

网络信息检索技术及 搜索引擎系统开发

高 凯 郭立炜 许云峰 编著



科学出版社
www.sciencep.com

47

网络信息检索技术及 搜索引擎系统开发

高 凯 郭立炜 许云峰 编著

随着网络信息技术的飞速发展，网络信息检索技术已成为人们获取信息的重要手段。搜索引擎系统作为网络信息检索的核心技术，在信息检索领域发挥着越来越重要的作用。本书系统地论述了网络信息检索技术的基本原理、搜索引擎系统开发及智能化技术的研究现状。本书共分4章，第1章介绍网络信息检索技术的基本概念、搜索引擎系统的组成及分类；第2章介绍搜索引擎系统的开发技术，包括搜索引擎系统的架构设计、索引技术、检索技术等；第3章介绍搜索引擎系统的智能化技术，包括自然语言处理、机器学习、数据挖掘等技术；第4章介绍搜索引擎系统的性能优化技术，包括搜索引擎系统的性能测试、性能优化技术等。本书可作为高等院校计算机专业及相关专业的教材，也可供从事搜索引擎系统开发及智能化技术研究的工程技术人员参考。

科学出版社

北京

很好的作用。本书第4章介绍了搜索引擎系统的性能优化技术，包括搜索引擎系统的性能测试、性能优化技术等。本书可作为高等院校计算机专业及相关专业的教材，也可供从事搜索引擎系统开发及智能化技术研究的工程技术人员参考。

内 容 简 介

本书全面、系统地讲述了网络信息检索技术的基本原理，并阐述了其在搜索引擎系统开发及其智能化实现中的应用。在全面介绍了网络信息检索技术、标引与索引、检索结果处理、中英文分词、网络信息获取及预处理之后，本书对信息采集中的网页去重与相似网页聚类、信息的动态采集、基于自然语言理解的检索处理、相关概念反馈、检索纠错、检索结果排序、基于用户浏览历史的网页预取技术等多个方面进行了较深入的研究与分析。

全书体系完整，内容新颖，条理清晰，组织合理，可为高校相关专业学生的学习和科研工作提供帮助，也可为从事搜索引擎技术开发的工程技术人员、希望了解搜索引擎技术的爱好者等提供参考。

图书在版编目(CIP)数据

网络信息检索技术及搜索引擎系统开发/高凯，郭立伟，许云峰编著。

—北京：科学出版社，2010

ISBN 978-7-03-026143-4

I. 网… II. ①高…②郭…③许… III. ①计算机网络—情报检索
②计算机网络—程序设计 IV. G354.4 TP393.09

中国版本图书馆 CIP 数据核字(2009)第 220299 号

责任编辑：赖文华 陈晓萍/责任校对：赵 燕

责任印制：吕春珉/封面设计：耕者设计工作室

科学出版社出版

北京东黄城根北街 16 号

邮政编码：100717

<http://www.sciencep.com>

骏杰印刷厂印刷

科学出版社发行 各地新华书店经销

*

2010 年 2 月第 一 版 开本：B5 (720×1000)

2010 年 2 月第一次印刷 印张：16

印数：1—3 000 字数：323 000

定价：32.00 元

(如有印装质量问题，我社负责调换〈环伟〉)

