经

典

原

版

书

库

现代信息检索

(英文版)

Modern Information Retrieval

Ricardo Baeza-Yates Berthier Ribeiro-Neto



Ricardo Baeza-Yates Berthier Ribeiro-Neto

等著

现代信息检索

(英文版)

Modern Information Retrieval

本书介绍了有关信息检索方面的所有新变化,而且其组织(包括支持本书的主页 www.dcc.ufmg.br/irbook)使读者既可以对现代信息检索有一个全面的了解,又可以获取现代信息检索所有关键主题的详细知识。本书的主要内容由信息检索领域的代表人物 Baeza-Yates 和 Ribeiro-Neto 编写,对于那些希望深入研究关键领域的读者,书中还提供了由其他主要研究人员编写的关于特殊主题的发展现状的内容:

- 并行和分布式信息检索——算法和体系结构。
- 用户界面和可视化──查询组织和结果可视化的主要界面范型。
- 多媒体信息检索:模型与语言——包括 MULTOS 和 SQL3。
- 索引和搜索——R树、GEMINI和QBIC。
- 图书馆和图书目录系统——在线系统和公共访问目录。
- 数字图书馆——有效部署面临的主要挑战。
- 文本信息检索──所有主要的信息检索模型、查询操作、文本操作、索引和搜索。
- Web 挑战、方法和模型、搜索引擎、目录、查询语言、元搜索及趋势。

本书可以作为信息检索专业必修课程及相关专业研究生课程的教材。同时,本书对于计算机科学、信息科学和图书馆科学专业的学生,以及相关产品的程序员及分析人员,也是非常有价值的。

作者简介

Ricardo Baeza-Yates于加拿大滑铁卢大学获得计算机科学博士学位。曾担任智利计算机科学学会 (SCCC) 主席,现任智利大学计算机科学系全职教授,同时也是世界上多所大学的客座教授,还是ACM、AMS、EATCS、IEEE、SCCC及 SIAM 会员。他的主要研究方向为算法与数据结构、文本检索、图形界面以及可视化在数据库中的应用。

Berthier Ribeiro-Neto 于加利福尼亚大学洛杉矶分校获得计算机科学博士学位。现任巴西 Minas Gerais 联合大学计算机科学系副教授,同时也是 ACM、ASIS 及 IEEE 会员。他的主要研究方向是信息检索系统、数字图书馆、Web 界面及视频点播。





☑ 网上购书: www.china-pub.com

北京市西城区百万庄南街 1 号 100037 读者服务热线: (010)68995259, 68995264 读者服务信箱: hzedu@hzbook.com http://www.hzbook.com

限中国大陆地区销售

ISBN 7-111-13704-3/TP · 3402 定价: 49.00 元

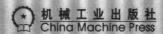
现代信息检索

(英文版)

Modern Information Retrieval

Ricardo Baeza-Yates Berthier Ribeiro-Neto

等著



Ricardo Baeza-Yates, Berthier Ribeiro-Neto, et al.: Modern Information Retrieval (ISBN: 0-201-39829-X).

Copyright © 1999 by the ACM press, A Division of the Association for Computing Machinary, Inc. (ACM).

This edition of Modern Information Retrieval is published by arrangement with Pearson Education Limited. Licensed for sale in the mainland territory of the People's Republic of China only, excluding Hong Kong, Macau, and Taiwan.

本书英文影印版由英国Pearson Education培生教育出版集团授权出版。未经出版者书面许可,不得以任何方式复制或抄袭本书内容。

此影印版只限在中国大陆地区销售(不包括香港、澳门、台湾地区)版权所有,侵权必究。

本书版权登记号:图字:01-2003-8980

图书在版编目(CIP)数据

现代信息检索(英文版)/ 巴伊赞-耶茨(Baeza-Yates, R.)等著.-北京: 机械工业出版社, 2004.2

(经典原版书库)

书名原文: Modern Information Retrieval ISBN 7-111-13704-3

I. 现… II. 巴… III. 情报检索 - 英文 IV. G252.7

中国版本图书馆CIP数据核字(2003)第121060号

机械工业出版社(北京市西域区百万庄大街22号 邮政编码 100037)

