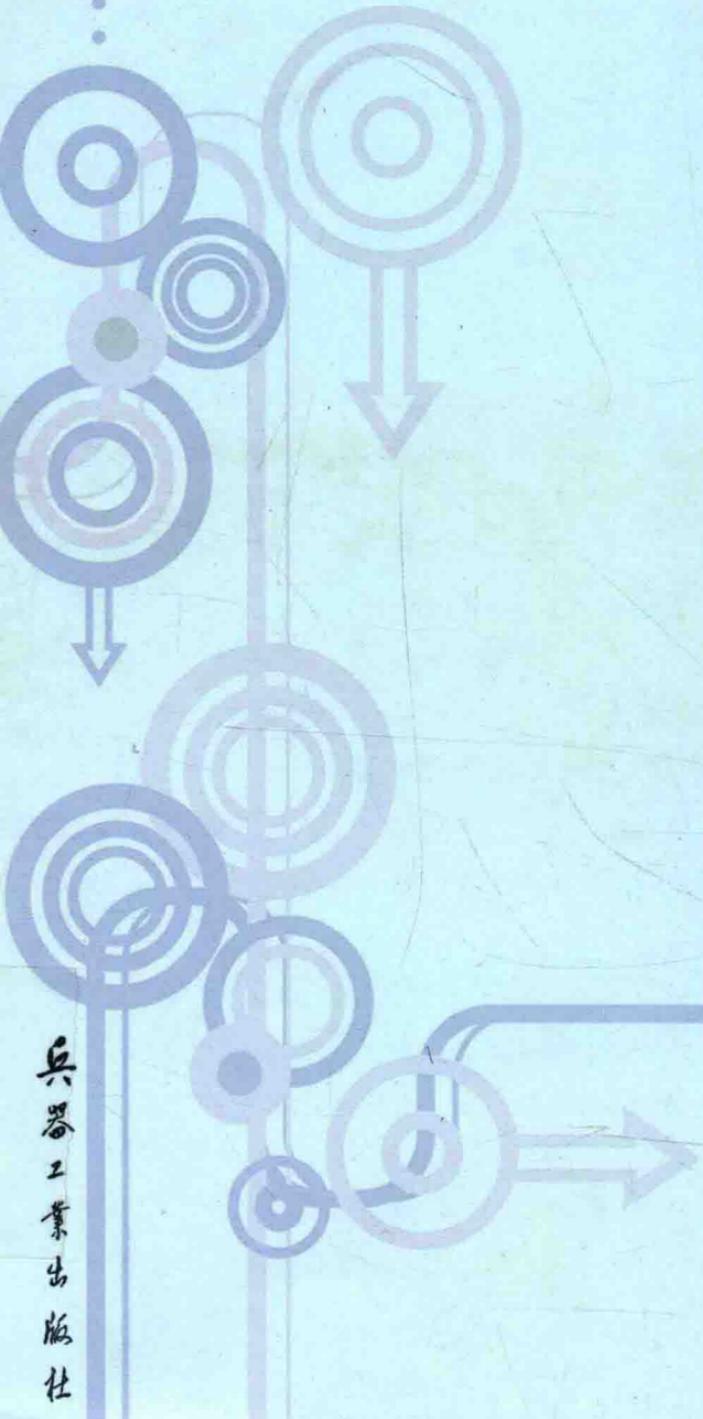


关联规则挖掘研究

关联规则挖掘是数据挖掘的一个重要组成部分，频集挖掘是关联规则挖掘的关键步骤，它在很大程度上决定了关联规则挖掘的效率。对于强规则、相关分析、时间序列、频集挖掘也有着十分重要的意义。

史月美 宗春梅 ◎著



当今社会，企业要想在激烈的市场竞争中立于不败之地，就必须学会运用各种先进的管理方法。其中，关联规则挖掘技术就是一种非常有效的方法。它可以帮助企业发现隐藏在大量数据中的有用信息，从而为企业决策提供有力的支持。

