

<<数据挖掘>>

图书基本信息

书名：<<数据挖掘>>

13位ISBN编号：9787111391401

10位ISBN编号：7111391403

出版时间：2012-8

出版时间：机械工业出版社

作者：（美） Jiawei Han,（加） Micheline Kamber,（加） Jian Pei

页数：468

译者：范明,孟小峰

版权说明：本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问：<http://www.tushu007.com>

<<数据挖掘>>

前言

前言：社会的计算机化显著地增强了我们产生和收集数据的能力。

大量数据从我们生活的每个角落涌出。

存储的或瞬态的数据的爆炸性增长已激起对新技术和自动工具的需求，以帮助我们智能地将海量数据转换成有用的信息和知识。

这导致称做数据挖掘的一个计算机科学前沿学科的产生，这是一个充满希望和欣欣向荣并具有广泛应用的学科。

数据挖掘通常又称为数据中的知识发现（KDD），是自动地或方便地提取代表知识的模式；这些模式隐藏在大型数据库、数据仓库、Web、其他大量信息库或数据流中。

本书考察知识发现和数据挖掘的基本概念和技术。

作为一个多学科领域，数据挖掘从多个学科汲取营养。

这些学科包括统计学、机器学习、模式识别、数据库技术、信息检索、网络科学、知识库系统、人工智能、高性能计算和数据可视化。

我们提供发现隐藏在大型数据集中的模式的技术，关注可行性、有用性、有效性和可伸缩性问题。

因此，本书不打算作为数据库系统、机器学习、统计学或其他某领域的导论，尽管我们确实提供了这些领域的必要背景材料，以便读者理解它们各自在数据挖掘中的作用。

本书是对数据挖掘的全面介绍。

对于计算科学的学生、应用开发人员、行业专业人员以及涉及以上列举的学科的研究人员，本书应当是有用的。

数据挖掘出现于20世纪80年代后期，20世纪90年代有了突飞猛进的发展，并可望在新千年继续繁荣。

本书全面展示该领域，介绍有趣的数据挖掘技术和系统，并讨论数据挖掘的应用和研究方向。

写本书的重要动机是需要建立一个学习数据挖掘的有组织的框架——由于这个快速发展领域的多学科特点，这是一项具有挑战性的任务。

我们希望本书有助于具有不同背景和经验丰富的人交换关于数据挖掘的见解，为进一步促进这个令人激动的、不断发展的领域的成长做出贡献。

本书的组织 自本书第1版、第2版出版以来，数据挖掘领域已经取得了重大进展，开发出了许多新的数据挖掘方法、系统和应用，特别是对于处理包括信息网络、图、复杂结构和数据流，以及文本、Web、多媒体、时间序列、时间空间数据在内的新的数据类型。

这种快速发展、新技术不断涌现使得在一本书中涵盖整个领域的广泛内容非常困难。

因此，我们决定与其继续扩大本书的涵盖面，还不如让本书以足够的广度和深度涵盖该领域的核心内容，而把复杂数据类型的处理留给另一本即将面世的书。

第3版对本书的前两版做了全面修订，加强和重新组织了全书的技术内容，显著地扩充和加强处理一般数据类型挖掘的核心技术。

第2版中讨论特定主题的章节（例如，数据预处理、频繁模式挖掘、分类和聚类）在这一版都被扩充，每章都分成两章。

对于这些主题，一章囊括基本概念和技术，而另一章提供高级概念和方法。

第2版关于复杂数据类型的章节（例如，流数据、序列数据、图结构数据、社会网络数据和多重关系数据，以及文本、Web、多媒体和时间空间数据）现在保留给专门介绍数据挖掘的高级课题的新书。为了支持读者学习这些高级课题，我们把第2版的相关章节的电子版放在本书的网站上，作为第3版的配套材料。

第3版各章的简要内容如下（重点介绍新的内容）：第1章提供关于数据挖掘的多学科领域的导论。

该章讨论导致需要数据挖掘的数据库技术的发展历程和数据挖掘应用的重要性。

该章考察挖掘的数据类型，包括关系的、事务的和数据仓库数据，以及复杂的数据类型，如时间序列、序列、数据流、时间空间数据、多媒体数据、文本数据、图、社会网络和Web数据。