责任编辑: 迟振春

北京中加印刷有限公司印刷・新华书店北京发行所发行

2004年2月第1版第1次印刷

787mm×1092mm 1/16 · 33.75 印张

印数: 0001-3000册

定价: 49.00元

凡购本书,如有倒页、脱页、缺页,由本社发行部调换本社购书热线:(010)68326294

出版者的话

文艺复兴以降,源远流长的科学精神和逐步形成的学术规范,使西方国家在自然科学的各个领域取得了垄断性的优势;也正是这样的传统,使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中,美国的产业界与教育界越来越紧密地结合,计算机学科中的许多泰山北斗同时身处科研和教学的最前线,由此而产生的经典科学著作,不仅擘划了研究的范畴,还揭橥了学术的源变,既遵循学术规范,又自有学者个性,其价值并不会因年月的流逝而减退。

近年,在全球信息化大潮的推动下,我国的计算机产业发展迅猛,对专业人才的需求日益 迫切。这对计算机教育界和出版界都既是机遇,也是挑战;而专业教材的建设在教育战略上显 得举足轻重。在我国信息技术发展时间较短、从业人员较少的现状下,美国等发达国家在其计 算机科学发展的几十年间积淀的经典教材仍有许多值得借鉴之处。因此,引进一批国外优秀计 算机教材将对我国计算机教育事业的发展起积极的推动作用,也是与世界接轨、建设真正的世 界一流大学的必由之路。

机械工业出版社华章图文信息有限公司较早意识到"出版要为教育服务"。自1998年开始,华章公司就将工作重点放在了遴选、移译国外优秀教材上。经过几年的不懈努力,我们与Prentice Hall, Addison-Wesley, McGraw-Hill, Morgan Kaufmann等世界著名出版公司建立了良好的合作关系,从它们现有的数百种教材中甄选出Tanenbaum, Stroustrup, Kernighan, Jim Gray等大师名家的一批经典作品,以"计算机科学从书"为总称出版,供读者学习、研究及庋藏。大理石纹理的封面,也正体现了这套从书的品位和格调。

"计算机科学丛书"的出版工作得到了国内外学者的鼎力襄助,国内的专家不仅提供了中肯的选题指导,还不辞劳苦地担任了翻译和审校的工作;而原书的作者也相当关注其作品在中国的传播,有的还专诚为其书的中译本作序。迄今,"计算机科学丛书"已经出版了近百个品种,这些书籍在读者中树立了良好的口碑,并被许多高校采用为正式教材和参考书籍,为进一步推广与发展打下了坚实的基础。

随着学科建设的初步完善和教材改革的逐渐深化,教育界对国外计算机教材的需求和应用都步入一个新的阶段。为此,华章公司将加大引进教材的力度,在"华章教育"的总规划之下出版三个系列的计算机教材:除"计算机科学从书"之外,对影印版的教材,则单独开辟出"经典原版书库";同时,引进全美通行的教学辅导书"Schaum's Outlines"系列组成"全美经典学习指导系列"。为了保证这三套丛书的权威性,同时也为了更好地为学校和老师们服务,华章公司聘请了中国科学院、北京大学、清华大学、国防科技大学、复旦大学、上海交通大学、南京大学、浙江大学、中国科技大学、哈尔滨工业大学、西安交通大学、中国人民大学、北京航空航天大学、北京邮电大学、中山大学、解放军理工大学、郑州大学、湖北工学院、中国国

家信息安全测评认证中心等国内重点大学和科研机构在计算机的各个领域的著名学者组成"专家指导委员会",为我们提供选题意见和出版监督。

这三套丛书是响应教育部提出的使用外版教材的号召,为国内高校的计算机及相关专业的教学度身订造的。其中许多教材均已为M. I. T., Stanford, U.C. Berkeley, C. M. U. 等世界名牌大学所采用。不仅涵盖了程序设计、数据结构、操作系统、计算机体系结构、数据库、编译原理、软件工程、图形学、通信与网络、离散数学等国内大学计算机专业普遍开设的核心课程,而且各具特色——有的出自语言设计者之手、有的历经三十年而不衰、有的已被全世界的几百所高校采用。在这些圆熟通博的名师大作的指引之下,读者必将在计算机科学的宫殿中由登堂而入室。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑,这些因素使我们的图书有了质量的保证,但我们的目标是尽善尽美,而反馈的意见正是我们达到这一终极目标的重要帮助。教材的出版只是我们的后续服务的起点。华章公司欢迎老师和读者对我们的工作提出建议或给予指正,我们的联系方法如下:

电子邮件: hzedu@hzbook.com 联系电话: (010) 68995264

联系地址:北京市西城区百万庄南街1号

邮政编码: 100037

专家指导委员会

(按姓氏笔画顺序)

