

世界著名计算机教材精选

# Web数据挖掘

Bing Liu 著

俞勇 薛贵荣  
韩定一 译



## WEB DATA MINING

清华大学出版社

世界著名计算机教材精选

**Web Data Mining**

# Web 数据挖掘

Bing Liu 著  
俞 勇  
薛贵荣 译  
韩定一

清华大学出版社  
北京

English reprint edition copyright © 2009 by Springer-Verlag and TSINGHUAUNIVERSITY PRESS.

Original English language title from Proprietor's edition of the Work.

Original English language title: Web Data Mining by Bing Liu. Copyright © 2009

All Rights Reserved.

This edition has been authorized by Springer-Verlag (Berlin/Heidelberg/New York) for sale in the People's Republic of China only and not for export therefrom.

本书中文翻译版由 Springer-Verlag 授权给清华大学出版社出版发行。

北京市版权局著作权合同登记号 图字 01-2008-0564 号

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

Web 数据挖掘/(美)刘兵(Liu, B.)著;俞勇等译.—北京: 清华大学出版社, 2009.4  
(世界著名计算机教材精选)

书名原文: Web Data Mining

ISBN 978-7-302-19338-8

I. W… II. ①刘… ②俞… III. 数据采集—教材 IV. TP311.13

中国版本图书馆 CIP 数据核字(2009)第 010593 号

责任编辑: 龙啟铭

责任校对: 徐俊伟

责任印制: 李红英

出版发行: 清华大学出版社 地址: 北京清华大学学研大厦 A 座

<http://www.tup.com.cn> 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

质 量 反 馈: 010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者: 北京密云胶印厂

装 订 者: 三河市兴旺装订有限公司

经 销: 全国新华书店

开 本: 185×260 印 张: 24.5 字 数: 594 千字

版 次: 2009 年 4 月第 1 版 印 次: 2009 年 4 月第 1 次印刷

印 数: 1~3000

定 价: 49.00 元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系调换。

联系电话: 010-62770177 转 3103 产品编号: 027154-01

## 译者序

作为互联网上最重要的应用之一,Web(万维网)提供了便捷的文档发布与获取机制,并逐步成为各类信息资源的聚集地。据 Google 于 2008 年发布的官方报告,它们已经在互联网上发现超过 1 万亿个 Web 文档,而且这个数字还在以每天几十亿的速度持续增长。面对如此巨大的信息量,普通 Web 用户往往迷失其中,他们迫切需要一种机制快速定位到所需信息。Web 数据挖掘便应运而生,并且伴随 Web 的发展而备受关注。Web 数据挖掘它建立在信息检索、数据挖掘以及知识管理等技术的基础上,通过对大量 Web 文档进行分析来获得隐含的知识和模式,从而帮助人们更好地进行信息搜索和决策制定。反过来,可以说,也正是 Web 挖掘技术的不断进展,推动了 Web 的进一步蓬勃发展。

目前 Web 数据挖掘已经引起了学术界、工业界、社会学家的广泛关注,也吸引了众多研究人员与开发人员投身其中。国内外很多大学与研究机构先后开设了 Web 挖掘课程。但长期以来并没有专门针对 Web 挖掘的教材与专著。刘兵教授出版的这本著作填补了该领域的空白。该教材针对 Web 挖掘中众多关键主题进行了深入分析。清华大学出版社独具慧眼,决定将该书翻译成中文版在国内出版,这必将对我国 Web 挖掘的教学与研究产生积极的推动作用,有幸承担该书的翻译工作,我们感到十分荣幸。

本书是由伊利诺伊大学芝加哥分校(UIC)的刘兵(Bing Liu)教授历经一年的时间所著的"Web Data Mining"的翻译版。刘兵教授是 Web 挖掘研究领域的国际知名专家,曾担任多个国际期刊的编辑,也是多个国际学术会议(如 WWW、KDD 与 AAAI 等)的程序委员会委员。刘兵教授在 Web 内容挖掘、互联网观点挖掘、数据挖掘等领域有非常高的造诣。他先后在国际著名学术期刊与重要国际学术会议上发表论文一百多篇。本教材中的部分章节也融入了刘兵教授从事 Web 挖掘研究多年的心血。

