

TURING

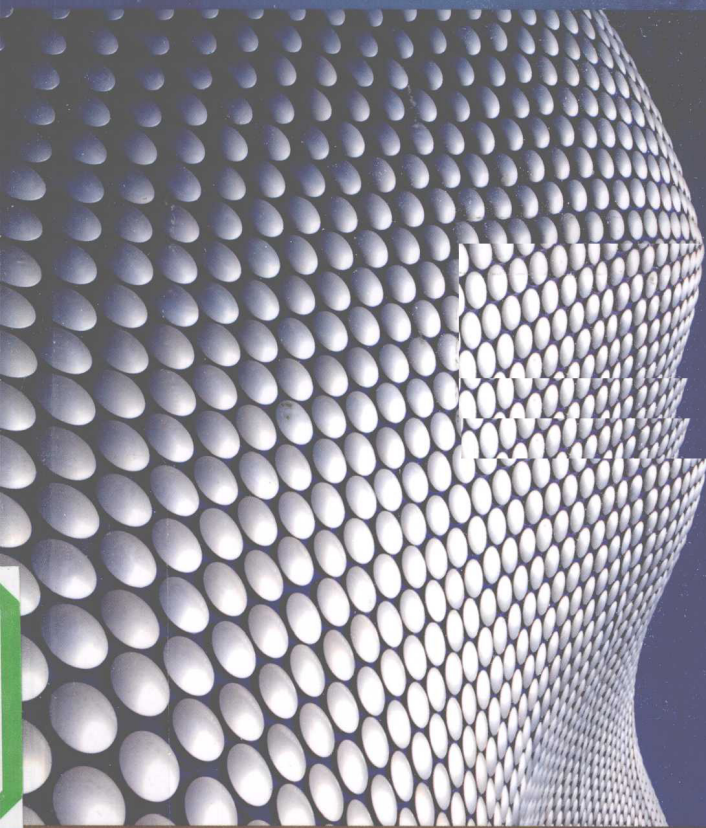
图灵计算机科学丛书

CAMBRIDGE

# 信息检索导论

Introduction to Information Retrieval

[美] Christopher D. Manning  
[美] Prabhakar Raghavan 著 王斌 译  
[德] Hinrich Schütze



人民邮电出版社  
POSTS & TELECOM PRESS

# 信息检索导论

第二版

清华大学出版社

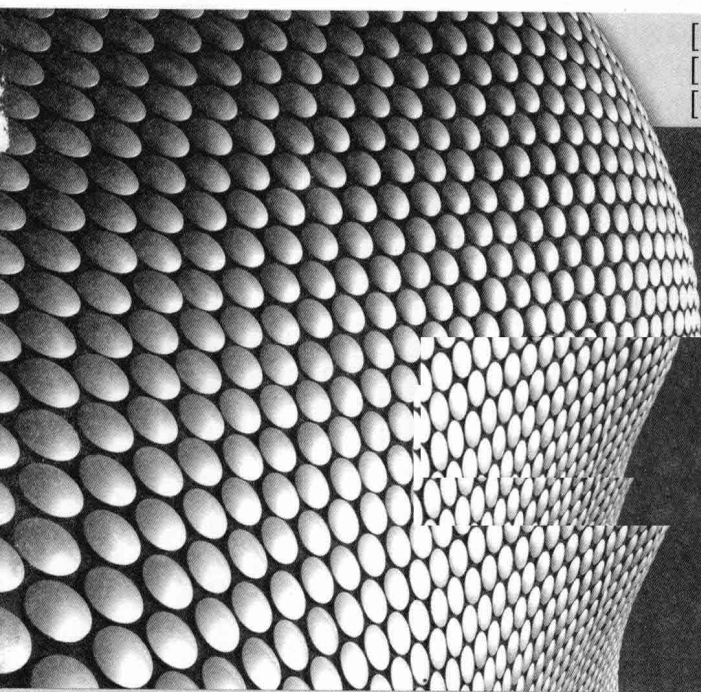
TURING

图灵计算机科学丛书

# 信息检索导论

Introduction to Information Retrieval

[美] Christopher D. Manning 著 王斌 译  
[美] Prabhakar Raghavan  
[德] Hinrich Schütze



人民邮电出版社  
北京

## 图书在版编目 (CIP) 数据

信息检索导论 / (美) 曼宁, (美) 拉哈万, (德) 舒策著; 王斌译. — 北京: 人民邮电出版社, 2010.9

(图灵计算机科学丛书)

书名原文: Introduction to Information Retrieval

ISBN 978-7-115-23424-7

I. ①信… II. ①曼… ②拉… ③舒… ④王… III. ①情报检索—教材 IV. ①G252.7

中国版本图书馆CIP数据核字(2010)第142051号

## 内 容 提 要

本书是一本讲授信息检索的经典教材。全书共 21 章, 前 8 章详述了信息检索的基础知识, 包括倒排索引、布尔检索及词项权重计算和评分算法等, 后 13 章介绍了一些高级话题, 如基于语言建模的信息检索模型、基于机器学习的排序方法和 Web 搜索技术等。另外, 本书还着重讨论了文本聚类技术这一信息检索中不可或缺的组成部分。全书语言流畅, 由浅入深, 一气呵成。

本书适合作为高等院校相关专业高年级本科生和研究生的课程教材, 也可供信息检索领域的研究人员和专业人士参考。

图灵计算机科学丛书

## 信息检索导论

◆ 著 [美] Christopher D. Manning  
[美] Prabhakar Raghavan  
[德] Hinrich Schütze

译 王 斌

责任编辑 杨海玲

执行编辑 罗词亮 陈 潇

◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街14号  
邮编 100061 电子函件 315@ptpress.com.cn  
网址 <http://www.ptpress.com.cn>  
三河市海波印务有限公司印刷

◆ 开本: 787×1092 1/16

印张: 24.25

字数: 620千字

印数: 1-3 000册

2010年9月第1版

2010年9月河北第1次印刷

著作权合同登记号 图字: 01-2009-7281号

ISBN 978-7-115-23424-7

定价: 69.00 元

读者服务热线: (010)51095186 印装质量热线: (010)67129223

反盗版热线: (010)67171154

# 版 权 声 明

*Introduction to Information Retrieval* (978-0-521-86571-5) by Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze first published by Cambridge University Press 2008.