尤晋元	王 珊	冯博琴	史忠植	史美林
石教英	吕 建	孙玉芳	吴世忠	吴时霖
张立昂	李伟琴	李师贤	李建中	杨冬青
邵维忠	陆丽娜	陆鑫达	陈向群	周伯生
周立柱	周克定	周傲英	孟小峰	岳丽华
范 明	郑国梁	施伯.乐	钟玉琢	唐世渭
袁崇义	高传善	梅宏	程 旭	程时端
谢希仁	裘宗燕	戴葵		

秘 书 组

武卫东 温莉芳 刘 江 杨海玲

Preface

Information retrieval (IR) has changed considerably in recent years with the expansion of the World Wide Web and the advent of modern and inexpensive graphical user interfaces and mass storage devices. As a result, traditional IR textbooks have become quite out of date and this has led to the introduction of new IR books. Nevertheless, we believe that there is still great need for a book that approaches the field in a rigorous and complete way from a computer-science perspective (as opposed to a user-centered perspective). This book is an effort to partially fulfill this gap and should be useful for a first course on information retrieval as well as for a graduate course on the topic.

The book comprises two portions which complement and balance each other. The core portion includes nine chapters authored or coauthored by the designers of the book. The second portion, which is fully integrated with the first, is formed by six state-of-the-art chapters written by leading researchers in their fields. The same notation and glossary are used in all the chapters. Thus, despite the fact that several people have contributed to the text, this book is really much more a textbook than an edited collection of chapters written by separate authors. Furthermore, unlike a collection of chapters, we have carefully designed the contents and organization of the book to present a cohesive view of all the important aspects of modern information retrieval.

From IR models to indexing text, from IR visual tools and interfaces to the Web, from IR multimedia to digital libraries, the book provides both breadth of coverage and richness of detail. It is our hope that, given the now clear relevance and significance of information retrieval to modern society, the book will contribute to further disseminate the study of the discipline at information science, computer science, and library science departments throughout the world.

To Helena, Rosa, and our children

Amo los libros exploradores, libros con bosque o nieve, profundidad o cielo

de Oda al Libro (I),

Pablo Neruda

I love books that explore, books with a forest or snow, depth or sky

from Ode to the Book (I),

Pablo Neruda

território de homens livres que será nosso país e será pátria de todos. Irmãos, cantai ese mundo que não verei, mas virá um dia, dentro de mil anos, talvez mais... não tenho pressa.

de Cidade Prevista no livro A Rosa do Povo, 1945

Carlos Drummond de Andrade

territory of free men that will be our country and will be the nation of all Brothers, sing this world which I'll not see, but which will come one day, in a thousand years, maybe more...no hurry.

from Prevised City in the book The Rose of the People, 1945

Carlos Drummond de Andrade

Acknowledgements

We would like to deeply thank the various people who, during the several months in which this endeavor lasted, provided us with useful and helpful assistance. Without their care and consideration, this book would likely not have matured.

First, we would like to thank all the chapter contributors, for their dedication and interest. To Elisa Bertino, Eric Brown, Barbara Catania, Christos Faloutsos, Elena Ferrari, Ed Fox, Marti Hearst, Gonzalo Navarro, Edie Rasmussen, Ohm Sornil, and Nivio Ziviani, who contributed with writings that reflect expertise we certainly do not fully profess ourselves. And for all their patience throughout an editing and cross-reviewing process which constitutes a rather difficult balancing act.

Second, we would like to thank all the people who demonstrated interest in publishing this book, particularly Scott Delman and Doug Sery.

Third, we would like to commend the interest, encouragement, and great job done by Addison Wesley Longman throughout the overall process, represented by Keith Mansfield, Karen Sutherland, Bridget Allen, David Harison, Sheila Chatten, Helen Hodge and Lisa Talbot. The reviewers they contacted read an early (and rather preliminary) proposal of this book and provided us with good feedback and invaluable insights. The chapter on Parallel and Distributed IR was moved from the part on Applications of IR (where it did not fit well) to the part on Text IR due to the objective argument of an unknown referee. A separate chapter on Retrieval Evaluation was only included after another zealous referee strongly made the case for the importance of this subject.

