



普通高等教育“十一五”国家级规划教材

21世纪高等学校计算机专业核心课程规划教材

# 数据挖掘原理与算法 (第3版)

毛国君 段立娟 编著



清华大学出版社



普通高等教育“十一五”国家级规划教材

21世纪高等学校计算机专业核心课程规划教材

# 数据挖掘原理与算法

(第3版)

毛国君 段立娟 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,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们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。

全书共分 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,而且也起到对本章内容相关文献的注释性索引功能。

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性教学。为了更好地让教师进行选择性教学,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 13701121933

### 图书在版编目(CIP)数据

数据挖掘原理与算法/毛国君,段立娟编著.--3 版.--北京: 清华大学出版社,2016

21 世纪高等学校计算机专业核心课程规划教材

ISBN 978-7-302-41581-7

I. ①数… II. ①毛… ②段… 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(2015)第 220315 号

责任编辑: 刘 星

封面设计: 杨 兮

责任校对: 焦丽丽

责任印制: 刘海龙

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](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈: 010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

课 件 下 载: <http://www.tup.com.cn>, 010-62795954

印 刷 者: 北京富博印刷有限公司

装 订 者: 北京市密云县京文制本装订厂

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 21.75 字 数: 504 千字

版 次: 2005 年 6 月第 1 版 2016 年 1 月第 3 版 印 次: 2016 年 1 月第 1 次印刷

印 数: 17801~19800

定 价: 39.50 元

## FOREWORD

# 前言

《数据挖掘原理与算法》经过第1版和第2版,历经十几年的历程,得到了研究者、教师、学生及计算机从业者的肯定和鼓励,在此表示衷心的感谢。据不完全统计,前两版已经被国内二十多所高校作为研究生或者本科生教材使用。在使用过程中,许多人也对第2版中的文字错误、内容编排等提出一些很好的建议。加之数据挖掘技术本身的发展对再次改版提出了强烈需求。第3版除了对必要的文字等错误进行修正外,重点增加了大数据挖掘等新的数据挖掘的需求和技术分析,对Web挖掘的内容进行了重新编排,并增加了必要的新方法。这样,第3版的内容及其编排更趋合理,近年来出现的公认的典型算法和技术也得到加强,使之很好地适应读者在教学或者学习中的新需求。

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。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、人工智能等在内的相关学科的发展,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。可以预见,随着大数据概念的提出和应用,数据挖掘也必将是支撑大数据分析的最重要和最核心的技术之一。

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;另一方面,作为比较新的交叉研究领域,不同背景的研究人员(数据库、人工智能、数学等)可能提供不同的视点,而且本身仍在发展中。本书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。由于没有合适的教材可用,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,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参与本书的编写。本书的许多内容是作者们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、高年级本科生的选修教材,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性学习或查阅。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,这些算法完全可以在理解的基础上进行改进或改造应用到实际工作中。

全书共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解。

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,作为作者的导师,他们在作者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和中央财经大学一些研究生的心血,他们在本书算法实例整理和验证等方面做了很多工作,在此就不一一列举了。此外,也感谢使用第2版图书的教师和学生,他们的使用给予我们进一步编好该书的动力,同时提出的许多意见也提升了第3版内容编排的质量。

作 者

2015年9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
1.1.3 大数据时代的数据挖掘技术需求分析 .....	5
1.2 数据挖掘研究的发展趋势 .....	7
1.3 数据挖掘概念 .....	10
1.3.1 从商业角度看数据挖掘技术 .....	10
1.3.2 数据挖掘的技术含义 .....	10
1.3.3 数据挖掘研究的理论基础 .....	12
1.4 数据挖掘技术的分类问题 .....	13
1.5 数据挖掘常用的知识表示模式与方法 .....	15
1.5.1 广义知识挖掘 .....	15
1.5.2 关联知识挖掘 .....	17
1.5.3 类知识挖掘 .....	17
1.5.4 预测型知识挖掘 .....	22
1.5.5 特异型知识挖掘 .....	23
1.6 不同数据存储形式下的数据挖掘问题 .....	24
1.6.1 事务数据库中的数据挖掘 .....	24
1.6.2 关系型数据库中的数据挖掘 .....	25
1.6.3 数据仓库中的数据挖掘 .....	26
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 .....	27
1.6.5 面向应用的新型数据源中的数据挖掘 .....	27
1.6.6 Web 数据源中的数据挖掘 .....	27
1.7 粗糙集方法及其在数据挖掘中的应用 .....	29

