

TURING

图灵计算机科学丛书

PEARSON

Introduction to Data Mining

数据挖掘导论（完整版）

Pang-Ning Tan

[美] Michael Steinbach 著

Vipin Kumar

范明 范宏建 等译



人民邮电出版社
POSTS & TELECOM PRESS

TURING

图灵计算机科学丛书

Introduction to Data Mining

数据挖掘导论

(完整版)

人民邮电出版社

北京

图书在版编目 (CIP) 数据

数据挖掘导论：完整版 / (美) 陈封能, (美) 斯坦巴赫 (Steinbach, M.), (美) 库玛尔 (Kumar, V.) 著; 范明等译. -- 2版. -- 北京: 人民邮电出版社, 2011. 1

(图灵计算机科学丛书)
ISBN 978-7-115-24100-9

I. ①数… II. ①陈… ②斯… ③库… ④范… III. ①数据采集 IV. ①TP274

中国版本图书馆CIP数据核字(2010)第209213号

内 容 提 要

本书全面介绍了数据挖掘的理论和方法,旨在为读者提供将数据挖掘应用于实际问题所必需的知识。本书涵盖五个主题:数据、分类、关联分析、聚类和异常检测。除异常检测外,每个主题都包含两章:前面一章讲述基本概念、代表性算法和评估技术,后面一章较深入地讨论高级概念和算法。目的是使读者在透彻地理解数据挖掘基础的同时,还能了解更多重要的高级主题。此外,书中还提供了大量示例、图表和习题。

本书适合作为相关专业高年级本科生和研究生数据挖掘课程的教材,同时也可作为数据挖掘研究和应用开发人员的参考书。

图灵计算机科学丛书

数据挖掘导论 (完整版)

- ◆ 著 [美] Pang-Ning Tan Michael Steinbach Vipin Kumar
译 范明 范宏建等
责任编辑 杨海玲
执行编辑 丁晓昀
- ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街14号
邮编 100061 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京鑫正大印刷有限公司印刷
- ◆ 开本: 787×1092 1/16
印张: 30
字数: 787千字 2011年1月第2版
印数: 10 001-13 000册 2011年1月北京第1次印刷
著作权合同登记号 图字: 01-2005-5236号

ISBN 978-7-115-24100-9

定价: 69.00元

读者服务热线: (010)51095186 印装质量热线: (010)67129223

反盗版热线: (010)67171154

版 权 声 明

Authorized translation from the English language edition, entitled *Introduction to Data Mining*, 0321321367 by Pang-Ning Tan, Michael Steinbach, and Vipin Kumar, published by Pearson Education, Inc., publishing as Addison Wesley, Copyright © 2006 by Pearson Education, Inc.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from Pearson Education, Inc.

CHINESE SIMPLIFIED language edition published by PEARSON EDUCATION ASIA LTD. and POSTS & TELECOM PRESS Copyright © 2011.

本书中文简体字版由 Pearson Education Asia Ltd. 授权人民邮电出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

本书封面贴有 Pearson Education (培生教育出版集团) 激光防伪标签，无标签者不得销售。版权所有，侵权必究。

Preface of the Chinese Edition

It is with great pleasure that we welcome the Chinese translation of our book by Professors Fan, Dr. Fan, *et al.*, who have previously translated several well-known statistics and data mining texts. Data mining is an area in computer science that aims to analyze the rapidly increasing amounts of business, scientific, and engineering data for knowledge and other profitable uses. The field has seen tremendous growth and development, with the great influx of scholars and researchers, not only from the Western countries but also from the Far East. We thank Professors Fan and Dr. Fan for their effort in doing the translation, which allows the book to reach a much broader audience among those students and researchers who are well-versed in the Chinese language. We hope that the readers of our book will find it to be both useful and engaging, and wish them the greatest success.

Pang-Ning Tan, Michael Steinbach, and Vipin Kumar
Michigan State University and the University of Minnesota, December 2005

中文版序

我们非常欢迎由范明教授和范宏建博士等人将我们的书翻译成中文，他们在此之前翻译了几本关于统计学和数据挖掘方面的著名教材。数据挖掘是计算机科学的一个领域，其目的是通过分析快速增长的商业、科学和工程数据来获取知识和其他利益。我们已经目睹了这个领域的迅猛增长和发展，学者和研究人员大量涌入其中，他们不仅来自西方国家，而且来自远东地区。我们感谢范明教授和范宏建博士，他们的翻译成果使本书得以传播到更广的读者群，包括那些精通中文的学生和研究人员。我们期望读者会发现这是一部有用的和引人入胜的书籍。祝你们成功！

Pang-Ning Tan
Michael Steinbach
Vipin Kumar

2005年12月于密歇根州立大学和明尼苏达大学

完整版译者序

图灵教育已经走过了 5 年。在图灵公司成立之初，我受其之托，翻译 P. Tan、M. Steinbach 和 V. Kumar 的力作 *Introduction to Data Mining*。2006 年 5 月，该书的中文版《数据挖掘导论》正式与读者见面。