全书主要包括前言和 12 个章节。本书的翻译和审校由俞勇、薛贵荣和韩定一共同完成。其中,俞勇负责前言、第 1 章和第 2 章,薛贵荣负责第 3~7 章,韩定一负责第 8~12 章。参加翻译工作的还有韩定一(前言、第 1 章、第 8 章)、徐生良(第 2 章)、凌霄(第 3 章)、郭晋文(第 4 章、第 5 章)、王亮(第 6 章)、陈林虎(第 7 章)、傅临云(第 9 章)、第 7 张迪(第 10 章)、包胜华(第 11

## 译者序

章)和王乐天(第 12 章)等。上海交通大学 APEX 数据和知识管理实验室的全体同学参加了本书的校对工作。

在本书的翻译过程中,得到了刘兵教授的大力支持。他向译者提供了全文书稿的最终版本,并对翻译工作提出了指导性建议。同时,感谢微软亚洲研究院李航博士的引荐,使我们有机会学习和翻译此书。最后,感谢清华大学出版社的龙启铭编辑,是他的远见,使得本书能够尽快与读者见面。

由于本书所涉及到内容非常广泛,许多术语目前尚无固定译法,翻译难度相对较大。尽管我们对某些术语进行了推敲,但仍然可能出现词不达意的地方。此外,由于译者水平有限,译文中不当之处也在所难免。我们也真诚地希望同行与读者朋友们不吝赐教。如果您能将您的意见与建议发往 yyu@apex.sjtu.edu.cn,我们将不胜感激。

由于时间仓促,书中存在许多不足,敬请各位读者批评指正。译者在此向所有关心和支持本书的读者表示衷心的感谢!

# 序言

过去几十年里,Web 的迅速发展使其成为世界上规模最大的公共数据源。紧跟字典源。Web 数据挖掘的目标是从 Web 超链接、网页内容和使用日志中探寻有用的信息。依据在挖掘过程中使用的数据类别,Web 挖掘的任务可以被划分为三种主要类型:Web 结构挖掘、Web 内容挖掘和 Web 使用挖掘。Web 结构挖掘从表征 Web 结构的超链接中寻找知识。Web 内容挖掘从网页内容中抽取有用的信息和知识。而 Web 使用挖掘则从记录每位用户点击情况的使用日志中挖掘用户的访问模式。

本书旨在讲述这些任务以及它们的核心挖掘算法;尽可能涵盖每个话题的广泛内容,给出足够多的细节,以便读者无须借助额外的阅读,即可获得相对完整的关于算法和技术的知识。其中结构化数据的抽取、信息整合、观点挖掘和 Web 使用挖掘等 4 章是本书的特色,这些内容在已有书籍中没有提及,但它们在 Web 数据挖掘中却占有非常重要的地位。当然,传统的 Web 挖掘主题,如搜索、页面爬取和资源探索以及链接分析在书中也作了详细描述。

本书尽管题为“Web 数据挖掘”,却依然涵盖了数据挖掘和信息检索的核心主题;因为 Web 挖掘大量使用了它们的算法和技术。数据挖掘部分主要由关联规则和序列模式、监督学习(分类)、无监督学习(聚类)这三大最重要的数据挖掘任务,以及半监督学习这个相对深入的主题组成。而信息检索对于 Web 挖掘而言最重要的核心主题都有所阐述。因此,本书自然地分为两大部分。第一部分,包括第 2~5 章,介绍数据挖掘的基础。第二部分,本章包括第 6~12 章,介绍 Web 相关的挖掘任务。

