

新编高等院校信息管理与信息系统专业核心教材

网络信息挖掘

Network Information Mining

黄晓斌 编著



电子工业出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

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

新编高等院校信息管理与信息系统专业核心教材

网络信息挖掘

Network Information Mining

黄晓斌 编著

電子工業出版社

Publishing House of Electronics Industry

北京 · BEIJING

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有，侵权必究。

图书在版编目（CIP）数据

网络信息挖掘/黄晓斌编著. —北京：电子工业出版社，2005.1

新编高等院校信息管理与信息系统专业核心教材

ISBN 7-121-00520-4

I . 网… II . 黄… III . 计算机网络—情报检索—高等学校—教材 IV . G354.4

中国版本图书馆 CIP 数据核字（2004）第 123433 号

责任编辑：刘宪兰 何 雄

印 刷：北京牛山世兴印刷厂

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

经 销：各地新华书店

开 本：787×980 1/16 印张：16 字数：410 千字

印 次：2005 年 1 月第 1 次印刷

印 数：5 000 册 定价：24.00 元

凡购买电子工业出版社的图书，如有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系。联系电话：(010) 68279077。质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

新编高等院校信息管理与信息系统

专业核心教材顾问

(按姓氏笔画排序)

马费成 陈禹 黄梯云

新编高等院校信息管理与信息系统

专业核心教材编委会

(按姓氏笔画排序)

马费成 王要武 叶继元

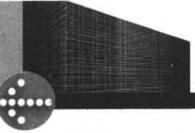
李一军 汪玉凯 陈京民

吴玲达 张维明 张基温

赵国俊 高阳 徐孝凯

黄晓斌 戴宗坤

执行主编：张基温



新编高等院校
信息与管理
信息系统
专业核心教材

顾问（按姓氏笔画排序）

马费成 陈禹 黄梯云

编委会（按姓氏笔画排序）

马费成 王要武 叶继元
李一军 汪玉凯 陈京民
吴玲达 张维明 张基温
赵国俊 高阳 徐孝凯
黄晓斌 戴宗坤

执行主编

张基温

新编高等院校信息管理与信息系统
专业核心教材

书 目

信息网络技术原理

计算机系统原理

数据仓库与数据挖掘技术

信息系统安全导论

管理信息系统

信息检索导论

信息系统工程

多媒体技术

信息资源开发与管理

数据库技术与应用

电子政务

电子商务原理

数据结构

网络信息挖掘

总序

Z O N G X U

20世纪70年代，当强大的信息化巨潮还蕴藏在大洋深处，我们的陆地只有一阵微风吹来之时，有识之士们就开始推动信息化专业人才的培养计划，为迎接即将到来的信息化巨潮扩军备战。他们一方面推动着信息技术的普及；一方面根据不同领域的需求，从不同的角度创办了不同类型的信息化专业，这就是管理信息系统专业、经济信息管理专业、科技信息管理专业、医学信息管理专业、林业信息管理专业、农业信息管理专业……实际上，这些专业培养目标可以概括为：为各行业、各部门培养以CIO为目标的信息化专门人才。从这一点上看，这些专业的课程设置应当具有相当大的共同性。1996年，出于多种考虑，教育部将这些专业合并为一个——信息管理与信息系统专业。

以CIO为目标的信息化专门人才是一类管理人才。但是他们所管理的主要对象是信息。这样的知识需求，将信息管理与信息系统专业定位于管理学科，与信息学、经济学、法学等学科交叉。这样的学科特点，给课程建设和教材建设带来不少困难。近30年来，尽管我们与许多的同行已经进行了不懈的努力，把信息管理与信息系统专业的课程建设和教材建设向前推进了一大步，但是仍然不尽人意，许多课程和教材还没有体现信息管理专业的特色和需要。在多次有关的研讨会上，大家一致呼吁编写一套真正体现信息管理与信息系统专业特色的教材。

新编和出版一套专业教材是要冒风险的，而编写和出版一套以瞬息万变的信息和信息技术为管理对象的专业教材就要冒更大的风险。国内信息业界著名的出版商——电子工业出版社，以超人的胆略愿意同我们一道承担这一风险，组织编写出版一套新的信息管理与信息系统专业核心教材。这套教材冠以“新编”二字，是试图在其体系上能比已有教材更体现信息专业的特色，同时在内容上要能反映最新信息技术的进步以及最新信息管理思想和方法。