在过去的 4 年多时间里，这本书受到了许多读者的关注，众多高校和研究所把它作为研究生和高年级本科生的数据挖掘相关课程的首选教材和参考书之一。热心的读者对译文提出了许多有益的意见和建议。在此，我们表示衷心感谢！

在迎接图灵公司成立五周年之际，出版社决定对 2006 年的版本进行修订，并补充翻译原著的附录，出版该书的中文完整版。

如一些读者所愿，完整版包含了原著的 5 个附录，涉及线性代数、维度归约、概率统计、回归和优化。这些内容是数据挖掘的数学基础，许多读者都已经从相关的数学专著和教科书中获得了这些知识。原著包含它们是希望这本书是自包含的，使未系统学习过这些数学知识，或虽然学过但有点淡忘的读者在阅读本书时不必四处翻阅数学文献。用作教材时，可以根据学生的情况，选择这些附录作为预备知识提前讲述。完整版包含这些附录也有尊重原著作者的考虑。

完整版对原译文进行了勘误（包括原著作者的勘误），修订了一些翻译生硬的句子，希望能够增强译著的可读性。此外，对于个别术语的译法也做了适当的调整。例如，*maximal frequent itemset* 改译为“极大频繁项集”，*closed frequent itemset* 改译为“闭频繁项集”。

修订和附录 A~E 的翻译均由我本人完成。

Introduction to Data Mining 自出版以来受到了广泛欢迎，已经成为数据挖掘领域的经典文献，希望中文版也能受到更多读者的青睐。

范 明

2010 年 10 月于郑州大学

译者序

自从我与孟小峰等人翻译 J. Han 和 M. Kamber 的《数据挖掘：概念与技术》以来，我们高兴地看到数据挖掘的研究正在我国蓬勃开展。许多学者和研究人员都对这个新兴的学科领域表现出了极大的兴趣，他们之中不仅有来自数据库领域的专家，而且不乏统计学、人工智能和模式识别、机器学习等领域的研究者。国内的学者和研究者在数据挖掘方面的研究已经取得了一些令人鼓舞的成果，并且正在逐渐与国际学术界同步。

数据挖掘的产生和发展一直是分析和理解数据的实际需求推动的。数据挖掘研究的进展也正是在于一直重视与其他领域研究者的合作。数据挖掘从工业、农业、医疗卫生和商业的需求中获得动力，从统计学、机器学习等领域的长期研究与发展中汲取营养。我们相信，只要有理解数据的需求，就有推动数据挖掘研究与应用发展的动力；只要依靠多学科的团队，就能应对新的数据分析任务带来的挑战。

P. Tan、M. Steinbach 和 V. Kumar 编写的这本《数据挖掘导论》是继《数据挖掘：概念与技术》一书之后的另一本重要的数据挖掘著作。三位作者都从事数据挖掘研究多年，其中 Vipin Kumar 教授是数据挖掘和高性能计算领域的国际知名学者。本书原版在正式出版之前就已经被斯坦福大学、得克萨斯大学奥斯汀分校等众多名校采用。J. Han 教授也高度评价该书：“这是一本全新数据挖掘的教材，值得大力推荐。它将成为我们的主要参考书。”

本书不需要读者具备数据库背景，只需要少量统计学或数学背景知识，而且取材涉及的学科和应用领域较多，实用性强，因此适合的读者面较广。本书强调如何用数据挖掘知识解决各种实际问题，强调所挖掘的知识模式的评估。例如，就像我们能够从天空中的白云想象出各种动物和物体一样，每个聚类算法能够从几乎所有的数据集中发现聚类。如果数据集中根本不存在自然的簇，所产生的聚类很难说具有实际意义。

全书共分 10 章。范明负责第 1~8 章的翻译，范宏建负责第 9 章和第 10 章的翻译。蒋宏杰、贾玉祥、许红涛和温箫笛也参加本书的最初翻译工作。全书的译文由范明负责统一定稿。在翻译的过程中，对发现的错误进行了更正，并得到原书作者的确认。

感谢 P. Tan、M. Steinbach 和 V. Kumar 为中文版撰写序言。感谢人民邮电出版社图灵公司的编辑们，他们在第一时间内引进本书，并组织翻译，使得中文版能够如此之快地与读者见面。

译文中的错误和不当之处，敬请读者朋友指正。意见和建议请发往 mfan@zzu.edu.cn。希望读者喜欢这本译著，希望这本译著有助于推动我国的数据挖掘研究与应用的深入开展。

