

<<搜索引擎>>

图书基本信息

书名：<<搜索引擎>>

13位ISBN编号：9787111288084

10位ISBN编号：7111288084

出版时间：2010-6

出版时间：机械工业出版社

作者：W.Bruce Croft,Donald Metzler,Trevor Strohman

页数：309

译者：刘挺,秦兵,张宇,车万翔

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>>

前言

本书综述了信息检索中的重要问题，并介绍了这些问题如何对搜索引擎的设计与实现产生影响。本书并不是按照相同的详细程度描述每个主题，相反，我们侧重于那些对于实现搜索引擎组件以及组件背后的信息检索模型最重要的部分。

网络搜索引擎显然是一个重要的话题，我们主要覆盖了在网络上使用的搜索技术，但搜索引擎在其他场合中也有应用，这就是为什么我们重点强调各种搜索引擎背后的信息检索理论与概念的原因。

本书的目标读者群主要是计算机科学或计算机工程的本科生，但研究生也会发现本书是有用的，此外，本书也适合多数情报科学专业的学生。

最后，无论读者是什么背景，通过阅读本书都可以对他们动手开发搜索引擎有所帮助。

本书中涉及数学知识，但并不深奥。

书中也有代码和程序设计的练习，但对于那些已经完成了基本计算机科学与程序设计课程的人来说，完全可以掌握。

每章末尾的练习使用了被称为Galago的基于Java的开源搜索引擎。

Galago既是为本书所设计的，也借鉴了Lemur和Indri项目的经验。

换句话说，这是一个功能齐全的能够支持真正应用的搜索引擎。

许多编程练习都是针对Galago组件的使用、修改和扩展。

内容在第1章，我们对信息检索及它与搜索引擎的关系做了一个高层次的回顾。

在第2章，我们描述了搜索引擎的架构，这一章全面介绍搜索引擎的各个组件，但没有涉及细节问题。

在第3章，我们关注于爬取、文档信息源和其他用于获取被检索信息的技术。

第4章描述了文本的统计特征，以及用来处理和识别重要特征的技术，并为建立索引做准备。

第5章描述了怎样为有效的搜索建立索引，以及怎样利用索引处理查询。

在第6章，我们描述了怎样处理查询，并把它们转换为更好的形式，以表达用户的信息需求。

<<搜索引擎>>

内容概要

本书介绍了信息检索（IR）中的关键问题，以及这些问题如何影响搜索引擎的设计与实现，并且用数学模型强化了重要的概念。

对于网络搜索引擎这一重要的话题，书中主要涵盖了在网络上广泛使用的搜索技术。

本书适用于高等院校计算机科学或计算机工程专业的本科生、研究生，对于专业人士而言，本书也不失为一本理想的入门教材。

作者简介

作者：（美国）克罗夫特（W.Bruce Croft）（美国）Donald Metzler（美国）Trevor Strohman 译者：刘挺 秦兵 张宇 等
克罗夫特（W.Bruce Croft），马萨诸塞大学阿默斯特分校计算机科学特聘教授、ACM会士。