Fourth, we would like to thank all the people who discussed this project with us. Doug Oard provided us with an early critique of the proposal. Gary Marchionini was an earlier supporter and provided us with useful contacts during the process. Bruce Croft encouraged our efforts from the beginning. Alberto Mendelzon provided us with an initial proposal and a compilation of references for the chapter on searching the Web. Ed Fox found time in a rather busy schedule to provide us with an insightful review of the introduction (which resulted in a great improvement) and a thorough review of the chapter on Modeling. Marti Hearst demonstrated interest in our proposal early on, provided assistance throughout the editing process, and has been an enthusiastic supporter and partner.

x ACKNOWLEDGEMENTS

Fifth, we thank the support of our institutions, the Departments of Computer Science of the University of Chile and of the Federal University of Minas Gerais, as well as the funding provided by national research agencies (CNPq in Brazil and CONICYT in Chile) and international collaboration projects, in particular CYTED project VII.13 AMYRI (Environment for Information Managing and Retrieval in the World Wide Web) and Finep project SIAM (Information Systems for Mobile Computers) under the Pronex program.

Most important, to Helena, Rosa, and our children, who put up with a string of trips abroad, lost weekends, and odd working hours.

List of Trademarks

Alta Vista is a trademark of Compaq Computer Corporation
FrameMaker is a trademark of Adobe Systems Incorporated
IBM SP2 is a trademark of International Business Machines Corporation
Netscape Communicator is a trademark of Netscape Communications Corporation
Solaris, Sun 3/50 and Sun UltraSparc-1 are trademarks of Sun Microsystems, Inc.
Thinking Machines CM-2 is a trademark of Thinking Machines Corporation
Unix is licensed through X/Open Company Ltd
Word is a trademark of Microsoft Corporation
WordPerfect is a trademark of of Corel Corporation

Biographies

Biographies of Main Authors

Ricardo Baeza-Yates received a bachelor degree in Computer Science in 1983 from the University of Chile. Later, he received an MSc in Computer Science (1985), a professional title in electrical engineering (1985), and an MEng in EE (1986) from the same university. He received his PhD in Computer Science from the University of Waterloo, Canada, in 1989. He has been the president of the Chilean Computer Science Society (SCCC) from 1992 to 1995 and from 1997 to 1998. During 1993, he received the Organization of the American States award for young researchers in exact sciences. Currently, he is a full professor at the Computer Science Department of the University of Chile, where he was the chairperson in the period 1993 to 1995. He is coauthor of the second edition of the Handbook of Algorithms and Data Structures, Addison-Wesley, 1991; and coeditor of Information Retrieval: Algorithms and Data Structures, Prentice Hall, 1992. He has also contributed several papers to journals published by professional organizations such as ACM, IEEE, and SIAM.

His research interests include algorithms and data structures, text retrieval, graphical interfaces, and visualization applied to databases. He currently coordinates an IberoAmerican project on models and techniques for searching the Web financed by the Spanish agency Cyted. He has been a visiting professor or an invited speaker at several conferences and universities around the world, as well as referee for several journals, conferences, NSF, etc. He is a member of the ACM, AMS, EATCS, IEEE, SCCC, and SIAM.

Berthier Ribeiro-Neto received a bachelor degree in Math, a BS degree in Electrical Engineering, and an MS degree in Computer Science, all from the Federal University of Minas Gerais, Brazil. In 1995, he was awarded a Ph.D. in Computer Science from the University of California at Los Angeles. Since then, he has been with the Computer Science Department of the Federal University of Minas Gerais where he is an Associate Professor.

His main interests are information retrieval systems, digital libraries, interfaces for the Web, and video on demand. He has been involved in a number

of research projects financed through Brazilian national agencies such as the Ministry of Science and Technology (MCT) and the National Research Council (CNPq). From the projects currently underway, the two main ones deal with wireless information systems (project SIAM financed within program PRONEX) and video on demand (project ALMADEM financed within program PROTEM III). Dr Ribeiro-Neto is also involved with an IberoAmerican project on information systems for the Web coordinated by Professor Ricardo Baeza-Yates. He was the chair of SPIRE'98 (String Processing and Information Retrieval South American Symposium), is the chair of SBBD'99 (Brazilian Symposium on Databases), and has been on the committee of several conferences in Brazil, in South America and in the USA. He is a member of ACM, ASIS, and IEEE.

