

<<数据挖掘>>

图书基本信息

书名：<<数据挖掘>>

13位ISBN编号：9787302307143

10位ISBN编号：7302307148

出版时间：2013-1

出版时间：清华大学出版社

作者：坎塔尔季奇

译者：王晓海,吴志刚

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>>

内容概要

随着数据集规模和复杂度的持续上升，分析员必须利用更高级的软件工具来执行间接的、自动的智能化数据分析。

《数据挖掘：概念、模型、方法和算法(第2版)》介绍了通过分析高维数据空间中的海量原始数据来提取用于决策的新信息的尖端技术和方法。

本书开篇阐述数据挖掘原理，此后在示例的引导下详细讲解起源于统计学、机器学习、神经网络、模糊逻辑和演化计算等学科的具有代表性的、最前沿的挖掘方法和算法。

本书还着重描述如何恰当地选择方法和数据分析软件并合理地调整参数。

每章末尾附有复习题。

本书主要用作计算机科学、计算机工程和计算机信息系统专业的研究生数据挖掘教材，高年级本科生或具备同等教育背景的读者也完全可以理解本书的所有主题。

介绍支持向量机(SVM)和Kohonen映射

讲解DBSCAN、BIRCH和分布式DBSCAN聚类算法

介绍贝叶斯网络，讨论图形中的Betweenness和Centrality参数测量算法

分析在建立决策树时使用的CART算法和基尼指数

介绍Bagging & Boosting集成学习方法，并详述AdaBoost算法

讨论Relief以及PageRank算法

讨论文本挖掘的潜在语义分析(LSA)，并分析如何测定文本文档之间的语义相似性

讲解时态、空间、Web、文本、并行和分布式数据挖掘等新主题

更详细地讲解数据挖掘技术商业、隐私、安全和法律方面的内容

<<数据挖掘>>

作者简介

作者：(美)坎塔尔季奇 (Mehmed Kantardzic) 译者：王晓海 吴志刚 王晓海，曾任总参某部应用研发中心副主任、信息服务中心主任，高级工程师，长期从事数据库应用系统的科研开发工作，负责主持多个大型数据库系统的开发和维护，荣获多项军队科技进步奖，享受军队优秀人才岗位津贴，出版多部论(译)著，在数据库挖掘、数据库应用开发、数据安全保护、数据恢复与数据去密等领域具有丰富的实践经验。

已出版的论著和译著：(《Oracle Streams 11g数据复制》，2012年，清华大学出版社)《SQL Server 2000管理、开发及应用实例详解》，2006年，人民邮电出版社、《空时编码技术》，2004年，机械工业出版社、《远程通信网络基础》，1996年，电子工业出版社。

吴志刚，工学博士，北京邮电大学副教授，长期从事网络与信息安全技术、数据库技术等领域的学术与科研工作，作为负责人主持过上述领域多项国家863计划、发改委产业化示范项目和国家级重大项目，获得技术专利2项，已在国内外学术期刊和国际会议上发表20余篇学术论文。

<<数据挖掘>>

书籍目录

第1章数据挖掘的概念 1.1概述 1.2数据挖掘的起源 1.3数据挖掘过程 1.4大型数据集 1.5数据仓库 1.6数据挖掘的商业方面：为什么数据挖掘项目会失败 1.7本书结构安排 1.8复习题 1.9参考书目 第2章数据准备 2.1原始数据的表述 2.2原始数据的特性 2.3原始数据的转换 2.3.1标准化 2.3.2数据平整 2.3.3 差值和比率 2.4丢失数据 2.5时间相关数据 2.6异常点分析 2.7复习题 2.8参考书目 第3章数据归约 3.1大型数据集的维度 3.2特征归约 3.2.1特征选择 3.2.2特征提取 3.3 Relief算法 3.4特征排列的熵度量 3.5主成分分析 3.6值归约 3.7特征离散化：ChiMerge技术 3.8案例归约 3.9复习题 3.10参考书目 第4章从数据中学习 4.1 学习机器 4.2统计学习原理 4.3学习方法的类型 4.4常见的学习任务 4.5支持向量机 4.6 KNN：最近邻分类器 4.7模型选择与泛化 4.8模型的评估 4.9 90%准确的情形 4.9.1保险欺诈检测 4.9.2改进心脏护理 4.10复习题 4.11参考书目 第5章统计方法 5.1统计推断 5.2评测数据集的差异 5.3贝叶斯定理 5.4预测回归 5.5方差分析 5.6对数回归 5.7对数—线性模型 5.8线性判别分析 5.9复习题 5.10参考书目 第6章决策树和决策规则 6.1决策树 6.2 C4.5算法：生成决策树 6.3未知属性值 6.4修剪决策树 6.5 C4.5算法：生成决策规则 6.6 CART算法和Gini指标 6.7决策树和决策规则的局限性 6.8复习题 6.9参考书目 第7章人工神经网络 7.1人工神经元的模型 7.2人工神经网络的结构 7.3 学习过程 7.4使用ANN完成的学习任务 7.4.1模式联想 7.4.2模式识别 7.5多层感知机 7.6竞争网络和竞争学习 7.7 SOM 7.8复习题 7.9参考书目 第8章集成学习 8.1集成学习方法论 8.2多学习器组合方案 8.3 bagging和boosting 8.4 AdaBoost算法 8.5复习题 8.6参考书目 第9章聚类分析 9.1聚类的概念 9.2相似度的度量 9.3凝聚层次聚类 9.4分区聚类 9.5增量聚类 9.6 DBSCAN算法 9.7 BIRCH算法 9.8聚类验证 9.9复习题 9.10参考书目 第10章关联规则 10.1购物篮分析 10.2 Apriori算法 10.3从频繁项集中得到关联规则 10.4提高Apriori算法的效率 10.5 FP增长方法 10.6关联分类方法 10.7多维关联规则挖掘 10.8复习题 10.9参考书目 第11章Web挖掘和文本挖掘 11.1 Web挖掘 11.2 Web内容、结构与使用挖掘 11.3 HITS和LOGSOM算法 11.4挖掘路径遍历模式 11.5 PageRank算法 11.6文本挖掘 11.7潜在语义分析 11.8复习题 11.9参考书目 第12章数据挖掘高级技术 12.1 图挖掘 12.2时态数据挖掘 12.2.1时态数据表示 12.2.2序列之间的相似性度量 12.2.3时态数据模型 12.2.4数据挖掘 12.3空间数据挖掘 (SDM) 12.4分布式数据挖掘 (DDM) 12.5关联并不意味着存在因果关系 12.6数据挖掘的隐私、安全及法律问题 12.7复习题 12.8参考书目 第13章遗传算法 13.1遗传算法的基本原理 13.2用遗传算法进行优化 13.2.1编码方案和初始化 13.2.2适合度估计 13.2.3选择 13.2.4交叉 13.2.5突变 13.3遗传算法的简单例证 13.3.1表述 13.3.2初始群体 13.3.3评价 13.3.4交替 13.3.5遗传算子 13.3.6评价 (第二次迭代) 13.4图式 13.5旅行推销员问题 13.6使用遗传算法的机器学习 13.6.1规则交换 13.6.2规则概化 13.6.3规则特化 13.6.4规则分割 13.7遗传算法用于聚类 13.8复习题 13.9参考书目 第14章模糊集和模糊逻辑 14.1 模糊集 14.2模糊集的运算 14.3扩展原理和模糊关系 14.4模糊逻辑和模糊推理系统 14.5多因子评价 14.6从数据中提取模糊模型 14.7数据挖掘和模糊集 14.8 复习题 14.9参考书目 第15章可视化方法 15.1感知和可视化 15.2科学可视化和信息可视化 15.3平行坐标 15.4放射性可视化 15.5使用自组织映射进行可视化 15.6数据挖掘的可视化系统 15.7复习题 15.8参考书目 附录A数据挖掘工具 附录B数据挖掘应用