有两大指导性原则贯穿本书始末。其一,本书的基础内容适合本科生阅读,但也包括足够的深度资料,以满足打算在 Web 数据挖掘和相关领域研读博士学位的研究生。书中对读者的预备知识几乎没有作任何要求,任何对算法和概率知识稍有理解的人都应当能够顺利地读完本书。其二,本书从实践的角度来审视 Web 挖掘的技术。这一点非常重要,因为大多数 Web 挖掘任务都在现实世界中有所应用。在过去的几年中,我有幸直接或间接地与许多研究人员和工程人员一起工作,他们来自于多个搜索引擎、电子商务公司,甚至是利用 Web 信息感兴趣的传统公司。在这个时代,对前者

## 序言

过程中,我获得了许多现实世界问题的实践经历和第一手知识。我尽量将其中非机密的信息和知识通过本书传递给读者,因此本书能在理论和实践中有所平衡。我希望本书不仅能够成为学生的教科书,也能成为 Web 挖掘研究人员和实践人员获取知识、信息,甚至是创新想法的一个有效渠道。

## 致 谢

在撰写本书的过程中,许多研究人员都给予了我无私的帮助;没有他们的帮助,这本书也许永远无法成为现实。我最深切的感谢要给予 Filippo Menczer 和 Bamshad Mobasher,他们热情地撰写了本书中重要的两个章节,他们也是相关领域的专家。Filippo 负责 Web 爬取这一章,Bamshad 负责 Web 使用挖掘这一章。我还要感谢 Wee Sun Lee(李伟上),他帮助完成第 5 章半监督学习的很大一部分。

Jian Pei(裴健)帮助撰写了第 2 章中 PrefixSpan 算法,并且检查了 MS-PS 算法。Eduard Dragut 帮助撰写了第 10 章的最后一节,并且多次阅读并修改这一整章。Yuanlin Zhang 对第 9 章提出很多意见。我对他们所有人都有所亏欠。

还有许多研究人员以各种方式提供了帮助。Yang Dai(戴阳)和 Rudy Setiono 在支持向量机(SVM)上提供帮助。Chris Ding(丁宏强)对链接分析提供了帮助。Clement Yu(余德)和 ChengXiang Zhai(翟成祥)阅读了第 6 章。Amy Langville 阅读了第 7 章。Kevin C.-C. Chang(张振川)、Ji-Rong Wen(文继荣)和 Clement Yu(余德)帮助了第 10 章的许多方面。Justin Zobel 帮助理清了索引压缩的许多议题。Ion Muslea 帮助理清了包裹简介的一些议题。Divy Agrawal、Yunbo Cao(曹云波)、Edward Fox、Hang Li(李航)、Xiaoli Li(李晓黎)、Zhaojun Tan、Dell Zhang(张德)和 Zijian Zheng 帮助检查了各个章节。在此对他们表示感谢!

和许多研究人员的讨论也帮助本书的成形。这些人包括 Amir Ashkenazi、Imran Aziz、Roberto Bayardo、Wendell Baker、Ling Bao、Jeffrey Benkler、AnHai Doan、Byron Dom、Michael Gamon、Robert Grossman、Jiawei Han(韩家炜)、Wynne Hsu、Ronny Kohavi、David D. Lewis、Ian McAllister、Wei-Ying Ma(马维英)、Marco Maggini、Llew Mason、Kamel Nigan、Julian Qian、Yan Qu、Thomas M. Tirpak、Andrew Tomkins、Alexander Tuzhilin、Weimin Xiao、Gu Xu(徐谷)、Philip S. Yu 和 Mohammed Zaki。

我的学生们(不论已毕业或是在读)检查了许多算法的正确性并且作出了许多修正。他们包括 Gao Cong(从高)、Minqing Hu、Nitin Jindal、Xin Li、Yiming Ma、Yanhong Zhai 和 Kaidi Zhao。本书中一些章节是我在伊利诺伊斯大学芝加哥分校的研究生课程讲义。我要感谢这些课程的学生帮我实现了一部分算法。他们提出的问题在某些情况下也帮助我修正算法。在这里我不可能完全列出他们的名字,但我要特别感谢 John Castano、Xiaowen Ding、Murthy Ganapathibhotla、Cynthia Kersey、Hari Prasad、Divyakotti、Ravikanth Turlapati、Srikanth Tadikonda、Makio Tamura、Haisheng Wang 和 Chad Williams,他们指出讲义中文本、举例或算法的错误。来自德保罗大学的 Michael Bombyk 也指出了不少笔误。