<<数据挖掘>>

该章根据所挖掘的知识类型、所使用的技术以及目标应用的类型，对数据挖掘任务进行了一般分类。最后讨论该领域的主要挑战。

第2章介绍一般数据特征。

该章首先讨论数据对象和属性类型，然后介绍基本统计数据描述的典型度量。

该章概述各种类型数据的数据可视化技术。

除了数值数据的可视化方法外，还介绍文本、标签、图和多维数据的可视化方法。

第2章还介绍度量各种类型数据的相似性和相异性的方法。

第3章介绍数据预处理技术。

该章首先介绍数据质量的概念，然后讨论数据清理、数据集成、数据归约、数据变换和数据离散化的方法。

第4章和第5章是数据仓库、OLAP（联机分析处理）和数据立方体技术的引论。

第4章介绍数据仓库和OLAP的基本概念、建模、结构、一般实现，以及数据仓库和其他数据泛化的关系。

第5章更深入地考察数据立方体技术，详细地研究数据立方体的计算方法，包括Star-Cubing和高维OLAP方法。

该章还讨论数据立方体和OLAP技术的进一步研究，如抽样立方体、排序立方体、预测立方体、用于复杂数据挖掘查询的多特征立方体和发现驱动的数据立方体的探查。

第6章和第7章介绍挖掘大型数据集中的频繁模式、关联和相关性的方法。

第6章介绍基本概念，如购物篮分析，还有条理地提供了许多频繁项集挖掘技术。

这些涵盖从基本Apriori算法和它的变形，到改进性能的更高级的方法，包括频繁模式增长方法，使用数据的垂直形式的频繁模式挖掘，挖掘闭频繁项集和极大频繁项集。

该章还讨论模式评估方法并介绍挖掘相关模式的度量。

第7章介绍高级模式挖掘方法。

该章讨论多层和多维空间中的模式挖掘，挖掘稀有和负模式，挖掘巨型模式和高维空间数据，基于约束的模式挖掘和挖掘压缩或近似模式。

该章还介绍模式探查和应用的方法，包括频繁模式的语义注解。

第8章和第9章介绍数据分类方法。

由于分类方法的重要性和多样性，内容被划分成两章。

第8章介绍分类的基本概念和方法，包括决策树归纳、贝叶斯分类和基于规则的分类。

该章还讨论模型评估和选择方法，以及提高分类准确率的方法，包括组合方法和处理不平衡数据。

第9章讨论分类的高级方法，包括贝叶斯信念网络、后向传播的神经网络技术、支持向量机、使用频繁模式的分类、k-最邻近分类、基于案例的推理、遗传算法、粗糙集理论和模糊集方法。

附加的主题包括多类分类、半监督分类、主动学习和迁移学习。

聚类分析是第10章和第11章的主题。

第10章介绍数据聚类的基本概念和方法，包括基本聚类分析方法的概述、划分方法、层次方法、基于密度的方法和基于网格的方法。

该章还介绍聚类评估方法。

第11章讨论聚类的高级方法，包括基于概率模型的聚类、聚类高维数据、聚类图和网络数据，以及基于约束的聚类。

第12章专门讨论离群点检测。

本章介绍离群点的基本概念和离群点分析，并从各种监督力度（监督的、半监督的和无监督的）以及方法角度（统计学方法、基于邻近性的方法、基于聚类的方法和基于分类的方法）讨论离群点检测方法。