All rights reserved.

This simplified Chinese edition for the People's Republic of China is published by arrangement with the Press Syndicate of the University of Cambridge, Cambridge, United Kingdom.

© Cambridge University Press & Posts & Telecom Press 2010.

This book is in copyright. No reproduction of any part may take place without the written permission of Cambridge University Press and Posts & Telecom Press.

This edition is for sale in the People's Republic of China (excluding Hong Kong SAR, Macao SAR and Taiwan Province) only.

此版本仅限在中华人民共和国境内（不包括香港、澳门特别行政区及台湾地区）销售。

# 译者序

第一次见到这本书的电子版是在 2007 年的年底,当时北京大学的闫宏飞博士向我推荐了这本书。从网上下载书稿的电子版之后,我便迫不及待地在一周时间内通读了这本书。读完之后便萌发了翻译这本书的冲动,随后我就联系作者、联系剑桥大学出版社并通过朋友寻找获得授权的国内出版社。辗转数月之后,我被告知该书已经交由其他学者翻译,很快便可出版。听到这个消息,虽然我有些遗憾,但也算是心里的一块石头得以落地。所以,当去年 8 月人民邮电出版社突然联系并询问我是否有意翻译这本著作时,我心里的惊讶可想而知。当然,惊讶之余我毫不犹豫地接受了这份邀请,并从此开始了长达数月的翻译历程。

之所以愿意翻译这本书不仅仅是由于该书的作者都是学术界甚至业界鼎鼎大名的人物,更主要的是因为本书在内容和组织上都有独到之处。之前也有很多信息检索方面的教材,但是其中很多内容已经过时。信息检索是一门不断发展并和其他领域、技术不断融合的学科。这本书补充了一些近年来受到广泛关注的新内容。比如:基于语言建模的信息检索模型、基于机器学习的排序方法、检索结果的 Snippet 生成、聚类标签生成、XML 检索、搜索广告、网页作弊,等等。除此之外,本书每章末尾的“参考文献及补充读物”一节也给出了相关技术的最新进展。本书在内容上与传统教材的另一个显著不同之处是加大了文本分类/聚类技术的介绍篇幅,实际上这些技术已经成为当代信息检索不可分割的一部分。另一方面,本书在深度上超过了大部分传统教材。在介绍信息检索技术的同时,本书深入介绍了其背后所依赖的原理。因此,本书不仅可以用作信息检索领域的入门教材,还能满足对该领域进行深入研究的需要。另外,本书给出了很多实际当中的运行算法和实施细节,这些内容对于信息检索技术的实际应用有很好的参考价值。最后值得一提的是,本书在结构上也进行了巧妙构思。首先通过一个例子引出基本技术,然后通过基本技术的不断增强来介绍信息检索的其他技术。全书浑然一体,读起来也有一气呵成的感觉。

这么一本优秀的著作在给译者的翻译带来无穷动力的同时,无疑也给翻译带来了无形的压力。为了尽量保证每章译稿的质量并保持译文的前后一致性,整本书的初译工作全部由译者本人独立完成,在翻译过程中译者也阅读了大量相关的教材和论文,并前后进行了六次自我校对。在校对过程中,有很多学术界同仁也提出了很多宝贵的意见和建议。他们包括:中科院研究生院的朱廷劲教授、中科院自动化所的赵军研究员、中科院软件所的孙乐研究员、复旦大学的黄萱菁教授、江西师范大学的王明文教授、江西财经大学的刘德喜博士、北京大学的闫宏飞博士、何靖博士、清华大学的张敏博士、北京语言大学的徐燕博士等。译者所在的中科院计算所信息检索课题组及选修研究生院《现代信息检索》课程的部分学生也提出了大量修改建议,他们是:郎皓、李亚楠、顾智宇、李鹏、李锐、马宏远、张爱华、蒋在帆、沈沉、史亮、卫冰洁、崔雅超、赵琴琴、李恒训、袁平广、邱泳钦、李丹、鲁凯、徐飞、张帅、张启龙、廖凤、钟进文、朱亮、赵娟等。对于他们无私的帮助,我表示由衷的感谢。感谢我所在的前瞻研究实验室主任

李锦涛老师对我的翻译工作给予的支持和肯定。当然，本书的翻译工作得以顺利完成，还要感谢人民邮电出版社众多工作人员特别是责任编辑杨海玲女士在各方面的支持和帮助。另一个需要感谢的是我的妻子，在前前后后近八个月当中，除上班时间完成自己的科研工作外，我几乎所有的业余时间都用在翻译和校对上，而她却默默地承担起两岁的儿子的所有抚育责任。

翻译的过程中，我还有幸与原文的第二作者 Prabhakar Raghavan 教授进行了当面交流，他对我翻译工作给予了极大鼓励。在与原文作者的邮件交流中，我也澄清了一些理解上的误区，并修正了原书中的多处错误。

虽然得到了众人的帮助，自己也算认真努力，但由于本人专业水平、理解能力和写作功底都十分有限，加上时间上仍显仓促，最后的译稿中一定存在不少理解上的偏差，译文也会有许多生硬之处。希望读者能不吝提出修改的意见和建议，以便对现有译稿不断改进，直至为国内信息检索领域的读者真正造福为止。来信请联系 [wbxjj2008@gmail.com](mailto:wbxjj2008@gmail.com)，对译稿的修改结果也会及时公布在网站 <http://ir.ict.ac.cn/~wangbin/iir-book/> 上。原书的初稿电子版、相关课件、勘误表、论坛等信息也可以从网站 <http://nlp.stanford.edu/IR-book/information-retrieval-book.html> 下载。