关联规则挖掘研究

史月美 宗春梅 著

兵器工业出版社

内容简介

本书从数据挖掘领域引入,介绍了关联规则的经典算法,重点从多种角度阐述作者以及该领域研究者们在关联规则挖掘方面的最新研究成果,包括最新研究思想、最新算法以及最新应用等。本书在各部分都提供了算法所涉及的必要的背景知识,以便读者无需查阅其他资料就可更好地理解。

图书在版编目(CIP)数据

关联规则挖掘研究/史月美,宗春梅著. - 北京:
兵器工业出版社, 2016重印

ISBN 978 - 7 - 80248 - 502 - 0

I. ①关… II. ①史…②宗… III. ①数据采集—研
究 IV. ①TP274

中国版本图书馆 CIP 数据核字(2010)第 056854 号

出版发行:兵器工业出版社

责任编辑:周宜今

发行电话:010 - 68962596, 68962591

封面设计:揽胜视觉

邮 编:100089

责任校对:郭 芳

社 址:北京市海淀区车道沟 10 号

责任印制:赵春云

经 销:各地新华书店

开 本:787 × 1092 1/16

印 刷:北京毅峰迅捷印刷有限公司

印 张:30.25

版 次:2016年 5月第1版第4次印刷

字 数:355 千字

定 价:65.00 元

(版权所有 翻印必究 印装有误 负责调换)

前　　言

随着产生和收集数据的能力不断提高,数据集的规模越来越大。存储数据的几何级数的爆炸性增长激起业界人士对信息新技术和智能化工具的需求,欲将数据转换为知识,以帮助决策者英明决策。

数据挖掘是一个多学科、知识交叉的领域。而关联规则挖掘是数据挖掘的知识模式中的一种主要形式,因其形式与大多数人对知识模式的理解最为相似,使其成为近几年研究的热点。

本书从数据挖掘领域引入,介绍了关联规则的经典算法,重点从多种角度阐述作者以及该领域研究者们在关联规则挖掘方面的最新研究成果,包括最新研究思想、最新算法以及最新应用等。本书在各部分都提供了算法所涉及的必要的背景知识,以便读者无需查阅其他资料就可更好地理解。

全书共 8 章。第 1 章主要介绍数据挖掘的基础知识,使读者有一个研究关联规则挖掘的背景知识;第 2 章介绍关联规则的基本知识,重点介绍关联规则挖掘的经典算法 Apriori 算法及 FP 算法,以及它们的改进算法;第 3 章介绍粗糙集的基本理论及将粗糙集理论应用于关联规则挖掘的研究成果;第 4 章简要介绍遗传算法的基本理论及将遗传算法应用于关联规则挖掘的相关算法;第 5 章介绍蚁群算法的基本原理及将蚁群算法与关联规则挖掘相结合的研究成果;第 6 章阐述分布式环境下关联规则挖掘的研究进展及其应用;第 7 章介绍基于时序和极大团

理论的关联规则挖掘相关算法;第8章说明关联规则挖掘在相关领域的实际应用。

本书第1、3、6、8章由史月美编写,第2、4、5、7章由宗春梅编写,由赵青杉副教授认真仔细审查了全书,并提出了许多宝贵意见,在此一并致谢。

本书可作为高等院校计算机科学与技术专业研究人员、技术人员的参考书,也可作为从事信息系统建设人员的参考资料。

由于作者水平有限,本书难免有不妥之处,恳请读者及专家批评指正,或提出建设性意见。敬请通过 sym_2291@sina.com 与作者联系,不胜感激。

作 者

2009年11月

目 录

第1章 数据挖掘概述	1
1.1 数据挖掘的发展	1
1.2 数据挖掘的过程	3
1.3 数据预处理	4
1.3.1 数据清洗	4
1.3.2 数据集成	9
1.3.3 数据变换	12
1.3.4 数据归约	15
1.3.5 离散化和概念分层	27
1.4 数据挖掘的分类	37
1.5 数据挖掘的方法	39
1.6 数据挖掘的应用	40
第2章 关联规则挖掘基础理论	43
2.1 关联规则的概念	43
2.1.1 关联规则的形式化描述	44
2.1.2 关联规则的分类	45
2.2 关联规则挖掘的 Apriori 算法	47
2.2.1 关联规则的挖掘过程	47