范明

2006年2月于郑州大学

译者简介



范明 郑州大学信息工程学院教授，中国计算机学会数据库专业委员会委员、人工智能与模式识别专业委员会委员，长期从事计算机软件与理论教学和研究。主要讲授的课程包括程序设计、计算机操作系统、数据库系统原理、知识库系统原理、数据挖掘与数据仓库等。当前感兴趣的研究方向包括数据挖掘、数据仓库和机器学习。1989~1990年曾访问加拿大 Simon Fraser 大学计算机科学系，从事演绎数据库研究。1999年曾访问美国 Wright State 大学计算机科学与工程系，从事数据挖掘研究。先后发表论文40余篇。除本书外，近年来主持翻译的数据挖掘方面的著作还有 Jiawei Han 和 Micheline Kamber 的《数据挖掘：概念与技术》，Trevor Hastie、Robert Tibshirani 和 Jerome Friedman 的《统计学习基础：数据挖掘、推理与预测》。



范宏建 1999年毕业于郑州大学计算机科学系，同年进入中国科学院软件研究所攻读硕士学位，次年赴澳大利亚墨尔本大学攻读博士学位，师从澳大利亚科学院院士 Kotagiri Ramamohanarao 教授，2004年获计算机科学博士学位。先后在 WWW、PAKDD、RSFDGrC、IEEE GrC 和 Australian AI 等国际学术会议和 *IEEE Transactions on Knowledge and Data Engineering* 发表论文 10 余篇。目前是澳大利亚 AUSTRAC 的高级分析师，主要从事利用数据挖掘和机器学习技术进行金融数据分析的工作。

前 言

数据生成和收集技术的进步促使商业和科研领域产生了海量数据集。数据仓库能够存储多种数据，如：企业销售和运作的详细情况，地球轨道卫星发送回地球的高分辨率图像和遥感数据，对越来越多的有机体进行的基因组实验产生的序列、结构和机能数据。收集和存储数据变得轻松简便，已经完全改变了人们对数据分析的态度，人们开始尽可能地收集各个时期和各种来源的数据。人们相信收集的数据肯定会有价值，或者当初收集它就有明确的目的，或者只是先收集起来再说。

传统数据分析技术在应对这些新型数据集提出的挑战时存在种种局限性，而数据挖掘技术突破了这些局限。数据挖掘并不是要取代其他分析领域，而是以它们为基础。尽管数据挖掘的某些主题（如关联分析）是其独有的，但是，还有许多主题（如聚类、分类和异常检测）则建立在其他领域长期工作的基础之上。事实上，数据挖掘研究者们主动利用已有技术对增强和拓展这个领域以及推动它的快速发展起到了促进作用。

该领域一直强调与其他领域的研究者合作，因而充满了活力。要迎接新类型数据分析的挑战，抛开理解数据的人和数据所处的领域而简单地使用数据分析技术是不可行的。通常，能否组建好多学科研究团队，已经成为数据挖掘项目（如创建新的独创性算法）成败的决定因素。正如历史上统计学的许多进展都是由农业、工业、医疗卫生和商业需求推动的一样，今天，数据挖掘的许多进展也正在被这些领域的需求所推动。

自 1998 年春季开始，我们在明尼苏达大学为高年级本科生和研究生开设了数据挖掘课程。为这些课程准备的演示幻灯片和习题随着时间不断积累，成为本书的基础。数据挖掘的聚类技术综述最初是为该领域的某项研究而写的，它也成为本书第 8 章的雏形。随着时间的推移，又增加了关于数据、分类、关联分析和异常检测的几章。本书定稿后已在作者所在的学校（明尼苏达大学和密歇根州立大学）以及其他一些大学作为教材试用。

在此期间，出现了许多数据挖掘方面的书籍，但是都不能完全满足我们学生的需要——他们主要是计算机科学专业的研究生和本科生，也包括来自工科和其他专业的学生。他们的数学和计算机背景差异很大，但是都有一个共同目标：尽可能直接地学习数据挖掘，尽快地将其应用到各自的领域。因此，要求较多数学和统计学预备知识的书对他们中的许多人没有吸引力，需要坚实的数据库背景的书也有同样的问题。为了满足这些学生需求而逐渐写成的本书，现在的完稿使用了大量例子、习题并用简洁的语言描述了关键算法，尽可能直接把重点放在数据挖掘的主要概念上。

概述

具体而言，本书全面介绍了数据挖掘，方便学生、教师、研究人员和专业人士理解有关概念和技术。本书所涵盖的领域包括数据预处理、可视化、预测建模、关联分析、聚类和异常检测。

目标是讲述每个主题的基本概念和算法，从而为读者提供将数据挖掘应用于实际问题所需的必要背景。此外，本书也为有志于从事数据挖掘和相关领域研究的读者提供一个起点。

本书涵盖五个主题：数据、分类、关联分析、聚类和异常检测。除异常检测外，每个主题都分两章讲述。对于分类、关联分析和聚类，前面一章讲述基本概念、代表性算法和评估技术，后面一章深入讨论高级概念和算法。这样做的目的是使读者透彻地理解数据挖掘的基础，同时论述更多重要的高级主题。由于这种安排，本书既可用作为教材又可用作参考书。