与 Springer 出版社的员工一起工作是一段令人愉快的经历。感谢编辑 Ralf Gerstner 在 2005 年初征询我对撰写一本有关 Web 挖掘的书籍是否感兴趣。从那以后,我们一直保持着愉快的合作经历。我还要感谢校对 Mike Nugent 提高了本书内容的表达质量,以及制

## 序言

作编辑 Michael Reinfarth 引导我顺利完成了本书的出版过程。还有两位匿名评审也给出不少有见解的评论。伊利诺伊斯大学芝加哥分校计算机科学系对本项目提供了计算资源和工作环境的支持。

最后,我要感谢我的父母和兄弟姐妹,他们给予我一贯的支持和鼓励。我将最深刻的情感给予我自己的家庭成员:Yue、Shelley 和 Kate。他们也在许多方面给予支持和帮助。尽管 Shelley 和 Kate 还年幼,但他们阅读了本书的绝大部分,并且找出了不少笔误。我的妻子将家里一切事情打理得秩序井然,使我可以将充分的时间和精力用在这本书上。谨以此书献给他们!

Bing Liu(刘兵)

# 目录

## 第一部分 数据挖掘基础

|                  |    |
|------------------|----|
| 第1章 概述           | 3  |
| 1.1 什么是万维网       | 3  |
| 1.2 万维网和互联网的历史简述 | 4  |
| 1.3 Web数据挖掘      | 5  |
| 1.3.1 什么是数据挖掘    | 6  |
| 1.3.2 什么是Web数据挖掘 | 7  |
| 1.4 各章概要         | 8  |
| 1.5 如何阅读本书       | 10 |
| 文献评注             | 10 |

|               |    |
|---------------|----|
| 第2章 关联规则和序列模式 | 12 |
|---------------|----|

|                      |    |
|----------------------|----|
| 2.1 关联规则的基本概念        | 12 |
| 2.2 Apriori 算法       | 14 |
| 2.2.1 频繁项目集生成        | 14 |
| 2.2.2 关联规则生成         | 17 |
| 2.3 关联规则挖掘的数据格式      | 19 |
| 2.4 多最小支持度的关联规则挖掘    | 20 |
| 2.4.1 扩展模型           | 21 |
| 2.4.2 挖掘算法           | 22 |
| 2.4.3 规则生成           | 26 |
| 2.5 分类关联规则挖掘         | 27 |
| 2.5.1 问题描述           | 27 |
| 2.5.2 挖掘算法           | 28 |
| 2.5.3 多最小支持度分类关联规则挖掘 | 31 |
| 2.6 序列模式的基本概念        | 31 |
| 2.7 基于GSP挖掘序列模式      | 32 |
| 2.7.1 GSP算法          | 33 |