目前，国内开设信息管理与信息系统专业的高等院校已经超过200所。这样一个数字一方面表明信息化已经深入人心，信息化队伍的规模正在急速扩大，信息化队伍

的素质正在不断提高；另一方面，也给我们增加了巨大的压力，使我们深感责任重大。好在国内本领域的三位知名学者——黄梯云、陈禹、马费成以及其他一批著名专家和后起之秀愿意与我们共担风险，鼓舞了我们挑起这副重担的勇气。同时，我们也把这套教材的不断精化寄希望于广大的同仁，愿我们把这套教材越改越好，永改永新。

新编高等院校信息管理与信息系统
专业核心教材编委会

前 言

Q I A N Y A N

互联网信息资源越来越庞大，一方面为信息传播开辟了新的途径，另一方面也给用户的利用带来新的挑战。网络信息挖掘是利用现代信息技术与方法，通过构建数据挖掘和知识发现系统，对网络信息进行深入的分析，形成各种有用的知识产品，提供给用户使用，从而实现信息价值的增值。网络信息挖掘是信息管理研究的一个重要内容。通过对网络信息资源的开发、利用和分析研究，可以发现其形成和存在的规律及用户的行为规律，以便合理配置和优化资源，提高网络信息资源的建设和利用水平，从而把信息资源转化为生产力。

本书根据目前网络信息资源管理的发展趋势和完善信息管理专业学生知识结构的需要，对当前网络信息开发和利用技术进行系统的总结和分析，目的是研究网络信息挖掘与利用技术的理论基础，探索有关技术和方法问题，分析在具体应用中的可行性，探讨网络信息挖掘和开发利用的具体途径。本书有如下几个特点：

- 选题新颖，网络信息挖掘是计算机和信息管理领域的前沿课题，书中许多内容都是目前国内外最新的研究成果；
- 内容充实，本书系统分析网络信息挖掘的有关技术原理和规律，重点对网络信息挖掘的方法和技巧等问题进行了探索性研究，并有一定的案例分析；
- 材料比较丰富，书中具有许多实例，并包括各种有关资料和参考书目；
- 适用范围广，本书主要作为信息管理和计算机类本科生和研究生的专业课程教材，也可供广大的信息管理领域科研工作者及其爱好者阅读或参考。

本书按照原理、技术、方法和应用的体例进行编写，共分 9 章：第 1 章介绍了知识发现的基本定义、相关研究的发展和网络信息挖掘与知识发现的意义；第 2 章分析了基于网络的信息挖掘系统的基本要求和功能，知识发现的过程和方法、数据挖掘语言、数据挖掘系统的结构和评价；第 3 章论述了网络信息集成的特点、作用和方式；第 4 章对网络信息的结构挖掘进行了分析；第 5 章对网络信息的内容挖掘进行分析；第 6 章对网络信息的使用记录的挖掘方式方法进行研究；第 7 章讨论了常见的网络信息挖掘的策略；第 8 章介绍了网络信息的挖掘应用的一些实例；第 9 章着重分析网络信息挖掘存在的一些难点和知识发现的研究方向，并对中文网络信息挖掘和知识发现问题提出了一些看法。

本书在编写过程中参阅了国内外一些资料，特向作者表示感谢。电子工业出版社对

本书的出版给予了大力支持，责任编辑刘宪兰和何雄老师对作者给予了极大的鼓励和帮助，并提出了不少意见和建议，其他有关同志也为本书的出版付出了辛勤的劳动，在此表示衷心的感谢。

由于网络信息挖掘是一个新的课题，有许多问题在不断发展变化之中，加上本人水平有限，本书难免会有错漏之处，敬请广大读者批评指正。

黄晓斌
于中山大学康乐园

目 录

第1章 网络信息挖掘概论	1
1.1 知识发现的基本概念	2
1.1.1 知识发现的基本定义	2
1.1.2 知识发现的类型	2
1.1.3 网络信息知识发现的特点	3
1.1.4 相关概念的辨析	4
1.2 相关研究的发展	6
1.2.1 知识发现的研究背景	6
1.2.2 知识发现的产生和发展	7
1.2.3 网络信息知识发现的研究现状	8
1.3 网络信息知识发现的意义	10
1.3.1 实际意义	10
1.3.2 理论意义	13
本章小结	15
思考题 1	15
本章参考文献	16
第2章 网络信息的挖掘系统	17
2.1 网络信息知识发现系统的基本要求	18
2.2 网络信息知识发现的基本功能	19
2.2.1 知识发现的知识类型	19
2.2.2 网络信息知识发现的主要任务	20
2.3 知识发现的过程	21
2.3.1 知识发现的过程模型	21
2.3.2 知识发现的实现过程	21
2.4 数据挖掘的基本方法	23
2.4.1 数据挖掘的基本方法及其特点	23
2.4.2 网络信息知识发现方法的适用性分析	25
2.5 数据挖掘语言	27
2.5.1 数据挖掘语言的意义	27