销售部电话 010-62134988 编辑部电话 010-62135120-8003

版权所有，侵权必究

举报电话：010-64030229；010-64034315；13501151303

前 言

在社会信息化的今天，伴随着因特网的迅速普及和应用，网络已成为人们获取信息的重要渠道。为了帮助人们方便、高效地利用网络信息，搜索引擎应运而生。由于网上信息的海量性、冗余性及用户需求的多样性等，迄今搜索引擎在信息采集与处理、检索、个性化等智能化方面尚不能很好地满足用户的需求。如何借助搜索引擎来帮助人们方便、高效地利用网络信息已成为当前 IT 业的研究热点之一。虽然目前搜索引擎在一定程度上满足了人们的检索需求，但随着信息检索的对象从相对封闭、稳定一致、由独立数据库集中管理的内容扩展到开放、海量、冗余、更新快、分布广泛的 Web 内容，同时由于 Web 内容的特殊性（如网页质量参差不齐、镜像网页的存在等），又由于使用者由原来的专业检索人员扩展到包括务工人员、管理人士、教师、学生等在内的普通大众，人们对网络信息检索技术也就提出了新的、更高的要求。目前，网络信息检索在信息采集与处理、检索处理、个性化处理等智能化方面尚不能很好地满足用户的需求。因此，研究网络信息检索技术，并利用包括开源架构等在内的工具构建自己的搜索引擎系统，进而提高搜索引擎的智能化水平，是一项很有意义的工作。

本书较系统地论述了网络信息检索技术的基本原理，并进一步阐述了其在搜索引擎系统开发及其智能化实现中的应用。本书分为三部分。第一部分是基础知识和相关背景介绍部分，包括从第 1 章到第 3 章的内容。其中，第 1 章概要介绍了信息检索的起源和发展、信息检索模型及方法、网络信息检索的过程、网络信息检索性能评价、网络信息智能化处理、网络信息检索技术的未来发展等问题；第 2 章简要介绍了 Web 信息下载、页面分析与信息抽取方法、基于链接分析的网页相关性算法、检索结果排序、自然语言处理等问题；第 3 章就搜索引擎的发展、分类、功能、资源等进行了介绍。第二部分为利用开源资源实现搜索引擎系统的部分，包括从第 4 章到第 8 章的内容，主要介绍如何利用 Lucene 等开源资源来构建自己的搜索引擎。作为开源项目中的一朵奇葩，Lucene 提供了强大的全文索引和检索功能，并在搜索引擎、桌面检索系统、网站站内搜索、企业级内部文档管理与检索、情报分析系统、知识管理系统、数字图书馆检索系统中发挥了很好的作用。本书第 4 章介绍了 Lucene 的索引与检索机制及其应用、开发平台的搭建与配置等；第 5 章介绍 Lucene 中的中英文分词处理及其效果；第 6 章介绍检索结果排序及处理技术；第 7 章介绍如何利用开源资源来获取网络信息；第 8 章介绍如何对常见格式的网络资源进行解析与预处理。第三部分为搜索引擎智能化的研究与实现部分，包括从第 9 章到第 14 章的内容。其中，第 9 章对信息采集中的网页去重与相关网页聚类进行了研究；第 10 章讨论了信息的动态采集

与更新策略，以期搜索引擎能根据网站及其更新速度的不同，动态调整其信息采集与更新的频度；第 11 章则是面向自然语言提问的理解与处理，提供面向大众的支持自然语言提问的智能检索接口不仅能使人机交互更加人性化，还能促进搜索引擎的应用普及；第 12 章则给出一种参照多数用户在检索类似问题时的经验，为用户提供一些关联性和扩展性的相关概念反馈的方法；第 13 章给出一种相近检索与检索结果排序方法；第 14 章阐述了一种基于用户浏览兴趣的网页预取策略。

全书理论联系实际，涉及面广，体系完整，内容新颖，条理清晰，组织合理，图例丰富，说明详细，既可作为高等院校计算机应用技术专业和图书馆等相关专业的教材，也可作为工程技术人员的参考资料。本书由高凯、郭立炜、许云峰合作编著。编著工作分工如下：高凯组稿并提出写作大纲，郭立炜撰写第 1 章、第 2 章，许云峰撰写第 3 章，高凯完成了其他章节的撰写以及最后的审订和统稿工作。