1.7.1 粗糙集的一些重要概念 .....	29
1.7.2 粗糙集应用举例 .....	31
1.7.3 粗糙集方法在 KDD 中的应用范围 .....	32
1.8 数据挖掘的应用分析 .....	33
1.8.1 数据挖掘与 CRM .....	33
1.8.2 数据挖掘与社会网络 .....	34
1.8.3 数据挖掘应用的成功案例分析 .....	35
1.9 本章小结和文献注释 .....	37
习题 1 .....	42
<b>第 2 章 知识发现过程与应用结构 .....</b>	<b>44</b>
2.1 知识发现的基本过程 .....	44
2.1.1 数据抽取与集成技术要点 .....	46
2.1.2 数据清洗与预处理技术要点 .....	46
2.1.3 数据的选择与整理技术要点 .....	47
2.1.4 数据挖掘技术要点 .....	47
2.1.5 模式评估技术要点 .....	47
2.2 数据库中的知识发现处理过程模型 .....	48
2.2.1 阶梯处理过程模型 .....	48
2.2.2 螺旋处理过程模型 .....	49
2.2.3 以用户为中心的处理模型 .....	50
2.2.4 联机 KDD 模型 .....	52
2.2.5 支持多数据源多知识模式的 KDD 处理模型 .....	54
2.3 知识发现软件或工具的发展 .....	57
2.3.1 独立的知识发现软件 .....	57
2.3.2 横向的知识发现工具集 .....	57
2.3.3 纵向的知识发现解决方案 .....	58
2.3.4 KDD 系统介绍 .....	58
2.4 知识发现项目的过程化管理 .....	60
2.5 数据挖掘语言介绍 .....	62
2.5.1 数据挖掘语言的分类 .....	62
2.5.2 数据挖掘查询语言 .....	63
2.5.3 数据挖掘建模语言 .....	64
2.5.4 通用数据挖掘语言 .....	65
2.5.5 DMQL 挖掘查询语言介绍 .....	66
2.6 本章小结和文献注释 .....	69
习题 2 .....	71

第3章 关联规则挖掘理论和算法 .....	72
3.1 基本概念与解决方法 .....	72
3.2 经典的频繁项目集生成算法分析 .....	73
3.2.1 项目集空间理论 .....	73
3.2.2 经典的发现频繁项目集算法 .....	74
3.2.3 关联规则生成算法 .....	76
3.3 Apriori 算法的性能瓶颈问题 .....	78
3.4 Apriori 的改进算法 .....	79
3.4.1 基于数据分割的方法 .....	79
3.4.2 基于散列的方法 .....	80
3.4.3 基于采样的方法 .....	81
3.5 项目集空间理论的发展 .....	82
3.5.1 Close 算法 .....	83
3.5.2 FP-tree 算法 .....	87
3.6 项目集格空间和它的操作 .....	90
3.7 基于项目集操作的关联规则挖掘算法 .....	92
3.7.1 关联规则挖掘空间 .....	92
3.7.2 三个实用算子 .....	92
3.7.3 最大频繁项目集格的生成算法 .....	94
3.7.4 ISS-DM 算法执行示例 .....	94
3.8 改善关联规则挖掘质量问题 .....	95
3.8.1 用户主观层面 .....	95
3.8.2 系统客观层面 .....	96
3.9 约束数据挖掘问题 .....	96
3.9.1 约束在数据挖掘中的作用 .....	96
3.9.2 约束的类型 .....	97
3.10 时态约束关联规则挖掘 .....	100
3.11 关联规则挖掘中的一些更深入的问题 .....	103
3.11.1 多层次关联规则挖掘 .....	103
3.11.2 多维关联规则挖掘 .....	104
3.11.3 数量关联规则挖掘 .....	105
3.12 数量关联规则挖掘方法 .....	106
3.12.1 数量关联规则挖掘问题 .....	106
3.12.2 数量关联规则的分类 .....	107
3.12.3 数量关联规则挖掘的一般步骤 .....	108
3.12.4 数值属性离散化问题及算法 .....	111
3.13 本章小结和文献注释 .....	114