该章还讨论挖掘情境离群点和集体离群点，以及高维数据中的离群点检测。

最后，在第13章我们讨论数据挖掘的趋势、应用和研究前沿。

我们简略地介绍挖掘复杂数据类型，包括挖掘序列数据（例如，时间序列、符号序列和生物学序列），挖掘图和网络，以及挖掘空间、多媒体、文本和Web数据。

<<数据挖掘>>

这些数据挖掘方法的深入讨论留给正在撰写的数据挖掘高级课题一书。

然后, 该章转向讨论其他数据挖掘方法学, 包括统计学数据挖掘、数据挖掘基础、可视和听觉数据挖掘, 以及数据挖掘的应用。

讨论数据挖掘在金融数据分析、零售和电信产业、科学与工程, 以及入侵检测和预防方面的应用。

该章还讨论数据挖掘与推荐系统的联系。

由于数据挖掘出现在我们日常生活的方方面面, 所以我们讨论数据挖掘与社会, 包括无处不在和无形的数据挖掘, 以及隐私、安全和数据挖掘对社会的影响。

我们用考察数据挖掘的发展趋势结束本书。

书中楷体字用于强调定义的术语, 而黑体字用于突出主要思想。

本书与其他数据挖掘教材相比具有一些显著特点: 它广泛、深入地讨论了数据挖掘原理。

各章尽可能是自包含的, 使得读者可以按自己感兴趣的次序阅读。

高级章节提供了更大的视野, 感兴趣的读者可以选读。

本书提供了数据挖掘的所有主要方法, 还提供了关于多维OLAP分析等数据挖掘的重要主题, 这些主题在其他书中常常被忽略或很少提及。

本书还维护了一个网站, 其中包含大量在线资源, 为教师、学生和该领域的专业人员提供支持。

这些将在下面介绍。

致教师 本书旨在提供数据挖掘领域的一个广泛而深入的概览, 可以作为高年级本科生或一年级研究生的数据挖掘导论。

除了讲稿、教师指南和阅读材料列表等教学资源之外, 本书网站 (www.cs.uiuc.edu/~hanj/bk3 或 www.booksite.mkp.com/datamining3e) 还提供了一个样本课程安排。

根据授课学时、学生的背景和你的兴趣, 你可以选取章节的子集, 以不同的顺序进行讲授。

例如, 如果你只打算给学生讲授数据挖掘入门导论, 可以按照图P.1的建议。

注意, 根据需要, 必要时可以省略其中某些节或某些小节。

图P.1 入门导论课程的建议章节序列 根据学时和讲授范围, 你可以有选择地把更多的章节增加到这个基本序列中。

例如, 对高级分类方法更感兴趣的教师可以首先增加“第9章 分类: 高级方法”; 对模式挖掘更感兴趣的教师可以选择包括“第7章 高级模式挖掘”; 而对OLAP和数据立方体技术感兴趣的教师可以增加“第4章 数据仓库与联机分析处理”和“第5章 数据立方体技术”。

或者, 你可以选择在两个学期的系列课程中讲授整本书, 包括本书的所有章节, 时间允许的话, 加上图和网络挖掘这样的高级课题。

这些高级课题可以从本书网站提供的配套材料选择, 辅以挑选的研究论文。

本书的每一章都可以用做自学材料, 或者用做数据库系统、机器学习、模式识别和数据智能分析等相关课程的专题。

每章后面都有一些习题, 适合作为家庭作业。

这些习题或者是用于测验对内容的掌握情况的小问题, 或者是需要分析思考的大问题, 或者是实现设计。

有些习题也可以用做研究讨论课题。

每章后面的文献注释可以用来查找包含正文中提供的概念和方法的来源、相关课题的深入讨论和可能的扩展的研究文献。

致学生 我们希望本书将激发你对年青, 但正在快速发展的数据挖掘领域的兴趣。

我们试图以清晰的方式提供材料, 仔细地解释所涵盖的主题。

每一章后面都附有一个小结, 总结要点。

全书包含了许多图和解释, 以便使本书更加有趣和便于阅读。

尽管本书是作为教材编写的, 但是我们也试图把它组织成一本有用的参考书或手册, 以有助于你今后在数据挖掘方面进行深入研究和求职。

为阅读本书, 你需要知道什么?

你应当具有关于统计学、数据库系统和机器学习的概念和术语方面的知识。

<<数据挖掘>>

然而，我们尽力提供这些基础知识的足够背景，以便在读者对这些领域不太熟悉或者记忆有些淡忘时，也能够理解本书的讨论。

你应当具有一些程序设计经验。

特别是你应当能够阅读伪代码，能够理解像多维数组这样的简单数据结构。

致专业人员 本书旨在涵盖数据挖掘领域的广泛主题。

因此，本书是关于该主题的一本优秀手册。

由于每一章的编写都尽可能独立，所以读者可以关注自己最感兴趣的课题。

希望学习数据挖掘关键思想的应用程序和信息服务管理人员可以使用本书。

对于有兴趣使用数据挖掘技术解决其业务问题的银行、保险、医药和零售业的数据分析人员，本书也是有用的。

此外，本书也可以作为数据挖掘领域的全面综述，有助于研究人员提升数据挖掘技巧，扩展数据挖掘的应用范围。

本书所提供的技术和算法是实用的，介绍的算法适合于发现隐藏在大型、现实数据集中的模式和知识，而不是挑选在小型“玩具”数据库上运行良好的算法。

本书提供的每个算法都用伪代码解释。

伪代码类似于程序设计语言C，但也精心加以策划，使得不熟悉C或C++的程序员易于理解。

如果你想实现算法，你会发现将我们的伪代码转换成选定的程序设计语言程序是一项非常简单的任务。

本书资源网站 本书网站的地址是www.cs.uiuc.edu/~hanj/bk3，另一个是Morgan Kaufmann出版社的网站www.booksite.mkp.com/datamining3e。