在本书的写作与相关科研课题的研究工作中，得到了多方面的支持与帮助。自动标引和自动文摘模块分别是采用第一作者所在原 OA 实验室标引课题组和自动文摘课题组提供的 DLL，相关科研课题的研究以及项目开发得到王永成教授和李明禄教授的指导，得到国家高科技研究发展计划（863）项目子课题“教育资讯搜索引擎系统”及上海市信息化专项资金项目“智能中文新闻搜索引擎”的支持，得到上海市软件评测中心、上海市远程教育集团等相关单位对项目研发的支持与协作。相关工作中，课题组宋聚平实现了新闻搜索引擎系统的基本体系架构，李刚实现了新闻分类与摘要的动态显示，龙宇巍对国内搜索引擎信息覆盖面进行了调研并初步确定了新闻搜索引擎网页下载种子集，许欢庆实现了股票检索原型系统架构，陈肯完成了对简单逻辑词的初步分析和获取检索项的拼音码、初步的相关概念反馈处理等工作，课题组肖君、丘振华、杨威、钱兵、刘杰等也给予了大力的协作。宗宝琴、张林利协助完成后续的部分文字校对、参考文献整理工作。另外，中国科学院龙星计划、北京大学搜索引擎与互联网信息挖掘组以及闫宏飞副教授均为部分研究工作提供了帮助，国内外众多的信息检索与搜索引擎智能化方面的研究和相关网站亦为本书提供了良好的基础，本书的顺利完成也得益于参阅了大量的相关工作及研究成果，在此谨向这些文献的作者以及为本书提供帮助的老师、同仁和课题组成员致以诚挚的谢意和崇高的敬意。本书亦得到 2009 年度河北省教育厅科学研究计划（编号：2009435）的资助。在本书写作过程中，也得到了科学出版社赖文华、陈晓萍等的大力支持和帮助，在此一并表示衷心感谢。

由于我们的学识、水平有限，书中不妥之处在所难免，恳请广大读者批评指正。

高凯 郭立炜 许云峰

2010 年 1 月

目 录

3.1	3.1.1 权威教材	40
3.2	3.1.2 国际著名研究机构	41
3.3	3.1.3 著名国际会议	42
3.4	本章小结	42
3.5	参考文献	42
4	第1章 绪论	1
4.1	1.1 引言	1
4.2	1.2 信息检索的起源和发展	4
4.3	1.2.1 手工检索	4
4.4	1.2.2 脱机批处理检索	5
4.5	1.2.3 联机检索	5
4.6	1.2.4 光盘检索	5
4.7	1.2.5 网络信息检索	5
4.8	1.3 信息检索模型及方法	6
4.9	1.3.1 传统布尔检索与扩展布尔检索模型	6
4.10	1.3.2 向量空间模型	9
4.11	1.3.3 概率检索模型	10
4.12	1.3.4 模糊检索模型	10
4.13	1.3.5 逻辑检索模型	10
4.14	1.3.6 概念检索	11
4.15	1.3.7 案例检索	12
4.16	1.4 网络信息检索的过程	12
4.17	1.4.1 网络信息获取	13
4.18	1.4.2 信息加工	13
4.19	1.4.3 信息检索与结果提供	13
4.20	1.5 网络信息检索的性能评价	13
4.21	1.6 网络信息智能化处理	15
4.22	1.7 网络信息检索技术的未来	16
4.23	1.7.1 以智能化技术为核心的智能检索	16
4.24	1.7.2 多媒体信息检索	17
4.25	1.7.3 跨语言检索	17
4.26	1.7.4 个性化检索	18
4.27	本章小结	18
4.28	参考文献	18
5	第2章 网络信息处理	21
5.1	2.1 网络信息采集	21

2.2	网络信息抽取	23
2.3	网络信息的标引与索引	24
2.3.1	标引	25
2.3.2	索引	25
2.4	基于链接分析的网页相关性算法及检索结果排序	26
2.4.1	链接分析	26
2.4.2	HITS 算法	27
2.4.3	PageRank 算法及网页相关性评价	28
2.4.4	HITS 算法和 PageRank 算法的比较	29
2.5	基于自然语言处理的检索	29
2.5.1	自然语言理解的发展	30
2.5.2	基于规则分析的方法	31
2.5.3	基于统计分析的方法	31
2.5.4	自然语言检索	31
	本章小结	32
	参考文献	32
第 3 章	搜索引擎	34
3.1	搜索引擎概述	34
3.2	搜索引擎的发展历程	35
3.3	搜索引擎的分类	36
3.3.1	目录索引式搜索引擎	36
3.3.2	自动式搜索引擎	36
3.3.3	元搜索引擎	37
3.3.4	分布式搜索引擎	37
3.4	搜索引擎开发平台简介	38
3.4.1	Lucene	38
3.4.2	Lemur	38
3.4.3	LIUS	38
3.4.4	Egothor	38
3.4.5	Xapian	39
3.5	开源的 Web 搜索引擎系统简介	39
3.5.1	Nutch	39
3.5.2	YaCy	39
3.5.3	Compass	40
3.6	相关资源	40