## 目录

|   |           |
|---|-----------|
| 2.7.2 多最小支持度挖掘 .....                              | 34        |
| 2.8 基于 PrefixSpan 算法的序列模式挖掘 .....                 | 37        |
| 2.8.1 PrefixSpan 算法 .....                         | 38        |
| 2.8.2 多最小支持度挖掘 .....                              | 39        |
| 2.9 从序列模式中产生规则 .....                              | 41        |
| 2.9.1 序列规则 .....                                  | 41        |
| 2.9.2 标签序列规则 .....                                | 41        |
| 2.9.3 分类序列规则 .....                                | 42        |
| 文献评注 .....  | 42        |
| <b>第3章 监督学习 .....</b>                             | <b>45</b> |
| 3.1 基本概念 .....                                    | 45        |
| 3.2 决策树推理 .....                                   | 48        |
| 3.2.1 学习算法 .....                                  | 49        |
| 3.2.2 混杂度函数 .....                                 | 50        |
| 3.2.3 处理连续属性 .....                                | 53        |
| 3.2.4 其他一些问题 .....                                | 54        |
| 3.3 评估分类器 .....                                   | 56        |
| 3.3.1 评估方法 .....                                  | 56        |
| 3.3.2 查准率、查全率、F-score 和平衡点(Breakeven Point) ..... | 57        |
| 3.4 规则推理 .....                                    | 59        |
| 3.4.1 序列化覆盖 .....                                 | 59        |
| 3.4.2 规则学习：Learn-One-Rule 函数 .....                | 61        |
| 3.4.3 讨论 .....                                    | 63        |
| 3.5 基于关联规则的分类 .....                               | 63        |
| 3.5.1 使用类关联规则进行分类 .....                           | 64        |
| 3.5.2 使用类关联规则作为分类属性 .....                         | 66        |
| 3.5.3 使用古典的关联规则分类 .....                           | 66        |
| 3.6 朴素贝叶斯分类 .....                                 | 67        |
| 3.7 朴素贝叶斯文本分类 .....                               | 70        |
| 3.7.1 概率框架 .....                                  | 70        |
| 3.7.2 朴素贝叶斯模型 .....                               | 71        |
| 3.7.3 讨论 .....                                    | 73        |
| 3.8 支持向量机 .....                                   | 73        |
| 3.8.1 线性支持向量机：可分的情况 .....                         | 74        |
| 3.8.2 线性支持向量机：数据不可分的情况 .....                      | 78        |
| 3.8.3 非线性支持向量机：核方法 .....                          | 80        |
| 3.9 k-近邻学习 .....                                  | 82        |
| 3.10 分类器的集成 .....                                 | 83        |

|  |            |
|--|------------|
| 3.10.1 Bagging .....                                 | 83         |
| 3.10.2 Boosting .....                                | 84         |
| 3.11 文献评注 .....                                      | 84         |
| <b>第4章 无监督学习 .....</b>                               | <b>87</b>  |
| 4.1 基本概念 .....                                       | 87         |
| 4.2 k-均值聚类 .....                                     | 89         |
| 4.2.1 k-均值算法 .....                                   | 89         |
| 4.2.2 k-均值算法的硬盘版本 .....                              | 91         |
| 4.2.3 优势和劣势 .....                                    | 92         |
| 4.3 聚类的表示 .....                                      | 95         |
| 4.3.1 聚类的一般表示方法 .....                                | 95         |
| 4.3.2 任意形状的聚类 .....                                  | 95         |
| 4.4 层次聚类 .....                                       | 96         |
| 4.4.1 单链接方法 .....                                    | 97         |
| 4.4.2 全链接方法 .....                                    | 98         |
| 4.4.3 平均链接方法 .....                                   | 98         |
| 4.4.4 优势和劣势 .....                                    | 98         |
| 4.5 距离函数 .....                                       | 99         |
| 4.5.1 数值的属性(Numeric Attributes) .....                | 99         |
| 4.5.2 布尔属性和符号属性(Binary and Nominal Attributes) ..... | 99         |
| 4.5.3 文本文档 .....                                     | 101        |
| 4.6 数据标准化 .....                                      | 101        |
| 4.7 混合属性的处理 .....                                    | 103        |
| 4.8 采用哪种聚类算法 .....                                   | 104        |
| 4.9 聚类的评估 .....                                      | 104        |
| 4.10 发现数据区域和数据空洞 .....                               | 106        |
| 文献评注 .....   | 108        |
| <b>第5章 部分监督学习 .....</b>                              | <b>110</b> |
| 5.1 从已标注数据和无标注数据中学习 .....                            | 110        |
| 5.1.1 使用朴素贝叶斯分类器的 EM 算法 .....                        | 111        |
| 5.1.2 Co-Training .....                              | 114        |
| 5.1.3 自学习 .....                                      | 115        |
| 5.1.4 直推式支持向量机 .....                                 | 116        |
| 5.1.5 基于图的方法 .....                                   | 117        |
| 5.1.6 讨论 .....                                       | 119        |
| 5.2 从正例和无标注数据中学习 .....                               | 119        |
| 5.2.1 PU 学习的应用 .....                                 | 120        |