2.5.2 数据挖掘语言的设计原则	27
2.5.3 数据挖掘语言的类型	27
2.5.4 基于 Web 的挖掘语言	28
2.6 网络信息知识发现系统的结构	31
2.6.1 知识发现系统的一般结构	32
2.6.2 基于网络的知识发现系统	34
2.7 网络信息知识发现系统的评价	36
本章小结	37
思考题 2	37
本章参考文献	38
第 3 章 网络信息的集成	41
3.1 网络信息集成的基本问题	42
3.1.1 网络信息的特点	42
3.1.2 网络信息集成的作用	43
3.2 基于虚拟数据库的网络信息集成	45
3.2.1 虚拟数据库的含义	45
3.2.2 虚拟数据库的特征	45
3.2.3 虚拟数据库的体系结构	46
3.2.4 构建网络信息的虚拟数据库	48
3.3 基于 XML 的网络信息集成	49
3.3.1 XML 的数据集成意义	49
3.3.2 利用 XML 进行异构数据集成	50
3.3.3 XML 文档与数据库的数据交换	51
3.4 基于 Web 数据仓库的网络信息集成	52
3.4.1 Web 数据仓库的特征	52
3.4.2 基于 Web 数据仓库的体系结构	53
3.4.3 多层次的 Web 数据仓库	54
3.5 基于智能代理的网络信息集成	56
3.5.1 智能代理的特点	56
3.5.2 移动智能代理和多智能代理系统	57
3.5.3 基于多智能代理的网络信息集成	58
3.5.4 智能代理网络信息集成的特点	59
本章小结	59
思考题 3	60

本章参考文献	60
第4章 网络信息的结构挖掘	63
4.1 超文本结构的特点	64
4.1.1 超文本的构成	64
4.1.2 超文本链接的基本方式	65
4.1.3 XML 链接与 HTML 链接的比较	66
4.2 Web 结构挖掘	67
4.2.1 Web 结构挖掘的含义	67
4.2.2 相关研究分析	68
4.2.3 Web 链接挖掘研究的意义	69
4.2.4 Web 链接机制分析的局限性	70
4.3 网页排序挖掘法	71
4.3.1 PageRank 算法	71
4.3.2 HITS 算法	72
4.3.3 PageRank 和 HITS 的比较	73
4.4 基于链接挖掘的迷路问题解决方法	74
4.4.1 迷路问题的原因	74
4.4.2 解决迷路问题的主要途径	74
4.5 基于链接挖掘的超文本结构优化	77
4.5.1 网站内链接的结构优化	77
4.5.2 网站外链接的结构优化	78
4.5.3 超文本链接的动态优化	79
本章小结	80
思考题 4	80
本章参考文献	81
第5章 网络信息的内容挖掘	83
5.1 半结构化数据的挖掘	84
5.1.1 半结构化数据的特点	84
5.1.2 半结构化数据模型	85
5.2 基于 HTML 的数据挖掘	87
5.2.1 HTML 的主要特点	87
5.2.2 HTML 网页内容的抽取	88
5.3 基于 XML 的数据挖掘	90
5.3.1 XML 的主要特点	90

5.3.2 XML 在 Web 数据挖掘中的应用	92
5.4 HTML 向 XML 的转换	93
5.4.1 转换的必要性	93
5.4.2 转换的原理	93
5.4.3 转换的方法	94
5.4.4 转换的工具	94
5.5 非结构化数据的挖掘	95
5.5.1 非结构化数据	95
5.5.2 非结构化数据库的特点及其在信息资源数字化中的应用	96
5.5.3 非结构化数据的挖掘	98
5.6 文本挖掘	99
5.6.1 文本挖掘概述	99
5.6.2 文本挖掘的内容和方法	99
5.6.3 文本挖掘的工具	106
5.7 多媒体数据的挖掘	108
5.7.1 多媒体数据挖掘的特点	108
5.7.2 多媒体数据的特征提取	109
5.7.3 多媒体数据挖掘系统的功能模块	111
5.7.4 多媒体数据挖掘的过程	111
5.7.5 多媒体数据的挖掘方式	112
本章小结	114
思考题 5	114
本章参考文献	115
第 6 章 网络信息的使用记录挖掘	117
6.1 使用记录挖掘的特点	118
6.2 使用记录挖掘的作用	118
6.3 使用记录挖掘的方式	120
6.4 使用记录挖掘的数据源	121
6.4.1 Web 服务器日志	121
6.4.2 注册信息	123
6.4.3 曲奇 (Cookie) 数据记录	123
6.5 使用记录挖掘的过程	124
6.5.1 数据预处理阶段	124
6.5.2 模式识别阶段	125

