

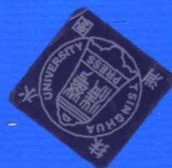
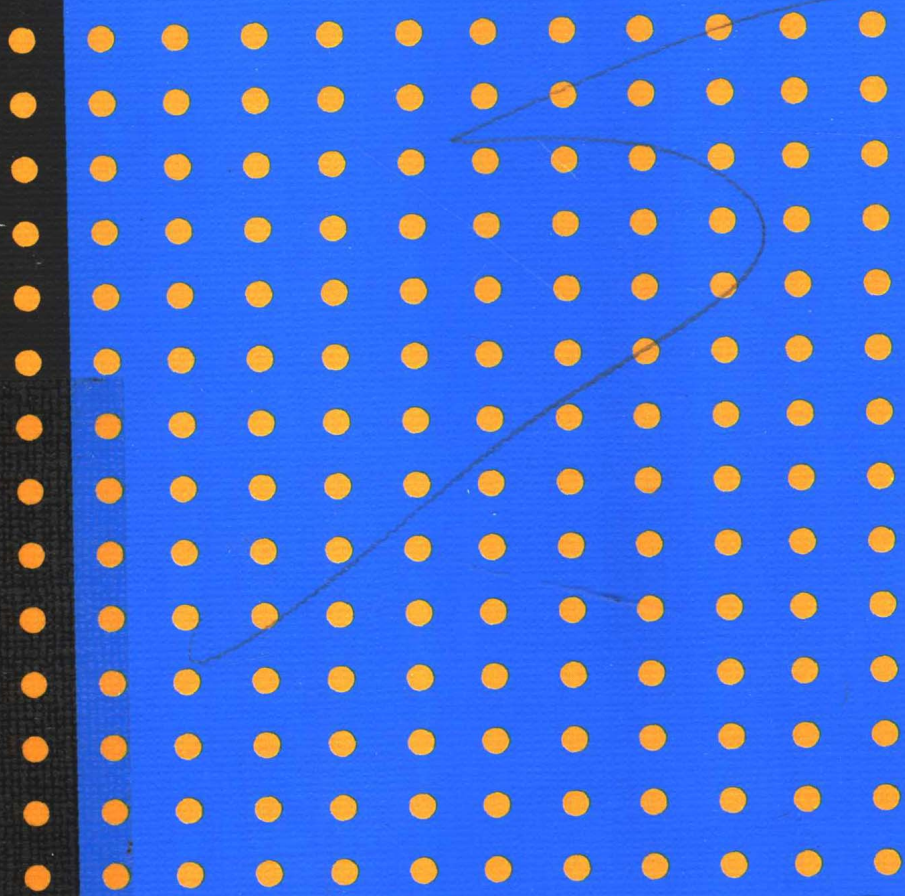


普通高等教育“十一五”国家级规划教材

重点大学计算机专业系列教材

数据挖掘原理与算法 (第二版)

毛国君 段立娟 王实 石云 编著



清华大学出版社



普通高等教育“十一五”国家规划教材

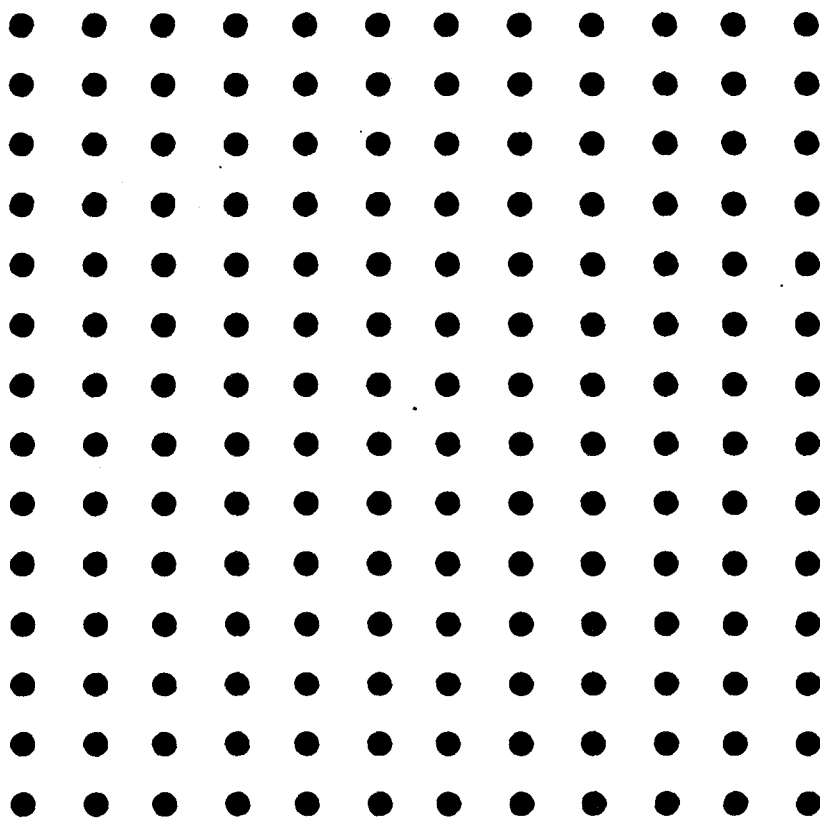
TP311.13/199=2

2007

重点大学计算机专业系列教材

数据挖掘原理与算法 (第二版)