<<数据挖掘>>

章节摘录

版权页：插图：12.4分布式数据挖掘（DDM）海量数据的涌现使得利用分布式系统对海量数据开展跨地理区域的分析的需求不断增长。

为海量数据驱动的知识发现，以及潜在的科学与商业理解带来了史无前例的发展机会。

在高性能分布式计算平台上（而不是集中式计算模型上）实现数据挖掘，其驱动力来自于技术和组织两个因素。

某些情况下，集中处理方式难以实现，因为需要长距离传输将大量的T级数据。

另外，集中方法违背了隐私规则，暴露了商业秘密，并带来其他一些社会问题。

这些问题的典型实例常见于医疗行业，其相关数据往往存在于多个组织商业机构中，例如制药公司、医院、政府实体（如美国食品和药物管理局）和非政府组织（如慈善和公共健康组织）。

每个组织都具有法律限制，例如隐私法规，有关专利信息的公司需求会给竞争对手带来巨大的商业利益。

因此既需要开发算法、工具、服务和基础结构用于实现分布式跨组织的数据挖掘，同时也需要考虑隐私保护问题。

这样一种朝着分布式、复杂环境发展的变化扩大了数据挖掘挑战的范围。

分布式数据所带来的新问题明显增加了数据挖掘过程的复杂性。

通过有线和无线网络，许多分布式计算环境，在计算和通信方面获得了进展。

这样的处理环境多数都涉及包含大量数据的分布式数据源、多个计算节点和分布式用户社区。

对这些分布式数据源进行监视和分析需要新的用于分布式应用的数据挖掘技术。

DDM领域处理这些问题——通过细致分析分布式源挖掘分布式数据源。

除数据分布外，网络的发展产生了大量复杂数据，包括自然语言文本、图像、时间序列、传感器数据、多关系及对象数据类型。

更复杂的是，包含分布式流数据的系统需要增量或在线挖掘工具，无论何时当底层数据发生变化时，需要完整地处理过程。

由于系统变化频繁，应用于如此复杂环境的数据挖掘技术必须适应巨大的动态变化，否则将会对系统的性能带来不良影响。

对所有这些特性提供支持的DDM系统需要有创新的解决方案。

Web架构（包含分层协议和服务）提供了合理的框架用于支持DDM。

新框架接受“融合通信和计算”的新趋势。

DDM接受数据可能自然地分布于不同的松耦合节点上的事实，这些分布的数据往往是通过网络连接起来的异构数据。

DDM提供用于通过分布式数据分析和使用最小数据通信建模发现新知识的技术。

同时，分布式系统交互需要以可靠、稳定、可扩展的方式实现。

最后，系统必须向用户隐藏技术方面的复杂性。

目前，能够通过e—services处理的商品不仅仅局限于类似电器、家具、机票等实体。

Intcmet及WWW的发展包含了软件、计算能力或有用的数据集这类资源。

这些新资源能够通过网络以服务的形式售卖或租赁给网络用户。

直观上看，数据挖掘适于作为一种e—service发布，因为该方法减少了高昂的用于支持该方法的基础架构的设置和维护开销。

<<数据挖掘>>

编辑推荐

《国外计算机科学经典教材:数据挖掘:概念、模型、方法和算法(第2版)》主要用作计算机科学、计算机工程和计算机信息系统专业的研究生数据挖掘教材,高年级本科生或具备同等教育背景的读者也完全可以理解《国外计算机科学经典教材:数据挖掘:概念、模型、方法和算法(第2版)》的所有主题。

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>