2.2.2 Apriori 算法	48
2.2.3 频集算法的改进策略	51
2.2.4 基于布尔矩阵的 Apriori 改进算法	54
2.2.5 基于十字链表的 Apriori 改进算法	59
2.3 FP 算法	63
2.3.1 FP - 增长算法	63
2.3.2 基于约束的 FP - 增长改进算法	66
2.3.3 一种无须子集检查的闭合频繁集挖掘算法	71
2.3.4 基于模式矩阵的 FP - 增长改进算法	87
2.4 多层关联规则挖掘	93
2.4.1 多层关联规则概念	93
2.4.2 多层关联规则的挖掘方法	94
2.4.3 基于 FP_tree 的多层关联规则挖掘算法	98
2.5 多维关联规则挖掘	102
2.5.1 多维关联规则的概念	102
2.5.2 多维关联规则的挖掘方法	103
2.5.3 基于频繁模式图的多维关联规则挖掘算法	104
2.5.4 基于 LR - RCEP 的多维关联规则挖掘算法	110
2.6 负关联规则挖掘	116
2.6.1 负关联规则相关概念	117
2.6.2 支持度 - 有效度框架	118
2.6.3 改进的挖掘算法	120
2.6.4 负关联规则挖掘中的频繁项集爆炸问题研究	122

第3章 基于粗糙集理论的关联规则挖掘研究	129
3.1 粗糙集理论	129
3.1.1 粗糙集理论的概念	129
3.1.2 粗糙集理论中的知识表示	134
3.1.3 属性约简与核	142
3.1.4 粗糙集对属性约简的一般方法	143
3.2 基于模板的关联规则挖掘算法	146
3.2.1 信息系统的 α -约简	146
3.2.2 信息系统中的模板	146
3.2.3 关联规则挖掘启发式算法	148
3.3 基于粗糙集的分类关联规则挖掘	154
3.3.1 属性的重要度	155
3.3.2 启发式属性约简算法	156
3.3.3 分类关联规则挖掘算法	157
3.4 基于 Rough Set 带结论域的关联规则挖掘	159
3.4.1 修改后的关联规则评价指标	159
3.4.2 带结论域的关联规则挖掘算法	160
3.5 基于粗糙集的多维关联规则挖掘	164
3.5.1 算法原理	165
3.5.2 表达用户个性化需要的挖掘语言	166
3.5.3 基于粗糙集的多维关联规则挖掘算法	167
3.6 基于 Rough Set 的 Web 日志挖掘研究	170
3.6.1 Web 挖掘研究现状	171
3.6.2 Web 日志的基本概念	172
3.6.3 基于 Rough set 的 Web 日志挖掘	172

第4章 基于遗传算法的关联规则挖掘研究	175
4.1 遗传算法基本理论	175
4.1.1 遗传算法的特点	175
4.1.2 遗传算法基本术语	177
4.1.3 基本遗传算法	179
4.1.4 遗传算法中的常用技术	183
4.1.5 遗传算法的关键参数确定	193
4.2 基于小生境遗传算法的连续属性关联规则挖掘	
.....	193
4.2.1 小生境遗传算法原理	194
4.2.2 算法描述	198
4.2.3 实验及分析	198
4.3 优化相关关联规则的发现	200
4.3.1 分类规则与关联规则的区别	200
4.3.2 优化关联规则发现算法基础	202
4.3.3 基于遗传算法的优化关联规则挖掘	208
4.3.4 实验及分析	213
4.4 基于遗传算法的频繁项挖掘算法	217
4.4.1 基于优化模型的频繁项挖掘问题表述	217
4.4.2 频繁项挖掘的非线性优化模型	218
4.4.3 基于遗传算法的频集挖掘方法	220
4.4.4 实验及分析	224
4.5 遗传优化模糊约束的频繁项集挖掘	226
4.5.1 带约束的频繁项集挖掘和模糊集	226
4.5.2 模糊约束	227
4.5.3 遗传寻优模糊集	230

