

**Broadview®**  
www.broadview.com.cn

- 首次揭示商业级中文搜索实现秘密
- 业内知名开发团队倾情奉献
- 引领Lucene开发技术升级

# 自己动手写搜索引擎

罗刚 编著



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

# 自己动手写搜索引擎

罗刚 编著

電子工業出版社

Publishing House of Electronics Industry

北京•BEIJING



## 内 容 简 介

本书是猎兔企业搜索开发团队的软件产品研发和项目实践的经验汇总。本书全方位展现出一个商用级别的 Lucene 搜索解决方案，主要包括爬虫、自然语言处理和搜索实现部分。

爬虫部分介绍了网页遍历方法和从网页提取主要内容的方法。

自然语言处理部分包括了中文分词从理论到实现以及在搜索引擎中的实用等细节。

其他自然语言处理的经典问题与实现包括：文档排重、文本分类、自动聚类、语法解析树、拼写检查、拼音转换等理论与实现方法。

在实现搜索方面，本书用简单的例子介绍了完整的搜索实现过程，覆盖了从索引库的设计和索引库与数据库的同步到搜索用户界面设计与实现。搜索用户界面包括实现布尔逻辑查询、按区间范围查询、搜索结果按日期排序等。本书还进一步介绍了搜索排序的优化方法。

最后以基于 Lucene 的搜索服务器 Solr 为例，展示了 Lucene 的最新应用方法。

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

自己动手写搜索引擎 / 罗刚编著. —北京：电子工业出版社，2009.12

ISBN 978-7-121-09640-2

I. 自… II. 罗… III. 互联网络—情报检索 IV. G354.4

中国版本图书馆 CIP 数据核字（2009）第 178728 号

责任编辑：江 立

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：23 字数：390 千字

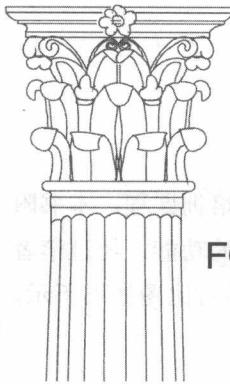
印 次：2009 年 12 月第 1 次印刷

印 数：4000 册 定价：55.00 元（含 DVD 光盘 1 张）

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](mailto:zlts@phei.com.cn)，盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

服务热线：(010) 88258888。



## FOREWORD

# 前　　言

**15** 在中国，随着互联网从城市到农村的普及，搜索引擎对日常生活产生越来越大的影响。例如，笔者自己一般每天就有 15 个左右的问题需要求助于搜索引擎。从 04 年开始笔者也从数据库相关软件开发转入搜索引擎相关开发工作。

**Google** 20 世纪末，在美国国家科学基金会的支持下，斯坦福大学的两个学生在他们的教授指导下开始了一个数字图书馆项目。后来，他们创建了 Google 公司，开创了通过互联网搜索技术共享人类信息的新纪元。Google 通过网络广告取得了巨大的商业回报，到现在仍然是世界 500 强企业中赢利能力最强的公司之一。NASDAQ 证券交易市场的最高股价是 Google 公司的股票。搜索引擎开发成为一项极有含金量的技术。

**Web** 开始写作《自己动手写搜索引擎》这本书以前，已经有一些介绍搜索理论或者搜索开发工具的图书，但是往往表现出来的是纯粹的理论推导和公式定理，或者仅仅是现成开源软件的介绍、分析和使用，并没有介绍其理论依据。有的读者是数学专业的博士，对于相关的数学模型一看就明白，但对于算法实现可能仍然缺少经验。有的读者是培训学校毕业的学生，可能对 Web 开发框架和软件工具的使用很熟悉，但缺少理论基础和深入创新的能力。本书的一个特点在于前面是原理介绍，接着是具体的代码实现。不仅讲解抽象的知识，更重要的是把知识转化成具体软件应用的过程也展示出来。