---

## 译者简介

---



**王斌** 博士，中国科学院计算技术研究所前瞻研究实验室信息检索课题组组长，副研究员，博士生导师。主持国家 973、863、国家自然科学基金、国际合作基金、部委及企业合作等课题 20 余项，在包括 SIGIR、CIKM、EMNLP 等在内的会议和刊物上发表学术论文 100 余篇。担任 CIKM、AIRS、CCIR、SEWM 等国际国内会议的程序委员会委员，同时是 ACM 和 IEEE 会员、中国计算机学会高级会员、中国中文信息学会会员、中文信息学会信息检索专业委员会委员及《中文信息学报》编委。自 2006 年起在中国科学院研究生院讲授《现代信息检索》研究生课程，迄今培养博士、硕士研究生近 30 名。

# 前 言

研究表明，直到 20 世纪 90 年代，大多数人还是首选通过别人而不是使用信息检索系统来获取信息。当然，那时候大多数人也往往通过旅行社来安排自己的行程。然而，在过去的十年中，信息检索效果的不断优化已经使 Web 搜索引擎的质量达到了一个新的水平，大多数用户在大部分情况下都对搜索的结果感到满意。Web 搜索引擎已经成为用户发现和获取信息的常规和首选渠道。以统计数据为证，2004 年美国 Pew 研究中心的一项因特网调查（Fallows 2004）结果表明，有 92% 的因特网用户认为因特网是人们获取日常信息的良好渠道。令很多人惊讶的是，信息检索也从一个以学术研究为主的领域，摇身一变而成为人们赖以获取日常信息的工具背后的基础学科。本书主要介绍该学科的核心理论基础，既考虑研究生科研的需求，也兼顾了高年级本科生学习的需求。

但是，信息检索并非始于 Web。在应对信息存取的各种挑战的过程中，信息检索逐渐发展成为一门给各种形式的内容搜索提供原理性方法的学科。信息检索起初主要面向科学文献和馆藏记录，但是很快就扩展到其他形式的内容，特别是新闻记者、律师、医生等特定领域专业人士所需的信息内容。信息检索中的很多学术研究都围绕上述内容展开，而其实践方面则主要是为公司或政府部门提供非结构化信息的获取服务，这些领域的研究和实践构成了本书的主要内容。

然而，近年来信息检索革新的主要推动力却来自万维网，因为网络上聚集了数以千万计的网络用户发布的内容。如果这些内容不能及时发现、标注和分析，并为有需求的人们提供相关的、全面的信息，那么它们的存在将毫无意义。到 20 世纪 90 年代末，很多人逐渐意识到，由于 Web 的规模呈指数级增长，继续给整个 Web 建立索引很快会变得毫无可能。但是，卓越的科学创新、一流的工程水平、日益低廉的计算机硬件价格及 Web 搜索商业化基础的壮大等一系列因素，促成了当今主流搜索引擎的产生与成长。这些搜索引擎一天之内能够完成对数十亿网页的数亿次搜索请求，并且每次搜索都能够在亚秒级时间内返回高质量的结果。

---

## 本书组织结构及课程设计

---

本书是在我们于斯坦福大学和斯图加特大学所讲授的一系列课程的教学成果总结。这些课程持续的时间从四分之一学期、半学期到一学期不等，主要面向低年级计算机专业的研究生，也曾用于高年级计算机专业的本科生和法律、医学信息学、统计、语言学及其他工程学科背景学生的教学。因此，本书主要的写作原则是提供一个学期的信息检索研究生课程，并尽量覆盖信息检索的学科重点。另一个原则是尽量让每章的内容能在约 75~90 分钟内讲授完。

本书前 8 章介绍信息检索的基础知识，特别是搜索引擎的核心理论。这八章对于任何信息检索课程来说都是核心部分。第 1 章主要介绍倒排索引，并说明如何通过这种索引实现简单的布尔查询。第 2 章介绍索引之前的文档预处理过程，并讨论在不同的功能和速度要求下对倒排索引



进行改进的方法。第 3 章主要介绍词典搜索的数据结构,并给出查询存在拼写错误或者与被搜索的文档中的词汇不能精确匹配时的处理方法。第 4 章主要介绍基于文本集合构建倒排索引的几个算法,并着重介绍具有高扩展性的分布式算法,这类算法适用于大规模文档集的索引构建。第 5 章介绍词典和倒排索引的压缩技术,这些技术对于实现大型搜索引擎的亚秒级查询响应十分关键。第 1~5 章中介绍的索引和查询仅针对布尔检索 (Boolean retrieval),即一篇文档和查询要么匹配,要么不匹配。那么,如何度量查询和文档的匹配程度,或者说如何根据文档和查询的匹配情况对结果打分呢?对这个问题的回答构成了第 6、第 7 章词项权重计算和评分算法的主要内容。也就是说,给定查询,我们可以利用这两章介绍的技术,按照文档评分的结果次序输出结果列表。第 8 章主要介绍信息检索系统的评价技术,即根据检索系统返回结果的相关性对不同系统进行评价,从而可以在基准文档集和查询上对不同系统的性能进行比较。

在前 8 章的基础上,本书的第 9~21 章涵盖了信息检索的一些高级话题。第 9 章介绍了相关反馈和查询扩展技术,其目的在于增加相关文档返回的可能性。第 10 章介绍了采用 XML 和 HTML 等标记语言的结构化文档的检索,这其中我们将结构化文档的检索进行约简,并采用第 6 章所介绍的向量空间模型进行求解。第 11 章和第 12 章介绍基于概率论的信息检索模型。其中,第 11 章介绍传统的概率检索模型,它提供了一个相关度计算框架,在给定一系列查询词项时,能够计算一篇文档与查询相关的概率。这个概率显然可以用于文档的评分和排序。第 12 章给出了另一种方法,即对文档集中的每篇文档建立一个语言模型,然后在每个模型下估计查询生成的概率。这个概率也显然可以用于文档的评分和排序。