他创建了智能信息检索研究中心，发表了200余篇论文，多次获奖，其中包括2003年由ACM SIGIR颁发的Gerard Salton奖。

Donald Metzler，马萨诸塞大学阿默斯特分校博士，是位于加州Santa Clara的雅虎研究中心搜索与计算机广告组的研究科学家。

Trevor Strohman，马萨诸塞大学阿默斯特分校博士。

他开发了Galago搜索引擎，也是Indri搜索引擎的主要开发者。

<<搜索引擎>>

书籍目录

出版者的话 译者序 前言 第1章 搜索引擎和信息检索 1.1 什么是信息检索 1.2 重要问题 1.3 搜索引擎 1.4 搜索工程师 参考文献和深入阅读 练习 第2章 搜索引擎的架构 2.1 什么是软件架构 2.2 基本的构件 2.3 组件及其功能 2.3.1 文本采集 2.3.2 文本转换 2.3.3 索引的创建 2.3.4 用户交互 2.3.5 排序 2.3.6 评价 2.4 搜索引擎是如何工作的 参考文献和深入阅读 练习 第3章 信息采集和信息源 3.1 确定搜索的内容 3.2 网络信息爬取 3.2.1 抓取网页 3.2.2 网络爬虫 3.2.3 时新性 3.2.4 面向主题的信息采集 3.2.5 深层网络 3.2.6 网站地图 3.2.7 分布式信息采集 3.3 文档和电子邮件的信息采集 3.4 文档信息源 3.5 转换问题 3.6 存储文档 3.6.1 使用数据库系统 3.6.2 随机存取 3.6.3 压缩和大规模文件 3.6.4 更新 3.6.5 BigTable 3.7 重复检测 3.8 去除噪声 参考文献和深入阅读 练习 第4章 文本处理 4.1 从词到词项 4.2 文本统计 4.2.1 词表增长 4.2.2 估计数据集和结果集大小 4.3 文档解析 4.3.1 概述 4.3.2 词素切分 4.3.3 停用词去除 4.3.4 词干提取 4.3.5 短语和n元串 4.4 文档结构和标记 4.5 链接分析 4.5.1 锚文本 4.5.2 PageRank 4.5.3 链接质量 4.6 信息抽取 4.7 国际化 参考文献和深入阅读 练习 第5章 基于索引的相关排序 5.1 概述 5.2 抽象的相关排序模型 5.3 倒排索引 5.3.1 文档 5.3.2 计数 5.3.3 位置 5.3.4 域与范围 5.3.5 分数 5.3.6 排列 5.4 压缩 5.4.1 熵与歧义 5.4.2 Delta编码 5.4.3 位对齐码 5.4.4 字节对齐码 5.4.5 实际应用中的压缩 5.4.6 展望 5.4.7 跳转和跳转指针 5.5 辅助结构 5.6 索引构建 5.6.1 简单构建 5.6.2 融合 5.6.3 并行与分布式 5.6.4 更新 5.7 查询处理 5.7.1 document-at-a-time评价 5.7.2 term-at-a-time评价 5.7.3 优化技术 5.7.4 结构化查询 5.7.5 分布式的评价 5.7.6 缓存 参考文献和深入阅读 练习 第6章 查询与界面 6.1 信息需求与查询 6.2 查询转换与提炼 6.2.1 停用词去除和词干提取 6.2.2 拼写检查和建议 6.2.3 查询扩展 6.2.4 相关反馈 6.2.5 上下文和个性化 6.3 搜索结果显示 6.3.1 搜索结果页面与页面摘要 6.3.2 广告与搜索 6.3.3 结果聚类 6.4 跨语言搜索 参考文献和深入阅读 练习 第7章 检索模型 7.1 检索模型概述 7.1.1 布尔检索 7.1.2 向量空间模型 7.2 概率模型 7.2.1 将信息检索作为分类问题 7.2.2 BM25排序算法 7.3 基于排序的语言模型 7.3.1 查询项似然排序 7.3.2 相关性模型和伪相关反馈 7.4 复杂查询和证据整合 7.4.1 推理网络模型 7.4.2 Galago查询语言 7.5 网络搜索 7.6 机器学习和信息检索 7.6.1 排序学习 7.6.2 主题模型和词汇不匹配 7.7 基于应用的模型 参考文献和深入阅读 练习 第8章 搜索引擎评价 8.1 搜索引擎评价的意义 8.2 评价语料 8.3 日志 8.4 效果评价 8.4.1 召回率和准确率 8.4.2 平均化和插值 8.4.3 关注排序靠前的文档 8.4.4 使用用户偏好 8.5 效率评价 8.6 训练、测试和统计 8.6.1 显著性检验 8.6.2 设置参数值 8.6.3 在线测试 8.7 基本要点 参考文献和深入阅读 练习 第9章 分类和聚类 9.1 分类 9.1.1 朴素贝叶斯 9.1.2 支持向量机 9.1.3 评价 9.1.4 分类器和特征选择 9.1.5 垃圾、情感及在线广告 9.2 聚类 9.2.1 层次聚类和K均值聚类 9.2.2 K近邻聚类 9.2.3 评价 9.2.4 如何选择K 9.2.5 聚类和搜索 参考文献和深入阅读 练习 第10章 社会化搜索 10.1 什么是社会化搜索 10.2 用户标签和人工索引 10.2.1 搜索标签 10.2.2 推测缺失的标签 10.2.3 浏览和标签云 10.3 社区内搜索 10.3.1 什么是社区 10.3.2 社区发现 10.3.3 基于社区的问答 10.3.4 协同搜索 10.4 过滤和推荐 10.4.1 文档过滤 10.4.2 协同过滤 10.5 P2P搜索和元搜索 10.5.1 分布式搜索 10.5.2 P2P网络 参考文献和深入阅读 练习 第11章 超越词袋 11.1 概述 11.2 基于特征的检索模型 11.3 词项依赖模型 11.4 再谈结构化 11.4.1 XML检索 11.4.2 实体搜索 11.5 问题越长,答案越好 11.6 词语、图片和音乐 11.7 搜索能否适用于所有情况 参考文献和深入阅读 练习 参考文献

章节摘录

插图：2.查询转换查询转换组件包括一系列的技术，这些技术用于在生成排好序的文档之前和之后改善初始查询。

最简单的处理涉及一些对文档进行文本转换的技术。

在查询文本上，需要进行词素切分、停用词去除和词干提取这些工作，以生成与文档词项具有可比性的索引词。

拼写检查（spell checking）和查询建议（query suggestion）是查询转换中的技术，生成与用户初始查询相似的输出。

在这两种情况下，向用户提供初始查询的一些候选查询，这些候选查询可能纠正了拼写错误或者是对用户所需信息的更规范描述。

这些技术通常会导致为网络应用搜集大量的查询日志（query log）。

查询扩展（query expansion）技术是对查询进行推荐或者增加一些额外的词项，但通常都是在对文档中词项的出现情况分析的基础上进行的。

该分析通常是用不同的信息源，如整个文档集合、检索到的文档或者用户计算机上的文档。

相关反馈（relevance feedback）是一种查询扩展技术，利用用户认为相关的文档中出现的词项对查询进行扩展。

3.结果输出结果输出组件负责对相关组件得到的排好序的文档的结果进行显示。

可能包含的任务有生成网页摘要（snippets）来对检索到的文档内容进行概括；强调（highlighting）文档中重要的词和段落；对输出结果聚类以找到文档相关的类别；以及将相应的广告增加到结果显示中。

在涉及多种语言的应用系统中，结果可能会被翻译成同一种的语言。

编辑推荐

《搜索引擎:信息检索实践》：计算机科学丛书

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>