**Lucene** 《自己动手写搜索引擎》是猎兔企业搜索开发团队的软件产品研发和项目实践的经验汇总。感谢 Lucene，它把搜索引擎开发工作变成了广大程序员都能够参与的游戏。所以本书选用 Lucene 来全方位展现一个商用级别的搜索解决方案。中文分词当前仍然是实现中文搜索的热门话题之一。本书重点介绍了中文分词的相关理论和代码实现，以及在搜索引擎中实用中文分词等细节。本书用简单的例子介绍了搜索引擎完整的实现过程，同时也没有忽略一些经典的算法实现。

该书适合需要具体实现搜索引擎的程序员使用，对于自然语言处理等相关研究人员也

有一定参考价值，同时猎兔搜索团队也已经开发出以本书为基础的专门培训课程。本书附带光盘中的代码经过了详细的注释。为了帮助初学者更容易地了解程序的功能，经过笔者的精心整理后，每个主要变量和每行主要的执行程序都加上了注释，前后对比图如下所示。

```
public int wordSegment(String Sentence)
{
    int senLen = Sentence.length();
    int i = 0, j = 0;
    int M = 12;
    String word;
    boolean bFind = false;
    while (i < senLen)
    {
        int N = i + M < senLen ? i + M : senLen;
        bFind = false;
        for (j = N; j > i; j--)
        {
            word = Sentence.substring(i, j);
            if (dic.Find(word))
            {
                System.out.print(word + " ");
                bFind = true;
                i = j;
                break;
            }
        }
        if (bFind == false)
        {
            word = Sentence.substring(i, i + 1);
            System.out.print(word + " ");
            ++i;
        }
    }
    System.out.println();
    return 1;
}
```

```
public int wordSegment(String Sentence)// 传入一个字符串作为要处理的对象。
{
    int senLen = Sentence.length(); // 首先计算出传入的这句话的字符长度
    int i = 0, j = 0; // i是用来控制截取的起始位置的变量，j是用来控制截取的结束位置的变量
    int M = 12; // 所取得的词组的最大值不超过12。
    String word; // 我们需要与词库中做对比的词。
    boolean bFind = false; // 用来判断是否是词库中的词。
    while (i < senLen) // 如果i的大小于此句话的长度就进入循环，如果i的大小等于了这句话的长度也就证明了匹配过程已经结束了
    {
        int N = i + M < senLen ? i + M : senLen; // 如果i+M<senLen为真就把i+M的值赋给N，如果其为假就把senLen的结果赋值给N
        // 我们首先做的事情就是从一句话中取出一个处理范围
        // 以上这句话就是为了界定我们所要处理的字或者是词组的大小。
        bFind = false; // 假设这个我们取出来的词组不是词库中的词
        for (j = N; j > i; j--) // 向正最大匹配
        {
            word = Sentence.substring(i, j); // 截取我们所需要的字符串。
            // 如果没有匹陪到合适的词的话就会一直在的在这里循环直到出循环。
            if (dic.Find(word))
            {
                System.out.print(word + " "); // 如果这个词是词库中的那么就打印出来
                bFind = true;
                i = j; // 如果你在词库中找到了这个词那么就截取的末尾的位置j赋给i
                // System.out.println(i);
                break; // 跳出for循环
            }
        }
        if (bFind == false) // 如果在我们所处理的范围内一直都没有找到匹配上的词，就一个字一个字的打印出来。
        {
            word = Sentence.substring(i, i + 1);
            System.out.print(word + " ");
            ++i; // 这句话其实也就是控制读的顺序
        }
    }
    System.out.println();
    return 1;
}
```

**3年** 笔者花了3年时间编写本书，但限于时间和水平，其中很多内容笔者仍然正在改写中，希望读者不吝赐教，使本书在重印时内容能更全面。如果条件允许，笔者也很愿意探索一些非比寻常的方式来介绍技术。

**API** 笔者把写书的过程看成是笔者自身和作者所组建的技术团队修炼和提高的过程。

笔者相信技术高手并非天生，而是来源于对卓越品质的追求和自我的不断提升。起点低不是问题，笔者所在的搜索开发团队就包括一些高中学历的程序员，不愿意学习和深入钻研问题才是最大的问题。从企业知识管理的角度，这本书也是有益的。团队里每次有新人进来，我们都先把相关文档发给他们，这样内部交流技术就多了一些基础，不用每次都从Lucene最基本的API等这些基本知识开始介绍。