第 13~18 章介绍了信息检索中各种形式的机器学习和数值方法。第 13~15 章主要关注文档分类的问题,即在给定一系列文档及其归属类别的前提下,将新文档分配到某个或者某几个类别中去。第 13 章首先指出统计分类是一个成功的搜索引擎所必需的关键技术之一,接着介绍了朴素贝叶斯算法(该算法概念虽然简单,但是文本分类的效率很高),最后给出了文本分类的评价技术。第 14 章将第 6 章所讲述的向量空间模型应用于文本分类,介绍了几种基于向量空间模型的分类方法,主要包括 Rocchio 和  $k$ NN ( $k$  nearest neighbor) 两种分类算法。本章最后给出了用于分类方法选择的偏差-方差折中准则,而偏差-方差折中也是学习问题的一个重要特点。第 15 章介绍了支持向量机,这是目前公认的效果最好的文本分类算法。另外,本章还将分类问题和一些看上去与文本分类无关的问题(比如如何从给定的训练集中推导出检索的评分函数)联系起来。

第 16~18 章主要介绍文档的聚类技术。第 16 章在概述信息检索中的一些重要聚类应用的基础上,主要介绍了两个扁平聚类算法: $K$ -均值算法和 EM 算法。前者是一个效率很高并被广泛应用的算法;后者虽然计算复杂度高一些,但是灵活性更好。第 17 章介绍信息检索对层次聚类(而非扁平聚类)的应用需求,并介绍了一些能产生层次簇结构的聚类算法。这一章还探讨了自动生成聚类标签的难题。第 18 章介绍了一些线性代数方法,它们是对聚类方法的扩展,并且为线性代数方法在信息检索的应用提供了极具吸引力的前景,其中最具有代表性的方法是隐性语义索引。

第 19~21 章主要介绍 Web 搜索这个具体的应用。第 19 章概述了 Web 搜索所面临的基本挑战,并给出了 Web 信息检索中的一些普遍使用的技术。第 20 章介绍了一个基本网络采集器的体系结构和必要需求。最后,第 21 章讨论了链接分析在 Web 搜索中的作用,其中用到了线性代数和高级概率论中的方法。

本书并没有囊括信息检索的所有主题，因为有些主题超出了信息检索入门课程的范围。当然，感兴趣的读者可以参见如下参考书籍。

*Cross-language IR* (跨语言检索): Grossman and Frieder 2004, 第 4 章; Oard and Dorr 1996。  
*Image and multimedia IR* (图像和多媒体检索): Grossman and Frieder 2004, 第 4 章; Baeza-Yates and Ribeiro-Neto 1999, 第 6、11 和 12 章; del Bimbo 1999; Lew 2001; Smeulders et al.2000。

*Speech retrieval* (语音检索): Coden et al.2002。

*Music retrieval* (音乐检索): Downie 2006 及网站 <http://www.ismir.net/>。

*User interfaces for IR* (信息检索中的用户界面): Baeza-Yates and Ribeiro-Neto 1999, 第 10 章。

*Parallel and peer-to-peer IR* (并行和 p2p 检索): Grossman and Frieder 2004, 第 7 章; Baeza-Yates and Ribeiro-Neto 1999, 第 9 章; Aberer 2001。

*Digital libraries* (数字图书馆): Baeza-Yates and Ribeiro-Neto 1999, 第 15 章; Lesk 2004。

*Information science perspective* (基于信息科学视角的信息检索): Korfhage 1997; Meadow et al.1999; Ingwersen and Järvelin 2005。

*Logic-based approaches to IR* (基于逻辑的信息检索): van Rijsbergen 1989。

*Natural language processing techniques* (自然语言处理技术): Manning and Schütze 1999; Jurafsky and Martin 2008; Lewis and Jones 1996。

---

## 预备知识

所有 21 章都需要数据结构和算法、线性代数以及概率论的基本知识。为方便读者和教师使用本书，各章更具体的要求如下。

第 1~5 章需要数据结构和算法的基本知识。第 6 章和第 7 章还另外需要线性代数的知识，包括向量和内积的基本概念。第 8~10 章不需要其他的预备知识。第 11 章需要概率论的基础知识，其中，11.1 节简单介绍了第 11~13 章所需要的基本概念。第 15 章假定读者熟悉非线性优化的基本概念，当然如果读者对非线性优化没有深入的了解，在阅读上也不会有太多的问题。第 18 章需要线性代数的基本知识，包括矩阵的秩、特征向量等概念，18.1 节对这些概念有一个简单的介绍。第 21 章还需要了解特征值和特征向量的概念。

---

## 书中的标记符号

本书正文中用铅笔符号 (✎) 标记例子，高级或者难度较大的章节用剪刀符号 (✂) 标记，习题用问号 (?) 标记，习题中分别用[\*]、[\*\*]、[\*\*\*]符号标识“容易”、“难度适中”和“难度大”的习题。

---

## 致谢

首先感谢剑桥大学出版社允许本书的电子样稿在网上公布，各种反馈促进了本书的写作过

程。感谢 Lauren Cowles，她是名出色的编辑，在本书样式、组织、覆盖面等方面提出了非常好的建议。如果说这本书最终达到了我们的写作预期的话，那么很大程度上应归功于她。

我们对那些在写作过程中对样稿提出建议和错误纠正意见的人们表示衷心的感谢。他们是 Cheryl Aasheim、Josh Attenberg、Luc Bélanger、Tom Breuel、Daniel Burckhardt、Georg Buscher、Fazli Can、Dinquan Chen、Ernest Davis、Pedro Domingos、Rodrigo Panchiniak Fernandes、Paolo Ferragina、Norbert Fuhr、Vignesh Ganapathy、Elmer Garduno、Xiubo Geng、David Gondek、Sergio Govoni、Corinna Habets、Ben Handy、Donna Harman、Benjamin Haskell、Thomas Hühn、Deepak Jain、Ralf Jankowitsch、Dinakar Jayarajan、Vinay Kakade、Mei Kobayashi、Wessel Kraaij、Rick Lafleur、Florian Laws、Hang Li、David Mann、Ennio Masi、Frank McCown、Paul McNamee、Sven Meyer zu Eissen、Alexander Murzaku、Gonzalo Navarro、Scott Olsson、Daniel Paiva、Tao Qin、Megha Raghava、Ghulam Raza、Michal Rosen-Zvi、Klaus Rothenhäusler、Kenyu L. Runner、Alexander Salamanca、Grigory Sapunov、Tobias Scheffer、Nico Schlaefer、Evgeny Shadchnev、Ian Soboroff、Benno Stein、Marcin Sydow、Andrew Turner、Jason Utt、Huey Vo、Travis Wade、Mike Walsh、Changliang Wang、Renjing Wang 及 Thomas Zeume。