毛国君 段立娟 王实 石云 编著



清华大学出版社
北京

内 容 简 介

本书是一本全面介绍数据挖掘和知识发现技术的专业书籍,它系统地阐述了数据挖掘和知识发现技术的产生、发展、应用以及相关概念、原理和算法,对数据挖掘中的主要技术分支,包括关联规则、分类、聚类、序列、空间以及 Web 挖掘等进行了理论剖析和算法描述。本书的许多内容是作者在攻读博士学位期间的工作总结,一方面,对于相关概念和技术的阐述尽量先从理论分析入手,在此基础上进行技术归纳;另一方面,为了保证技术的系统性,所有的挖掘模型和算法描述都在统一的技术归纳框架下进行。同时,为了避免抽象算法描述给读者带来的理解困难,本书的所有典型算法都通过具体跟踪执行实例来进一步说明。

本书共分 8 章,各章相对独立成篇,以利于读者选择性学习。在每章后面都设置专门一节来对本章内容和文献引用情况进行归纳,它不仅可以帮助读者对相关内容进行整理,而且也起到对本内容相关文献的注释性索引功能。第 1 章是绪论,系统地介绍了数据挖掘产生的商业和技术背景,从不同侧面剖析了数据挖掘的概念和应用价值;第 2 章给出了知识发现的过程分析和应用体系结构设计;第 3 章对关联规则挖掘的原理和算法进行全面阐述;第 4 章给出分类的主要理论和算法描述;第 5 章讨论聚类的常用技术和算法;第 6 章对时间序列分析技术和序列挖掘算法进行论述;第 7 章系统地介绍了 Web 挖掘的主要研究领域和相关技术及算法;第 8 章是对空间数据挖掘技术和算法的分析和讲述。

本书可作为计算机专业研究生或高年级本科生教材,也可以作为从事计算机研究和开发人员的参考资料。作为教材,教师可以根据课时安排进行选择教学。为了更好地让教师进行选择教学,本书配有专门的教师用书,对内容的重点、难点和课时分配给出了对应的建议,对重要的和难度较大的习题进行了分析和解答。对于研究人员,本书是一本高参考价值的专业书籍。对于软件技术人员,可以把它当作提高用书或参考资料,一些算法可以通过改造用于实际的应用系统中。

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 13501256678 13801310933

图书在版编目(CIP)数据

数据挖掘原理与算法/毛国君等编著. —2 版. —北京:清华大学出版社,2007. 12
(重点大学计算机专业系列教材)
ISBN 978-7-302-15876-9

I. 数… II. 毛… III. 数据采集 IV. TP274

中国版本图书馆 CIP 数据核字(2007)第 119309 号

责任编辑:丁 岭 李 晔

责任校对:时翠兰

责任印制:何 芊

出版发行:清华大学出版社 地 址:北京清华大学学研大厦 A 座

<http://www.tup.com.cn> 邮 编:100084

c-service@tup.tsinghua.edu.cn

社 总 机:010-62770175 邮购热线:010-62786544

投稿咨询:010-62772015 客户服务:010-62776969

印 刷 者:北京国马印刷厂

装 订 者:三河市溧源装订厂

经 销:全国新华书店

开 本:185×260 印 张:21.75 字 数:497 千字

版 次:2007 年 12 月第 2 版 印 次:2007 年 12 月第 1 次印刷

印 数:1~4000

定 价:29.00 元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系调换。联系电话:(010)62770177 转 3103 产品编号:027329-01

出版说明

随着国家信息化步伐的加快和高等教育规模的扩大, 社会对计算机专业人才的需求不仅体现在数量的增加上, 而且体现在质量要求的提高上, 培养具有研究和实践能力的高层次的计算机专业人才已成为许多重点大学计算机专业教育的主要目标。目前, 我国共有 16 个国家重点学科、20 个博士点一级学科、28 个博士点二级学科集中在教育部部属重点大学, 这些高校在计算机教学和科研方面具有一定优势, 并且大多以国际著名大学计算机教育为参照系, 具有系统完善的教学课程体系、教学实验体系、教学质量保证体系和人才培养评估体系等综合体系, 形成了培养一流人才的教学和科研环境。