这些网站为本书的读者和对数据挖掘感兴趣的人提供了一些附加材料，资源包括：每章的幻灯片。提供了用微软的PowerPoint制作的每章教案。

高级数据挖掘的配套章节。

本书第2版的第8~10章涵盖了挖掘复杂的数据类型，这超出了本书的主题，对这些高级主题感兴趣的读者可从网站上获取。

教师手册。

本书习题的完整答案通过出版社的网站只向教师提供。

课程提纲和教学计划。

使用本书和幻灯片用于数据挖掘导论课程和高级教程的本科生和研究生，可以获取这些资源。

带超链接的辅助阅读文献列表。

补充读物的原创性文章按章组织。

到数据挖掘数据集和软件的链接。

我们将提供到数据挖掘数据集和某些包含有趣的数据挖掘软件包的站点的链接，如到伊利诺伊大学厄巴纳-尚佩恩分校IlliMine的链接(<http://illimine.cs.uiuc.edu>)。

作业、考试和课程设计样本。

一组作业、考试和课程设计样本将在出版社的网站上向教师提供。

本书的插图。

这可能有助于你制作自己的课堂教学幻灯片。

本书目录。

PDF格式。

本书不同印次的勘误表。

欢迎读者指出本书中的错误。

一旦错误被证实，我们将更新勘误表，并对你的贡献致谢。

评论或建议请发往hanj@cs.uiuc.edu。

我们很高兴听到你的建议。

<<数据挖掘>>

内容概要

本书完整全面地讲述数据挖掘的概念、方法、技术和最新研究进展。本书对前两版做了全面修订，加强和重新组织了全书的技术内容，重点论述了数据预处理、频繁模式挖掘、分类和聚类等内容，还全面讲述了OLAP和离群点检测，并研讨了挖掘网络、复杂数据类型以及重要应用领域。

本书是数据挖掘和知识发现领域内的所有教师、研究人员、开发人员和用户都必读的参考书，是一本适用于数据分析、数据挖掘和知识发现课程的优秀教材，可以用做高年级本科生或者一年级研究生的数据挖掘导论教材。

<<数据挖掘>>

作者简介

Jiawei

Han (韩家炜), 是伊利诺伊大学厄巴纳-尚佩恩分校计算机科学系的Bliss教授。他因知识发现和数据挖掘研究方面的贡献而获得许多奖励, 包括ACM SIGKDD创新奖(2004)、IEEE计算机学会技术成就奖(2005)和IEEE W.Wallace McDowell奖(2009)。

他是ACM和IEEE会士。

他还担任《ACM Transactions on Knowledge

Discovery from Data》的执行主编(2006—2011)和许多杂志的编委, 包括《IEEE Transactions on Knowledge and Data Engineering》和《Data Mining Knowledge Discovery》。