很多人主动或应我们的要求对每个章节提出了非常细致的意见。就这点，我们要特别感谢如下这些人：James Allan、Omar Alonso、Ismail Sengor Altingovde、Vo NgocAnh、Roi Blanco、Eric Breck、Eric Brown、Mark Carman、Carlos Castillo、Junghoo Cho、Aron Culotta、Doug Cutting、Meghana Deodhar、Susan Dumais、Johannes Fürnkranz、Andreas Heß、Djoerd Hiemstra、David Hull、Thorsten Joachims、Siddharth Jonathan J. B.、Jaap Kamps、Mounia Lalmas、Amy Langville、Nicholas Lester、Dave Lewis、Stephen Liu、Daniel Lowd、Yosi Mass、Jeff Michels、Alessandro Moschitti、Amir Najmi、Marc Najork、Giorgio Maria Di Nunzio、Paul Ogilvie、Priyank Patel、Jan Pedersen、Kathryn Pedings、Vassilis Plachouras、Daniel Ramage、Stefan Riezler、Michael Schiehlen、Helmut Schmid、Falk Nicolas Scholer、Sabine Schulte im Walde、Fabrizio Sebastiani、Sarabjeet Singh、Alexander Strehl、John Tait、Shivakumar Vaithyanathan、Ellen Voorhees、Gerhard Weikum、Dawid Weiss、Yiming Yang、Yisong Yue、Jian Zhang 及 Justin Zobel。

最后，我们还要感谢书稿的审阅人员，他们为本书提出了大量高质量的建议。感谢他们为本书的内容和结构提出了重要建议，我们对他们表示深深的谢意，他们是：Pavel Berkhin、Stefan Büttcher、Jamie Callan、Byron Dom、Torsten Suel 及 Andrew Trotman。

第 13、14 和 15 章的部分初稿基于 Ray Mooney 慷慨提供的报告幻灯片。尽管后来的内容做了大量的修改，但是我们对 Ray Mooney 为这 3 章的贡献表示由衷的感谢，特别在各种分类器算法的时间复杂度分析方面，本书基本沿用了 Ray Mooney 的工作。

上述感谢的名单并不完整，我们仍然在不断整理来自各方面的反馈。当然，和其他作者一样，我们不一定留意到每一条建议，这点还请大家谅解。另外需要指出的是，本书的出版版本文责自负。

作者还要感谢斯坦福大学和斯图加特大学提供的优良的学术环境。在此环境中，大家可以自由交流思想，并通过授课来促进本书的写作和完善。C. D. Manning 感谢他的家人给他时间投入到本书的写作上，并希望明年能更多地利用周末的时间陪伴家人。P. Raghavan 感谢他的家人默默提供的支持，并感谢 Yahoo! 公司提供了一个良好的写作环境。H. Schütze 要感谢他的父母、家人和朋友在他写作期间给予他的支持。

## 网站和联系方式

---

与本书英文原版配套的网站地址是 <http://informationretrieval.org>。该网站不仅收录了多个相关资源的链接，还提供了每章的教学课件供大家下载使用。我们也欢迎大家将更多的反馈意见和建议发送至 [informationretrieval@yahogroups.com](mailto:informationretrieval@yahogroups.com)。

# 符号对照表

| 符 号                            | 原 书 | 含 义   |
|--------------------------------|-----|---|
| $\gamma$                       | 90  | $\gamma$ 编码   |
| $\gamma$                       | 237 | $\gamma(d)$ 表示分类或者聚类函数： $\gamma(d)$ 是 $d$ 所属的类或者簇                                     |
| $\Gamma$                       | 237 | 第13、14章中的有监督学习方法： $\Gamma(\mathbb{D})$ 是从训练集 $\mathbb{D}$ 上学到的分类函数 $\gamma$           |
| $\lambda$                      | 370 | 特征值   |
| $\bar{\mu}(\cdot)$             | 269 | 类质心（在Rocchio分类中）或簇质心（在 $K$ -均值和质心聚类中）   |
| $\Phi$                         | 105 | 训练样本  |
| $\sigma$                       | 374 | 奇异值   |
| $\Theta(\cdot)$                | 10  | 算法复杂度的紧上界   |
| $\omega, \omega_k$             | 328 | 聚类结果中的一个簇   |
| $\Omega$                       | 328 | 聚类结果或簇集合 $\{\omega_1, \dots, \omega_K\}$  |
| $\operatorname{argmax}_x f(x)$ | 164 | 使函数 $f$ 取最大值的 $x$ 的值  |
| $\operatorname{argmin}_x f(x)$ | 164 | 使函数 $f$ 取最小值的 $x$ 的值  |
| $c, c_j$                       | 237 | 分类中的一个类别  |
| $cf_i$                         | 82  | 词项 $i$ 的文档集频率（该词项在整个文档集中出现的总次数）   |
| $\mathbb{C}$                   | 237 | 类别集合 $\{c_1, \dots, c_J\}$  |
| $C$                            | 248 | 取值为类别集合 $\mathbb{C}$ 中元素的随机变量   |
| $C$                            | 369 | 词项-文档矩阵   |
| $d$                            | 4   | 文档集 $D$ 中的第 $d$ 篇文档的索引号   |
| $d$                            | 65  | 一篇文档  |
| $\vec{d}, \vec{q}$             | 163 | 文档向量及查询向量   |
| $D$                            | 326 | 所有文档的集合 $\{d_1, \dots, d_N\}$   |
| $D_c$                          | 269 | 类别 $c$ 中的文档集  |
| $\mathbb{D}$                   | 237 | 第13~15章中的已标记文档集 $\{\langle d_1, c_1 \rangle, \dots, \langle d_N, c_N \rangle\}$ ，即训练集 |
| $df_i$                         | 108 | 词项 $i$ 的文档频率（文档集中出现 $i$ 的文档数目）  |
| $H$                            | 91  | 熵   |
| $H_M$                          | 93  | 第 $M$ 个调和数  |
| $I(X; Y)$                      | 252 | 随机变量 $X$ 和 $Y$ 的互信息   |
| $\operatorname{idf}_i$         | 108 | 词项 $i$ 的逆文档频率   |
| $J$                            | 237 | 类别数目  |
| $k$                            | 267 | 集合中排名前 $k$ 的元素，如kNN中的前 $k$ 个邻居、检索文档的前 $k$ 个结果以及词汇表 $V$ 中选出的前 $k$ 个特征                  |
| $k$                            | 50  | $k$ 个字符组成的序列  |
| $K$                            | 326 | 簇的个数  |
| $L_d$                          | 214 | 文档 $d$ 的长度（以词条为单位计数）  |
| $L_a$                          | 242 | 测试文档或应用文档的长度（以词条为单位计数）  |
| $L_{\text{ave}}$               | 64  | 文档的平均长度（以词条为单位计数）   |