**搜索** 从学习的过程来看，有人会问：“你没有去过\*\*公司，怎么知道开发搜索引擎？”

笔者没有去过大搜索公司，所以只能从公开的论文、资料和代码等开始学习。感谢搜索引擎技术本身的发展，这方面的学习材料越来越容易通过搜索获得。笔者没有在课堂上学习过自然语言处理等方面的课程，刚开始做中文分词的时候，吴春尧博士问：“你知不知道语言模型？”笔者当时对这些概念一无所知。但是通过了解中文分词程序每一行代码的作用，阅读相关的论文，终于了解了其中一些常用的数学方法。

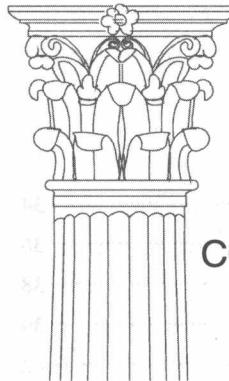
**Solr** 从开发过程来看，首先由客户或者其他开发人员提出问题，解决问题，重新审查整个过程并寻找更好的解决方法。比如实现搜索结果中分类统计，最开始用QueryFilter方法可以实现基本的功能，后来Solr的作者实现了一个更快的DocSetHitCollector方法，如果用查表法来改进这个DocSetHitCollector方法可以达到比Solr官方版本更快的计算速度。搜索软件开发的过程借用日本产品研发大师稻盛和夫的话与大家共勉——“对产品的现场，重新进行审视、体察、贴心、倾听。这样的话，就可以听见神灵的声音。现场产品传来喃喃细语告知我们解决问题的诀窍——这样试试如何？”

**准确** 从成书过程来看，这本书也找过一些内部的开发人员请他们挑问题使得本书的表达更准确。最后，电子工业出版社的胡辛征老师对本书提出了很多具体而必要的建议，

江立老师反复做了细致的修订工作，他们丰富的经验和认真负责的态度使得本书从粗糙的讲稿变成了书本的形式。

## 感谢

最后笔者衷心感谢家人、关心我的老师和朋友们、创业伙伴、搜索培训班的学员以及选择猎兔搜索软件的客户多年来的支持。第一次开办搜索培训班时，由于讲稿本身准备得不够充分，学员有诸多的反馈意见。经过不断的修改后，新学员的负面反馈越来越少。他们的意见也在时刻提醒笔者要保持谦虚的态度。猎兔搜索的客户对于完善搜索技术提供了重要的支持，客户公司包括：天下互联公司、北京电视台、港澳资讯公司、三星鹏泰公司、中国自动化网、中国工控网、中国机电贸易网、易车网、敦煌网、G宝盘、消息树网站等。



## CONTENTS

## 目 录

<b>第 1 章 遍历搜索引擎技术 .....</b>	<b>1</b>
1.1 30 分钟实现的搜索引擎 .....	1
1.1.1 准备工作环境（10 分钟） .....	1
1.1.2 编写代码（15 分钟） .....	3
1.1.3 发布运行（5 分钟） .....	5
1.2 Google 神话 .....	9
1.3 体验搜索引擎 .....	9
1.4 搜索语法 .....	10
1.5 你也可以做搜索引擎 .....	13
1.6 搜索引擎基本技术 .....	14
1.6.1 网络蜘蛛 .....	14
1.6.2 全文索引结构 .....	14
1.6.3 Lucene 全文检索引擎 .....	15
1.6.4 Nutch 网络搜索软件 .....	16
1.6.5 用户界面 .....	17
1.7 商业搜索引擎技术介绍 .....	19
1.7.1 通用搜索 .....	19
1.7.2 垂直搜索 .....	20
1.7.3 站内搜索 .....	21
1.7.4 桌面搜索 .....	23
1.8 本章小结 .....	24
<b>第 2 章 获得海量数据 .....</b>	<b>25</b>
2.1 自己的网络蜘蛛 .....	25
2.1.1 抓取网页 .....	25
2.1.2 网络蜘蛛遍历与实现 .....	26
2.1.3 改进网络蜘蛛 .....	30