第5章 基于蚁群算法的关联规则挖掘研究	233
5.1 蚁群算法基本理论	234
5.1.1 蚁群算法的概念	234
5.1.2 蚁群算法的原理分析	241
5.1.3 蚁群算法的算法描述	244
5.1.4 蚁群算法的特征	245
5.1.5 蚁群算法与其他仿生算法的比较	250
5.1.6 蚁群算法的发展与研究现状	254
5.2 基于蚁群算法的关联规则挖掘	257
5.2.1 研究背景	257
5.2.2 基于蚁群优化的关联规则挖掘算法	258
5.2.3 算法实现	261
5.2.4 实验结果与分析	263
5.3 基于多态蚁群算法的关联规则挖掘	264
5.3.1 自适应调整挥发系数的逆向蚁群算法	265
5.3.2 基于模拟退火算法的多道逆向蚁群算法	268
5.3.3 基于信息素扩散的多态蚁群算法	270
5.3.4 改进的多态蚁群算法在迷宫最短路径问题中的应用	274
5.4 基于混合蚁群算法的关联规则挖掘	279
5.4.1 研究背景	279
5.4.2 基于混合蚁群算法的关联规则算法 (Gaaa-miner)	281
5.4.3 基于混合蚁群算法的关联规则算法 (Antga-miner)	295
5.5 基于时间模型的蚁群算法的关联规则挖掘	301
5.5.1 利用基于时间模型的蚁群算法挖掘	302

5.5.2 天才时间蚁群算法	304
5.5.3 其他时间蚁群算法	307
5.6 基于遗传——蚂蚁的多维关联规则挖掘	308
5.6.1 算法设计	308
5.6.2 仿真实验	313
第6章 基于分布式数据库的关联规则挖掘研究	315
6.1 基于分布式数据库关联规则挖掘的经典算法	315
6.1.1 分布式关联规则挖掘的基本原理和方法	316
6.1.2 分布式关联规则挖掘的 FDM 算法	320
6.2 星形结构下的分布式关联规则挖掘方法 CDMA	338
6.2.1 CDMA 的基本思想	338
6.2.2 CDMA 关联规则挖掘算法	341
6.2.3 CDMA 算法分析	345
6.2.4 树形结构的关联规则分布式挖掘算法	347
6.3 隐私保护的分布式关联规则挖掘算法研究	349
6.3.1 隐私保护的概念	350
6.3.2 算法的改进策略	352
6.3.3 密码学及加密算法概述	354
6.3.4 P_ODMA 算法的设计及分析	361
6.3.5 P_ODMA 算法描述	375
6.4 分布式数据库关联规则更新算法	388
6.4.1 相关概念	389
6.4.2 IUAAR 算法基础	391
6.4.3 全局频繁项目集的维护算法	394

6.4.4 IUAAR 算法步骤	395
6.5 分布式数据库约束性关联规则的快速挖掘	397
6.5.1 相关概念	399
6.5.2 候选项集的生成函数	400
6.5.3 约束性频繁项集分布式挖掘算法 DCAR	403
6.5.4 算法实验与性能比较	406
6.6 基于分布式数据库采样的关联规则挖掘算法	408
6.6.1 相关知识	408
6.6.2 基于采样的关联规则挖掘算法 SMA	410
6.6.3 算法实验分析与比较	413
 第 7 章 基于时序和极大团的关联规则挖掘研究	416
7.1 时序逻辑及其模式	418
7.1.1 数据库的预处理	419
7.1.2 基本概念和描述	419
7.2 极大团及其算法研究	426
7.2.1 基本概念	426
7.2.2 F2setT 算法	428
7.2.3 MaxCliqueT 算法	430
7.3 基于时序逻辑的概率理论研究	432
7.3.1 基本概念	432
7.3.2 等概率的数学模型	433
7.3.3 不等概率的数学模型	434
7.4 基于极大频繁项目集的关联规则的生成算法	436