Micheline

Kamber, 由加拿大魁北克蒙特利尔Concordia大学获计算机科学(人工智能专业)硕士学位。

她曾是NSERC学者, 作为研究者在McGill大学、西蒙-弗雷泽大学和瑞士工作。

她的数据挖掘背景和以易于理解的形式写作的热情使得本书更受专业人员、教师和学生的欢迎。

Jian Pei (裴健), 现在是西蒙-弗雷泽大学计算机科学学院教授。

他在Jiawei

Han的指导下, 于2002年获西蒙-弗雷泽大学计算科学博士学位。

他在数据挖掘、数据库、Web搜索和信息检索的主要学术论坛发表了大量文章, 并积极服务于学术团体。

他的文章被引用数千次, 并获多次荣誉奖。

他是多种数据挖掘和数据分析杂志的助理编辑。

<<数据挖掘>>

书籍目录

出版者的话

中文版序

译者序

译者简介

第3版序

第2版序

前言

致谢

作者简介

第1章 引论

1.1 为什么进行数据挖掘

1.1.1 迈向信息时代

1.1.2 数据挖掘是信息技术的进化

1.2 什么是数据挖掘

1.3 可以挖掘什么类型的数据

1.3.1 数据库数据

1.3.2 数据仓库

1.3.3 事务数据

1.3.4 其他类型的数据

1.4 可以挖掘什么类型的模式

1.4.1 类/概念描述：特征化与区分

1.4.2 挖掘频繁模式、关联和相关性

1.4.3 用于预测分析的分类与回归

1.4.4 聚类分析

1.4.5 离群点分析

1.4.6 所有模式都是有趣的吗

1.5 使用什么技术

1.5.1 统计学

1.5.2 机器学习

1.5.3 数据库系统与数据仓库

1.5.4 信息检索

1.6 面向什么类型的应用

1.6.1 商务智能

1.6.2 Web搜索引擎

1.7 数据挖掘的主要问题

1.7.1 挖掘方法

1.7.2 用户界面

1.7.3 有效性和可伸缩性

1.7.4 数据库类型的多样性

1.7.5 数据挖掘与社会

1.8 小结

1.9 习题

1.10 文献注释

第2章 认识数据

2.1 数据对象与属性类型

<<数据挖掘>>

- 2.1.1 什么是属性
 - 2.1.2 标称属性
 - 2.1.3 二元属性
 - 2.1.4 序数属性
 - 2.1.5 数值属性
 - 2.1.6 离散属性与连续属性
 - 2.2 数据的基本统计描述
 - 2.2.1 中心趋势度量：均值、中位数和众数
 - 2.2.2 度量数据散布：极差、四分位数、方差、标准差和四分位数极差
 - 2.2.3 数据的基本统计描述的图形显示
 - 2.3 数据可视化
 - 2.3.1 基于像素的可视化技术
 - 2.3.2 几何投影可视化技术
 - 2.3.3 基于图符的可视化技术
 - 2.3.4 层次可视化技术
 - 2.3.5 可视化复杂对象和关系
 - 2.4 度量数据的相似性和相异性
 - 2.4.1 数据矩阵与相异性矩阵
 - 2.4.2 标称属性的邻近性度量
 - 2.4.3 二元属性的邻近性度量
 - 2.4.4 数值属性的相异性：闵可夫斯基距离
 - 2.4.5 序数属性的邻近性度量
 - 2.4.6 混合类型属性的相异性
 - 2.4.7 余弦相似性
 - 2.5 小结
 - 2.6 习题
 - 2.7 文献注释
- ### 第3章 数据预处理
- 3.1 数据预处理：概述
 - 3.1.1 数据质量：为什么要对数据预处理
 - 3.1.2 数据预处理的主要任务
 - 3.2 数据清理
 - 3.2.1 缺失值
 - 3.2.2 噪声数据
 - 3.2.3 数据清理作为一个过程
 - 3.3 数据集成
 - 3.3.1 实体识别问题
 - 3.3.2 冗余和相关分析
 - 3.3.3 元组重复
 - 3.3.4 数据值冲突的检测与处理
 - 3.4 数据归约
 - 3.4.1 数据归约策略概述
 - 3.4.2 小波变换
 - 3.4.3 主成分分析
 - 3.4.4 属性子集选择
 - 3.4.5 回归和对数线性模型：参数化数据归约
 - 3.4.6 直方图