3.6.1	权威教材	40
3.6.2	国际著名研究机构	41
3.6.3	著名国际会议	42
	本章小结	42
	参考文献	42
第 4 章	Lucene 的索引与检索机制及其应用	43
4.1	Lucene 简介	43
4.2	Lucene 的下载、安装与部署	44
4.2.1	下载 Lucene	44
4.2.2	配置环境变量	45
4.2.3	对 Lucene Demo 的测试	45
4.3	Lucene 的索引与检索机制概述	48
4.3.1	文本分析	48
4.3.2	Lucene 的索引方式	48
4.3.3	Lucene 索引文件的构成	50
4.3.4	Lucene 的检索	51
4.3.5	Lucene 的索引和检索主要流程	52
4.4	管理和操作索引	53
4.4.1	设定增量索引	53
4.4.2	更新索引	53
4.4.3	优化索引	56
4.4.4	管理索引	56
4.5	Lucene 的检索	57
4.5.1	构建检索	57
4.5.2	完成检索的主要步骤	58
4.6	根据用户提交的检索词构造查询	59
4.6.1	对单一域字段检索	60
4.6.2	对逻辑关系检索	60
4.6.3	对范围的检索	61
4.6.4	对前缀通配的检索	62
4.6.5	对 Query 的前缀和后缀通配的检索	63
4.6.6	模糊检索的实现	63
4.6.7	对多关键词的检索	64
4.6.8	通过 Query 的 SpanNearQuery 方式完成近似检索	65
4.7	基于 Lucene 应用程序: 开源搜索引擎系统 Nutch	66

4.7.1	Nutch 简介	66
4.7.2	在 Eclipse 中加载 Nutch	68
	本章小结	76
	参考文献	76
第 5 章	分词处理	77
5.1	概述	77
5.1.1	基于词典匹配的中文分词	77
5.1.2	基于词频统计的无词典中文分词	78
5.1.3	Lucene 的分析器	78
5.2	常用的中英文分词器及分词效果	79
5.2.1	停用词分析器	79
5.2.2	标准分析器	81
5.2.3	简单分析器	82
5.2.4	空格分析器	83
5.2.5	关键词分析器	84
5.2.6	ChineseAnalyzer	85
5.2.7	CJKAnalyzer	87
5.2.8	第三方分词工具 ICTCLAS	89
5.2.9	第三方分析软件 JE	90
5.2.10	第三方分析软件 IK_CAnalyzer	91
5.2.11	第三方分析软件 MIK_Canalyzer	93
	本章小结	94
	参考文献	94
第 6 章	检索结果排序及处理	95
6.1	检索结果集 Hits	95
6.2	检索结果的排序及控制	96
6.2.1	Lucene 的排序机制	96
6.2.2	通过改变文档的 Boost 因子来改变排序结果	98
6.2.3	使用 Lucene 的 Sort 类定制排序结果	100
6.2.4	对多个指定 Field 进行综合排序	101
6.3	检索结果的分页	102
6.4	检索结果的高亮显示	104
6.5	检索日志处理	107
6.5.1	下载及配置 Log4J	107
6.5.2	配置信息	107