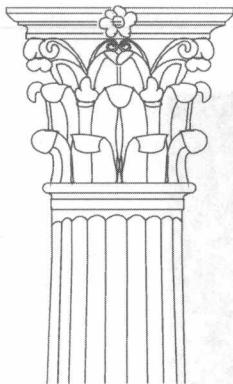
2.1.4 MP3 抓取 .....	34
2.1.5 RSS 抓取 .....	36
2.1.6 图片抓取 .....	38
2.1.7 垂直行业抓取 .....	39
2.2 抓取数据库中的内容 .....	42
2.2.1 建立数据视图 .....	42
2.2.2 JDBC 数据库连接 .....	43
2.2.3 增量抓取 .....	45
2.3 抓取本地硬盘上的文件 .....	47
2.4 本章小结 .....	49
<b>第3章 提取文档中的文本内容 .....</b>	<b>50</b>
3.1 从 HTML 文件中提取文本 .....	50
3.1.1 HtmlParser 介绍 .....	53
3.1.2 结构化信息提取 .....	63
3.1.3 查看网页的 DOM 结构 .....	68
3.1.4 正文提取的工具 NekoHTML .....	71
3.1.5 网页去噪 .....	73
3.1.6 网页结构相似度计算 .....	76
3.1.7 网站风格树去除文档噪声 .....	80
3.1.8 正文提取 .....	92
3.2 从非 HTML 文件中提取文本 .....	98
3.2.1 TEXT 文件 .....	98
3.2.2 PDF 文件 .....	98
3.2.3 Word 文件 .....	105
3.2.4 RTF 文件 .....	106
3.2.5 Excel 文件 .....	107
3.2.6 PowerPoint 文件 .....	108
3.3 流媒体内容提取 .....	109
3.3.1 音频流内容提取 .....	109
3.3.2 视频流内容提取 .....	111
3.4 抓取限制应对方法 .....	113
3.5 本章小结 .....	114

第4章 中文分词	115
4.1 Lucene 中的中文分词	115
4.2 Lietu 中文分词的使用	116
4.3 中文分词的原理	117
4.4 查找词典算法	118
4.5 最大概率分词方法	123
4.6 新词发现	127
4.7 词性标注	129
4.8 本章小结	139
第5章 自然语言处理	140
5.1 语法解析树	140
5.2 文档排重	141
5.3 中文关键词提取	142
5.3.1 关键词提取的基本方法	142
5.3.2 从网页中提取关键词	145
5.4 相关搜索	145
5.5 拼写检查	148
5.5.1 英文拼写检查	148
5.5.2 中文拼写检查	149
5.6 自动摘要	153
5.6.1 自动摘要技术	153
5.6.2 自动摘要的设计	154
5.6.3 Lucene 中的动态摘要	162
5.7 自动分类	163
5.7.1 Classifier4J	164
5.7.2 自动分类的接口定义	165
5.7.3 自动分类的 SVM 方法实现	166
5.7.4 多级分类	167
5.8 自动聚类	170
5.8.1 聚类的定义	170
5.8.2 K 均值聚类方法	170

5.8.3 K 均值实现	173
5.9 拼音转换	179
5.10 语义搜索	180
5.11 跨语言搜索	186
5.12 本章小结	188
<b>第 6 章 创建索引库</b>	<b>189</b>
6.1 设计索引库结构	190
6.1.1 理解 Lucene 的索引库结构	190
6.1.2 设计一个简单的索引库	192
6.2 创建和维护索引库	193
6.2.1 创建索引库	193
6.2.2 向索引库中添加索引文档	194
6.2.3 删除索引库中的索引文档	196
6.2.4 更新索引库中的索引文档	197
6.2.5 索引的合并	197
6.2.6 索引的定时更新	197
6.2.7 索引的备份和恢复	198
6.2.8 修复索引	199
6.3 读写并发控制	200
6.4 优化使用 Lucene	200
6.4.1 索引优化	201
6.4.2 查询优化	202
6.4.3 实现时间加权排序	206
6.4.4 实现字词混合索引	207
6.4.5 定制 Similarity	214
6.4.6 定制 Tokenizer	215
6.5 查询大容量索引	217
6.6 本章小结	218
<b>第 7 章 用户界面设计与实现</b>	<b>219</b>
7.1 Lucene 搜索接口（search 代码）	219
7.2 搜索页面设计	221