重点大学计算机学科的教学与科研氛围是培养一流计算机人才的基础, 其中专业教材的使用和建设则是这种氛围的重要组成部分, 一批具有学科方向特色优势的计算机专业教材作为各重点大学的重点建设项目成果得到肯定。为了展示和发扬各重点大学在计算机专业教育上的优势, 特别是专业教材建设上的优势, 同时配合各重点大学的计算机学科建设和专业课程教学需要, 在教育部相关教学指导委员会专家的建议和各重点大学的大力支持下, 清华大学出版社规划并出版本系列教材。本系列教材的建设旨在“汇聚学科精英、引领学科建设、培育专业英才”, 同时以教材示范各重点大学的优秀教学理念、教学方法、教学手段和教学内容等。

本系列教材在规划过程中体现了如下一些基本组织原则和特点。

1. 面向学科发展的前沿, 适应当前社会对计算机专业高级人才的培养需求。教材内容以基本理论为基础, 反映基本理论和原理的综合应用, 重视实践和应用环节。

2. 反映教学需要, 促进教学发展。教材要能适应多样化的教学需要, 正确把握教学内容和课程体系的改革方向。在选择教材内容和编写体系时注意体现素质教育、创新能力与实践能力的培养, 为学生知识、能力、素质协调发展创造条件。

3. 实施精品战略,突出重点,保证质量。规划教材建设的重点依然是专业基础课和专业主干课;特别注意选择并安排了一部分原来基础比较好的优秀教材或讲义修订再版,逐步形成精品教材;提倡并鼓励编写体现重点大学计算机专业教学内容和课程体系改革成果的教材。

4. 主张一纲多本,合理配套。专业基础课和专业主干课教材要配套,同一门课程可以有多本具有不同内容特点的教材。处理好教材统一性与多样化的关系;基本教材与辅助教材以及教学参考书的关系;文字教材与软件教材的关系,实现教材系列资源配套。

5. 依靠专家,择优落实。在制订教材规划时要依靠各课程专家在调查研究本课程教材建设现状的基础上提出规划选题。在落实主编人选时,要引入竞争机制,通过申报、评审确定主编。书稿完成后要认真实行审稿程序,确保出书质量。

繁荣教材出版事业,提高教材质量的关键是教师。建立一支高水平的以老带新的教材编写队伍才能保证教材的编写质量,希望有志于教材建设的教师能够加入到我们的编写队伍中来。

教材编委会

前言

众所周知,数据库技术从 20 世纪 80 年代开始,已经得到广泛的应用。随着数据库容量的膨胀,特别是数据仓库以及 Web 等新型数据源的日益普及,人们面临的主要问题不再是缺乏足够的信息可以使用,而是面对浩瀚的数据海洋如何有效地利用这些数据。面对这一挑战,数据挖掘和知识发现技术应运而生,并显示出强大的生命力。数据挖掘和知识发现使数据处理技术进入了一个更高级的阶段。它不仅能对过去的数据进行查询,而且能够找出过去数据之间的潜在联系,进行更高层次的分析,以便更好地解决决策、预测等问题。历经了十几年的发展,数据挖掘技术本身已经积累了一批有价值的理论和技术成果。同时,包括统计学、人工智能等在内相关学科的发展从某种程度上对数据挖掘技术的发展起到了极大的推动作用。根据麻省理工学院的《科技评论》评估,“数据挖掘”技术是对未来人类产生重大影响的十大新兴技术之一。毫不夸张地说,如今的“数据挖掘”已经成为计算机、信息科学以及相关领域的一个时髦名词,而且在诸如银行、电信、保险、交通、零售(如超级市场)以及天文学、分子生物学等领域得到应用。

诚然,要真正理解数据挖掘技术并不是一件容易的事。一方面,数据挖掘技术覆盖范围很广泛,需要从理论到应用、从概念到算法的完整过程来理解;另一方面,作为年轻的交叉研究领域,不同背景的研究人员(数据库、人工智能、数学等)可能提供不同的视角,而且本身仍在发展中。第一作者长期从事相关方面的教学工作,其中面临的问题之一就是教材的选择。由于目前相关书籍较少,而且侧重点不同,内容的完整性和科学性有待商榷。由于没有合适教材可用,在教学的初期不得不通过指定大量参考书或文献来解决,之后也采用补充讲义的形式来扩充。同时,对于一些软件工程师或工程硕士、在职硕士进修班等要求提高实践能力的人员来说,也需要在科学的理论(原理)框架下理解和掌握数据挖掘技术。基于这样的要求,本书的第一作者在多年各类教学和软件工程的实践基础上,对积累的素材进行了整理和加工,并且邀请段立娟博士、王实博士和石云博士