习题 3 .....	116
<b>第 4 章 分类方法.....</b>	<b>119</b>
4.1 分类的基本概念与步骤 .....	120
4.2 基于距离的分类算法 .....	122
4.3 决策树分类方法 .....	125
4.3.1 决策树基本算法概述.....	126
4.3.2 ID3 算法 .....	128
4.3.3 C4.5 算法 .....	133
4.4 贝叶斯分类 .....	138
4.4.1 贝叶斯定理.....	138
4.4.2 朴素贝叶斯分类.....	138
4.4.3 EM 算法 .....	141
4.5 规则归纳 .....	145
4.5.1 AQ 算法 .....	145
4.5.2 CN2 算法 .....	149
4.5.3 FOIL 算法 .....	156
4.6 与分类有关的其他问题 .....	160
4.6.1 分类数据预处理.....	160
4.6.2 分类器性能的表示与评估.....	161
4.7 本章小结和文献注释 .....	163
习题 4 .....	165
<b>第 5 章 聚类方法.....</b>	<b>169</b>
5.1 概述 .....	169
5.1.1 聚类分析在数据挖掘中的应用.....	171
5.1.2 聚类分析算法的概念与基本分类.....	171
5.1.3 距离与相似性的度量.....	174
5.2 划分聚类方法 .....	177
5.2.1 $k$ -平均算法 .....	177
5.2.2 PAM .....	180
5.2.3 其他方法.....	184
5.3 层次聚类方法 .....	184
5.3.1 AGNES 算法 .....	185
5.3.2 DIANA 算法 .....	186
5.3.3 其他聚类方法.....	188
5.4 密度聚类方法 .....	189
5.5 其他聚类方法 .....	193

5.5.1 STING 算法 .....	193
5.5.2 SOM 算法 .....	194
5.5.3 COBWEB 算法 .....	194
5.5.4 模糊聚类算法 FCM .....	195
5.6 本章小结和文献注释 .....	195
习题 5 .....	197
<b>第 6 章 时间序列和序列模式挖掘 .....</b>	<b>199</b>
6.1 时间序列及其应用 .....	199
6.2 时间序列预测的常用方法 .....	200
6.2.1 确定性时间序列预测方法 .....	200
6.2.2 随机时间序列预测方法 .....	201
6.2.3 其他方法 .....	201
6.3 基于 ARMA 模型的序列匹配方法 .....	201
6.3.1 基本概念 .....	201
6.3.2 利用基本概念建立模型 .....	202
6.3.3 构造判别函数 .....	203
6.4 基于离散傅里叶变换的时间序列相似性查找 .....	204
6.4.1 完全匹配 .....	205
6.4.2 子序列匹配 .....	206
6.5 基于规范变换的查找方法 .....	208
6.5.1 基本概念 .....	209
6.5.2 查找方法 .....	209
6.6 序列挖掘 .....	211
6.6.1 基本概念 .....	212
6.6.2 数据源的形式 .....	212
6.6.3 序列模式挖掘的一般步骤 .....	214
6.7 AprioriAll 算法 .....	215
6.8 AprioriSome 算法 .....	218
6.9 GSP 算法 .....	222
6.10 本章小结和文献注释 .....	224
习题 6 .....	227
<b>第 7 章 Web 挖掘技术 .....</b>	<b>229</b>
7.1 Web 挖掘的意义 .....	229
7.2 Web 挖掘的分类 .....	230
7.3 Web 挖掘的含义 .....	232
7.3.1 Web 挖掘与信息检索 .....	232

7.3.2 Web 挖掘与信息抽取 .....	232
7.4 Web 挖掘的数据来源 .....	233
7.4.1 服务器日志数据 .....	233
7.4.2 在线市场数据 .....	234
7.4.3 Web 页面 .....	234
7.4.4 Web 页面超链接关系 .....	235
7.4.5 其他信息 .....	235
7.5 Web 内容挖掘方法 .....	235
7.5.1 爬虫与 Web 内容挖掘 .....	236
7.5.2 虚拟的 Web 视图 .....	236
7.5.3 个性化与 Web 内容挖掘 .....	237
7.5.4 对 Web 页面内文本信息的挖掘 .....	237
7.5.5 对 Web 页面内多媒体信息挖掘 .....	238
7.5.6 Web 页面内容的预处理 .....	238
7.6 Web 访问信息挖掘方法 .....	239
7.6.1 Web 访问信息挖掘的特点 .....	239
7.6.2 Web 访问信息挖掘的意义 .....	241
7.6.3 Web 访问信息挖掘的数据源 .....	242
7.6.4 Web 访问信息挖掘的一般过程 .....	245
7.6.5 Web 访问信息挖掘的数据清理 .....	246
7.6.6 用户识别方法 .....	247
7.6.7 会话识别方法 .....	249
7.6.8 其他预处理技术 .....	252
7.6.9 Web 访问挖掘的应用方法 .....	252
7.6.10 Web 访问信息挖掘的要素构成 .....	254
7.6.11 Web 访问信息挖掘应用 .....	255
7.7 Web 结构挖掘方法 .....	264
7.7.1 页面等级(分级)的评价方法 .....	264
7.7.2 PageRank 算法 .....	265
7.7.3 权威页面和中心页面 .....	268
7.7.4 Web 站点结构的预处理 .....	269
7.8 本章小结和文献注释 .....	271
习题 7 .....	275
<b>第 8 章 空间挖掘 .....</b>	<b>277</b>
8.1 引言 .....	277
8.2 空间数据概要 .....	278
8.2.1 空间数据的复杂性特征 .....	278

