.1 大数据时代的数据挖掘	1
1.1.1 数据挖掘	1
1.1.2 从数据挖掘应用的角度看大数据	3
1.2 数据挖掘技术发展和历史	3
1.3 十大数据挖掘算法简介	4
第二章 系统日志和事件的挖掘	8
2.1 摘要	8
2.2 系统日志分析的目的	8
2.2.1 系统问题诊断	8
2.2.2 调试与优化	9
2.2.3 系统安全维护	9
2.3 日志数据分析管理系统的架构	10
2.3.1 日志数据的收集和预处理	11
2.3.2 历史日志数据存储	11
2.3.3 日志事件数据的分析和结果展示以及使用	11
2.4 系统日志的数据形式	11
2.4.1 无结构的日志数据	12
2.4.2 结构化与半结构化的日志数据	13
2.4.3 非结构化数据的转换	14
2.5 基于日志数据的异常检测	15
2.5.1 基于监督学习的异常检测	15
2.5.2 基于无监督学习的异常检测	18

2.6 系统故障根源跟踪	21
2.6.1 日志事件的依赖性挖掘	22
2.6.2 基于依赖关系的系统故障追踪	30
2.7 日志事件总结	31
2.7.1 事件总结算法基本要求及相关工作	31
2.7.2 基于事件发生频率变迁描述的事件总结	32
2.7.3 基于马尔科夫模型描述的事件总结	32
2.7.4 基于事件关系网络描述的事件总结	33
2.8 本章小结	34
2.9 中英文术语对照表	34
参考文献	35
第三章 数据挖掘在云计算中的应用	38
3.1 摘要	38
3.2 云计算背景介绍	38
3.3 数据挖掘在云计算中的应用	39
3.4 案例介绍及困难分析：容量规划与虚拟机储备	41
3.4.1 问题背景	41
3.4.2 问题抽象与描述	42
3.4.3 预测结果评估	43
3.4.4 预测的困难性	44
3.5 案例具体分析及解决	44
3.5.1 预测困难性的体现	44
3.5.2 资源预测解决方案	46
3.5.3 数据预处理问题	47
3.5.4 预测评估标准选择	50
3.5.5 集成学习策略	52
3.6 案例分析结果	53
3.6.1 资源请求时间序列预测结果分析	53
3.6.2 资源销毁时间序列预测结果分析	54
3.6.3 虚拟机储备时间序列预测结果分析	55

3.7 本章小结.....	56
3.8 附录：时间序列分析模型介绍.....	57
3.8.1 滑动窗口平均数预测.....	57
3.8.2 自回归预测.....	57
3.8.3 人工神经网络.....	58
3.8.4 支持向量回归机.....	59
3.8.5 基因表达式编程.....	60
3.9 术语解释.....	61
参考文献.....	63
第四章 恶意软件智能检测.....	65
4.1 摘要	65
4.2 应用背景	65
4.2.1 互联网安全现状.....	65
4.2.2 “云安全”计划.....	66
4.2.3 数据挖掘在恶意软件智能检测中的应用.....	66
4.3 数据采集与预处理	67
4.3.1 恶意软件的定义	67
4.3.2 恶意软件的分类及特点	67
4.3.3 恶意软件的特征表达	68
4.4 数据挖掘的算法与实现	73
4.4.1 数据挖掘的任务	73
4.4.2 分类学习方法在恶意软件检测中的算法与实现	73
4.4.3 分类集成学习在恶意软件检测中的算法与实现	79
4.4.4 聚类及聚类融合在恶意软件检测中的算法与实现	81
4.5 系统实现	87
4.5.1 系统架构	87
4.5.2 系统实际应用效果与分析	88
4.6 本章小结	90
4.7 中英文对照表	91
参考文献	92

第五章 社交媒体挖掘	95
5.1 摘要	95
5.2 社交媒体数据挖掘简介	95
5.2.1 社交媒体分析的特点综述	96
5.2.2 社交媒体典型应用	97
5.3 社交网络数据	97
5.4 数据挖掘在社交媒体热点问题上的应用	98
5.4.1 社交媒体数据挖掘需求	99
5.4.2 信息扩散分析 (Information Diffusion)	99
5.4.3 链接的预测 (Link Prediction)	102
5.4.4 专家与关键人物的挖掘	106
5.4.5 搜索	111
5.4.6 信任 (Trust)	115
5.4.7 社交网络的内容与情感挖掘	118
5.5 本章小结	118
5.6 术语解释	119
参考文献	120
第六章 推荐系统	123
6.1 摘要	123
6.2 个性化推荐系统概述	123
6.3 推荐技术	125
6.3.1 基于内容的推荐方法	127
6.3.2 基于协同过滤的推荐方法	130
6.3.3 基于混合过滤的推荐方法	133
6.3.4 小结	135
6.4 推荐系统评测	135
6.4.1 实验环境	135
6.4.2 评测指标	138
6.4.3 小结	142

6.5 推荐系统实例	142
6.5.1 新闻推荐	142
6.5.2 人才推荐	148
6.6 推荐系统前景展望	154
6.6.1 多维度推荐	155
6.6.2 推荐中的时间动态性	156
6.7 本章小结	156
6.8 术语解释	157
参考文献	159
第七章 智能广告	163
7.1 摘要	163
7.2 引言	163
7.3 计算广告产业链介绍	164
7.3.1 广告计价模式	166
7.3.2 广告竞价模式	167
7.4 计算广告系统介绍	167
7.4.1 离线分析平台	167
7.4.2 实时投放平台	169
7.4.3 广告系统评估标准	171
7.5 搜索广告	171
7.5.1 广告索引	173
7.5.2 广告匹配模型	174
7.5.3 CTR预测与广告投放	175
7.5.4 拍卖策略	176
7.6 上下文广告	177
7.6.1 广告匹配	178
7.6.2 关键字提取	180
7.6.3 广告排序模型	180

