

<<自己动手写搜索引擎>>

图书基本信息

书名：<<自己动手写搜索引擎>>

13位ISBN编号：9787121096402

10位ISBN编号：7121096404

出版时间：2009-11

出版时间：电子工业出版社

作者：罗刚

页数：353

版权说明：本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问：<http://www.tushu007.com>

## <<自己动手写搜索引擎>>

### 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。

Google20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。

## <<自己动手写搜索引擎>>

### 内容概要

《自己动手写搜索引擎》是猎兔企业搜索开发团队的产品研发和项目实践的经验汇总。

《自己动手写搜索引擎》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，《自己动手写搜索引擎》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。

《自己动手写搜索引擎》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。

## <<自己动手写搜索引擎>>

### 作者简介

罗刚，猎兔搜索（<http://www.lietu.com>）创始人，当前猎兔搜索在北京和上海均设有研发部。带领猎兔搜索技术开发团队先后开发出猎兔中文分词系统、猎兔智能垂直搜索系统以及网络信息监测系统，实现互联网信息的采集、过滤、搜索和实时监测。

## &lt;&lt;自己动手写搜索引擎&gt;&gt;

## 书籍目录

第1章 遍历搜索引擎技术/11.1 30分钟实现的搜索引擎/11.1.1 准备工作环境(10分钟)/11.1.2 编写代码(15分钟)/31.1.3 发布运行(5分钟)/51.2 Google神话/91.3 体验搜索引擎/91.4 搜索语法/101.5 你也可以做搜索引擎/131.6 搜索引擎基本技术/141.6.1 网络蜘蛛/141.6.2 全文索引结构/141.6.3 Lucene全文检索引擎/151.6.4 Nutch网络搜索软件/161.6.5 用户界面/171.7 商业搜索引擎技术介绍/191.7.1 通用搜索/191.7.2 垂直搜索/201.7.3 站内搜索/211.7.4 桌面搜索/231.8 本章小结/24第2章 获得海量数据/252.1 自己的网络蜘蛛/252.1.1 抓取网页/252.1.2 网络蜘蛛遍历与实现/262.1.3 改进网络蜘蛛/302.1.4 MP3抓取/342.1.5 RSS抓取/362.1.6 图片抓取/382.1.7 垂直行业抓取/392.2 抓取数据库中的内容/422.2.1 建立数据视图/422.2.2 JDBC数据库连接/432.2.3 增量抓取/452.3 抓取本地硬盘上的文件/472.4 本章小结/49第3章 提取文档中的文本内容/503.1 从HTML文件中提取文本/503.1.1 HtmlParser介绍/533.1.2 结构化信息提取/633.1.3 查看网页的DOM结构/683.1.4 正文提取的工具NekoHTML/713.1.5 网页去噪/733.1.6 网页结构相似度计算/763.1.7 网站风格树去除文档噪声/803.1.8 正文提取/923.2 从非HTML文件中提取文本/983.2.1 TEXT文件/983.2.2 PDF文件/983.2.3 Word文件/1053.2.4 RTF文件/1063.2.5 Excel文件/1073.2.6 PowerPoint文件/1083.3 流媒体内容提取/1093.3.1 音频流内容提取/1093.3.2 视频流内容提取/1113.4 抓取限制应对方法/1133.5 本章小结/114第4章 中文分词/1154.1 Lucene中的中文分词/1154.2 Lietu中文分词的使用/1164.3 中文分词的原理/1174.4 查找词典算法/1184.5 最大概率分词方法/1234.6 新词发现/1274.7 词性标注/1294.8 本章小结/139第5章 自然语言处理/1405.1 语法解析树/1405.2 文档排重/1415.3 中文关键词提取/1425.3.1 关键词提取的基本方法/1425.3.2 从网页中提取关键词/1455.4 相关搜索/1455.5 拼写检查/1485.5.1 英文拼写检查/1485.5.2 中文拼写检查/1495.6 自动摘要/1535.6.1 自动摘要技术/1535.6.2 自动摘要的设计/1545.6.3 Lucene中的动态摘要/1625.7 自动分类/1635.7.1 Classifier4J/1645.7.2 自动分类的接口定义/1655.7.3 自动分类的SVM方法实现/1665.7.4 多级分类/1675.8 自动聚类/1705.8.1 聚类的定义/1705.8.2 K均值聚类方法/1705.8.3 K均值实现/1735.9 拼音转换/1795.10 语义搜索/1805.11 跨语言搜索/1865.12 本章小结/188第6章 创建索引库/1896.1 设计索引库结构/1906.1.1 理解Lucene的索引库结构/1906.1.2 设计一个简单的索引库/1926.2 创建和维护索引库/1936.2.1 创建索引库/1936.2.2 向索引库中添加索引文档/1946.2.3 删除索引库中的索引文档/1966.2.4 更新索引库中的索引文档/1976.2.5 索引的合并/1976.2.6 索引的定时更新/1976.2.7 索引的备份和恢复/1986.2.8 修复索引/1996.3 读写并发控制/2006.4 优化使用Lucene/2006.4.1 索引优化/2016.4.2 查询优化/2026.4.3 实现时间加权排序/2066.4.4 实现字词混合索引/2076.4.5 定制Similarity/2146.4.6 定制Tokenizer/2156.5 查询大容量索引/2176.6 本章小结/218第7章 用户界面设计与实现/2197.1 Lucene搜索接口(search代码)/2197.2 搜索页面设计/2217.2.1 用于显示搜索结果的taglib/2217.2.2 用于搜索结果分页的taglib/2237.2.3 设计一个简单的搜索页面/2257.3 实现搜索接口/2277.3.1 布尔搜索/2277.3.2 指定范围搜索/2287.3.3 搜索结果排序/2337.3.4 搜索页面的索引缓存与更新/2347.4 实现关键词高亮显示/2367.5 实现分类统计视图/2397.6 实现相似文档搜索/2447.7 实现AJAX自动完成/2467.7.1 总体结构/2477.7.2 服务器端处理/2477.7.3 浏览器端处理/2497.7.4 服务器端改进/2507.7.5 部署总结/2617.8 jQuery实现的自动完成/2627.9 集成其他功能/2677.9.1 拼写检查/2677.9.2 分类统计/2677.9.3 相关搜索/2717.9.4 再次查找/2747.9.5 搜索日志/2757.10 搜索日志分析/2767.11 本章小结/280第8章 其他高级主题/2818.1 使用Solr实现分布式搜索/2818.1.1 Solr服务器端的配置与中文支持/2828.1.2 把数据放进Solr/2878.1.3 删除数据/2898.1.4 客户端搜索界面/2908.1.5 Solr索引库的查找/2928.1.6 索引分发/2948.1.7 Solr搜索优化/2988.1.8 Solr中字词混合索引/3028.1.9 相关检索/3048.1.10 搜索结果去重/3078.1.11 分布式搜索/3118.1.12 SolrJ查询分析器/3158.1.13 扩展SolrJ/3258.1.14 扩展Solr/3278.1.15 Solr的.NET客户端/3338.1.16 Solr的PHP客户端/3348.2 图像的OCR识别/3368.3 竞价排名/3438.4 Web图分析/3448.5 使用并行程序分析数据/3508.6 RSS搜索/3518.7 本章小结/353参考资料/354