6.5.3 模式的分析	126
6.6 使用记录挖掘的方法.....	126
6.6.1 相关研究分析	126
6.6.2 序列模式挖掘法	127
6.6.3 文本挖掘法	128
6.6.4 概率分布分析法	128
6.6.5 关联规则分析法	129
6.6.6 聚类算法	130
本章小结	132
思考题 6	132
本章参考文献	133
第 7 章 网络信息的挖掘策略	135
7.1 元数据的挖掘	136
7.1.1 元数据的基本问题	136
7.1.2 元数据的挖掘意义	140
7.2 引文数据库的挖掘	141
7.2.1 引文分析的特点	141
7.2.2 引文数据挖掘的数据源	141
7.2.3 引文数据的挖掘策略	145
7.2.4 引文数据挖掘应注意的问题	148
7.3 网络电子出版物的挖掘	149
7.3.1 网络电子出版物的类型	149
7.3.2 网络电子出版物的特点	149
7.3.3 网络电子出版物的挖掘策略	151
7.4 数字图书馆的挖掘	151
7.4.1 数字图书馆的特点	151
7.4.2 数字图书馆知识发现的内容	153
7.4.3 数字图书馆的挖掘的特点	155
本章小结	156
思考题 7	156
本章参考文献	157
第 8 章 网络信息挖掘的应用	159
8.1 网络信息挖掘在电子商务中的应用	160
8.1.1 网络信息挖掘在电子商务中应用的必要性	160

8.1.2 网络信息挖掘在电子商务中的主要应用	162
8.1.3 电子商务中网络信息挖掘的主要方式	164
8.1.4 实例：网络信息挖掘在网上书店的应用[3]	167
8.2 网络信息挖掘在网络广告分析中的应用	170
8.2.1 网络广告的优势	171
8.2.2 网络广告的分类	172
8.2.3 网络广告的发展趋势	174
8.2.4 网络信息挖掘在网络广告中的作用	175
8.2.5 网络广告传播效果的挖掘分析	176
8.3 网络信息挖掘在客户关系管理的应用	177
8.3.1 客户关系管理的含义与意义	177
8.3.2 网络数据挖掘在客户关系管理中的应用	178
8.3.3 客户档案的数据挖掘方法	183
8.4 网络信息挖掘在电子政务中的应用	185
8.4.1 电子政务概述	185
8.4.2 电子政务信息管理与开发利用的意义	186
8.4.3 网络信息挖掘在电子政务中的应用	186
8.4.4 民意信息的挖掘分析	188
8.4.5 政府公共服务信息的挖掘分析	189
8.5 网络信息挖掘在网络信息管理中的应用	189
8.5.1 在电子邮件管理中的应用	189
8.5.2 在 BBS 管理中的应用	191
8.5.3 在搜索引擎中的应用	192
8.5.4 在网络知识检索与管理中的应用	193
8.5.5 在网络入侵检测中的应用	195
8.5.6 在网络个性化服务中的运用	198
8.6 网络信息挖掘在竞争情报工作中的应用	200
8.6.1 网络信息挖掘在竞争情报工作中的作用	200
8.6.2 网络竞争情报信息的特点	201
8.6.3 网络信息挖掘在竞争情报搜集和处理中的应用	202
8.6.4 竞争信息管理系统实现策略	203
8.6.5 竞争情报软件	204
8.6.6 实例：TRS 网络信息雷达系统	206
8.6.7 基于网络专利信息的竞争情报挖掘分析	209

本章小结	211
思考题 8	212
本章参考文献	213
第 9 章 网络信息挖掘的研究方向	215
9.1 网络信息挖掘存在的问题	216
9.2 知识发现的发展趋势	218
9.3 网络信息挖掘的研究方向	221
9.4 关于我国网络信息知识发现问题的思考	224
9.4.1 我国网络信息资源建设的现状	224
9.4.2 对我国网络信息知识发现问题的思考	225
本章小结	226
思考题 9	227
本章参考文献	227
参考文献	228