7.2.1	用于显示搜索结果的 taglib	221
7.2.2	用于搜索结果分页的 taglib	223
7.2.3	设计一个简单的搜索页面	225
7.3	实现搜索接口	227
7.3.1	布尔搜索	227
7.3.2	指定范围搜索	228
7.3.3	搜索结果排序	233
7.3.4	搜索页面的索引缓存与更新	234
7.4	实现关键词高亮显示	236
7.5	实现分类统计视图	239
7.6	实现相似文档搜索	244
7.7	实现 AJAX 自动完成	246
7.7.1	总体结构	247
7.7.2	服务器端处理	247
7.7.3	浏览器端处理	249
7.7.4	服务器端改进	250
7.7.5	部署总结	261
7.8	jQuery 实现的自动完成	262
7.9	集成其他功能	267
7.9.1	拼写检查	267
7.9.2	分类统计	267
7.9.3	相关搜索	271
7.9.4	再次查找	274
7.9.5	搜索日志	275
7.10	搜索日志分析	276
7.11	本章小结	280
第 8 章	其他高级主题	281
8.1	使用 Solr 实现分布式搜索	281
8.1.1	Solr 服务器端的配置与中文支持	282
8.1.2	把数据放进 Solr	287
8.1.3	删除数据	289
8.1.4	客户端搜索界面	290

8.1.5	Solr 索引库的查找.....	292
8.1.6	索引分发.....	294
8.1.7	Solr 搜索优化 .....	298
8.1.8	Solr 中字词混合索引.....	302
8.1.9	相关检索.....	304
8.1.10	搜索结果去重.....	307
8.1.11	分布式搜索 .....	311
8.1.12	SolrJ 查询分析器.....	315
8.1.13	扩展 SolrJ.....	325
8.1.14	扩展 Solr .....	327
8.1.15	Solr 的.NET 客户端.....	333
8.1.16	Solr 的 PHP 客户端.....	334
8.2	图像的 OCR 识别.....	336
8.3	竞价排名 .....	343
8.4	Web 图分析 .....	344
8.5	使用并行程序分析数据.....	350
8.6	RSS 搜索 .....	351
8.7	本章小结 .....	353
	参考资料 .....	354



# 第1章 遍历搜索引擎技术

搜索引擎是我们每天上网经常使用的功能，本书介绍的搜索技术需要 Java 编程语言基础。本章从快速实现基本的搜索入手，然后深入展开分析搜索的基本技术。



## 1.1 30分钟实现的搜索引擎

首先从一个简单的搜索引擎入手，实现一个简单的指定目录文件的搜索引擎。实现之前需要读者具有 Java 开发方面的基础知识。

### 1.1.1 准备工作环境（10分钟）

首先要准备一个 Java 的开发环境。当前可以使用 JDK 1.6。JDK 1.6 可以从 Sun 的官方网站 [java.sun.com](http://java.sun.com) 下载得到。使用默认方式安装即可。

然后要使用的是一个用来管理搜索引擎索引库的 jar 包，叫做 Lucene。目前可以从 <http://lucene.apache.org/java/docs/index.html> 下载到最新的 Lucene，当前的版本是 2.3。另外，使用集成开发环境 Eclipse，其开发界面如图 1-1 所示。