<<数据挖掘>>

- 3.4.7 聚类
 - 3.4.8 抽样
 - 3.4.9 数据立方体聚集
 - 3.5 数据变换与数据离散化
 - 3.5.1 数据变换策略概述
 - 3.5.2 通过规范化变换数据
 - 3.5.3 通过分箱离散化
 - 3.5.4 通过直方图分析离散化
 - 3.5.5 通过聚类、决策树和相关分析离散化
 - 3.5.6 标称数据的概念分层产生
 - 3.6 小结
 - 3.7 习题
 - 3.8 文献注释
- 第4章 数据仓库与联机分析处理
- 4.1 数据仓库：基本概念
 - 4.1.1 什么是数据仓库
 - 4.1.2 操作数据库系统与数据仓库的区别
 - 4.1.3 为什么需要分离的数据仓库
 - 4.1.4 数据仓库：一种多层体系结构
 - 4.1.5 数据仓库模型：企业仓库、数据集市和虚拟仓库
 - 4.1.6 数据提取、变换和装入
 - 4.1.7 元数据库
 - 4.2 数据仓库建模：数据立方体与OLAP
 - 4.2.1 数据立方体：一种多维数据模型
 - 4.2.2 星形、雪花形和事实星座：多维数据模型的模式
 - 4.2.3 维：概念分层的作用
 - 4.2.4 度量的分类和计算
 - 4.2.5 典型的OLAP操作
 - 4.2.6 查询多维数据库的星网查询模型
 - 4.3 数据仓库的设计与使用
 - 4.3.1 数据仓库的设计的商务分析框架
 - 4.3.2 数据仓库的设计过程
 - 4.3.3 数据仓库用于信息处理
 - 4.3.4 从联机分析处理到多维数据挖掘
 - 4.4 数据仓库的实现
 - 4.4.1 数据立方体的有效计算：概述
 - 4.4.2 索引OLAP数据：位图索引和连接索引
 - 4.4.3 OLAP查询的有效处理
 - 4.4.4 OLAP服务器结构：ROLAP、MOLAP、HOLAP的比较
 - 4.5 数据泛化：面向属性的归纳
 - 4.5.1 数据特征的面向属性的归纳
 - 4.5.2 面向属性归纳的有效实现
 - 4.5.3 类比较的面向属性归纳
 - 4.6 小结
 - 4.7 习题
 - 4.8 文献注释
- 第5章 数据立方体技术