进行本书的编写。本书的许多内容是几位作者在攻读博士学位期间的工作总结。这些保证了本书的系统性、先进性和实用性。

本书可作为计算机专业研究生教材、高年级本科生选修教材,也可以作为从事计算机研究和开发人员的参考资料。为了保证内容的先进性和深度,对重点内容进行了重点阐述。本书内容相对全面,各章之间耦合度小。作为教材,教师可以根据学生类型、学时安排等进行选择性教学。作为参考书,读者可以根据自己的基础进行选择学习或查阅。由于在每章后面都设置专门一节来对本章内容和文献引用情况进行归纳,它不仅可以帮助读者对相关内容进行整理,而且对读者,特别是研究人员,也起到文献的注释性索引功能。本书的所有典型算法都通过具体跟踪执行实例来进一步说明,这对于读者正确理解和应用算法是有益的。对于工程技术人员来说,这些算法完全可以通过理解基础上的改进或改造应用到实际工作中。

本书共分8章。第1章是绪论,系统地介绍了数据挖掘的概念、产生背景以及应用价值;第2章给出了知识发现的过程分析和应用体系结构设计,并对数据挖掘应用系统的主要功能部件和关键步骤进行了较为详尽的剖析;第3章全面阐述了关联规则挖掘的原理和算法,并对一些新的焦点问题(如多维、数量、约束关联规则挖掘)的最新成果尽可能地加以介绍;第4章给出分类的主要理论和算法描述;第5章讨论聚类的常用技术和算法;第6章对时间序列分析技术和序列挖掘算法进行论述;第7章系统地介绍了Web挖掘的主要研究领域和相关技术及算法;第8章是对空间数据挖掘技术和算法的分析和讲述。本书的第1~3章由毛国君执笔,第4~6章由段立娟执笔,第7章由王实执笔,第8章由石云执笔,全书由毛国君统稿。

特别感谢北京工业大学刘椿年教授和中国科学院高文和孙玉方研究员,因为他们作为作者们的导师,在作者们攻读博士学位期间对本书素材的积累提供了极大的帮助。本书也凝聚了北京工业大学硕士研究生徐启贵、鲁杰、尤春梅、邱洪君、罗春雨、孙岳、刘旭、杨霞玲、孙晓希、韩连华和一些本科学生的心血,他们在本书算法实例整理和验证等方面进行了很好的工作。此外,第一作者也感谢北京工业大学参加过相关课程学习的各类学生,它们的许多意见和文字更正,提高了本书的内容编排质量。同时第一作者也感谢所有作者及其家人,我们的合作是愉快的,家人的支持是保证本书顺利出版的必要前提。相信通过我们出色而有成效的工作将为读者提供一本有价值的专业书籍。本书对应的讲义和2005版,被许多高校作为研究生和本科生教材使用,经过近2年的尝试,使用者发现一些问题,并且提出了许多宝贵的意见。另外,本书得到了北京工业大学研究生课程建设经费的支持,这对本书的改版起到了重要的作用。对此,作者一并表示感谢。

作 者

2007年5月于北京

目录

第 1 章 绪论	1
1.1 数据挖掘技术的产生与发展	2
1.1.1 数据挖掘技术的商业需求分析	2
1.1.2 数据挖掘产生的技术背景分析	3
1.2 数据挖掘研究的发展趋势	5
1.3 数据挖掘概念	7
1.3.1 从商业角度看数据挖掘技术	7
1.3.2 数据挖掘的技术含义	7
1.3.3 数据挖掘研究的理论基础	9
1.4 数据挖掘技术的分类问题	10
1.5 数据挖掘常用的知识表示模式与方法	12
1.5.1 广义知识挖掘	12
1.5.2 关联知识挖掘	14
1.5.3 类知识挖掘	14
1.5.4 预测型知识挖掘	19
1.5.5 特异型知识挖掘	20
1.6 不同数据存储形式下的数据挖掘问题	21
1.6.1 事务数据库中的数据挖掘	21
1.6.2 关系型数据库中的数据挖掘	22
1.6.3 数据仓库中的数据挖掘	23
1.6.4 在关系模型基础上发展的新型数据库中的数 据挖掘	24
1.6.5 面向应用的新型数据源中的数据挖掘	24
1.6.6 Web 数据源中的数据挖掘	24
1.7 粗糙集方法及其在数据挖掘中的应用	26
1.7.1 粗糙集的一些重要概念	27