6.5.3	Servlet 启动文件	108
6.5.4	测试	112
	本章小结	113
	参考文献	113
第 7 章	网络信息获取	114
7.1	网络蜘蛛的工作原理	114
7.2	开源网络蜘蛛简介	115
7.2.1	Weblech	115
7.2.2	J-spider	117
7.3	Nutch 网络蜘蛛的工作机制及其使用	118
7.3.1	确定种子集	118
7.3.2	下载网页	119
7.3.3	查阅爬行日志	120
7.3.4	修改配置文件	122
	本章小结	127
	参考文献	127
第 8 章	网络信息预处理	128
8.1	使用 PDFBOX 预处理 PDF 文档	129
8.2	使用 JACOB 预处理 WORD 文档	132
8.3	使用 HTMLParser 预处理 HTML 文档	134
8.4	使用 POI 处理 OFFICE 文档	138
8.4.1	处理 EXCEL 文档	138
8.4.2	处理 WORD 文档	139
8.5	使用 Lucene 处理 SQL Server 数据表	142
	本章小结	148
	参考文献	149
第 9 章	信息采集中的网页去重与相似网页聚类	150
9.1	概述	150
9.2	相关工作	152
9.3	对同源网页的去重	153
9.4	同源网页去重性能评测	155
9.4.1	测试数据集与测试环境	155
9.4.2	同源网页去重算法性能比较与分析	156
9.5	相似网页聚类	157
9.5.1	网页主题概念的自动标引	158

9.5.2	主题概念权值的确定	159
9.5.3	主题概念抽取的主要流程与示例	160
9.5.4	对主题概念标引过程中可能存在的问题的说明	162
9.5.5	网页间相似关系的度量与聚类处理	162
9.6	对内容雷同网页聚类的性能评测	163
9.6.1	应用环境	163
9.6.2	网页聚类示例	164
9.6.3	召回率与聚类准确率统计	168
9.6.4	可能存在的问题及改进计划	170
	本章小结	172
	附录	172
	参考文献	177
第 10 章	信息的动态采集与更新	179
10.1	概述	179
10.2	相关工作	180
10.3	泊松过程	181
10.4	用泊松过程描述更新事件	182
10.5	更新事件到达时间的条件分布	182
10.6	网页动态采集及调整策略	184
10.7	基于相关性的网页动态采集调整	187
10.8	网页动态采集实验结果与分析	190
10.8.1	网页更新事件的分布与统计	190
10.8.2	更新效果分析及对可能存在的问题的说明	192
10.8.3	系统资源利用分析	193
10.8.4	局限性及下一步的工作	194
	本章小结	194
	参考文献	195
第 11 章	面向自然语言提问的理解与处理	196
11.1	概述	196
11.2	相关工作	197
11.3	基于句模分析的自然语言提问处理	199
11.3.1	概述	199
11.3.2	句模	199
11.3.3	核心检索项的抽取	200
11.3.4	概念检索	202

11.4	核心检索项间逻辑关系的识别与处理	203
11.4.1	研究背景	203
11.4.2	对自然语言提问的形式化表示	203
11.4.3	基于产生式规则的归约	204
11.4.4	对二义性问题的处理	206
11.4.5	对语义的处理及其局限性	206
11.5	性能评测	207
11.5.1	对检索数量的定量分析	207
11.5.2	对检索项间逻辑关系处理的分析	209
11.5.3	查全率和查准率统计与分析	211
11.5.4	对尚存问题的说明	214
	本章小结	215
	参考文献	215
第 12 章	相关概念反馈	217
12.1	概述	217
12.2	相关工作	217
12.3	相关概念反馈的实现	218
12.3.1	基于用户检索提问的相关概念获取	218
12.3.2	基于 FPR 算法的相关概念获取	219
	本章小结	223
	参考文献	223
第 13 章	相近检索与检索结果排序	225
13.1	查询纠错与相近检索概述	225
13.2	性能测试与分析	226
13.3	可能存在的问题	229
13.4	有关检索结果排序的相关工作	230
13.5	检索结果排序策略	230
13.6	相关性权值的确定	231
13.7	检索效果示例及对可能存在问题的说明	232
	本章小结	233
	参考文献	233
第 14 章	基于用户浏览兴趣的网页预取	234
14.1	概述	234
14.2	相关工作	235
14.2.1	个性化技术	235