为了帮助读者理解书中概念，我们提供大量示例、图表和习题。每一章的结尾给出了文献注释，是为那些对更高级的主题、重要的历史文献和当前趋势感兴趣的读者提供的。

致教师

作为一本教材，本书广泛适合于高年级本科生和研究生。由于学习这门课程的学生背景不同，他们可能不具备广博的统计学和数据库知识，因此本书只要求最低限度的预备知识——不需要数据库知识，并假定读者只有一般的统计学或数学背景。本书尽可能自成一体。统计学、线性代数和机器学习的必要基础知识或者已经融入正文，或者包含在附录中。

由于讨论主要数据挖掘主题的各章也是自成一体的，因此主题的讲授次序相当灵活。核心题材在第 2、4、6、8 和 10 章介绍。数据导论（第 2 章）应当最先讨论，基本的分类、关联分析和聚类（分别是第 4、6、8 章）可以以任意次序讲述。由于异常处理（第 10 章）与分类（第 4 章）和聚类（第 8 章）有一定的关系，这两章应当在第 10 章之前讲述。还可以根据课程安排和师生的兴趣从高级的分类、关联分析和聚类（分别为第 5、7、9 章）中选讲一些主题。我们也建议教师用数据挖掘的实际项目和练习强化课程的教学。尽管这样做很耗费时间，但是实践性的作业可以大大提高这门课程的价值。

支持材料

本书的教辅材料可以在 Addison-Wesley 的网站（www.aw-bc.com/cssupport）上找到^①。提供给所有读者的支持材料如下。

- 课程幻灯片。
- 学生项目建议。
- 数据挖掘资源，如数据挖掘算法和数据集。
- 联机指南，使用实际的数据集和数据分析软件，为本书介绍的部分数据挖掘技术提供例子讲解。

其他支持材料（包括习题答案）只向采纳本书做教材的教师提供。意见和建议以及勘误请通过 dmbok@cs.unm.edu 发给作者。

致谢

许多人都为本书做出了贡献。我们首先向家人表示感谢，这本书是献给他们的。没有他们的耐心和支持，不可能写出本书。

我们要感谢明尼苏达大学和密歇根州立大学数据挖掘小组的学生所做的贡献。Eui-Hong

^① 相关材料也可以从图灵网站（www.turingbook.com）本书网页免费注册下载。——编者注

(Sam) Han 和 Mahesh Joshi 帮助我们准备了最初的数据挖掘课程。他们编制的某些习题和演示幻灯片已经收录在本书及其辅助幻灯片中。小组中的其他学生也为本书的初稿提出建议或以各种方式做出贡献，他们是 Shyam Boriah、Haibin Cheng、Varun Chandola、Eric Eilertson、Levent Ertöz、Jing Gao、Rohit Gupta、Sridhar Iyer、Jung-Eun Lee、Benjamin Mayer、Aysel Ozgur、Uygar Oztekin、Gaurav Pandey、Kashif Riaz、Jerry Scripps、Gyorgy Simon、Hui Xiong、Jieping Ye 和 Pusheng Zhang。我们还要感谢明尼苏达大学和密歇根州立大学选修数据挖掘课程的学生，他们使用了本书的初稿，并提供了极富价值的反馈。我们特别感谢 Bernardo Craemer、Arifin Ruslim、Jamshid Vayghan 和 Yu Wei 的有益的建议。

Joydeep Ghosh (得克萨斯大学) 和 Sanjay Ranka (佛罗里达大学) 试用了本书的初稿。我们也直接从得克萨斯大学下列学生那里获得了许多有用的建议：Pankaj Adhikari、Rajiv Bhatia、Frederic Bosche、Arindam Chakraborty、Meghana Deodhar、Chris Everson、David Gardner、Saad Godil、Todd Hay、Clint Jones、Ajay Joshi、Joonsoo Lee、Yue Luo、Anuj Nanavati、Tyler Olsen、Sunyoung Park、Aashish Phansalkar、Geoff Prewett、Michael Ryoo、Daryl Shannon 和 Mei Yang。

Ronald Kostoff (ONR) 阅读了聚类部分的初稿，并提出了许多建议。Musetta Steinbach 发现了图中的一些错误。

我们要感谢明尼苏达大学和密歇根州立大学的同事，他们帮助创建了良好的数据挖掘研究环境。他们是 Dan Boley、Joyce Chai、Anil Jain、Ravi Janardan、Rong Jin、George Karypis、Haesun Park、William F. Punch、Shashi Shekhar 和 Jaideep Srivastava。我们还要向我们的数据挖掘项目的合作者表示谢意，他们是 Ramesh Agrawal、Steve Cannon、Piet C. de Groen、Fran Hill、Yongdae Kim、Steve Klooster、Kerry Long、Nihar Mahapatra、Chris Potter、Jonathan Shapiro、Kevin Silverstein、Nevin Young 和 Zhi-Li Zhang。