| 符 号                     | 原 书 | 含 义  |
|-------------------------|-----|--|
| $M$                     | 4   | 词汇表大小 (即 $ V $ )   |
| $M_a$                   | 242 | 测试文档或应用文档的词汇量  |
| $M_{ave}$               | 71  | 文档集中每篇文档的平均词汇量   |
| $M_d$                   | 218 | 文档 $d$ 的模型   |
| $N$                     | 4   | 检索或训练文档集中的文档数目   |
| $N_c$                   | 240 | 类别 $c$ 中的文档数目  |
| $N(\omega)$             | 275 | 事件 $\omega$ 发生的次数  |
| $O(\cdot)$              | 10  | 算法复杂度的界  |
| $O(\cdot)$              | 203 | 事件的优势率   |
| $P$                     | 142 | 正确率  |
| $P(\cdot)$              | 202 | 概率   |
| $P$                     | 425 | 转移概率矩阵   |
| $q$                     | 55  | 查询   |
| $R$                     | 143 | 召回率  |
| $s_i$                   | 53  | 字符串  |
| $s_i$                   | 103 | 域评分布尔值   |
| $\text{sim}(d_1, d_2)$  | 111 | 文档 $d_1$ 和 $d_2$ 的相似度  |
| $T$                     | 40  | 文档集中所有词条的数目  |
| $T_c$                   | 240 | 词 $t$ 在 $c$ 类文档中的出现次数  |
| $t$                     | 4   | 词汇表 $V$ 中第 $t$ 个词项的索引号   |
| $t$                     | 56  | 词汇表中的一个词项  |
| $\text{tf}_{t,d}$       | 107 | 词项 $t$ 在文档 $d$ 中的出现频率 (即 $t$ 在 $d$ 中的出现次数)                     |
| $U_t$                   | 246 | 表示词项 $t$ 存在与否的随机变量, 当 $t$ 存在时, 值为1, 否则为0                       |
| $V$                     | 190 | 文档中的所有词项 $\{t_1, \dots, t_M\}$ 组成的词汇表 (也称为词典lexicon)           |
| $\bar{v}(d)$            | 111 | 文档 $d$ 经长度归一化后的文档向量  |
| $\tilde{v}(d)$          | 110 | 文档 $d$ 未经长度归一化的文档向量  |
| $w_{t,d}$               | 115 | 词项 $t$ 在文档 $d$ 中的权重  |
| $w$                     | 103 | 权重, 比如域的权重或者词项的权重  |
| $\bar{w}^T \bar{x} = b$ | 269 | 超平面方程: $\bar{w}$ 是超平面的法向量, $w_i$ 是 $\bar{w}$ 的第 $i$ 个分量        |
| $\bar{x}$               | 204 | 基于词项表示的文档向量 $\bar{x} = (x_1, \dots, x_M)$ , 更一般地说, 为文档的特征表示    |
| $X$                     | 246 | 取值为词汇表 $V$ 中元素的随机变量 (比如, 某个文档位置 $k$ 上的词)                       |
| $\mathbb{X}$            | 237 | 文本分类中的文档空间   |
| $ A $                   | 56  | 集合 $A$ 的势: 集合 $A$ 中的元素个数                                       |
| $ S $                   | 570 | 方阵 $S$ 的行列式  |
| $ s_i $                 | 53  | $s_i$ 的长度 (以字符计)   |
| $ \bar{x} $             | 128 | 向量 $\bar{x}$ 的大小   |
| $ \bar{x} - \bar{y} $   | 121 | 向量 $\bar{x}$ 、 $\bar{y}$ 的欧氏距离, 也即向量 $(\bar{x} - \bar{y})$ 的大小 |