第8章 关联规则挖掘的应用	439
8.1 关联规则挖掘在灾害天气预测中的应用	439
8.2 关联规则挖掘在CRM中的应用	441
8.3 关联规则挖掘在概念检索中的应用	444
8.4 关联规则挖掘在网络入侵检测中的应用	446
8.5 关联规则挖掘在煤矿安全预警系统中的应用 ..	448
8.6 关联规则挖掘在股票分析预测中的应用	450
8.7 关联规则挖掘在基因表达数据中的应用	457
8.8 关联规则挖掘在交通事故数据分析中的应用 ..	459
参考文献	467

第1章 数据挖掘概述

数据采集和存储技术的进步导致庞大的数据库日益增多，这已经发生在人类耕耘的几乎所有领域，从普通的超市业务数据、信用卡使用记录、电话呼叫清单以及政府的统计数据等领域，到不太普通的天体图像、分子数据库和医疗记录等领域。能否从这些数据中提取出对数据库拥有者有价值的信息呢？人们对这个问题的兴趣在不断增长，而且已经形成了致力于这个任务的一门学科——数据挖掘。

目前的数据库系统可以高效地实现数据的录入、查询、统计等功能，但无法发现数据中存在的关系和规则，无法根据现有的数据预测未来的发展趋势。缺乏挖掘数据背后隐藏知识的手段，导致了“数据爆炸但信息、知识贫乏”的现象。而我们要想在这些数据和信息的基础上进行迅速准确的决策就必须借助一些新的技术和自动工具，以便将海量的数据转化成有价值的信息和知识。这些都促使数据挖掘的出现并推动其发展。

1.1 数据挖掘的发展

数据挖掘是一个逐渐演变的过程。在电子化数据处理的初期，人们试图通过某些方式来实现自动决策支持，当时的机器学习成为人们关心的焦点，机器学习的过程就是将一些已知的并已被成功解决的问题作为范例输入计算机，机器通过学习这些

范例总结并生成相应的规则,这些规则具有通用性,使用它们可以解决某一类的问题。随后,随着神经网络技术的形成和发展,使人们的注意力转向知识工程。知识工程不同于机器学习那样给计算机输入范例,让其生成出规则,而是直接为计算机输入已被代码化的规则,计算机通过使用这些规则来解决某些问题。专家系统就是使用这种方法所得到的成果,但它有投资大、效果不甚理想等不足。20世纪80年代人们又在新的神经网络理论的指导下,重新回到机器学习的方法上,并将其成果应用于处理大型商业数据库,从而产生了一个新的术语——数据库中的知识发现(Knowledge Discovery in Database, KDD)。它首次出现在1989年8月在底特律举行的第十一届国际联合人工智能学术会议上。为了统一认识,在1996年出版的总结该领域进展的权威论文集《知识发现与数据进展》中,Fayyd, Piatetsky-Shapiro and Smyth给出了KDD和数据挖掘的最新定义,将二者加以区分。

KDD的定义为:KDD是从数据中辨别有效的、新颖的、潜在有用的、最终可理解的模式的过程。数据挖掘的定义为:数据挖掘是KDD中通过特定的算法在可接受的计算效率限制内生成特定模式的一个步骤。

由此可见,整个KDD过程是一个以知识使用者为中心、人机交互的探索过程。数据挖掘是数据库中知识发现的一个步骤,但又是最重要的一步。因此,往往可以不加区别地使用KDD和数据挖掘。一般在研究领域被称作数据库中的知识发现,在工程领域则称为数据挖掘。1989年举行了第一届专题讨论会后,1991年、1993年、1994年又连续举行了KDD专题讨论会。1995年讨论会开始发展为一年一次的国际学术大会。会议较全面地讨论了数据挖掘与知识发现(Data Mining and Knowledge Discovery, DMKD)的基础理论、新的发现算法、数据