1.7.2	粗糙集应用举例	28
1.7.3	粗糙集方法在 KDD 中的应用范围	29
1.8	数据挖掘的应用分析	30
1.8.1	数据挖掘与 CRM	30
1.8.2	数据挖掘应用的成功案例分析	31
1.9	本章小结和文献注释	33
	习题 1	37
第 2 章	知识发现过程与应用结构	39
2.1	知识发现的基本过程	39
2.1.1	数据抽取与集成技术要点	41
2.1.2	数据清洗与预处理技术要点	41
2.1.3	数据的选择与整理技术要点	42
2.1.4	数据挖掘技术要点	42
2.1.5	模式评估技术要点	42
2.2	数据库中的知识发现处理过程模型	43
2.2.1	阶梯处理过程模型	43
2.2.2	螺旋处理过程模型	44
2.2.3	以用户为中心的处理模型	45
2.2.4	联机 KDD 模型	47
2.2.5	支持多数据源多知识模式的 KDD 处理模型	49
2.3	知识发现软件或工具的发展	52
2.3.1	独立的知识发现软件	52
2.3.2	横向的知识发现工具集	52
2.3.3	纵向的知识发现解决方案	53
2.3.4	KDD 系统介绍	53
2.4	知识发现项目的过程化管理	55
2.5	数据挖掘语言介绍	57
2.5.1	数据挖掘语言的分类	57
2.5.2	数据挖掘查询语言	58
2.5.3	数据挖掘建模语言	59
2.5.4	通用数据挖掘语言	60
2.5.5	DMQL 挖掘查询语言介绍	61
2.6	本章小结和文献注释	64
	习题 2	66
第 3 章	关联规则挖掘理论和算法	67
3.1	基本概念与解决方法	67

3.2	经典的频繁项目集生成算法分析	68
3.2.1	项目集空间理论	68
3.2.2	经典的发现频繁项目集算法	69
3.2.3	关联规则生成算法	71
3.3	Apriori 算法的性能瓶颈问题	73
3.4	Apriori 的改进算法	74
3.4.1	基于数据分割(Partition)的方法	74
3.4.2	基于散列(Hash)的方法	75
3.4.3	基于采样(Sampling)的方法	76
3.5	对项目集空间理论的发展	77
3.5.1	Close 算法	78
3.5.2	FP-tree 算法	82
3.6	项目集格空间和它的操作	85
3.7	基于项目集操作的关联规则挖掘算法	87
3.7.1	关联规则挖掘空间	87
3.7.2	三个实用算子	87
3.7.3	最大频繁项目集格的生成算法	89
3.7.4	ISS-DM 算法执行示例	89
3.8	改善关联规则挖掘质量问题	90
3.8.1	用户主观层面	90
3.8.2	系统客观层面	91
3.9	约束数据挖掘问题	91
3.9.1	约束在数据挖掘中的作用	91
3.9.2	约束的类型	92
3.10	时态约束关联规则挖掘	95
3.11	关联规则挖掘中的一些更深入的问题	98
3.11.1	多层次关联规则挖掘	98
3.11.2	多维关联规则挖掘	99
3.11.3	数量关联规则挖掘	100
3.12	数量关联规则挖掘方法	101
3.12.1	数量关联规则挖掘问题	101
3.12.2	数量关联规则的分类	102
3.12.3	数量关联规则挖掘的一般步骤	103
3.12.4	数值属性离散化问题及算法	106
3.13	本章小结和文献注释	109
	习题 3	111

第 4 章 分类方法	114
4.1 分类的基本概念与步骤	115
4.2 基于距离的分类算法	117
4.3 决策树分类方法	120
4.3.1 决策树基本算法概述.....	120
4.3.2 ID3 算法	122
4.3.3 C4.5 算法	128
4.4 贝叶斯分类	132
4.4.1 贝叶斯定理.....	132
4.4.2 朴素贝叶斯分类.....	133
4.4.3 EM 算法	135
4.5 规则归纳	139
4.5.1 AQ 算法	140
4.5.2 CN2 算法	143
4.5.3 FOIL 算法	150
4.6 与分类有关的其他问题	155
4.6.1 分类数据预处理.....	155
4.6.2 分类器性能表示与评估.....	156
4.7 本章小结和文献注释	158
习题 4	160
第 5 章 聚类方法	164
5.1 概述	164
5.1.1 聚类分析在数据挖掘中的应用.....	166
5.1.2 聚类分析算法的概念与基本分类.....	166
5.1.3 距离与相似性的度量.....	169
5.2 划分聚类方法	172
5.2.1 k -平均算法.....	172
5.2.2 PAM	175
5.2.3 其他方法.....	179
5.3 层次聚类方法	179
5.3.1 AGNES 算法	180
5.3.2 DIANA 算法	181
5.3.3 其他聚类方法.....	183
5.4 密度聚类方法	184