<<数据挖掘>>

5.1 数据立方体计算：基本概念

5.1.1 立方体物化：完全立方体、冰山立方体、闭立方体和立方体外壳

5.1.2 数据立方体计算的一般策略

5.2 数据立方体计算方法

5.2.1 完全立方体计算的多路数组聚集

5.2.2 BUC：从顶点方体向下计算冰山立方体

5.2.3 Star-Cubing：使用动态星树结构计算冰山立方体

5.2.4 为快速高维OLAP预计算壳片段

5.3 使用探索立方体技术处理高级查询

5.3.1 抽样立方体：样本数据上基于OLAP的挖掘

5.3.2 排序立方体：top-k查询的有效计算

5.4 数据立方体空间的多维数据分析

5.4.1 预测立方体：立方体空间的预测挖掘

5.4.2 多特征立方体：多粒度上的复杂聚集

5.4.3 基于异常的、发现驱动的立方体空间探查

5.5 小结

5.6 习题

5.7 文献注释

第6章 挖掘频繁模式、关联和相关性：基本概念和方法

6.1 基本概念

6.1.1 购物篮分析：一个诱发例子

6.1.2 频繁项集、闭项集和关联规则

6.2 频繁项集挖掘方法

6.2.1 Apriori算法：通过限制候选产生发现频繁项集

6.2.2 由频繁项集产生关联规则

6.2.3 提高Apriori算法的效率

6.2.4 挖掘频繁项集的模式增长方法

6.2.5 使用垂直数据格式挖掘频繁项集

6.2.6 挖掘闭模式和极大模式

6.3 哪些模式是有趣的：模式评估方法

6.3.1 强规则不一定是有趣的

6.3.2 从关联分析到相关分析

6.3.3 模式评估度量比较

6.4 小结

6.5 习题

6.6 文献注释

第7章 高级模式挖掘

7.1 模式挖掘：一个路线图

7.2 多层、多维空间中的模式挖掘

7.2.1 挖掘多层关联规则

7.2.2 挖掘多维关联规则

7.2.3 挖掘量化关联规则

7.2.4 挖掘稀有模式和负模式

7.3 基于约束的频繁模式挖掘

7.3.1 关联规则的元规则制导挖掘

7.3.2 基于约束的模式产生：模式空间剪枝和数据空间剪枝

7.4 挖掘高维数据和巨型模式

<<数据挖掘>>

- 7.5 挖掘压缩或近似模式
 - 7.5.1 通过模式聚类挖掘压缩模式
 - 7.5.2 提取感知冗余的top-k模式
- 7.6 模式探索与应用
 - 7.6.1 频繁模式的语义注解
 - 7.6.2 模式挖掘的应用
- 7.7 小结
- 7.8 习题
- 7.9 文献注释
- 第8章 分类：基本概念
 - 8.1 基本概念
 - 8.1.1 什么是分类
 - 8.1.2 分类的一般方法
 - 8.2 决策树归纳
 - 8.2.1 决策树归纳
 - 8.2.2 属性选择度量
 - 8.2.3 树剪枝
 - 8.2.4 可伸缩性与决策树归纳
 - 8.2.5 决策树归纳的可视化挖掘
 - 8.3 贝叶斯分类方法
 - 8.3.1 贝叶斯定理
 - 8.3.2 朴素贝叶斯分类
 - 8.4 基于规则的分类
 - 8.4.1 使用IF-THEN规则分类
 - 8.4.2 由决策树提取规则
 - 8.4.3 使用顺序覆盖算法的规则归纳
 - 8.5 模型评估与选择
 - 8.5.1 评估分类器性能的度量
 - 8.5.2 保持方法和随机二次抽样
 - 8.5.3 交叉验证
 - 8.5.4 自助法
 - 8.5.5 使用统计显著性检验选择模型
 - 8.5.6 基于成本效益和ROC曲线比较分类器
 - 8.6 提高分类准确率的技术
 - 8.6.1 组合分类方法简介
 - 8.6.2 装袋
 - 8.6.3 提升和AdaBoost
 - 8.6.4 随机森林
 - 8.6.5 提高类不平衡数据的分类准确率
 - 8.7 小结
 - 8.8 习题
 - 8.9 文献注释
- 第9章 分类：高级方法
 - 9.1 贝叶斯信念网络
 - 9.1.1 概念和机制
 - 9.1.2 训练贝叶斯信念网络
 - 9.2 用后向传播分类

<<数据挖掘>>

- 9.2.1 多层前馈神经网络
- 9.2.2 定义网络拓扑
- 9.2.3 后向传播
- 9.2.4 黑盒内部：后向传播和可解释性
- 9.3 支持向量机
- 9.3.1 数据线性可分的情况
- 9.3.2 数据非线性可分的情况
- 9.4 使用频繁模式分类
- 9.4.1 关联分类
- 9.4.2 基于有区别力的频繁模式分类
- 9.5 惰性学习法(或从近邻学习)
- 9.5.1 k-最近邻分类
- 9.5.2 基于案例的推理
- 9.6 其他分类方法
- 9.6.1 遗传算法
- 9.6.2 粗糙集方法
- 9.6.3 模糊集方法
- 9.7 关于分类的其他问题
- 9.7.1 多类分类
- 9.7.2 半监督分类
- 9.7.3 主动学习
- 9.7.4 迁移学习
- 9.8 小结
- 9.9 习题
- 9.10 文献注释
- 第10章 聚类分析：基本概念和方法
- 10.1 聚类分析
- 10.1.1 什么是聚类分析
- 10.1.2 对聚类分析的要求
- 10.1.3 基本聚类方法概述
- 10.2 划分方法
- 10.2.1 k-均值：一种基于形心的技术
- 10.2.2 k-中心点：一种基于代表对象的技术
- 10.3 层次方法
- 10.3.1 凝聚的与分裂的层次聚类
- 10.3.2 算法方法的距离度量
- 10.3.3 BIRCH：使用聚类特征树的多阶段聚类
- 10.3.4 Chameleon：使用动态建模的多阶段层次聚类
- 10.3.5 概率层次聚类
- 10.4 基于密度的方法
- 10.4.1 DBSCAN：一种基于高密度连通区域的基于密度的聚类
- 10.4.2 OPTICS：通过点排序识别聚类结构
- 10.4.3 DENCLUE：基于密度分布函数的聚类
- 10.5 基于网格的方法
- 10.5.1 STING：统计信息网格
- 10.5.2 CLIQUE：一种类似于Apriori的子空间聚类方法
- 10.6 聚类评估

