

<<搜索引擎基础教程>>

图书基本信息

书名：<<搜索引擎基础教程>>

13位ISBN编号：9787302220497

10位ISBN编号：7302220492

出版时间：2010-7

出版时间：清华大学出版社

作者：袁津生，李群

页数：320

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>

<<搜索引擎基础教程>>

内容概要

本书从教学的角度出发，对搜索引擎的原理及开发技术进行了全面的介绍，内容包括搜索引擎的基本原理、网页抓取技术、信息预处理技术、信息索引技术、信息查询技术和多媒体信息检索技术。

另外，本书还对搜索引擎开发技术进行了详细的讨论。

本书适合高等院校计算机科学与技术专业及相关专业的高年级学生和研究生阅读参考，也适合相关领域的工程技术人员参阅。

<<搜索引擎基础教程>>

书籍目录

第1章 搜索引擎概述 1.1 搜索引擎的概念、原理及历史与发展 1.1.1 搜索引擎的概念 1.1.2 搜索引擎的原理 1.2 搜索引擎的历史与发展趋势 1.2.1 搜索引擎的发展史 1.2.2 搜索引擎的发展趋势 1.3 搜索引擎的分类 1.3.1 全文搜索引擎 1.3.2 目录索引搜索引擎 1.3.3 元搜索引擎 1.3.4 分布式搜索引擎 1.4 搜索引擎的关键技术 1.4.1 信息收集和存储技术 1.4.2 信息预处理技术 1.4.3 信息索引技术 1.5 主要搜索引擎介绍 1.5.1 谷歌搜索 1.5.2 雅虎搜索 1.5.3 百度搜索 1.5.4 北大天网搜索 1.6 小结 思考题第2章 搜索引擎基础 2.1 搜索引擎的体系结构 2.1.1 搜索器 2.1.2 索引器 2.1.3 检索器 2.1.4 用户接口 2.2 搜索引擎的工作原理 2.2.1 网页搜集 2.2.2 网页处理 2.2.3 查询服务 2.3 搜索引擎的数据结构 2.3.1 存储结构 2.3.2 信息库 2.3.3 文本索引 2.3.4 词典 2.3.5 采样表 2.3.6 前向索引 2.3.7 后向索引 2.4 元搜索引擎 2.4.1 元搜索引擎的基本构成 2.4.2 元搜索引擎的分类 2.4.3 常用元搜索引擎介绍 2.4.4 元搜索引擎的特点 2.4.5 主要技术指标 2.5 个性化搜索引擎 2.5.1 系统模块及其功能 2.5.2 个性化搜索引擎的关键技术 2.6 智能搜索引擎 2.6.1 智能搜索引擎特征 2.6.2 智能搜索引擎主要技术 2.7 小结 思考题第3章 网页抓取技术 3.1 搜索引擎爬虫 3.1.1 网络爬虫工作原理 3.1.2 开源网络爬虫简介 3.1.3 网页信息的抓取 3.2 搜索引擎爬虫的关键技术 3.2.1 网页抓取优先策略 3.2.2 深度优先策略 3.2.3 广度优先策略 3.2.4 最佳优先策略 3.2.5 不重复抓取策略 3.2.6 网页重访策略 3.2.7 网页抓取提速策略 3.2.8 Robots协议 3.3 小结 思考题第4章 网页信息预处理技术第5章 信息索引技术第6章 信息查询与评价技术第7章 多媒体信息检索技术第8章 搭建基于Lucene的搜索引擎第9章 搭建基于Nutch的搜索引擎参考文献

章节摘录

插图：在抓取过程中，可以使用多个爬虫来合作抓取，这样可以进一步降低每个爬虫的用于记录历史抓取情况的哈希表大小。

如果有 n 个爬虫，则可将哈希表继续压缩到原有大小的 $\frac{1}{n}$ 。

如果 n 个爬虫分别运行在不同的机器上，那么每个机器被哈希表占用的内存用量将非常少；通常保持抓取历史记录所需要的内存在百兆字节左右是恰当的。

通过不重复抓取的方法初步解决了死循环的问题，即抓过的不再抓。

然而实际操作中还有这样一个问题，如果任意两个网页存在链接，则链接它们的最短路径值为17。

这样，爬虫无论用何种遍历方法都不能保证一定会按照最佳路径抓取每一个网页，因为任何一个网页都可能从多个种子站点开始广度优先被遍历到。

为了防止爬虫无限制的广度优先抓取，必须在某个深度上进行限制。

到达这个深度后就应该停止抓取，这个深度的取值就是万维网直径长度。

当在最大深度上停止时，那些深度过大的未抓网页，总是期望可以从其他种子站点更加经济地到达。

例如，种子站点B和C在抓取到深度为17的时候，立即停止抓取，把抓取剩余网页的机会留给从种子站点A出发的进行抓取工作的爬虫。

此外，深度优先策略和广度优先策略的组合可以有效地保证抓取过程中的封闭性。

即在抓取过程（遍历路径）中总是在抓取相同域名下的网页，而很少出现其他域名下的网页。

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>