Biographies of Contributors

Elisa Bertino is Professor of Computer Science in the Department of Computer Science of the University of Milano where she heads the Database Systems Group. She has been a visiting researcher at the IBM Research Laboratory (now Almaden) in San Jose, at the Microelectronics and Computer Technology Corporation in Austin, Texas, and at Rutgers University in Newark, New Jersey. Her main research interests include object-oriented databases, distributed databases, deductive databases, multimedia databases, interoperability of heterogeneous systems, integration of artificial intelligence and database techniques, and database security. In those areas, Professor Bertino has published several papers in refereed journals, and in proceedings of international conferences and symposia. She is a coauthor of the books Object-Oriented Database Systems - Concepts and Architectures, Addison-Wesley 1993; Indexing Techniques for Advanced Database Systems, Kluwer 1997; and Intelligent Database Systems, Addison-Wesley forthcoming. She is or has been on the editorial boards of the following scientific journals: the IEEE Transactions on Knowledge and Data Engineering, the International Journal of Theory and Practice of Object Systems, the Very Large Database Systems (VLDB) Journal, the Parallel and Distributed Database Journal, the Journal of Computer Security, Data & Knowledge Engineering, and the International Journal of Information Technology.

Eric Brown has been a Research Staff Member at the IBM T.J. Watson Research Center in Yorktown Heights, NY, since 1995. Prior to that he was a Research Assistant at the Center for Intelligent Information Retrieval at the University of Massachusetts, Amherst. He holds a BSc from the University of Vermont and an MS and PhD from the University of Massachusetts, Amherst. Dr. Brown conducts research in large scale information retrieval systems, automatic text categorization, and hypermedia systems for digital libraries and knowledge management. He has published a number of papers in the field of information retrieval.

Barbara Catania is a researcher at the University of Milano, Italy. She received an MS degree in Information Sciences in 1993 from the University of

Genova and a PhD in Computer Science in 1998 from the University of Milano. She has also been a visiting researcher at the European Computer-Industry Research Center, Munich, Germany. Her main research interests include multimedia databases, constraint databases, deductive databases, and indexing techniques in object-oriented and constraint databases. In those areas, Dr Catania has published several papers in refereed journals, and in proceedings of international conferences and symposia. She is also a coauthor of the book *Indexing Techniques for Advanced Database Systems*, Kluwer 1997.

Christos Faloutsos received a BSc in Electrical Engineering (1981) from the National Technical University of Athens, Greece and an MSc and PhD in Computer Science from the University of Toronto, Canada. Professor Faloutsos is currently a faculty member at Carnegie Mellon University. Prior to joining CMU he was on the faculty of the Department of Computer Science at the University of Maryland, College Park. He has spent sabbaticals at IBM-Almaden and AT&T Bell Labs. He received the Presidential Young Investigator Award from the National Science Foundation in 1989, two 'best paper' awards (SIGMOD 94, VLDB 97), and three teaching awards. He has published over 70 refereed articles and one monograph, and has filed for three patents. His research interests include physical database design, searching methods for text, geographic information systems, indexing methods for multimedia databases, and data mining.

Elena Ferrari is an Assistant Professor at the Computer Science Department of the University of Milano, Italy. She received an MS in Information Sciences in 1992 and a PhD in Computer Science in 1998 from the University of Milano. Her main research interests include multimedia databases, temporal object-oriented data models, and database security. In those areas, Dr Ferrari has published several papers in refereed journals, and in proceedings of international conferences and symposia. She has been a visiting researcher at George Mason University in Fairfax, Virginia, and at Rutgers University in Newark, New Jersey.

Dr Edward A. Fox holds a PhD and MS in Computer Science from Cornell University, and a BS from MIT. Since 1983 he has been at Virginia Polytechnic Institute and State University (Virginia Tech), where he serves as Associate Director for Research at the Computing Center, Professor of Computer Science, Director of the Digital Library Research Laboratory, and Director of the Internet Technology Innovation Center. He served as vice chair and chair of ACM SIGIR from 1987 to 1995, helped found the ACM conferences on multimedia and digital libraries, and serves on a number of editorial boards. His research is focused on digital libraries, multimedia, information retrieval, WWW/Internet, educational technologies, and related areas.

Marti Hearst is an Assistant Professor at the University of California Berkeley in the School of Information Management and Systems. From 1994 to 1997 she was a Member of the Research Staff at Xerox PARC. She received her BA, MS, and PhD degrees in Computer Science from the University of California at Berkeley. Professor Hearst's research focuses on user interfaces and robust language analysis for information access systems, and on furthering the understanding of how people use and understand such systems.