明尼苏达大学和密歇根州立大学的计算机科学与工程系为本书写作及研究提供了计算资源和支持环境。ARDA、ARL、ARO、DOE、NASA 和 NSF 等机构为本书作者提供了研究资助。特别应该提到的是，Kamal Abdali、Dick Brackney、Jagdish Chandra、Joe Coughlan、Michael Coyle、Stephen Davis、Frederica Darema、Richard Hirsch、Chandrika Kamath、Raju Namburu、N. Radhakrishnan、James Sidoran、Bhavani Thuraisingham、Walt Tiernin、Maria Zemankova 和 Xiaodong Zhang 有力地支持了我们的数据挖掘和高性能计算研究。

与培生出版集团的工作人员的合作令人愉快。具体地，我们要感谢 Michelle Brown、Matt Goldstein、Katherine Harutunian、Marilyn Lloyd、Kathy Smith 和 Joyce Wells。我们还要感谢 George Nichols 帮助绘图，Paul Anagnostopoulos 提供 L^AT_EX 支持。我们感谢出版社邀请的审稿人：Chien-Chung Chan (阿克伦大学)、Zhengxin Chen (内布拉斯加大学奥马哈分校)、Chris Clifton (普度大学)、Joydeep Ghosh (得克萨斯大学奥斯汀分校)、Nazli Goharian (伊利诺伊理工学院)、J. Michael Hardin (阿拉巴马大学)、James Hearne (西华盛顿大学)、Hillol Kargupta (马里兰大学巴尔的摩分校和 Agnik 公司)、Eamonn Keogh (加利福尼亚大学里弗赛德分校)、Bing Liu (伊利诺伊大学芝加哥分校)、Mariofanna Milanova (阿肯色大学小石城分校)、Srinivasan Parthasarathy (俄亥俄州立大学)、Zbigniew W. Ras (北卡罗莱纳大学夏洛特分校)、Xintao Wu (北卡罗莱纳大学夏洛特分校) 和 Mohammed J. Zaki (伦斯勒理工学院)。

目 录

第 1 章 绪论	1	第 3 章 探索数据	59
1.1 什么是数据挖掘.....	2	3.1 鸢尾花数据集.....	59
1.2 数据挖掘要解决的问题.....	2	3.2 汇总统计.....	60
1.3 数据挖掘的起源.....	3	3.2.1 频率和众数.....	60
1.4 数据挖掘任务.....	4	3.2.2 百分位数.....	61
1.5 本书的内容与组织.....	7	3.2.3 位置度量: 均值和中位数.....	61
文献注释.....	7	3.2.4 散布度量: 极差和方差.....	62
参考文献.....	8	3.2.5 多元汇总统计.....	63
习题.....	10	3.2.6 汇总数据的其他方法.....	64
第 2 章 数据	13	3.3 可视化.....	64
2.1 数据类型.....	14	3.3.1 可视化的动机.....	64
2.1.1 属性与度量.....	15	3.3.2 一般概念.....	65
2.1.2 数据集的类型.....	18	3.3.3 技术.....	67
2.2 数据质量.....	22	3.3.4 可视化高维数据.....	75
2.2.1 测量和数据收集问题.....	22	3.3.5 注意事项.....	79
2.2.2 关于应用的问题.....	26	3.4 OLAP 和多维数据分析.....	79
2.3 数据预处理.....	27	3.4.1 用多维数组表示鸢尾花数据.....	80
2.3.1 聚集.....	27	3.4.2 多维数据: 一般情况.....	81
2.3.2 抽样.....	28	3.4.3 分析多维数据.....	82
2.3.3 维归约.....	30	3.4.4 关于多维数据分析的最后评述.....	84
2.3.4 特征子集选择.....	31	文献注释.....	84
2.3.5 特征创建.....	33	参考文献.....	85
2.3.6 离散化和二元化.....	34	习题.....	86
2.3.7 变量变换.....	38	第 4 章 分类: 基本概念、决策树与模型	
2.4 相似性和相异性的度量.....	38	评估	89
2.4.1 基础.....	39	4.1 预备知识.....	89
2.4.2 简单属性之间的相似度和相		4.2 解决分类问题的一般方法.....	90
异度.....	40	4.3 决策树归纳.....	92
2.4.3 数据对象之间的相异度.....	41	4.3.1 决策树的工作原理.....	92
2.4.4 数据对象之间的相似度.....	43	4.3.2 如何建立决策树.....	93
2.4.5 邻近性度量的例子.....	43	4.3.3 表示属性测试条件的方法.....	95
2.4.6 邻近度计算问题.....	48	4.3.4 选择最佳划分的度量.....	96
2.4.7 选取正确的邻近性度量.....	50	4.3.5 决策树归纳算法.....	101
文献注释.....	50	4.3.6 例子: Web 机器人检测.....	102
参考文献.....	52	4.3.7 决策树归纳的特点.....	103
习题.....	53	4.4 模型的过分拟合.....	106