7.7 显示广告.....	181
7.7.1 用户定位.....	182
7.7.2 CTR预测	183
7.8 本章小结.....	184
7.9 术语解释.....	184
参考文献.....	186
第八章 灾难信息管理	193
8.1 摘要.....	193
8.2 灾难管理的背景和目标.....	193
8.3 灾难管理应用中数据的特点和难点.....	194
8.4 灾难管理工作流程和工具.....	195
8.5 灾难管理数据流和功能模块.....	197
8.5.1 信息抽取（Information Extraction, IE）	197
8.5.2 信息检索（Information Retrieval, IR）	198
8.5.3 信息过滤（Information Filtering, IF）	198
8.5.4 决策支持（Decision Support, DS）	199
8.6 数据挖掘在灾难管理中的作用	199
8.7 案例分析.....	201
8.7.1 项目背景.....	201
8.7.2 数据资源	201
8.7.3 系统目标	203
8.7.4 系统实现及功能组件.....	203
8.8 算法分析和评价标准.....	205
8.8.1 定向爬虫（Focused Crawler）	205
8.8.2 信息提取（Information Extraction）	207
8.8.3 多文档文摘（Multi-Document Summarization）	208
8.8.4 动态查询（Dynamic Query Form）	208
8.8.5 动态展板（Dynamic Dashboard）	209
8.8.6 社区发现（Community Generation）	209
8.8.7 推荐（Recommendation）	210

8.9 本章小结	212
8.10 中英文对照表	212
参考文献	214
第九章 文本挖掘	216
9.1 摘要	216
9.2 文本表示(Text Representation)	216
9.3 话题挖掘(Topic Mining)	218
9.3.1 非负矩阵分解(NMF)	218
9.3.2 概率潜在语义分析(PLSA)	218
9.3.3 潜在狄利克雷分配模型(LDA)	219
9.3.4 分析与实例比较	221
9.4 多文档自动文摘	222
9.4.1 目标函数选择：句子重要性评价	222
9.4.2 优化方法	225
9.4.3 其他的自动文摘问题	226
9.4.4 实例分析	227
9.5 情感分析和摘要	229
9.5.1 基于频繁项集(frequent item set)的方法	229
9.5.2 实例分析	232
9.5.3 基于方面(Aspect-based)的话题模型分析方法	233
9.6 剧情摘要	237
9.6.1 连点成线方法(Connecting Dots)	237
9.6.2 有向施泰纳树扩展支配集方法	241
9.6.3 地铁网络模型(Metro Map)	244
9.7 本章小结	246
9.8 中英文对照表	247
参考文献	248

第十章 多媒体数据挖掘	251
10.1 摘要	251
10.2 多媒体基本概念	251
10.2.1 数字化	251
10.2.2 多样性	252
10.2.3 集成性	252
10.2.4 交互性	252
10.2.5 非线性	252
10.2.6 实时性	252
10.3 多媒体数据挖掘概述	253
10.3.1 背景	253
10.3.2 研究及应用现状	253
10.4 多媒体数据的特征抽取	254
10.4.1 文本特征抽取	254
10.4.2 图像特征表示	255
10.5 数据挖掘在图像检索中的应用	257
10.5.1 应用背景	257
10.5.2 数据集描述	258
10.5.3 数据挖掘在图像检索中的算法分析	259
10.5.4 图像检索案例	261
10.6 数据挖掘在多媒体信息融合中的应用	266
10.6.1 应用背景	266
10.6.2 数据集描述	267
10.6.3 数据挖掘在多媒体信息融合中的算法分析	268
10.6.4 多媒体信息融合案例	269
10.7 本章小结	282
10.8 中英文对照表	283
参考文献	285

第十一章 空间数据挖掘	288
11.1 简介	288
11.2 空间数据挖掘特点	288
11.3 空间位置预测	289
11.3.1 自回归模型	289
11.3.2 马尔可夫随机场模型	290
11.4 空间异常检测	290
11.5 空间同位规则挖掘	291
11.5.1 参照中心特征模型	293
11.5.2 中心窗口模型	294
11.5.3 中心事件模型	294
11.6 案例分析	294
11.6.1 TerryFly GeoCloud系统功能介绍	294
11.6.2 实际案例分析	297
11.7 空间数据挖掘最新研究方向	299
11.7.1 时空数据挖掘	301
11.7.2 移动对象数据挖掘与检索	302
11.8 本章小结	303
11.9 中英文对照表	303
参考文献	305
第十二章 生物信息学和健康医疗	308
12.1 摘要	308
12.2 生物学背景知识概述	308
12.3 数据挖掘在基因芯片数据处理中的应用	310
12.3.1 基因芯片技术概述	310
12.3.2 基因芯片的应用概述	311
12.3.3 基因表达谱芯片数据的采集与预处理	311
12.3.4 数据挖掘应用算法概述	312
12.3.5 下一代测序技术	315
12.3.6 多源生物数据融合	317