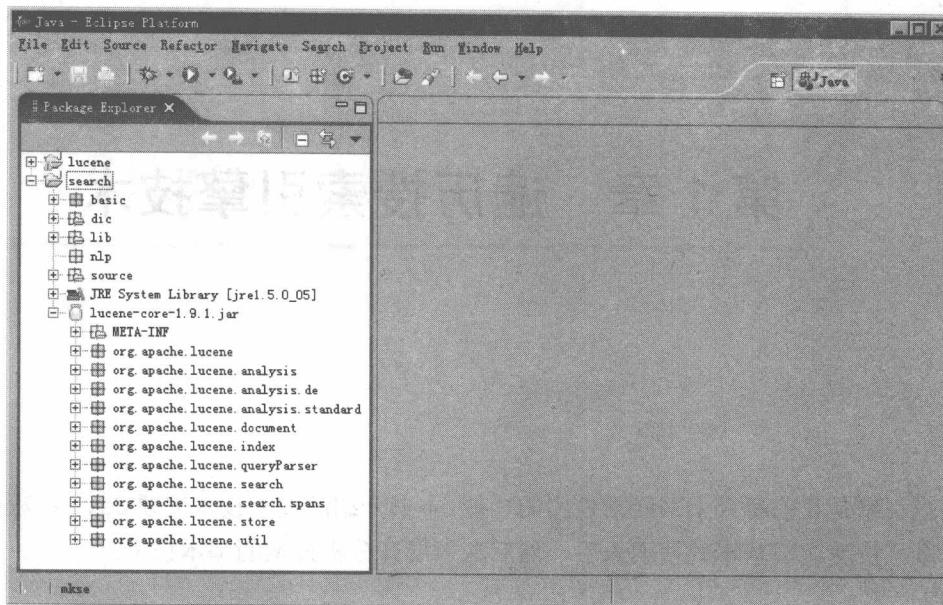


图 1-1 Java 开发界面

如果需要用 Web 界面搜索，还要下载 Tomcat，当前可以从 <http://tomcat.apache.org/> 下载到，推荐使用 Tomcat 5.5 以上的版本。

然后增加 Tomcat 的内存使用量，防止内存溢出。

如果是在 Linux 下，可以用 vi 编辑器修改./catalina.sh 文件，增加行：

```
JAVA_OPTS=-Xmx600m
```

如果是在 Windows 下，可以用文本编辑器“记事本”修改./catalina.bat 文件，增加行：

```
set JAVA_OPTS=-Xmx600m
```

**【作者提示】**如果不增加内存使用量，运行时可能会出现 `java.lang.OutOfMemoryError` 异常。

### 1.1.2 编写代码（15分钟）

搜索引擎的基础在于对全文索引库的管理，在 Lucene 中，通过 **IndexWriter** 来写入索引库。伪代码如下：

1. 创建 **IndexWriter**，准备写索引；
2. 遍历要索引的路径；
3. 优化索引。

下面是主要的实现代码：

```
public void go() throws Exception {
    long start = System.currentTimeMillis();
    if (verbose) {
        System.out.println("Creating index in: " + indexDir);
        //创建索引目录或者建立增量索引
        if (incremental) System.out.println("- using incremental mode");
    }
    Index = new IndexWriter(new File(indexDir), new StandardAnalyzer(),
                           !incremental); //打开或创建索引库, indexDir 是索引存放的路径

    File dir = new File(sSourceDir); //待索引的文件存放的路径
    indexDir(dir); //索引路径
    index.optimize(); //索引优化
    index.close(); //关闭索引库
    if (verbose)
        System.out.println("index complete in :" + (System.currentTimeMillis()
                           - start)/1000);
}
```

下面这段代码把文件内容加到索引库：

```
private void indexFile(File item) {
```



```
if (verbose) System.out.println("Adding FILE: " + item);
News news = loadFile(item); //把文件中的内容加载到news对象
if (news != null && news.body != null) {
    Document doc = new Document();
    //创建网址列
    Field f = new Field("url", news.URL,
        Field.Store.YES, Field.Index.UN_TOKENIZED,
        Field.TermVector.NO);
    doc.add(f);
    //创建标题列
    f = new Field("title", news.title,
        Field.Store.YES, Field.Index.TOKENIZED,
        Field.TermVector.WITH_POSITIONS_OFFSETS);
    doc.add(f);
    //创建内容列
    f = new Field("body", news.body.toString(),
        Field.Store.YES, Field.Index.TOKENIZED,
        Field.TermVector.WITH_POSITIONS_OFFSETS);
    doc.add(f);
    try{
        //文档增加到索引库
        index.addDocument(doc);
    }
    catch(Exception e)
    {
        e.printStackTrace();
        System.exit(-1);
    }
}
}
```

完整的代码可以在本书附带的光盘中找到。

运行以后，一般会生成以下三个索引文件：