8.2.2 空间查询问题.....	279
8.2.3 空间数据结构.....	280
8.2.4 专题地图.....	284
8.3 空间数据挖掘基础 .....	284
8.4 空间统计学 .....	286
8.5 泛化与特化 .....	287
8.5.1 逐步求精.....	287
8.5.2 泛化.....	287
8.5.3 最临近方法.....	289
8.5.4 统计信息网格方法.....	289
8.6 空间规则 .....	291
8.7 空间分类算法 .....	293
8.7.1 ID3 扩展 .....	293
8.7.2 空间决策树.....	293
8.8 空间聚类算法 .....	294
8.8.1 基于随机搜索的聚类方法 CLARANS 扩展 .....	295
8.8.2 大型空间数据库基于距离分布的聚类算法 DBCLASD .....	296
8.8.3 BANG .....	297
8.8.4 小波聚类.....	297
8.8.5 近似值.....	297
8.9 空间挖掘的其他问题 .....	299
8.10 空间数据挖掘原型系统介绍.....	302
8.11 空间数据挖掘的研究现状.....	304
8.12 空间数据挖掘的研究与发展方向.....	305
8.13 空间数据挖掘与相关学科的关系.....	307
8.13.1 空间数据挖掘与空间数据库.....	307
8.13.2 空间数据挖掘与空间数据仓库.....	308
8.13.3 空间数据挖掘与空间联机分析处理.....	308
8.13.4 空间数据挖掘与地理信息系统.....	309
8.14 数字地球.....	310
8.15 本章小结和文献注释.....	310
习题 8 .....	313
参考文献.....	314

# 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入的方向发展。

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。同时在目前的商业活动中,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

发展趋势或做出正确决策的。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,为企业提供能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,是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,而且现在看来仍有日益加速增长的趋势。随之而来的是硬盘、CPU等关键部件的价格大幅度下降,使得人们收集、存储和处理数据的能力和欲望不断提高。经过几十年的发展,计算机的体系结构,特别是并行处理技术已经逐渐成熟并获得普遍应用,而且成为支持大型数据处理应用的基础。计算机性能的提高和先进的体系结构的发展使数据挖掘技术的研究和应用成为可能。

历经了十几年的发展,包括基于统计学、人工智能等在内的理论与技术成果已经被成功地应用到商业处理和分析中。这些应用从某种程度上为数据挖掘技术的提出和发展起到了极大的推动作用。数据挖掘系统的核心模块技术和算法都离不开这些理论和技术的支持。从某种意义上讲,这些理论本身的发展和应用为数据挖掘提供了有价值的理论和应用积累。数理统计是一个有几百年发展历史的应用数学学科,至今仍然是应用数学中最重要、最活跃的学科之一。如今相当强大有效的数理统计方法和工具,已成为信息咨询业的基础。然而它和数据库技术的结合性研究应该说是近十几年才被重视。以前的基于数理统计方法的应用大多都是通过专用程序来实现的。我们知道,大多数的统计分析技术是基于严格的数学理论和高超的应用技巧的,这使得一般的用户很难从容地驾驭它。一旦人们有了从数据查询到知识发现、从数据演绎到数据归纳的要求,概率论和数理统计就获得了新的生命力。从这个意义上说,数据挖掘技术是数理统计分析应用的延伸和发展。假如人们利用数据库的方式从被动地查询变成了主动地发现知识的话,那么概率论和数理统计这一古老的学科可以为我们从数据归纳到知识发现提供理论基础。