## 目录

|                        |     |
|------------------------|-----|
| 5.2.2 理论基础.....        | 121 |
| 5.2.3 建立分类器：两步方法 ..... | 122 |
| 5.2.4 建立分类器：直接方法 ..... | 127 |
| 5.2.5 讨论.....          | 128 |
| 附录：朴素贝叶斯 EM 算法的推导..... | 129 |
| 文献评注.....              | 131 |

## 第二部分 Web 挖掘

|                          |     |
|--------------------------|-----|
| 第 6 章 信息检索与 Web 搜索 ..... | 135 |
|--------------------------|-----|

|                        |     |
|------------------------|-----|
| 6.1 信息检索中的基本概念 .....   | 136 |
| 6.2 信息检索模型 .....       | 138 |
| 6.2.1 布尔模型.....        | 138 |
| 6.2.2 向量空间模型.....      | 139 |
| 6.2.3 统计语言模型.....      | 141 |
| 6.3 关联性反馈 .....        | 142 |
| 6.4 评估标准 .....         | 143 |
| 6.5 文本和网页的预处理 .....    | 147 |
| 6.5.1 停用词移除.....       | 147 |
| 6.5.2 词干提取.....        | 147 |
| 6.5.3 其他文本预处理步骤.....   | 148 |
| 6.5.4 网页预处理步骤.....     | 148 |
| 6.5.5 副本探测.....        | 149 |
| 6.6 倒排索引及其压缩 .....     | 150 |
| 6.6.1 倒排索引.....        | 150 |
| 6.6.2 使用倒排索引搜索.....    | 151 |
| 6.6.3 索引的建立.....       | 152 |
| 6.6.4 索引的压缩.....       | 153 |
| 6.7 隐式语义索引 .....       | 157 |
| 6.7.1 奇异值分解.....       | 158 |
| 6.7.2 查询和检索.....       | 159 |
| 6.7.3 实例.....          | 160 |
| 6.7.4 讨论.....          | 163 |
| 6.8 Web 搜索 .....       | 163 |
| 6.9 元搜索引擎和组合多种排序 ..... | 165 |
| 6.9.1 使用相似度分数的合并.....  | 166 |
| 6.9.2 使用排名位置的合并.....   | 166 |
| 6.10 网络作弊.....         | 168 |
| 6.10.1 内容作弊.....       | 169 |

|                         |            |
|-------------------------|------------|
| 6.10.2 链接作弊             | 169        |
| 6.10.3 隐藏技术             | 170        |
| 6.10.4 抵制作弊             | 171        |
| 文献评注                    | 172        |
| <b>第7章 链接分析</b>         | <b>174</b> |
| 7.1 社会关系网分析             | 175        |
| 7.1.1 中心性               | 175        |
| 7.1.2 权威                | 177        |
| 7.2 同引分析和引文耦合           | 178        |
| 7.2.1 同引分析              | 178        |
| 7.2.2 引文耦合              | 179        |
| 7.3 PageRank            | 179        |
| 7.3.1 PageRank 算法       | 180        |
| 7.3.2 PageRank 算法的优点和缺点 | 185        |
| 7.3.3 Timed PageRank    | 185        |
| 7.4 HITS                | 186        |
| 7.4.1 HITS 算法           | 187        |
| 7.4.2 寻找其他的特征向量         | 189        |
| 7.4.3 同引分析和引文耦合的关系      | 189        |
| 7.4.4 HITS 算法的优点和缺点     | 189        |
| 7.5 社区发现                | 191        |
| 7.5.1 问题定义              | 191        |
| 7.5.2 二分核心社区            | 192        |
| 7.5.3 最大流社区             | 193        |
| 7.5.4 基于中介性的电子邮件社区      | 195        |
| 7.5.5 命名实体的重叠社区         | 196        |
| 文献评注                    | 197        |
| <b>第8章 Web 爬取</b>       | <b>199</b> |
| 8.1 一个简单爬虫算法            | 199        |
| 8.1.1 宽度优先爬虫            | 201        |
| 8.1.2 带偏好的爬虫            | 201        |
| 8.2 实现议题                | 202        |
| 8.2.1 网页获取              | 202        |
| 8.2.2 网页解析              | 202        |
| 8.2.3 删除无用词并提取词干        | 204        |
| 8.2.4 链接提取和规范化          | 204        |
| 8.2.5 爬虫陷阱              | 206        |

