

SEARCH ENGINE'S Principle·Practice·Application

SEARCH ENGINE'S

Principle·Practice·Application

搜索与引擎

原理、实践与应用

SEARCH

卢亮 张博文 编著



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

SEARCH ENGINE'S
Principle·Practice·Application

搜索与引擎
原理、实践与应用

卢亮 张博文 编著

電子工業出版社

Publishing House of Electronics Industry
北京·BEIJING

内 容 简 介

本书是搜索引擎业界资深的搜索引擎专家长久以来经验的积累与研究工作的心路历程。本书从搜索引擎的历史和现状开始展开，为广大读者展现了搜索引擎简单界面背后的复杂技术、原理和应用，从数据挖掘到搜索引擎的数据结构，从搜索引擎爬虫到分布式搜索引擎的设计均是作者精心研究的结果和过程，对研究搜索引擎的人士颇有实用价值和参考意义。

本书适合刚从事搜索引擎行业乃至互联网行业的从业人员、如网站设计者、程序员、个人网站的站长，本书还适合作为大中专学院相关专业及培训机构培训的参考书籍。

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

搜索引擎原理、实践与应用 / 卢亮，张博文编著. —北京：电子工业出版社，2007.9

ISBN 978-7-121-04810-4

I. 搜… II. ①卢… ②张… III. 互联网络—情报检索 IV. G354.4

中国版本图书馆 CIP 数据核字（2007）第 119283 号

责任编辑：孙学瑛

印 刷：北京智力达印刷有限公司

装 订：北京中新伟业印刷有限公司

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

开 本：787×980 1/16 印张：19.25 字数：210 千字

印 次：2007 年 9 月第 1 次印刷

印 数：4000 册 定价：42.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，
联系及邮购电话：(010) 88254888。

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。

服务热线：(010) 88258888。

前　　言

搜索引擎目前已经成为我们日常上网必备的工具之一，通过搜索引擎我们可以检索到需要的信息，检索可以说是通向互联网世界的大门。Google，百度等这些优秀的搜索引擎已经是我们耳熟能详的名字了。然而搜索引擎拥有较高的技术门槛，目前世界上也只有为数不多的公司能提供搜索引擎服务，不仅如此，在各行各业中，随着企业数据库规模的不断增涨，对搜索引擎技术的需求也日益旺盛。

为什么写作本书

搜索引擎技术是一门融合各方面知识于一体的技术，简单的搜索框背后隐藏着无限的复杂，作者致力于将搜索引擎技术平民化作为目标，撰写此书的目的就是使徘徊在搜索引擎技术大门外的读者入门。

关于本书作者

本书由卢亮先生与张博文先生合著，作者长久以来合作从事搜索引擎研究工作，并且是业界资深的搜索引擎专家，对搜索引擎产业有相当深入地了解，并且拥有丰富的研发经验。此书的灵魂便是作者长久以来经验的积累与研究工作的心路历程。

本书主要内容

第1章“搜索引擎的历史和现状”介绍了搜索引擎的历史背景，使读者能够对搜索引擎现今的发展有一个明确地了解。

第 2 章“数据挖掘”介绍了数据挖掘技术，数据挖掘技术是搜索引擎技术的基础。

第 3 章“搜索引擎的数据结构”介绍了搜索引擎在存储与索引方法方面的基本原理，并且介绍了检索的一般规则。

第 4 章“搜索引擎的基本结构”是本书重要的一个章节，它介绍了搜索引擎的基本结构与组成部分，使读者能够了解搜索引擎技术的全貌，接下来的章节都是围绕此章作深入地研究。

第 5 章“搜索引擎爬虫”是搜索引擎采集数据的手段，本章介绍了搜索引擎爬虫的两种基本策略，并且给出了一个搜索引擎爬虫实例。

第 6 章“搜索引擎索引系统”介绍了搜索引擎中最核心的部分，即索引技术，此章揭示了搜索引擎为何能瞬间检索到数据的原理。并给出了一个索引系统的实例。

第 7 章“分布式搜索引擎设计”介绍了搜索引擎在分布式环境下的设计策略，现今搜索引擎的技术门槛在于大数据量，大访问量的情况下处理问题的能力，所以本书也着力于研究分布式搜索引擎的解决方案，如果读者已经对搜索引擎的实现原理有所了解，那么分布式搜索方案中的一些观点能够对这部分读者起到良好的进阶作用。

第 8 章“Google 搜索引擎的结构”介绍了世界著名搜索引擎 Google 是如何设计的，Google 是我们学习的榜样，学习它的分布式设计对我们也有相当大的启迪作用。

第 9 章“中文分词”。中文分词我们单独作为一章来讲，因为中文分词是实现中文搜索引擎过程中最为基础最为重要的技术。本章介绍了几种常用的分词策略，并给出了一个分词实例供读者学习。