# 目 录

|                             |    |
|-----------------------------|----|
| 第 1 章 布尔检索                  | 1  |
| 1.1 一个信息检索的例子               | 2  |
| 1.2 构建倒排索引的初体验              | 5  |
| 1.3 布尔查询的处理                 | 8  |
| 1.4 对基本布尔操作的扩展及有序检索         | 11 |
| 1.5 参考文献及补充读物               | 13 |
| 第 2 章 词项词典及倒排记录表            | 14 |
| 2.1 文档分析及编码转换               | 14 |
| 2.1.1 字符序列的生成               | 14 |
| 2.1.2 文档单位的选择               | 16 |
| 2.2 词项集合的确定                 | 16 |
| 2.2.1 词条化                   | 16 |
| 2.2.2 去除停用词                 | 19 |
| 2.2.3 词项归一化                 | 20 |
| 2.2.4 词干还原和词形归并             | 23 |
| 2.3 基于跳表的倒排记录表快速合并算法        | 26 |
| 2.4 含位置信息的倒排记录表及短语查询        | 28 |
| 2.4.1 二元词索引                 | 28 |
| 2.4.2 位置信息索引                | 29 |
| 2.4.3 混合索引机制                | 31 |
| 2.5 参考文献及补充读物               | 32 |
| 第 3 章 词典及容错式检索              | 34 |
| 3.1 词典搜索的数据结构               | 34 |
| 3.2 通配符查询                   | 36 |
| 3.2.1 一般的通配符查询              | 37 |
| 3.2.2 支持通配符查询的 $k$ -gram 索引 | 38 |
| 3.3 拼写校正                    | 39 |
| 3.3.1 拼写校正的实现               | 39 |
| 3.3.2 拼写校正的方法               | 40 |
| 3.3.3 编辑距离                  | 40 |
| 3.3.4 拼写校正中的 $k$ -gram 索引   | 42 |
| 3.3.5 上下文敏感的拼写校正            | 43 |
| 3.4 基于发音的校正技术               | 44 |
| 3.5 参考文献及补充读物               | 45 |
| 第 4 章 索引构建                  | 46 |
| 4.1 硬件基础                    | 46 |
| 4.2 基于块的排序索引方法              | 47 |
| 4.3 内存式单遍扫描索引构建方法           | 50 |
| 4.4 分布式索引构建方法               | 51 |
| 4.5 动态索引构建方法                | 54 |
| 4.6 其他索引类型                  | 56 |
| 4.7 参考文献及补充读物               | 57 |
| 第 5 章 索引压缩                  | 59 |
| 5.1 信息检索中词项的统计特性            | 59 |
| 5.1.1 Heaps 定律: 词项数目的估计     | 61 |
| 5.1.2 Zipf 定律: 对词项的分布建模     | 62 |
| 5.2 词典压缩                    | 63 |
| 5.2.1 将词典看成单一字符串的压缩方法       | 63 |
| 5.2.2 按块存储                  | 64 |
| 5.3 倒排记录表的压缩                | 66 |
| 5.3.1 可变字节码                 | 67 |
| 5.3.2 $\gamma$ 编码           | 68 |
| 5.4 参考文献及补充读物               | 74 |
| 第 6 章 文档评分、词项权重计算及向量空间模型    | 76 |
| 6.1 参数化索引及域索引               | 76 |
| 6.1.1 域加权评分                 | 78 |
| 6.1.2 权重学习                  | 79 |
| 6.1.3 最优权重 $g$ 的计算          | 80 |
| 6.2 词项频率及权重计算               | 81 |
| 6.2.1 逆文档频率                 | 81 |
| 6.2.2 tf-idf 权重计算           | 82 |
| 6.3 向量空间模型                  | 83 |
| 6.3.1 内积                    | 83 |
| 6.3.2 查询向量                  | 86 |

|                           |     |                             |     |
|---------------------------|-----|-----------------------------|-----|
| 6.3.3 向量相似度计算             | 87  | 8.8 参考文献及补充读物               | 118 |
| 6.4 其他tf-idf权重计算方法        | 88  | <b>第9章 相关反馈及查询扩展</b>        | 120 |
| 6.4.1 tf的亚线性尺度变换方法        | 88  | 9.1 相关反馈及伪相关反馈              | 120 |
| 6.4.2 基于最大值的tf归一化         | 88  | 9.1.1 Rocchio相关反馈算法         | 122 |
| 6.4.3 文档权重和查询权重机制         | 89  | 9.1.2 基于概率的相关反馈方法           | 125 |
| 6.4.4 文档长度的回转归一化          | 89  | 9.1.3 相关反馈的作用时机             | 125 |
| 6.5 参考文献及补充读物             | 92  | 9.1.4 Web上的相关反馈             | 126 |
| <b>第7章 一个完整搜索系统中的评分计算</b> | 93  | 9.1.5 相关反馈策略的评价             | 127 |
| 7.1 快速评分及排序               | 93  | 9.1.6 伪相关反馈                 | 127 |
| 7.1.1 非精确返回前K篇文档的方法       | 94  | 9.1.7 间接相关反馈                | 128 |
| 7.1.2 索引去除技术              | 94  | 9.1.8 小结                    | 128 |
| 7.1.3 胜者表                 | 95  | 9.2 查询重构的全局方法               | 128 |
| 7.1.4 静态得分和排序             | 95  | 9.2.1 查询重构的词汇表工具            | 128 |
| 7.1.5 影响度排序               | 96  | 9.2.2 查询扩展                  | 129 |
| 7.1.6 簇剪枝方法               | 97  | 9.2.3 同义词词典的自动构建            | 130 |
| 7.2 信息检索系统的组成             | 98  | 9.3 参考文献及补充读物               | 131 |
| 7.2.1 层次型索引               | 98  | <b>第10章 XML检索</b>           | 133 |
| 7.2.2 查询词项的邻近性            | 98  | 10.1 XML的基本概念               | 134 |
| 7.2.3 查询分析及文档评分函数的设计      | 99  | 10.2 XML检索中的挑战性问题           | 137 |
| 7.2.4 搜索系统的组成             | 100 | 10.3 基于向量空间模型的XML检索         | 140 |
| 7.3 向量空间模型对各种查询操作的支持      | 101 | 10.4 XML检索的评价               | 144 |
| 7.3.1 布尔查询                | 101 | 10.5 XML检索:以文本为中心与以数据为中心的对比 | 146 |
| 7.3.2 通配符查询               | 102 | 10.6 参考文献及补充读物              | 148 |
| 7.3.3 短语查询                | 102 | <b>第11章 概率检索模型</b>          | 150 |
| 7.4 参考文献及补充读物             | 102 | 11.1 概率论基础知识                | 150 |
| <b>第8章 信息检索的评价</b>        | 103 | 11.2 概率排序原理                 | 151 |
| 8.1 信息检索系统的评价             | 103 | 11.2.1 1/0风险的情况             | 151 |
| 8.2 标准测试集                 | 104 | 11.2.2 基于检索代价的概率排序原理        | 152 |
| 8.3 无序检索结果集合的评价           | 105 | 11.3 二值独立模型                 | 152 |
| 8.4 有序检索结果的评价方法           | 108 | 11.3.1 排序函数的推导              | 153 |
| 8.5 相关性判定                 | 112 | 11.3.2 理论上的概率估计方法           | 155 |
| 8.6 更广的视角看评价:系统质量及用户效用    | 115 | 11.3.3 实际中的概率估计方法           | 156 |
| 8.6.1 系统相关问题              | 115 | 11.3.4 基于概率的相关反馈方法          | 157 |
| 8.6.2 用户效用                | 115 | 11.4 概率模型的相关评论及扩展           | 158 |
| 8.6.3 对已有系统的改进            | 116 | 11.4.1 概率模型的评论              | 158 |
| 8.7 结果片段                  | 116 | 11.4.2 词项之间的树型依赖            | 159 |
|                           |     | 11.4.3 Okapi BM25:一个非二值的模型  | 160 |