## <<自己动手写搜索引擎>>

### 章节摘录

1995年，两个年轻的学生Larry Page和Sergey Brin在一点上达成了共识——从大量数据中检索信息是计算系统面临的最大的挑战之一。

1996年，他们创建了一个叫做BackRub的搜索引擎。

这个搜索引擎后来叫做Google。

1998年，Page和Brin在Larry的大学宿舍创立Google公司的第一个数据中心。

2000年，Google开始成为全球最大的搜索引擎一直到现在。

2005年，Google股票市值超过1000亿美元。

公司首次登陆华尔街时，“让世界更美好”便是它们阐明的目标之一。

除了搜索引擎，集成AJAX技术的Gmail邮箱和Google地图创造了更好的Web界面视觉效果。

和Gmail邮箱集成的聊天工具Gtalk也得到广泛好评。

和MSN在客户端保存历史聊天记录不同，Gtalk聊天工具可以在服务器端保存并搜索聊天记录。

可以浏览整个地球的Google Earth也带给人们一种全新的体验。

你想起什么大脑没有记住的知识了吗？

用搜索引擎吧。

它往往不会让你失望。

事实上这正是Google的创始人设想的。

Google创始人Brin在一次用户大会上讲道：“如果你想搜就能搜，几乎像拥有第二个大脑，那就妙极了。

”就在Brin（一直是两人中站在前台的那个）说在兴头儿上的时候，Page在不声不响地露面之后，便带着“共谋者般的微笑”走出了房间。

## <<自己动手写搜索引擎>>

### 编辑推荐

首次揭示商业级搜索引擎实现秘密。

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》特色：根据猎兔搜索开发团队多年搜索和自然语言处理研发经验，将国内搜索工程项目实践结合当前流行理论实现搜索技术。

全方位展现出一个商用级别的Lucene搜索解决方案，主要包括网络爬虫、中文分词和搜索实现等

。

用深入浅出的方式介绍了隐马尔可夫模型等流行理论。

相关的算法全部利用Java语言实现。

#### 版权说明

本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问:<http://www.tushu007.com>