14.2.2	网页预取	236
14.3	基于 Session-tree 的网页预取	237
14.3.1	用户行为分析	237
14.3.2	Session-tree 结构及算法流程	238
14.4	性能分析及对可能存在问题的说明	240
	本章小结	241
	参考文献	242
11.2.1	概述	168
11.2.2	相关工作	170
11.2.3	基于句法分析的自然语言规则处理	172
11.2.4	核心检索项的抽取	174
11.2.5	原型的构造	176
11.2.6	原型的应用	178
11.2.7	本章小结	177
11.2.8	参考文献	179
11.3.1	概述	181
11.3.2	相关工作	181
11.3.3	基于句法分析的自然语言规则处理	181
11.3.4	核心检索项的抽取	181
11.3.5	原型的构造	181
11.3.6	原型的应用	181
11.3.7	本章小结	181
11.3.8	参考文献	181
11.4.1	概述	191
11.4.2	相关工作	191
11.4.3	基于句法分析的自然语言规则处理	191
11.4.4	核心检索项的抽取	191
11.4.5	原型的构造	191
11.4.6	原型的应用	191
11.4.7	本章小结	191
11.4.8	参考文献	191
11.5.1	概述	196
11.5.2	相关工作	196
11.5.3	基于句法分析的自然语言规则处理	196
11.5.4	核心检索项的抽取	196
11.5.5	原型的构造	196
11.5.6	原型的应用	196
11.5.7	本章小结	196
11.5.8	参考文献	196

第1章 绪 论

随着现代网络的飞速发展，中国的互联网普及实现再次飞跃，赶上并超过了全球平均水平。据中国网络信息中心 CNNIC 在 2009 年 1 月发布的统计数据显示，截至 2008 年底，中国网民规模达到 2.98 亿人，较 2007 年增长 41.9%（见图 1-1），因特网普及率达到 22.6%，略高于全球平均水平 21.9%（注：数据来源于 <http://www.internetworldstats.com>，对比的其他国家和地区因特网普及率为 2008 年 6 月底数据）。

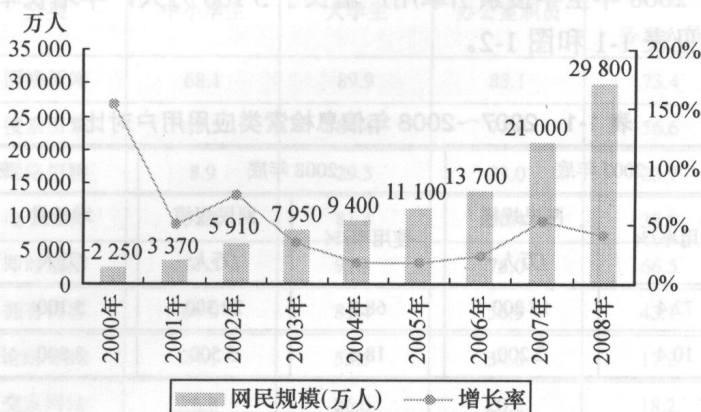


图 1-1 2000~2008 年中国的网民规模及增长率

伴随着网络应用的普及，网络信息也呈爆炸式增长。英国科学家詹姆斯·马丁认为，人类的知识在 19 世纪是每 50 年增加 1 倍，20 世纪中叶是每 10 年增加 1 倍，到 20 世纪 70 年代就已经缩短为每 5 年增加 1 倍^[1]。迄今，信息更如爆炸般产生，而且信息的生产能力已超过了人们对其处理和吸收的能力。正如美国作家奈斯比特在其著作《大趋势》一书中曾指出的那样：“我们虽淹没在信息的海洋中，但是却渴求所需的知识”。为什么会出现这种情况呢？主要原因之一是因为缺乏有效的信息检索与知识获取手段。因此，如何帮助人们快速、方便、准确地从信息海洋中找到所需信息已成为时代发展的迫切需要。