## 目录

|                           |            |
|---------------------------|------------|
| 8.2.6 网页库                 | 206        |
| 8.2.7 并发性                 | 207        |
| 8.3 通用爬虫                  | 208        |
| 8.3.1 可扩展性                | 208        |
| 8.3.2 覆盖度、新鲜度和重要度         | 209        |
| 8.4 限定爬虫                  | 210        |
| 8.5 主题爬虫                  | 212        |
| 8.5.1 主题本地性和线索            | 213        |
| 8.5.2 最优优先变种              | 217        |
| 8.5.3 自适应                 | 219        |
| 8.6 评价标准                  | 223        |
| 8.7 爬虫道德和冲突               | 226        |
| 8.8 最新进展                  | 228        |
| 文献评注                      | 230        |
| <b>第9章 结构化数据抽取: 包装器生成</b> | <b>231</b> |
| 9.1 预备知识                  | 231        |
| 9.1.1 两种富含数据的网页           | 232        |
| 9.1.2 数据模型                | 233        |
| 9.1.3 数据实例的 HTML 标记编码     | 235        |
| 9.2 包装器归纳                 | 236        |
| 9.2.1 从一张网页抽取             | 237        |
| 9.2.2 学习抽取规则              | 238        |
| 9.2.3 识别提供信息的样例           | 242        |
| 9.2.4 包装器维护               | 242        |
| 9.3 基于实例的包装器学习            | 243        |
| 9.4 自动包装器生成中的一些问题         | 245        |
| 9.4.1 两个抽取问题              | 246        |
| 9.4.2 作为正则表达式的模式          | 246        |
| 9.5 字符串匹配和树匹配             | 247        |
| 9.5.1 字符串编辑距离             | 247        |
| 9.5.2 树匹配                 | 249        |
| 9.6 多重对齐                  | 252        |
| 9.6.1 中星方法                | 252        |
| 9.6.2 部分树对齐               | 253        |
| 9.7 构建 DOM 树              | 257        |
| 9.8 基于列表页的抽取: 平坦数据记录      | 258        |
| 9.8.1 有关数据记录的两个观察结果       | 258        |
| 9.8.2 挖掘数据区域              | 259        |