4.4.1	噪声导致的过分拟合	107	5.5.4	非线性支持向量机	164
4.4.2	缺乏代表性样本导致的过分拟合	109	5.5.5	支持向量机的特征	168
4.4.3	过分拟合与多重比较过程	109	5.6	组合方法	168
4.4.4	泛化误差估计	110	5.6.1	组合方法的基本原理	168
4.4.5	处理决策树归纳中的过分拟合	113	5.6.2	构建组合分类器的方法	169
4.5	评估分类器的性能	114	5.6.3	偏倚-方差分解	171
4.5.1	保持方法	114	5.6.4	装袋	173
4.5.2	随机二次抽样	115	5.6.5	提升	175
4.5.3	交叉验证	115	5.6.6	随机森林	178
4.5.4	自助法	115	5.6.7	组合方法的实验比较	179
4.6	比较分类器的方法	116	5.7	不平衡类问题	180
4.6.1	估计准确度的置信区间	116	5.7.1	可选度量	180
4.6.2	比较两个模型的性能	117	5.7.2	接受者操作特征曲线	182
4.6.3	比较两种分类法的性能	118	5.7.3	代价敏感学习	184
文献注释		118	5.7.4	基于抽样的方法	186
参考文献		120	5.8	多类问题	187
习题		122	文献注释		189
参考文献		120	参考文献		190
习题		122	习题		193
第5章 分类：其他技术		127	第6章 关联分析：基本概念和算法		201
5.1	基于规则的分类器	127	6.1	问题定义	202
5.1.1	基于规则的分类器的工作原理	128	6.2	频繁项集的产生	204
5.1.2	规则的排序方案	129	6.2.1	先验原理	205
5.1.3	如何建立基于规则的分类器	130	6.2.2	Apriori 算法的频繁项集产生	206
5.1.4	规则提取的直接方法	130	6.2.3	候选的产生与剪枝	208
5.1.5	规则提取的间接方法	135	6.2.4	支持度计数	210
5.1.6	基于规则的分类器的特征	136	6.2.5	计算复杂度	213
5.2	最近邻分类器	137	6.3	规则产生	215
5.2.1	算法	138	6.3.1	基于置信度的剪枝	215
5.2.2	最近邻分类器的特征	138	6.3.2	Apriori 算法中规则的产生	215
5.3	贝叶斯分类器	139	6.3.3	例：美国国会投票记录	217
5.3.1	贝叶斯定理	139	6.4	频繁项集的紧凑表示	217
5.3.2	贝叶斯定理在分类中的应用	140	6.4.1	极大频繁项集	217
5.3.3	朴素贝叶斯分类器	141	6.4.2	闭频繁项集	219
5.3.4	贝叶斯误差率	145	6.5	产生频繁项集的其他方法	221
5.3.5	贝叶斯信念网络	147	6.6	FP 增长算法	223
5.4	人工神经网络	150	6.6.1	FP 树表示法	224
5.4.1	感知器	151	6.6.2	FP 增长算法的频繁项集产生	225
5.4.2	多层人工神经网络	153	6.7	关联模式的评估	228
5.4.3	人工神经网络的特点	155	6.7.1	兴趣度的客观度量	228
5.5	支持向量机	156	6.7.2	多个二元变量的度量	235
5.5.1	最大边缘超平面	156	6.7.3	辛普森悖论	236
5.5.2	线性支持向量机：可分情况	157	6.8	倾斜支持度分布的影响	237
5.5.3	线性支持向量机：不可分情况	162			