1.1 引 言

信息检索是伴随着人类社会的进步而发展起来的。从广义上来说，信息检索一般指用户为处理解决各种问题而查找、识别、获取相关的事实、数据、文献的

活动及过程，例如，用户在图书馆查询相关信息资料的行为就属于一种信息检索行为。而狭义的信息检索就是指用户在计算机信息检索系统上进行信息查询的行为^[1]。一般认为信息检索是经过了手工检索、脱机批处理检索、联机检索、光盘检索等多个阶段后，逐步发展到今天的网络信息检索。可以说，利用信息检索工具，特别是网络信息检索工具，已经成为现代商业社会和日常生活中不可或缺的一部分。在国外，“Google”已经不再单纯是一个名词，有时它充当的则是动词的角色；在国内，“百度一下”等流行语的出现也说明现在人们越来越依赖于搜索引擎。据 2009 年 1 月我国澳门特别行政区因特网使用现状统计报告公布的数据，在网上娱乐方面，利用因特网搜索引擎寻找信息的网民是最多的，占 78%。而在中国大陆，中国互联网络信息中心 2009 年 1 月发布的统计报告显示，搜索引擎是网民在因特网中获取所需信息的基础应用，其使用率为 68%，在各因特网应用中位列第 4，2008 年全年搜索引擎用户增长了 5 100 万人，年增长率达到 33.6%，具体统计数据见表 1-1 和图 1-2。

表 1-1 2007~2008 年信息检索类应用用户对比

应用	2007 年底		2008 年底		变化	
	使用率/%	网民规模 /万人	使用率/%	网民规模 /万人	增长量 /万人	增长率/%
搜索引擎	72.4	15 200	68.0	20 300	5 100	33.6
网络求职	10.4	2 200	18.6	5 500	3 300	150.0



图 1-2 搜索引擎用户规模

但搜索引擎的使用存在明显的城乡、年龄、学历、收入差异：城镇网民搜索引擎使用率明显高于农村；20~40 岁网民搜索引擎使用率明显高于其他人群；学历越高，搜索引擎使用率越高；收入越高，搜索引擎使用率越高。各网络应用在重点群体中的普及率见表 1-2（注：数据来源于中国互联网络信息中心）。可见，搜索引擎应用人群的特点决定了它在网络领域的高商业价值。虽然由于网民规模快速增长，新增网民中低学历网民比重增大，而该部分网民的搜索引擎的使用率

较低,导致搜索引擎的整体使用率下降,但这并不妨碍搜索引擎正在成为人们获取信息的重要方式。网络信息检索工具——搜索引擎——已经成为人们获取网络信息的重要手段,是载着人们在信息海洋中遨游的快艇。但由于人们对搜索引擎提出了新的、更高的要求^[2-4],又由于网上信息的海量性、冗余性,及用户需求的多样性等,迄今搜索引擎在信息采集与处理、检索、个性化等智能化方面尚不能很好地满足用户的需求。据中国互联网络信息中心 2003 年 7 月至 2005 年发布的几次统计报告显示,用户对搜索引擎性能感到非常满意的只有 23.4%、27.4%、26.9%、28.4% (注:2005 年 7 月后的调查报告中无此项统计数据)。可见,搜索引擎的性能仍有许多需改进之处^[5]。