第 10 章“分类与聚类”介绍了文本自动分类与自动聚类技术，这是数据挖掘技术中重要的两项应用，并且给出了一个分类模型的实验，目前互联网上较大的搜索引擎都采用了此两项技术。

第 11 章“内容消重和 SPAM 消除”阐述了屏蔽互联网上垃圾信息的方法，互联网上有大量的垃圾信息，垃圾信息太多会直接的影响搜索引擎的效率与用户的使用感受，并给出了一个典型的 SPAM 案例。

第 12 章“图像搜索引擎”将会介绍的一门较为新颖的检索技术——图像检索，本章将读者的视野从文字扩展到图像，对读者起到较好的进阶作用。

附录提供了本书中用到的系统化数据。

在一些重要章节之后，我们会介绍一些实例，这些实例能帮助作者更好的理解此章节的内容。比如在第 5 章“搜索引擎爬虫”章节中将给出一个实例 Larbin，在第 6 章“搜索引擎索引系统”章节中给出实例 Lucene(一个高效的索引库)和 Booso.com(基于 Lucene 的博客搜索系统)，在第 9 章“中文分词”环节中给出了一个中文自动分词的实例。

无论是初学者还是长久以来关注搜索引擎技术的读者，本书都会帮助你入门和提高。现在就让我们一起开启搜索引擎技术的大门吧！

由于作者水平有限，书中不足及错误之处在所难免，敬请专家和读者给予批评指正。

卢亮 张博文

2007 年 7 月

目 录

第 1 章 搜索引擎的历史和现状	1
1.1 搜索引擎的历史	1
1.2 搜索引擎的分类	4
1.3 搜索引擎的现状	6
1.4 世界主要搜索引擎介绍	7
1.4.1 Google	7
1.4.2 百度	10
1.4.3 Technorati	11
1.4.4 Alltheweb	12
1.4.5 Ask.com	13
1.4.6 MSN Search	14
小结	15
参考文献	15
第 2 章 数据挖掘	17
2.1 数据挖掘概览	17
2.1.1 为什么要做数据挖掘	19
2.1.2 数据挖掘的任务	20
2.1.3 关联挖掘和分析	21
2.1.4 分类	21
2.1.5 聚类分析	23
2.1.6 序列模式分析	24
2.1.7 偏差分析	24
2.2 数据挖掘的常用技术	25
2.2.1 人工神经网络	25

2.2.2 统计分析	26
小结	27
参考文献	27
第 3 章 搜索引擎的数据结构	29
3.1 存储结构	29
3.1.1 四种基本存储方法	30
3.2 索引	33
3.2.1 倒排索引	33
3.3 结构化查询语言	36
3.4 海量数据系统	38
小结	39
参考文献	39
第 4 章 搜索引擎的基本结构	40
4.1 网络爬虫	41
4.2 排序	43
4.2.1 网页的权重	43
4.2.2 向量空间模型 VSM	44
4.2.3 扩展相关信息	46
4.3 索引系统	49
4.4 缓存机制	55
小结	55
参考文献	56
第 5 章 搜索引擎爬虫	57
5.1 深度优先与广度优先	58
5.1.1 网页链接情况概览	58
5.1.2 取得有效的网页文字	59

5.1.3 深度优先	61
5.1.4 广度优先	63
5.2 实例——Larbin	65
5.2.1 Larbin 简介	65
5.2.2 Larbin 的作用	66
5.2.3 Larbin 的使用	66
5.2.4 Larbin 的性能特征	68
小结	68
参考文献	69
第 6 章 搜索引擎索引系统	70
6.1 索引结构	70
6.2 使用直接 I/O 提高文件扫描性能	75
6.3 倒排表合并策略	78
6.4 利用内存临时索引技术实现即时搜索	81
6.5 对正排文件进行压缩减少磁盘占用	82
6.6 实例——Lucene	83
6.6.1 Lucene 的原理	84
6.6.2 对多个字段进行联合检索	92
6.6.3 对索引进行优化	93
6.6.4 利用 RangeQuery 进行范围查询	94
6.6.5 组合查询	95
6.6.6 多个索引进行联合搜索	96
6.7 实例——Booso.com	98
小结	111
参考文献	111
第 7 章 分布式搜索引擎设计	112
7.1 分布式搜索引擎的核心问题	112