<<数据挖掘>>

- 10.6.1 估计聚类趋势
- 10.6.2 确定簇数
- 10.6.3 测定聚类质量
- 10.7 小结
- 10.8 习题
- 10.9 文献注释
- 第11章 高级聚类分析
 - 11.1 基于概率模型的聚类
 - 11.1.1 模糊簇
 - 11.1.2 基于概率模型的聚类
 - 11.1.3 期望最大化算法
 - 11.2 聚类高维数据
 - 11.2.1 聚类高维数据：问题、挑战和主要方法
 - 11.2.2 子空间聚类方法
 - 11.2.3 双聚类
 - 11.2.4 维归约方法和谱聚类
 - 11.3 聚类图和网络数据
 - 11.3.1 应用与挑战
 - 11.3.2 相似性度量
 - 11.3.3 图聚类方法
 - 11.4 具有约束的聚类
 - 11.4.1 约束的分类
 - 11.4.2 具有约束的聚类方法
 - 11.5 小结
 - 11.6 习题
 - 11.7 文献注释
- 第12章 离群点检测
 - 12.1 离群点和离群点分析
 - 12.1.1 什么是离群点
 - 12.1.2 离群点的类型
 - 12.1.3 离群点检测的挑战
 - 12.2 离群点检测方法
 - 12.2.1 监督、半监督和无监督方法
 - 12.2.2 统计方法、基于邻近性的方法和基于聚类的方法
 - 12.3 统计学方法
 - 12.3.1 参数方法
 - 12.3.2 非参数方法
 - 12.4 基于邻近性的方法
 - 12.4.1 基于距离的离群点检测和嵌套循环方法
 - 12.4.2 基于网格的方法
 - 12.4.3 基于密度的离群点检测
 - 12.5 基于聚类的方法
 - 12.6 基于分类的方法
 - 12.7 挖掘情境离群点和集体离群点
 - 12.7.1 把情境离群点检测转换成传统的离群点检测
 - 12.7.2 关于情境对正常行为建模
 - 12.7.3 挖掘集体离群点

<<数据挖掘>>

- 12.8 高维数据中的离群点检测
 - 12.8.1 扩充的传统离群点检测
 - 12.8.2 发现子空间中的离群点
 - 12.8.3 高维离群点建模
- 12.9 小结
- 12.10 习题
- 12.11 文献注释
- 第13章 数据挖掘的发展趋势和研究前沿
 - 13.1 挖掘复杂的数据类型
 - 13.1.1 挖掘序列数据：时间序列、符号序列和生物学序列
 - 13.1.2 挖掘图和网络
 - 13.1.3 挖掘其他类型的数据
 - 13.2 数据挖掘的其他方法
 - 13.2.1 统计学数据挖掘
 - 13.2.2 关于数据挖掘基础的观点
 - 13.2.3 可视和听觉数据挖掘
 - 13.3 数据挖掘应用
 - 13.3.1 金融数据分析的数据挖掘
 - 13.3.2 零售和电信业的数据挖掘
 - 13.3.3 科学与工程数据挖掘
 - 13.3.4 入侵检测和预防数据挖掘
 - 13.3.5 数据挖掘与推荐系统
 - 13.4 数据挖掘与社会
 - 13.4.1 普适的和无形的数据挖掘
 - 13.4.2 数据挖掘的隐私、安全和社会影响
 - 13.5 数据挖掘的发展趋势
 - 13.6 小结
 - 13.7 习题
 - 13.8 文献注释
- 参考文献
- 索引

章节摘录

版权页：插图：第二种技术称做广义关系阈值控制，为广义关系设置一个阈值。如果广义关系中不同元组的个数超过该阈值，则应当进行进一步泛化；否则，不再进一步泛化。这样的阈值也可以在数据挖掘系统中提供（通常取值范围为10~30），或者由专家或用户设置，并且允许调整。

例如，如果用户感到广义关系太小，则他可以加大该阈值；这意味着下钻。

否则，为进一步泛化关系，他可以减小该阈值；这意味着上卷。

这两种技术可以顺序使用：首先使用属性泛化阈值控制技术泛化每个属性，然后使用关系阈值控制进一步压缩广义关系。

无论使用哪种泛化控制技术，都应当允许用户调整泛化阈值，以便得到有趣的概念描述。

在许多面向数据库的归纳过程中，用户感兴趣的是在不同的抽象层得到数据的量化信息或统计信息。

因此，在归纳过程中收集计数和其他聚集值是非常重要的。

从概念上讲，这可以通过采用如下办法来实现。

聚集函数count（）与每个数据库元组相关联。

对于初始工作关系的每个元组，它的值被初始化为1。

通过删除属性和属性泛化，初始关系中的元组可能被泛化，导致相同的元组分组。

在这种情况下，形成一个组的所有相等元组应当合并成一个元组。

<<数据挖掘>>

编辑推荐

- 数据挖掘领域最具里程碑意义的经典著作
- 完整全面阐述该领域的重要知识和技术创新海报：

版权说明

本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问:<http://www.tushu007.com>