5.5	其他聚类方法	188
5.5.1	STING 算法	188
5.5.2	SOM 算法	189
5.5.3	COBWEB 算法	189
5.5.4	模糊聚类算法 FCM	190
5.6	本章小结和文献注释	190
	习题 5	192
第 6 章	时间序列和序列模式挖掘	194
6.1	时间序列及其应用	194
6.2	时间序列预测的常用方法	195
6.2.1	确定性时间序列预测方法	195
6.2.2	随机时间序列预测方法	196
6.2.3	其他方法	196
6.3	基于 ARMA 模型的序列匹配方法	196
6.3.1	基本概念	196
6.3.2	利用基本概念建立模型	197
6.3.3	构造判别函数	198
6.4	基于离散傅里叶变换的时间序列相似性查找	199
6.4.1	完全匹配	200
6.4.2	子序列匹配	201
6.5	基于规范变换的查找方法	203
6.5.1	基本概念	204
6.5.2	查找方法	204
6.6	序列挖掘	206
6.6.1	基本概念	207
6.6.2	数据源的形式	207
6.6.3	序列模式挖掘的一般步骤	209
6.7	AprioriAll 算法	210
6.8	AprioriSome 算法	213
6.9	GSP 算法	217
6.10	本章小结和文献注释	219
	习题 6	222
第 7 章	Web 挖掘技术	224
7.1	Web 挖掘的意义	224
7.2	Web 挖掘的分类	225
7.3	Web 挖掘的含义	227

7.3.1	Web 挖掘与信息检索	227
7.3.2	Web 挖掘与信息抽取	227
7.4	Web 挖掘的数据来源	228
7.4.1	服务器日志数据	228
7.4.2	在线市场数据	229
7.4.3	Web 页面	229
7.4.4	Web 页面超链接关系	230
7.4.5	其他信息	230
7.5	Web 内容挖掘方法	230
7.5.1	爬虫与 Web 内容挖掘	231
7.5.2	虚拟的 Web 视图	231
7.5.3	个性化与 Web 内容挖掘	232
7.5.4	对 Web 页面内文本信息的挖掘	232
7.5.5	对 Web 页面内多媒体信息挖掘	233
7.5.6	Web 页面内容的预处理	233
7.6	Web 访问信息挖掘方法	234
7.6.1	Web 访问信息挖掘的特点	234
7.6.2	Web 访问信息挖掘的意义	236
7.6.3	Web 访问信息挖掘的数据源	237
7.6.4	Web 访问信息挖掘的预处理	240
7.6.5	其他信息的预处理技术	244
7.6.6	在 Web 访问挖掘中的常用技术	246
7.6.7	Web 访问信息挖掘的要素构成	247
7.6.8	利用 Web 访问信息挖掘实现用户建模	248
7.6.9	利用 Web 访问信息挖掘发现导航模式	250
7.6.10	利用 Web 访问信息挖掘改进访问效率	252
7.6.11	利用 Web 访问信息挖掘进行个性化服务	253
7.6.12	利用 Web 访问信息挖掘进行商业智能发现	255
7.6.13	利用 Web 访问信息挖掘进行用户移动模式发现	256
7.6.14	利用协作推荐的方法实现实时个性化推荐的例子	257
7.7	Web 结构挖掘方法	260
7.7.1	页面重要性的评价方法	260
7.7.2	页面等级	261
7.7.3	权威页面和中心页面	261
7.7.4	Web 站点结构的预处理	262
7.8	本章小结和文献注释	264
	习题 7	268

第 8 章 空间挖掘	271
8.1 引言	271
8.2 空间数据概要	272
8.2.1 空间数据的复杂性特征.....	272
8.2.2 空间查询问题.....	273
8.2.3 空间数据结构.....	274
8.2.4 专题地图.....	278
8.3 空间数据挖掘基础	278
8.4 空间统计学	280
8.5 泛化与特化	281
8.5.1 逐步求精.....	281
8.5.2 泛化.....	281
8.5.3 最临近方法.....	283
8.5.4 统计信息网格方法 STING	283
8.6 空间规则	285
8.7 空间分类算法	287
8.7.1 ID3 扩展	287
8.7.2 空间决策树.....	287
8.8 空间聚类算法	288
8.8.1 基于随机搜索的聚类方法 CLARANS 扩展	289
8.8.2 大型空间数据库基于距离分布的聚类算法 DBCLASD	290
8.8.3 BANG	291
8.8.4 小波聚类.....	291
8.8.5 近似值.....	291
8.9 空间挖掘的其他问题	293
8.10 空间数据挖掘原型系统介绍.....	296
8.11 空间数据挖掘的研究现状.....	298
8.12 空间数据挖掘的研究与发展方向.....	299
8.13 空间数据挖掘与相关学科的关系.....	302
8.13.1 空间数据挖掘与空间数据库.....	302
8.13.2 空间数据挖掘与空间数据仓库.....	302
8.13.3 空间数据挖掘与空间联机分析处理.....	303
8.13.4 空间数据挖掘与地理信息系统.....	303
8.14 数字地球.....	304
8.15 本章小结和文献注释.....	304
习题 8	307
参考文献	309

