

<<信息检索>>

图书基本信息

书名：<<信息检索>>

13位ISBN编号：9787115212252

10位ISBN编号：7115212252

出版时间：2009-10

出版时间：人民邮电出版社

作者：（美）格罗斯曼，（美）弗里德 著

页数：332

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>

<<信息检索>>

内容概要

本书是“信息检索”课程的优秀教材，书中对信息检索的概念、原理和算法进行了详细介绍，内容主要包括检索策略、检索实用工具、跨语言信息检索、查询处理、集成结构化及数据和文本、并行信息检索以及分布式信息检索等，并给出了阐述算法的大量实例。

本书有一定的深度和广度，而且所有的内容都用当前的技术阐述，是高等院校计算机及信息管理等相关专业本科生和研究生的理想教材，对信息检索领域的科研和技术人员也是很好的参考书。

<<信息检索>>

作者简介

格罗斯曼 (David A.Grossman) , 佐治亚梅森大学博士。
现在伊利诺伊理工大学计算机系任教。
曾在美国政府部门高级技术服务中心和研究发展办公室担任项目经理。
主要研究领域包括信息检索、结构化与非结构化数据集成以及数据挖掘。

书籍目录

1. INTRODUCTION
2. RETRIEVAL STRATEGIES 2.1 Vector Space Model 2.2 Probabilistic Retrieval Strategies 2.3 Language Models 2.4 Inference Networks 2.5 Extended Boolean Retrieval 2.6 Latent Semantic Indexing 2.7 Neural Networks 2.8 Genetic Algorithms 2.9 Fuzzy Set Retrieval 2.10 Summary 2.11 Exercises
3. RETRIEVAL UTILITIES 3.1 Relevance Feedback 3.2 Clustering 3.3 Passage-based Retrieval 3.4 N-grams 3.5 Regression Analysis 3.6 Thesauri 3.7 Semantic Networks 3.8 Parsing 3.9 Summary 3.10 Exercises
4. CROSS-LANGUAGE INFORMATION RETRIEVAL 4.1 Introduction 4.2 Crossing the Language Barrier 4.3 Cross-Language Retrieval Strategies 4.4 Cross Language Utilities 4.5 Summary 4.6 Exercises
5. EFFICIENCY 5.1 Inverted Index 5.2 Query Processing 5.3 Signature Files 5.4 Duplicate Document Detection 5.5 Summary 5.6 Exercises
6. INTEGRATING STRUCTURED DATA AND TEXT 6.1 Review of the Relational Model 6.2 A Historical Progression 6.3 Information Retrieval as a Relational Application 6.4 Semi-Structured Search using a Relational Schema 6.5 Multi-dimensional Data Model 6.6 Mediators 6.7 Summary 6.8 Exercises
7. PARALLEL INFORMATION RETRIEVAL 7.1 Parallel Text Scanning 7.2 Parallel Indexing 7.3 Clustering and Classification 7.4 Large Parallel Systems 7.5 Summary 7.6 Exercises
8. DISTRIBUTED INFORMATION RETRIEVAL 8.1 A Theoretical Model of Distributed Retrieval 8.2 Web Search 8.3 Result Fusion 8.4 Peer-to-Peer Information Systems 8.5 Other Architectures 8.6 Summary 8.7 Exercises
9. SUMMARY AND FUTURE DIRECTIONS
References
Index

章节摘录

3.4.1 DAmore and Mah Initial information retrieval research focused on n-grams as presented in [DAmore and Mah, 1985]. The motivation behind their work was the fact that it is difficult to develop mathematical models for terms since the potential for a term that has not been seen before is infinite. With n-grams, only a fixed number of n-grams can exist for a given value of n. A mathematical model was developed to estimate the noise in indexing and to determine appropriate document similarity measures. DAmore and Mah's method replaces terms with n-grams in the vector space model. The only remaining issue is computing the weights for each n-gram. Instead of simply using n-gram frequencies, a scaling method is used to normalize the length of the document. DAmore and Mah's contention was that a large document contains more n-grams than a small document, so it should be scaled based on its length. To compute the weights for a given n-gram, DAmore and Mah estimated the number of occurrences of an n-gram in a document. The first simplifying assumption was that n-grams occur with equal likelihood and follow a binomial distribution. Hence, it was no more likely for n-gram "ABC" to occur than "DEE". The Zipfian distribution that is widely accepted for terms is not true for n-grams. DAmore and Mah noted that n-grams are not equally likely to occur, but the removal of frequently occurring terms from the document collection resulted in n-grams that follow a more binomial distribution than the terms. DAmore and Mah computed the expected number of occurrences of an n-gram in a particular document. This is the product of the number of n-grams in the document (the document length) and the probability that the n-gram occurs. The n-grams probability of occurrence is computed as the ratio of its number of occurrences to the total number of n-grams in the document. DAmore and Mah continued their application of the binomial distribution to derive an expected variance and, subsequently,

媒体关注与评论

“ 本书涉及最新的研究成果，语言经得起推敲，还精心准备了大量的实例说明，适合作为研究生和本科生信息检索课程的首选教材。

” ——美国马萨诸塞大学阿默斯特校区计算机系杰出教授 W.Bruce Croft “ 推荐把本书作为计算机专业学生的首选教材，同时也适用于SEO专业人员和Web开发者阅读，将搜索技术，算法和启发式方法运用于他们的项目中。

” ——信息技术与服务顾问 E.Garcia博士

<<信息检索>>

编辑推荐

随着Google、百度等搜索引擎公司的崛起，信息检索已经成为令人振奋的热门研究领域。

《信息检索：算法与启发式方法(英文版·第2版)》从发展的角度描述了ad hoc信息检索，讨论了用来实现大规模数据检索的最新算法，详细介绍了推理网络和系统的效率，并且对每种方法都给出了详细可行的实例。

此外，《信息检索：算法与启发式方法(英文版·第2版)》整合了结构化和非结构化数据的处理技术，这是其他教材所不具备的。

第2版新增加了IR语言模型和跨语言检索，还讨论了许多当前的热点话题，如XML、P2P信息检索、文本查重、文档并行聚类、不同检索策略的融合、信息中间表示等。

《信息检索：算法与启发式方法(英文版·第2版)》兼顾了学科广度和主题深度，把握了最新的发展趋势，是信息检索领域的一本名著，更为许多著名高校（如美国普林斯顿大学、罗格斯大学）采用为教材。

随着Google、百度等搜索引擎公司的崛起，信息检索已经成为令人振奋的热门研究领域。

《信息检索：算法与启发式方法(英文版·第2版)》从发展的角度描述了ad hoc信息检索，讨论了用来实现大规模数据检索的最新算法，详细介绍了推理网络和系统的效率，并且对每种方法都给出了详细可行的实例。

此外，《信息检索：算法与启发式方法(英文版·第2版)》整合了结构化和非结构化数据的处理技术。这是其他教材所不具备的。

第2版新增加了IR语言模型和跨语言检索，还讨论了许多当前的热点话题，如XML、P2P信息检索、文本查重、文档并行聚类、不同检索策略的融合、信息中间表示等。

《信息检索：算法与启发式方法(英文版·第2版)》兼顾了学科广度和主题深度，把握了最新的发展趋势，是信息检索领域的一本名著，更为许多著名高校（如美国普林斯顿大学、罗格斯大学）采用为教材。

<<信息检索>>

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>