人工智能是计算机科学研究中争议最多而又始终保持强大生命力的研究领域。专家系统曾经是人工智能研究工作者的骄傲。专家系统实质上是一个问题求解系统。领域专家长期以来面向一个特定领域的经验世界,通过人脑的思维活动积累了大量有用信息。在研制一个专家系统时,首先,知识工程师要从领域专家那里获取知识,这一过程是非常复杂的个人到个人之间的交流过程,有很强的个性和随机性。因此,知识获取成为专家系统研究中公认的瓶颈问题。其次,知识工程师在整理表达从领域专家那里获得的知识时,一般用if-then等规则表达,这种表达局限性太大,勉强抽象出来的规则有很强的工艺色彩,知识表示又成为一大难题。此外,即使某个领域的知识通过一定手段获取并表达了,但这样做成的专家系统对常识和百科知识出奇地贫乏,而人类专家的知识是以大量常识知识为基础的。人工智能学家Feigenbaum估计,一般人拥有的常识存入计算机大约有100万条事实和抽象经验法则,离开常识的专家系统有时会比傻子还傻。另外,由于专家系统是主观整理知识,因此这种机制不可避免地带有偏见和错误。以上诸多难题大大限制了专家系统的应用。数据挖掘继承了专家系统的高度实用性的特点,并且以数据为基本出发点,客观地挖掘知识。机器学习应该说是得到了充分的研究和发展,从事机器学习的科学家们,不再满足自己构造的小样本学习模式的象牙塔,开始正视现实生活中大量的、不完全的、有噪声的、模糊的、随机的大数据样本,进而也走上了数据挖掘的道路。因

此,可以说,数据挖掘研究在继承已有的人工智能相关领域的研究成果的基础上,摆脱了以前象牙塔式的研究模式,真正客观地开始从数据集中发现蕴藏的知识。

### 1.1.3 大数据时代的数据挖掘技术需求分析

大数据(Big Data)概念虽然最早是在 20 世纪 80 年代提出的,来自于 *Nature* 2008 年推出的 Big Data 专刊,但是真正受到广泛研究与应用探索应该算是在 2011 年,其中重要的标志是麦肯锡咨询公司发布的《大数据:下一个创新、竞争和生产率的前沿》报告。然而,大数据如此迅猛发展,形成的是目前“边研究边应用”的局面,当然也带来新的问题。

毋庸置疑,数据挖掘技术将是大数据分析的核心和骨干技术之一。当然,大数据时代也对数据挖掘技术的发展提出新的挑战性的问题。我们可以从大数据的发展历史、对应的概念演变入手,分析大数据时代对数据挖掘的技术需求。

大数据研究的发展可以粗略划分成三个阶段。

① 2000 年及以前,称为“大数据概念萌芽阶段”。就科技文献而言,从主流的学术引擎上(Google 学术),以 Big Data 为关键词,检索出的 2000 年及以前的学术论文不超过 50 篇。这足以说明该时期虽然有了大数据这个名词,但是并没有受到学术界和商界的广泛重视。当然该时期有了“大数据的萌芽”,其中一个重要标志是:开始将互联网作为大数据的一个重要来源。由于互联网数据的特别关注,使得计算机界早已提出的海量数据处理问题得到扩展。众所周知,海量数据处理起源于大容量科学数据计算需要,所以其研究和应用主要还是面向单一的数据结构(如数据库表)。

② 2001—2010 年,大数据概念得到广泛讨论,应用价值获得共识,我们把它称为“大数据概念探索阶段”。例如,这一阶段的年均科技文献已经超过 100 篇。除美国以外,包括中国在内的其他国家的论文数量也显著增长。说明大数据概念已经得到普遍认可。当然,该阶段大数据的概念还是在探索中,和我们今天认识的大数据还是有差距的。特别地,由于当时数据处理中的“数据挖掘”研究已经进入高峰时期,因此,提出和讨论的大数据概念总是有数据挖掘的影子。然而,从科学发展的历史的角度来说,任何概念总要经过初始提出和不断探讨才能越来越清楚,而这中间数据挖掘概念及其丰富的研究成果为大数据概念的逐步演化起到了关键的作用:提供了理论和方法上的储备。

③ 2011 年及以后,大数据的概念进一步深化,已经成为学术研究的焦点,成为许多应用的支撑概念。特别地,检索 2011—2013 年的论文,其年均论文数量已经接近 1000 篇。可谓百花齐放、齐头并进,可以看出大数据已经受到学界和商界的高度重视。

下面摘录了一些引用比较多的关于大数据概念的解释。

① 2011 年,麦肯锡公司发布的《大数据:下一个创新、竞争和生产率的前沿》报告认为:“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”。

② 2012 年,国际数据公司(IDC)则认为大数据有四个特征,即 4V 属性:数据规模大(Volume)、数据高速聚集(Velocity)、数据类型多样(Variety)、数据价值巨大(Value)。关于第四个属性,IBM 的研究报告也有真实性(Veracity)之说,所以很多学者认为,前三个属性是最关键的。

③ 2012 年,顾能(Gartner)公司技术报告则认为:“大数据是需要新处理模式才能具