绪 论

第 1 章

数据挖掘(Data Mining)是一个多学科交叉研究领域,它融合了数据库(Database)技术、人工智能(Artificial Intelligence)、机器学习(Machine Learning)、统计学(Statistics)、知识工程(Knowledge Engineering)、面向对象方法(Object-Oriented Method)、信息检索(Information Retrieval)、高性能计算(High-Performance Computing)以及数据可视化(Data Visualization)等最新技术的研究成果。经过十几年的研究,产生了许多新概念和新方法。特别是最近几年,一些基本概念和方法趋于清晰,它的研究正向着更深的方向发展。

数据挖掘之所以被称为未来信息处理的骨干技术之一,主要在于它以一种全新的概念改变着人类利用数据的方式。20世纪,数据库技术取得了决定性的成果并且已经得到广泛的应用。但是,数据库技术作为一种基本的信息存储和管理方式,仍然以联机事务处理(On-Line Transaction Processing, OLTP)为核心应用,缺少对决策、分析、预测等高级功能的支持机制。众所周知,随着数据库容量的膨胀,特别是数据仓库(Data Warehouse)以及 Web 等新型数据源的日益普及,联机分析处理(On-Line Analytic Processing, OLAP)、决策支持(Decision Support)以及分类(Classification)、聚类(Clustering)等复杂应用成为必然。面对这一挑战,数据挖掘和知识发现(Knowledge Discovery)技术应运而生,并显示出强大的生命力。数据挖掘和知识发现使数据处理技术进入了一个更高级的阶段。它不仅能对过去的数据进行查询,并且能够找出过去数据之间的潜在联系,进行更高层次的分析,以便更好地做出理想的决策、预测未来的发展趋势等。通过数据挖掘,有价值的知识、规则或高层次的信息就能从数据库的相关数据集合中抽取出来,从而使大型数据库作为一个丰富、可靠的资源为知识的提取服务。

特别需要指出的是,数据挖掘技术从一开始就是面向应用的。它不仅仅是面向特定数据库的简单检索查询应用,而是要对这些数据进行微观、中观乃至宏观的统计、分析、综合和推理,进而发现潜在的知识。这里所说的知识发现,不是要求发现放之四海而皆准的真理,也不是要去发现崭新

的自然科学定理和纯数学公式。所有发现的知识都是相对的,是面向特定领域的,同时还要能够易于被用户理解。

1.1 数据挖掘技术的产生与发展

1.1.1 数据挖掘技术的商业需求分析

数据挖掘之所以吸引专家学者的研究兴趣和引起商业厂家的广泛关注,主要在于大型数据系统的广泛使用和把数据转换成有用知识的迫切需要。20世纪60年代,为了适应信息的电子化要求,信息技术一直从简单的文件处理系统向有效的数据库系统变革。20世纪70年代,数据库系统的三个主要模式:层次、网络和关系型数据库的研究和开发取得了重要进展。20世纪80年代,关系型数据库及其相关的数据模型工具、数据索引及数据组织技术被广泛采用,并且成为了整个数据库市场的主导。从20世纪80年代中期开始,关系型数据库技术和新型技术的结合成为数据库研究和开发的重要标志。从数据模型上看,诸如扩展关系、面向对象、对象-关系(Object-Relation)以及演绎模型等被应用到数据库系统中。从应用的数据类型上看,包括空间、时态、多媒体以及Web等新型数据成为数据库应用的重要数据源。同时,事务数据库(Transaction Database)、主动数据库(Active Database)、知识库(Knowledge Base)、办公信息库(Information Base)等技术也得到蓬勃发展。从数据的分布角度看,分布式数据库(Distributed Database)及其透明性、并发控制、并行处理等成为必须面对的课题。进入90年代,分布式数据库理论上趋于成熟,分布式数据库技术得到了广泛应用。目前,由于各种新型技术与数据库技术的有机结合,使数据库领域中的新内容、新应用、新技术层出不穷,形成了庞大的数据库家族。但是,这些数据库的应用都是以实时查询处理技术为基础的。从本质上说,查询是对数据库的被动使用。由于简单查询只是数据库内容的选择性输出,因此它和人们期望的分析预测、决策支持等高级应用仍有很大距离。