挖掘与数据仓库及 OLAP 的结合、可视化技术、知识表示方法、Web 中的数据挖掘等。迄今为止,对关系数据库和事务数据库进行数据挖掘和知识发现的研究已经取得了一定的进展,最有影响的发现算法有:加拿大 Simon Fraser 大学 J. Han 教授的概念树提升算法^[2]、IBM 的 R. Agrawal 的关联算法^[3]、澳大利亚的 J. R. Quinlan 教授的分类算法^[4]、密西根州立大学 Erick Goodman 的遗传算法等。IBM、GTE、SAS、Microsoft、Silicon Graphics、Integral Solutions、Thinking Machines 等公司,相继开发出一些实用的 KDD 商业系统和原型系统,如市场分析用的 BehaviorScan、Explorer、MDT (Management Discovery Tool),金融投资领域的 Stock Selector、AI (Automated Investor),欺诈预警用的 Falcon、FAIS、Clone detector 等。与国外相比,国内对 DMKD 的研究稍晚,没有形成整体力量。1993 年国家自然科学基金首次支持该领域的研究项目。目前,国内的许多科研单位和高等院校开展了知识发现的基础理论及其应用研究,如清华大学、中科院计算技术研究所等。

1.2 数据挖掘的过程

数据挖掘(Data Mining)就是对观测到的数据集(经常是很庞大的)进行分析,目的是发现未知的关系和以数据拥有者可以理解并对其有价值的新颖方式来总结数据。这些数据可以存放在数据库、数据仓库或其他信息存储中。这是一个年轻的跨学科领域,源于诸如数据库系统、数据仓库、统计学、机器学习、数据可视化、信息检索和高性能计算。其他有贡献的领域包括神经网络、模式识别、空间数据分析、图像数据库、信号处理和许多应用领域,包括商务、经济学和生物信息学等。

数据挖掘也称为术语数据库中的知识发现。数据挖掘过程

一般分为三个阶段:数据准备、数据挖掘、结果评估,如图 1-1 所示。

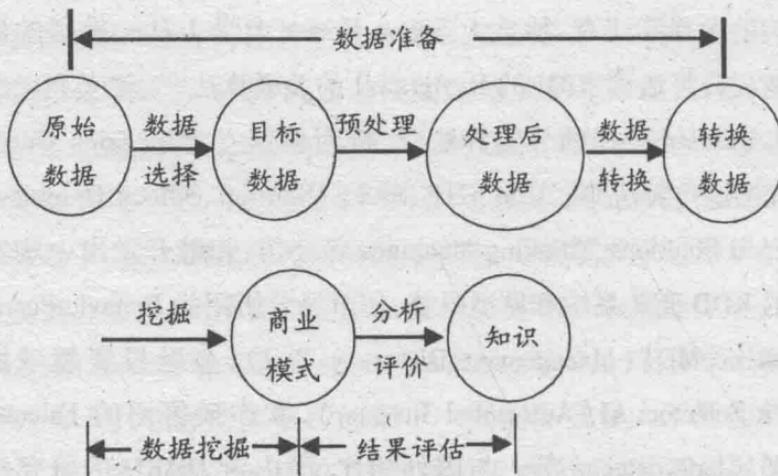


图 1-1 数据挖掘过程

1.3 数据预处理

现实世界中的原始数据经常存在不一致、重复、不完整、含噪声、维度高等问题,怎样预处理原始数据才能使得数据挖掘过程更加简便有效?数据预处理的目的就是为数据挖掘过程提供干净、准确、简洁的数据,提高数据挖掘效率和准确性,是数据挖掘中非常重要的环节。

常用的数据预处理的方法有数据清洗、数据集成、数据变换和数据规约。

1.3.1 数据清洗

数据清洗(Data Cleaning)就是清除数据噪声和与挖掘主题明显无关的数据。通常包括:添补遗漏数据、平滑噪声数据、识