|                            |            |
|----------------------------|------------|
| 9.8.3 从数据区域中识别数据记录.....    | 263        |
| 9.8.4 数据项对齐与抽取.....        | 263        |
| 9.8.5 利用视觉信息.....          | 264        |
| 9.8.6 一些其他技术.....          | 264        |
| 9.9 基于列表页的抽取：嵌套数据记录.....   | 265        |
| 9.10 基于多张网页的抽取.....        | 269        |
| 9.10.1 采用前几节中的技术.....      | 270        |
| 9.10.2 RoadRunner 算法.....  | 270        |
| 9.11 一些其他问题.....           | 271        |
| 9.11.1 从其他网页中抽取.....       | 271        |
| 9.11.2 析取还是可选.....         | 272        |
| 9.11.3 一个集合类型还是一个元组类型..... | 273        |
| 9.11.4 标注与整合.....          | 273        |
| 9.11.5 领域相关的抽取.....        | 273        |
| 9.12 讨论.....               | 274        |
| 文献评注.....                  | 274        |
| <b>第 10 章 信息集成 .....</b>   | <b>276</b> |
| 10.1 什么是样式表匹配.....         | 277        |
| 10.2 样式表匹配的预处理工作.....      | 278        |
| 10.3 样式表层次的匹配.....         | 279        |
| 10.3.1 基于语言学的算法.....       | 279        |
| 10.3.2 基于样式表中限制的算法.....    | 280        |
| 10.4 基于领域和实例层次的匹配.....     | 280        |
| 10.5 不同相似度的联合.....         | 282        |
| 10.6 1:m 匹配 .....          | 283        |
| 10.7 其他问题.....             | 284        |
| 10.7.1 重用以前的匹配结果.....      | 284        |
| 10.7.2 大量样式表的匹配.....       | 285        |
| 10.7.3 样式表匹配的结果.....       | 285        |
| 10.7.4 用户交互.....           | 285        |
| 10.8 Web 搜索界面的集成 .....     | 285        |
| 10.8.1 基于聚类的算法.....        | 287        |
| 10.8.2 基于互关系的方法.....       | 289        |
| 10.8.3 基于实例的方法.....        | 290        |
| 10.9 构建一个全局的搜索界面.....      | 292        |
| 10.9.1 结构上的正确性和合并算法.....   | 293        |
| 10.9.2 词汇的正确性.....         | 294        |
| 10.9.3 实例的正确性.....         | 295        |
| 文献评注.....                  | 295        |

## 目录

|                                   |     |
|-----------------------------------|-----|
| 第 11 章 观点挖掘 .....                 | 296 |
| 11.1 意见分类 .....                   | 297 |
| 11.1.1 基于意见短语的分类 .....            | 297 |
| 11.1.2 采用文本分类方法进行意见分类 .....       | 299 |
| 11.1.3 基于评分函数进行分类 .....           | 299 |
| 11.2 基于特征的观点挖掘和摘要 .....           | 300 |
| 11.2.1 问题定义 .....                 | 301 |
| 11.2.2 对象特征抽取 .....               | 305 |
| 11.2.3 格式 1 中正面和负面评价部分的特征抽取 ..... | 306 |
| 11.2.4 符合格式 2 和 3 的评审上的特征抽取 ..... | 308 |
| 11.2.5 观点倾向分类 .....               | 309 |
| 11.3 比较性句子和比较关系挖掘 .....           | 310 |
| 11.3.1 问题定义 .....                 | 311 |
| 11.3.2 等级比较性语句的识别 .....           | 312 |
| 11.3.3 比较关系的抽取 .....              | 314 |
| 11.4 观点搜索 .....                   | 315 |
| 11.5 观点欺诈 .....                   | 316 |
| 11.5.1 观点欺诈的目标和行为 .....           | 317 |
| 11.5.2 欺诈和欺诈者的种类 .....            | 317 |
| 11.5.3 隐藏技巧 .....                 | 318 |
| 11.5.4 欺诈检测 .....                 | 318 |
| 文献评注 .....                        | 320 |
| 第 12 章 Web 使用挖掘 .....             | 322 |
| 12.1 数据收集和预处理 .....               | 323 |
| 12.1.1 数据的来源和类型 .....             | 323 |
| 12.1.2 Web 使用记录数据预处理的关键元素 .....   | 326 |
| 12.2 Web 使用记录挖掘的数据建模 .....        | 331 |
| 12.3 Web 用法模式的发现和分析 .....         | 334 |
| 12.3.1 会话和访问者分析 .....             | 334 |
| 12.3.2 聚类分析和访问者分割 .....           | 334 |
| 12.3.3 关联及相关度分析 .....             | 337 |
| 12.3.4 序列和导航模式分析 .....            | 340 |
| 12.3.5 基于 Web 用户事务的分类和预测 .....    | 342 |
| 12.4 讨论和展望 .....                  | 343 |
| 文献评注 .....                        | 344 |
| 参考文献 .....                        | 345 |

# 第一部分

## 数据挖掘基础