7.2	分布式元搜索引擎	113
7.3	散列式分布搜索引擎	120
7.3.1	散列式分布搜索引擎原理	121
7.3.2	对散列式分布搜索引擎的改进	129
7.3.3	散列式分布搜索引擎的缺陷	131
7.4	索引与文档的分开存放	132
7.4.1	文档服务器的设计	134
7.4.2	文档服务器的分布式处理	135
7.5	对分布式结构建立缓存机制	137
7.6	混合分布式搜索引擎	139
7.7	分布式搜索引擎的扩展性	142
7.7.1	一种实用的节点动态调整方案	143
7.8	P2P 分布搜索引擎	146
7.9	局部遍历型搜索引擎	150
	小结	152
	参考文献	153
第 8 章 Google 搜索引擎的结构		155
8.1	Google 要解决的问题	155
8.2	Google 的分布式设计	157
8.3	Google 文件系统	161
8.4	MapReduce	166
8.5	BigTable	171
8.6	相关搜索的实现	176
	小结	176
	参考文献	177
第 9 章 中文分词		178
9.1	中文信息的特征	178

9.2 搜索引擎的分词	179
9.3 分词的方法	181
9.3.1 正向最大匹配分词	181
9.3.2 逆向最大匹配分词	185
9.4 基于统计的分词方法	187
9.5 其他分词系统	188
9.6 混合分词	191
9.7 对分不出来的词的处理	192
9.8 停止词训练方法	194
9.9 实例——分词程序	195
小结	207
参考文献	208
第 10 章 分类与聚类	209
10.1 分类与聚类介绍	209
10.1.1 自动分类	209
10.1.2 自动聚类	211
10.2 自动分类的原理	211
10.2.1 自动分类技术概览	211
10.2.2 矢量模型	214
10.2.3 在多文档情况下的矢量模型的修正	214
10.2.4 TF*IDF 的修正	215
10.2.5 基于位置的修正	215
10.3 文本信息的噪声模型	216
10.3.1 文本信息模型概览	216
10.3.2 噪声模型	216
10.3.3 噪声的提取	217
10.4 分类的实验	218
10.4.1 实施步骤	218

10.4.2 样本与类向量	219
10.4.3 分类实验的流程	220
10.4.4 分类结果的输出	221
10.5 利用层次聚类法实现文本自动聚类	222
10.5.1 层次聚类技术	222
10.5.2 实现步骤	223
10.5.3 自动聚类实例	224
小结	226
参考文献	226
第 11 章 内容消重和 SPAM 消除	228
11.1 信息指纹	229
11.2 内容消重	230
11.3 Spam 的识别和剔除	232
11.3.1 一种简单实用的识别 SPAM 方法	234
11.3.2 一个最大的 spam 案例	237
小结	238
参考文献	239
第 12 章 图像搜索引擎	241
12.1 简介	241
12.2 图像的收集过程	242
12.2.1 图像的发现	243
12.2.2 图像的文件形式	244
12.3 主题分类和索引	245
12.3.1 文本处理	245
12.3.2 图像的主题分类 (Taxonomy)	246
12.3.3 关键词表和目录名与主题的匹配	247
12.4 搜索、浏览、挖掘	248

12.5 基于图像内容的搜索技术	249
12.5.1 色阶图的相似性	250
12.5.2 自动图像归类	251
12.5.3 相关反馈	251
12.5.4 色阶图的调节	252
12.6 基于图像内容的搜索实例	252
12.6.1 Retrievr 系统	253
12.6.2 图像聚类的 Booso 算法	255
12.6.3 一个图像处理的实例	256
12.6.4 基于图像内容聚类的实现	259
小结	260
参考资料	260
附录	262

第1章

搜索引擎的历史和现状

01

长久以来，人类的发展离不开知识获取与发现的过程，进入互联网时代后，信息量开始爆炸性地增长，大量信息扑面而来，然而人们的接受能力却十分有限，此时人们急需一种技术手段，能够使信息的获取更加简单，准确，并且直接找到人们想要的信息，去除杂乱的干扰信息。在此情况下，搜索引擎诞生了，并且经过数十年的发展，目前已经成为人们日常生活中必不可少的工具。

1.1 搜索引擎的历史

众所周知，最早的互联网站点大部分是建立在大学校园里的。1990年，由于当时并未出现WWW（World Wide Web，万维网），所以FTP（File Transportation Protocol，文件传输协议）软件成为共享文件的主要工具。要共享文件，就必须建立一个FTP服务器。而要检索FTP数据的人就必须使用FTP客户端，这样就造成数据零零散散地分布在各个不同的地方，搜索引擎的祖先——Archie——就应运而生，它是由Montreal的McGill University的学生Alan

Emtage、Peter Deutsch、Bill Wheelan 发明的。人们使用 Archie 时必须输入精确的文件名进行搜索，Archie 才能告诉用户哪一个 FTP 地址可以下载该文件。