|  |     |   |     |
|--|-----|---|-----|
| 11.4.4 IR中的贝叶斯网络<br>方法 .....           | 161 | <b>第 15 章 支持向量机及文档机器学习<br/>方法</b> ..... | 221 |
| 11.5 参考文献及补充读物 .....                   | 162 | 15.1 二类线性可分条件下的支持向量机 .....              | 221 |
| <b>第 12 章 基于语言建模的信息检索<br/>模型</b> ..... | 163 | 15.2 支持向量机的扩展 .....                     | 226 |
| 12.1 语言模型 .....                        | 163 | 15.2.1 软间隔分类 .....                      | 226 |
| 12.1.1 有穷自动机和语言模型 .....                | 163 | 15.2.2 多类情况下的支持向量机 .....                | 228 |
| 12.1.2 语言模型的种类 .....                   | 165 | 15.2.3 非线性支持向量机 .....                   | 228 |
| 12.1.3 词的多项式分布 .....                   | 166 | 15.2.4 实验结果 .....                       | 230 |
| 12.2 查询似然模型 .....                      | 167 | 15.3 有关文本文档分类的考虑 .....                  | 231 |
| 12.2.1 IR中的查询似然模型 .....                | 167 | 15.3.1 分类器类型的选择 .....                   | 231 |
| 12.2.2 查询生成概率的估计 .....                 | 167 | 15.3.2 分类器效果的提高 .....                   | 233 |
| 12.2.3 Ponte和Croft进行的实验 .....          | 169 | 15.4 ad hoc检索中的机器学习方法 .....             | 236 |
| 12.3 语言建模的方法与其他检索方法的<br>比较 .....       | 171 | 15.4.1 基于机器学习评分的简单<br>例子 .....          | 236 |
| 12.4 扩展的LM方法 .....                     | 172 | 15.4.2 基于机器学习的检索结果<br>排序 .....          | 238 |
| 12.5 参考文献及补充读物 .....                   | 173 | 15.5 参考文献及补充读物 .....                    | 239 |
| <b>第 13 章 文本分类及朴素贝叶斯方法</b> .....       | 175 | <b>第 16 章 扁平聚类</b> .....                | 241 |
| 13.1 文本分类问题 .....                      | 177 | 16.1 信息检索中的聚类应用 .....                   | 242 |
| 13.2 朴素贝叶斯文本分类 .....                   | 178 | 16.2 问题描述 .....                         | 244 |
| 13.3 贝努利模型 .....                       | 182 | 16.3 聚类算法的评价 .....                      | 246 |
| 13.4 NB的性质 .....                       | 183 | 16.4 K-均值算法 .....                       | 248 |
| 13.5 特征选择 .....                        | 188 | 16.5 基于模型的聚类 .....                      | 254 |
| 13.5.1 互信息 .....                       | 188 | 16.6 参考文献及补充读物 .....                    | 258 |
| 13.5.2 $\chi^2$ 统计量 .....              | 191 | <b>第 17 章 层次聚类</b> .....                | 260 |
| 13.5.3 基于频率的特征选择方法 .....               | 192 | 17.1 凝聚式层次聚类 .....                      | 260 |
| 13.5.4 多类问题的特征选择方法 .....               | 193 | 17.2 单连接及全连接聚类算法 .....                  | 263 |
| 13.5.5 不同特征选择方法的比较 .....               | 193 | 17.3 组平均凝聚式聚类 .....                     | 268 |
| 13.6 文本分类的评价 .....                     | 194 | 17.4 质心聚类 .....                         | 269 |
| 13.7 参考文献及补充读物 .....                   | 199 | 17.5 层次凝聚式聚类的最优性 .....                  | 270 |
| <b>第 14 章 基于向量空间模型的文本<br/>分类</b> ..... | 200 | 17.6 分裂式聚类 .....                        | 272 |
| 14.1 文档表示及向量空间中的关联度<br>计算 .....        | 201 | 17.7 簇标签生成 .....                        | 273 |
| 14.2 Rocchio分类方法 .....                 | 202 | 17.8 实施中的注意事项 .....                     | 274 |
| 14.3 k近邻分类器 .....                      | 205 | 17.9 参考文献及补充读物 .....                    | 275 |
| 14.4 线性及非线性分类器 .....                   | 209 | <b>第 18 章 矩阵分解及隐性语义索引</b> .....         | 277 |
| 14.5 多类问题的分类 .....                     | 212 | 18.1 线性代数基础 .....                       | 277 |
| 14.6 偏差-方差折中准则 .....                   | 214 | 18.2 词项-文档矩阵及SVD .....                  | 280 |
| 14.7 参考文献及补充读物 .....                   | 219 | 18.3 低秩逼近 .....                         | 282 |
|  |     | 18.4 LSI .....                          | 284 |
|  |     | 18.5 参考文献及补充读物 .....                    | 287 |