文献注释	240	8.2.3 二分 K 均值	316
参考文献	244	8.2.4 K 均值和不同的簇类型	317
习题	250	8.2.5 优点与缺点	318
		8.2.6 K 均值作为优化问题	319
第 7 章 关联分析: 高级概念	259	8.3 凝聚层次聚类	320
7.1 处理分类属性	259	8.3.1 基本凝聚层次聚类算法	321
7.2 处理连续属性	261	8.3.2 特殊技术	322
7.2.1 基于离散化的方法	261	8.3.3 簇邻近度的 Lance-Williams	
7.2.2 基于统计学的方法	263	公式	325
7.2.3 非离散化方法	265	8.3.4 层次聚类的主要问题	326
7.3 处理概念分层	266	8.3.5 优点与缺点	327
7.4 序列模式	267	8.4 DBSCAN	327
7.4.1 问题描述	267	8.4.1 传统的密度: 基于中心的方法	327
7.4.2 序列模式发现	269	8.4.2 DBSCAN 算法	328
7.4.3 时限约束	271	8.4.3 优点与缺点	329
7.4.4 可选计数方案	274	8.5 簇评估	330
7.5 子图模式	275	8.5.1 概述	332
7.5.1 图与子图	276	8.5.2 非监督簇评估: 使用凝聚度和	
7.5.2 频繁子图挖掘	277	分离度	332
7.5.3 类 Apriori 方法	278	8.5.3 非监督簇评估: 使用邻近度	
7.5.4 候选产生	279	矩阵	336
7.5.5 候选剪枝	282	8.5.4 层次聚类的非监督评估	338
7.5.6 支持度计数	285	8.5.5 确定正确的簇个数	339
7.6 非频繁模式	285	8.5.6 聚类趋势	339
7.6.1 负模式	285	8.5.7 簇有效性的监督度量	340
7.6.2 负相关模式	286	8.5.8 评估簇有效性度量的显著性	343
7.6.3 非频繁模式、负模式和负相关		文献注释	344
模式比较	287	参考文献	345
7.6.4 挖掘有趣的非频繁模式的技术	288	习题	347
7.6.5 基于挖掘负模式的技术	288		
7.6.6 基于支持度期望的技术	290	第 9 章 聚类分析: 其他问题与算法	355
文献注释	292	9.1 数据、簇和聚类算法的特性	355
参考文献	293	9.1.1 例子: 比较 K 均值和	
习题	295	DBSCAN	355
		9.1.2 数据特性	356
第 8 章 聚类分析: 基本概念和算法	305	9.1.3 簇特性	357
8.1 概述	306	9.1.4 聚类算法的一般特性	358
8.1.1 什么是聚类分析	306	9.2 基于原型的聚类	359
8.1.2 不同的聚类类型	307	9.2.1 模糊聚类	359
8.1.3 不同的簇类型	308	9.2.2 使用混合模型的聚类	362
8.2 K 均值	310	9.2.3 自组织映射	369
8.2.1 基本 K 均值算法	310	9.3 基于密度的聚类	372
8.2.2 K 均值: 附加的问题	315	9.3.1 基于网格的聚类	372

9.3.2	子空间聚类	374	10.1.3	类标号的使用	405
9.3.3	DENCLUE: 基于密度聚类的 一种基于核的方案	377	10.1.4	问题	405
9.4	基于图的聚类	379	10.2	统计方法	406
9.4.1	稀疏化	379	10.2.1	检测一元正态分布中的 离群点	407
9.4.2	最小生成树聚类	380	10.2.2	多元正态分布的离群点	408
9.4.3	OPOSSUM: 使用 METIS 的 稀疏相似度最优划分	381	10.2.3	异常检测的混合模型方法	410
9.4.4	Chameleon: 使用动态建模的 层次聚类	381	10.2.4	优点与缺点	411
9.4.5	共享最近邻相似度	385	10.3	基于邻近度的离群点检测	411
9.4.6	Jarvis-Patrick 聚类算法	387	10.4	基于密度的离群点检测	412
9.4.7	SNN 密度	388	10.4.1	使用相对密度的离群点检测	413
9.4.8	基于 SNN 密度的聚类	389	10.4.2	优点与缺点	414
9.5	可伸缩的聚类算法	390	10.5	基于聚类的技术	414
9.5.1	可伸缩: 一般问题和方法	391	10.5.1	评估对象属于簇的程度	415
9.5.2	BIRCH	392	10.5.2	离群点对初始聚类的影响	416
9.5.3	CURE	393	10.5.3	使用簇的个数	416
9.6	使用哪种聚类算法	395	10.5.4	优点与缺点	416
	文献注释	397		文献注释	417
	参考文献	398		参考文献	418
	习题	400		习题	420
第 10 章	异常检测	403	附录 A	线性代数	423
10.1	预备知识	404	附录 B	维归约	433
10.1.1	异常的成因	404	附录 C	概率统计	445
10.1.2	异常检测方法	404	附录 D	回归	451
			附录 E	优化	457

绪论

数据收集和数据存储技术的快速进步使得各组织机构可以积累海量数据。然而，提取有用的信息已经成为巨大的挑战。通常，由于数据量太大，无法使用传统的分析工具和技术处理它们。有时，即使数据集相对较小，但由于数据本身具有一些非传统特点，也不能使用传统的方法处理。在另外一些情况下，面临的问题不能使用已有的数据分析技术来解决。这样，就需要开发新的方法。

数据挖掘是一种技术，它将传统的分析方法与处理大量数据的复杂算法相结合。数据挖掘为探查和分析新的数据类型以及用新方法分析旧有数据类型提供了令人振奋的机会。本章，我们概述数据挖掘，并列举本书所涵盖的关键主题。我们首先介绍需要新的数据分析技术的一些大家熟知的应用。

商务 借助 POS（销售点）数据收集技术[条码扫描器、射频识别（RFID）和智能卡技术]，零售商可以在其商店的收银台收集顾客购物的最新数据。零售商可以利用这些信息，加上电子商务网站的日志、电购中心的顾客服务记录等其他的重要商务数据，更好地理解顾客的需求，做出明智的商务决策。