很快，第一个互联网爬虫——Mathew Gray 所开发的 World Wide Web Wanderer——出现了，爬虫(Spider)是用于抓取互联网信息的程序。其实，他最初是想测量 Web 的增长速度而开发这个爬虫来计算 Web 上活动的站点的个数的。但很快，他又升级了这个爬虫来抓取实际的 URL。

1993 年 2 月，Excite 由 6 个斯坦福的学生创建，Excite 从 Architext 项目衍生而来。他们想使用静态统计的方法来分析词之间的关系来使搜索更具效率。

1993 年 10 月，Martijn Koster 创建了 ALIWEB，ALIWEB 允许用户提交他们自己的网页简介信息，这意味着 ALIWEB 不需要通过网络来“爬”数据，也不会使用大量的带宽。然而有许多人不知道如何用 ALIWEB 来提交他们的站点。

1993 年 12 月，三个羽翼丰满的反馈式搜索引擎出现了，分别是 JumpStation，World Wide Web Worm 和 RBSE。JumpStation 与 World Wide Web Worm 类似，基本都是收集页面的标题和头部信息，然后使用简单的线性查找来进行检索。然而随着 Web 规模的迅速增长，JumpStation 由于太慢而终止了。而 RBSE 是第一个能够索引 Html 文件正文的搜索引擎。

1994 年 4 月，David Fili 与 Jerry Yang 创办了 Yahoo，Yahoo 分门别类地收集了一些站点供人们查询，这些站点都是手工收录的。

1994 年 4 月 20 日，华盛顿大学的 Brian Pinkerton 创建了

WebCrawler。WebCrawler 是第一个抓取整个网页全文的爬虫程序。很快，它非常流行，最后 AOL 收购了 WebCrawler，并且让它运行在自己的网络上。1997 年，Excite 购买了 WebCrawler，并且 AOL 开始使用 Excite 来加强自己的 NetFind。WebCrawler 为随后的许多应用打开了大门，在它登场的一年后，Lycos、Infoseek 和 OpenText 便出现了。

Lycos 的出现是搜索引擎历史上重要的一步，它由 Carnegie Mellon 大学于 1994 年 7 月设计。Michale Mauldin 负责搜索引擎的工作，直到现在，他还担任 Lycos 公司的首席科学家。1994 年 7 月 20 日，Lycos 发布了条目数数量为 54 000 个文档的搜索引擎，Lycos 还提供了前缀匹配和字符相近限制，到 1994 年 8 月，Lycos 已经能够搜索 394 000 个文档，到 1995 年，已经有一百五十万的文档，而 1996 年，Lycos 就已经索引了六千万的文档。

Infoseek 是另外一个重要的搜索引擎，起初它与 Lycos、Yahoo 没有多大区别，但其随后增加了很多附加功能，而后又与 Netscape 合作，即为 Netscape 提供搜索服务，当用户点击 Netscape 上的搜索按钮时，直接利用 Infoseek 进行搜索，而以前一直是由 Yahoo 提供该服务的。Infoseek 拥有先进的搜索技术，并且它能够允许用户在 24 小时内任意添加或删除自己的 URL。

1998 年，在今天仍产生巨大影响的 Google 公司诞生了。Google 使用了先进 PageRank 技术。Google 随后开始非常流行，并开始为 AOL 和 Yahoo 提供搜索服务。MSN Search 也在同年诞生。

2000 年，Teoma 搜索引擎诞生，它能够将网站按照其主题的流行程度进行分类，在 2001 年，Ask Jeeves 购买了 Teoma。

2003 年，Google 发布了在线广告 AdSense，让人们把相关的

广告放置在相关页面的任何地方。随后，Google 又陆续推出了更多的服务，详情可见 1.4 节的介绍。【1, 2】

知名搜索引擎出现的时间表如表 1.1 所示，大家可对搜索引擎的发展历史一目了然。

表 1-1 知名搜索引擎出现的时间表

年 代	名 称
1990 年	Archie
1993 年 2 月	Excite
1993 年 10 月	ALIWEB
1993 年 12 月	JumpStation, World Wide Web Worm 和 RBSE
1994 年, 4 月	Yahoo
1994 年 7 月	Lycos
1998 年	Google
2000 年	Teoma
2001 年	Teoma 被 Ask Jeeves 收购
2003 年	Google 发布 AdSense

1.2 搜索引擎的分类

搜索引擎按其工作方式主要可分为三种：全文搜索引擎、目录索引类搜索引擎和元搜索引擎。

全文搜索引擎是目前流行的大型搜索引擎普遍采用的形式。所谓全文，就是用户可以去搜索一篇文章的任何部分，不论是标题，还是正文，用户得到了更大的自由度。Google（其界面如图 1-1 所示）就是优秀的全文检索引擎。

目录式搜索引擎并不是严格意义上的搜索引擎，它只是将网站