新的需求推动新的技术的诞生。随着信息技术的高速发展,数据库应用的规模、范围和深度不断扩大,已经从单台机器发展到网络环境。近年来由于数据采集技术的更新,如商业条码的推广、企业和政府利用计算机管理事务的能力增强,产生了大规模的数据。数以百万计的数据库系统在运行,而且每天都在增加。决策所面对的数据量在不断增长,即使像使用IC卡和打电话这样简单的事务也能产生大量的数据。随着数据的急剧增长,现有信息管理系统中的数据分析工具已无法适应新的需求。因为无论是查询、统计还是报表,其处理方式都是对指定的数据进行简单的数字处理,而不能对这些数据所包含的内在信息进行提取。人们希望能够提供更高层次的数据分析功能,自动和智能地将待处理的数据转化为有用的信息和知识。

数据挖掘的基础是数据分析方法。数据分析是科学研究的基础,许多科学研究都是建立在数据收集和分析基础上的。同时目前的商业活动中,数据分析总是和一些特殊的人群的高智商行为联系起来,因为并不是每个人都能从过去的销售情况预测将来的发展趋势或作出正确决策的。但是,随着一个企业或行业业务数据的不断积累,特别是由于数据库的普及,人工去整理和理解如此大的数据源已经存在效率、准确性等问题。因此,

探讨自动化的数据分析技术,为企业能提供带来商业利润的决策信息就成为了必然。

事实上,可以将数据(Data)、信息(Information)和知识(Knowledge)看作是广义数据表现的不同形式。毫不夸张地说,人们对于数据的拥有欲是贪婪的,特别是计算机存储技术和网络技术的发展加速了人们收集数据的范围和容量。这种贪婪的结果导致了“数据丰富而信息贫乏(Data Rich & Information Poor)”现象的产生。数据库是目前组织和存储数据的最有效方法之一,但是面对日益膨胀的数据,数据库查询技术已表现出它的局限性。直观上说,信息或称有效信息是指对人们有帮助的数据。例如,在现实社会中,如果人均日阅读时间在 30 分钟的话,一个人一天最快只能浏览一份 20 版左右的报纸。如果你订阅了 100 份报纸,其实你每天也不过只阅读了一份而已。面对计算机中的海量数据,人们也处于同样的尴尬境地,缺乏获取有效信息的手段。知识是一种概念、规则、模式和规律等,它不会像数据或信息那么具体,但是它却是人们一直不懈追求的目标。事实上,在我们的生活中,人们只是把数据看作是形成知识的源泉。我们是通过正面的或反面的数据或信息来形成和验证知识的,同时又不断地利用知识来获得新的信息。因此,随着数据的膨胀和技术环境的进步,人们对联机决策和分析等高级信息处理的要求越来越迫切。在强大的商业需求的驱动下,商家们开始注意到有效地解决大容量数据的利用问题具有巨大的商机。学者们开始思考如何从大容量数据集中获取有用信息和知识的方法。因此,在 20 世纪 80 年代后期,产生了数据仓库和数据挖掘等信息处理思想。

1.1.2 数据挖掘产生的技术背景分析

任何技术的产生总是有它的技术背景的。数据挖掘技术的提出和普遍接受是由于计算机及其相关技术的发展为其提供了研究和应用的技术基础。

归纳数据挖掘产生的技术背景,下面一些相关技术的发展起到了决定性的作用:

- 数据库、数据仓库和 Internet 等信息技术的发展。
- 计算机性能的提高和先进的体系结构的发展。
- 统计学和人工智能等方法在数据分析中的研究和应用。

数据库技术从 20 世纪 80 年代开始,已经得到广泛的普及和应用。在关系型数据库的研究和产品提升过程中,人们一直在探索组织大型数据和快速访问的相关技术。高性能关系型数据库引擎以及相关的分布式查询、并发控制等技术的使用,已经提升了数据库的应用能力。在数据的快速访问、集成与抽取等问题的解决上积累了经验。数据仓库作为一种新型的数据存储和处理手段,被数据库厂商普遍接受并且相关辅助建模和管理工具快速推向市场,成为多数据源集成的一种有效的技术支撑环境。另外,Internet 的普及也为人们提供了丰富的数据源。据说,在美国,电视普及达到 5000 万户大约用了 15 年,而 Internet 上网普及达到 5000 万户仅用了 4 年。而且 Internet 技术本身的发展,已经不光是简单的信息浏览,以 Web 计算为核心的信息处理技术可以处理 Internet 环境下的多种信息源。因此,人们已经具备利用多种方式存储海量数据的能力。只有这样,数据挖掘技术才能有它的用武之地。这些丰富多彩的数据存储、管理以及访问技术的发展,为数据挖掘技术的研究和应用提供了丰富的土壤。

计算机芯片技术的发展,使计算机的处理和存储能力日益提高。大家熟知的摩尔定律告诉我们,计算机硬件的关键指标大约以每 18 个月翻一番的速度在增长,而且现在看