数据挖掘技术可以用来支持广泛的商务智能应用，如顾客分析、定向营销、 workflow 管理、商店分布和欺诈检测等。数据挖掘还能帮助零售商回答一些重要的商务问题，如“谁是最有价值的顾客？”“什么产品可以交叉销售^①或提升销售^②？”“公司明年的收入前景如何？”这些问题催生了一种新的数据分析技术——关联分析（见第 6 章和第 7 章）。

医学、科学与工程 医学、科学技术界的研究者正在快速积累大量数据，这些数据对获得有价值的新发现至关重要。例如，为了更深入地理解地球的气候系统，NASA 已经部署了一系列的地球轨道卫星，不停地收集地表、海洋和大气的全球观测数据。然而，由于这些数据的规模和时空特性，传统的方法常常不适合分析这些数据集。数据挖掘开发的技术可以帮助地球科学家回答如下问题：“干旱和飓风等生态系统扰动的频度和强度与全球变暖之间有何联系？”“海洋表面温度对地表降水量和温度有何影响？”“如何准确地预测一个地区的生长季节的开始和结束？”

再举一个例子，分子生物学研究者希望利用当前收集的大量基因组数据，更好地理解基因的结构和功能。过去，传统方法只允许科学家在一个实验中每次研究少量基因。微阵列技术的最新突破已经能让科学家在多种情况下，比较数以千计的基因特性。这种比较有助于确定每个基因的作用，或许可以查出导致特定疾病的基因。然而，由于数据的噪声和高维性，需要新的数据分析

① cross-sell，指根据顾客的兴趣推荐或显示相关商品以增加销售机会。——译者注

② up-sell，指尝试向曾经购买的顾客销售价格更高的商品。——译者注

方法。除分析基因序列数据外，数据挖掘还能用来处理生物学的其他难题，如蛋白质结构预测、多序列校准、生物化学路径建模和种系发生学。

1.1 什么是数据挖掘

数据挖掘是在大型数据存储库中，自动地发现有有用信息的过程。数据挖掘技术用来探查大型数据库，发现先前未知的有用模式。数据挖掘还可以预测未来观测结果，例如，预测一位新的顾客是否会在一家百货公司消费 100 美元以上。

并非所有的信息发现任务都被视为数据挖掘。例如，使用数据库管理系统查找个别的记录，或通过因特网的搜索引擎查找特定的 Web 页面，则是信息检索（information retrieval）领域的任务。虽然这些任务非常重要，可能涉及使用复杂的算法和数据结构，但是它们主要依赖传统的计算机科学技术和数据的明显特征来创建索引结构，从而有效地组织和检索信息。尽管如此，人们也在利用数据挖掘技术增强信息检索系统的能力。

数据挖掘与知识发现

数据挖掘是数据库中知识发现（knowledge discovery in database, KDD）不可缺少的一部分，而 KDD 是将未加工的数据转换为有用信息的整个过程，如图 1-1 所示。该过程包括一系列转换步骤，从数据的预处理到数据挖掘结果的后处理。

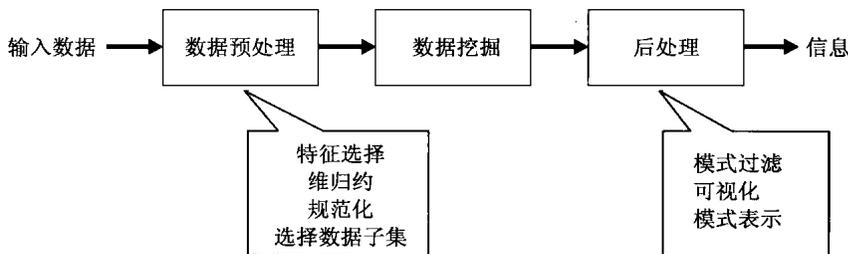


图 1-1 数据库中知识发现 (KDD) 过程

输入数据可以以各种形式存储（平展文件、电子数据表或关系表），并且可以驻留在集中的数据存储库中，或分布在多个站点上。数据预处理（preprocessing）的目的是将未加工的输入数据转换成适合分析的形式。数据预处理涉及的步骤包括融合来自多个数据源的数据，清洗数据以消除噪声和重复的观测值，选择与当前数据挖掘任务相关的记录和特征。由于收集和存储数据的方式多种多样，数据预处理可能是整个知识发现过程中最费力、最耗时的步骤。

“结束循环”（closing the loop）通常指将数据挖掘结果集成到决策支持系统的过程。例如，在商业应用中，数据挖掘的结果所揭示的规律可以结合商业活动管理工具，从而开展或测试有效的商品促销活动。这样的结合需要后处理（postprocessing）步骤，确保只将那些有效的和有用的结果集成到决策支持系统中。后处理的一个例子是可视化（见第 3 章），它使得数据分析者可以从各种不同的视角探查数据和数据挖掘结果。在后处理阶段，还能使用统计度量或假设检验，删除虚假的数据挖掘结果。

1.2 数据挖掘要解决的问题

前面提到，面临新的数据集带来的问题时，传统的数据分析技术常常遇到实际困难。下面是