表 1-2 各网络应用在重点群体中的普及率/%

应用	分类	中小学生	大学生	办公室职员	农村外出务工人员	总体
网络媒体	网络新闻	68.1	89.9	83.1	73.4	78.5
信息检索	搜索引擎	63.5	84.4	71.9	56.6	68.0
	网络招聘	8.9	29.5	23.0	23.7	18.6
网络通信	电子邮件	52.2	81.4	60.4	38.9	56.8
	即时通信	77.5	91.1	75.0	66.5	75.3
网络社区	拥有博客	64.0	81.4	50.9	43.1	54.3
	论坛/BBS	24.1	55.5	34.6	17.2	30.7
	交友网站	16.8	26.0	20.2	18.2	19.3
网络娱乐	网络音乐	86.9	94.0	83.4	78.2	83.7
	网络视频	67.4	84.4	68.1	57.3	67.7
	网络游戏	69.7	64.2	60.6	55.5	62.8
电子商务	网络购物	16.2	38.8	29.4	11.7	24.8
	网上卖物	2.1	5.2	4.4	0.8	3.7
	网上支付	9.6	30.5	22.4	7.9	17.6
	旅行预订	2.0	6.8	6.8	2.5	5.6
其他	网上银行	7.7	29.9	25.5	7.4	19.3
	网络炒股	4.7	4.7	15.5	4.1	11.4
	网上教育	16.2	25.6	17.3	7.8	16.5

近年来,随着 Internet 应用的普及,网络信息检索技术吸引了众多学者的目光,许多原先从事人工智能、数据挖掘、自然语言处理等领域的学者也先后加入到信息检索的研究与开发队伍中来,国际上也有专门针对这个研究领域的国际会

议,如 ACM SIGIR 等。纵观网络信息检索技术的发展,未来的搜索引擎将朝着智能化、多功能化、人机交互等方向发展^[6]。近年来,研究者投向用户的目光愈来愈多了,而搜索引擎的智能化技术则是当前研发的热点。智能化技术应用了相关领域的多项成果,在信息采集与处理、面向自然语言的理解与检索、个性化服务、信息自动分类等方面能够向用户提供更有价值的服务,因而能较好地满足用户的需要。单就如何提高搜索引擎的智能性来说,文献[7]的观点是:“全、快、好、准”应成为搜索引擎系统所追求的主要目标之一。“全”是指搜索引擎应尽量全面地采集所需信息。虽然目前尚无任何一个搜索引擎可以覆盖整个 Internet (据统计,搜索引擎索引网页数据量最大的也不超过整个 Internet 的 16%,所有搜索引擎索引网页数据量的总和大约是整个 Internet 网页量的 42%左右^[8,9]),但要力争采集到尽可能全面的信息;同时,由于目前 Internet 上有大量的重复及转载网页^[10,11],因此在尽量全面采集信息的同时,还要考虑到对重复内容的处理。“快”一般指信息采集快、加工处理快、检索快等。而“好”和“准”一般指便于用户方便地与系统进行交流、系统能提供相关的检索结果等,譬如,是否允许用户用自然语言进行提问、能否较为准确地理解用户的检索需求、能否进行相关概念反馈、能否提供信息自动分类等。

本书较全面地介绍了网络信息检索技术的原理技术、进展,并阐述了其在搜索引擎系统中的应用。全书共分三部分。第一部分(第1章~第3章)对信息检索技术、网络信息智能处理、搜索引擎技术等进行综述,力图使读者对网络信息检索技术和搜索引擎系统有一个全面认识;第二部分(第4章~第8章)介绍如何借助开源的 Lucene 来搭建自己的搜索引擎应用系统,介绍 Lucene 用于索引和检索 API 的使用、基于 Nutch 的搜索引擎系统的构建等;第三部分(第9章~第14章)则针对搜索引擎系统实现中涉及的部分智能化处理技术进行研究与实现,通过对信息采集过程中对内容雷同网页的去重与聚类处理、信息动态采集的研究与实现、面向自然语言提问的理解与检索、相关概念反馈技术、相关检索与查询纠错、检索结果排序、基于用户浏览历史的网页预取等的叙述,较全面地阐述了网络信息智能检索技术的研究与应用。

作为全书的绪论,本章首先从信息检索的起源、发展起步,介绍信息检索的基本原理,并给出网络信息检索的评价指标、研究内容等。

1.2 信息检索的起源和发展

1.2.1 手工检索

顾名思义,手工检索就是指人们以手工为主的方式进行信息检索。在长期的社会和生产实践中,手工检索曾经扮演了重要的角色,特别是在计算机出现以前