Gonzalo Navarro received his first degrees in Computer Science from ESLAI (Latin American Superior School of Informatics) in 1992 and from the University of La Plata (Argentina) in 1993. In 1995 he received his MSc in Computer Science from the University of Chile, obtaining a PhD in 1998. Between 1990 and 1993 he worked at IBM Argentina, on the development of interactive applications and on research on multimedia and hypermedia. Since 1994 he has worked in the Department of Computer Science of the University of Chile, doing research on design and analysis of algorithms, textual databases, and approximate search. He has published a number of papers and also served as referee on different journals (Algorithmica, TOCS, TOIS, etc.) and at conferences (SIGIR, CPM, ESA, etc.).

Edie Rasmussen is an Associate Professor in the School of Information Sciences, University of Pittsburgh. She has also held faculty appointments at institutions in Malaysia, Canada, and Singapore. Dr Rasmussen holds a BSc from the University of British Columbia and an MSc degree from McMaster University, both in Chemistry, an MLS degree from the University of Western Ontario, and a PhD in Information Studies from the University of Sheffield. Her current research interests include indexing and information retrieval in text and multimedia databases.

Ohm Sornil is currently a PhD candidate in the Department of Computer Science at Virginia Polytechnic and State University and a scholar of the Royal Thai Government. He received a BEng in Electrical Engineering from Kasetsart University, Thailand, in 1993 and an MS in Computer Science from Syracuse University in 1997. His research interests include information retrieval, digital libraries, communication networks, and hypermedia.

Nivio Ziviani is a Professor of Computer Science at the Federal University of Minas Gerais in Brazil, where he heads the laboratory for Treating Information. He received a BS in Mechanical Engineering from the Federal University of Minas Gerais in 1971, an MSc in Informatics from the Catholic University of Rio in 1976, and a PhD in Computer Science from the University of Waterloo, Canada, in 1982. He has obtained several research funds from the Brazilian Research Council (CNPq), Brazilian Agencies CAPES and FINEP, Spanish Agency CYTED (project AMYRI), and private institutions. He currently coordinates a four year project on Web and wireless information systems (called SIAM) financed by the Brazilian Ministry of Science and Technology. He is cofounder of the Miner Technology Group, owner of the Miner Family of agents to search the Web. He is the author of several papers in journals and conference proceedings covering topics in the areas of algorithms and data structures, information retrieval, text indexing, text searching, text compression, and related areas. Since January of 1998, he is the editor of the 'News from Latin America' section in the Bulletin of the European Association for Theoretical Computer Science. He has been chair and member of the program committee of several conferences and is a member of ACM, EATICS and SBC.

Contents

Preface			vii		
A	Acknowledgements				
Biographies					
1	Intro	oduction	1		
	1.1	Motivation	1		
		1.1.1 Information versus Data Retrieval	1		
		1.1.2 Information Retrieval at the Center of the Stage	2		
		1.1.3 Focus of the Book	3		
	1.2	Basic Concepts	3		
		1.2.1 The User Task	4		
		1.2.2 Logical View of the Documents	5		
	1.3	Past, Present, and Future	6		
		1.3.1 Early Developments	6		
		1.3.2 Information Retrieval in the Library	7		
		1.3.3 The Web and Digital Libraries	7		
		1.3.4 Practical Issues	8		
	1.4 The Retrieval Process				
	1.5	Organization of the Book	10		
		1.5.1 Book Topics	11		
		1.5.2 Book Chapters	12		
	1.6	How to Use this Book	15		
		1.6.1 Teaching Suggestions	15		
		1.6.2 The Book's Web Page	16		
	1.7	Bibliographic Discussion	17		
2	Mod	leling	19		
	2.1	Introduction	19		
	2.2	A Taxonomy of Information Retrieval Models	20		
	2.3	Retrieval: Ad hoc and Filtering	21		

XII	CONT	ENTS	
	2.4	A Formal Characterization of IR Models	23
	2.5	Classic Information Retrieval	24
		2.5.1 Basic Concepts	24
		2.5.2 Boolean Model	25
		2.5.3 Vector Model	27
		2.5.4 Probabilistic Model	36
		2.5.5 Brief Comparison of Classic Models	34
	2.6		34
		2.6.1 Fuzzy Set Model	34
		2.6.2 Extended Boolean Model	38
	2.7	Alternative Algebraic Models	41
			41
			44
			46
	2.8		48
			48
			49
			56
			59
			60
			61
	2.9		61
			62
			63
	2.10		65
			65
			66
		2.10.3 The Hypertext Model	66
	2.11	······································	69
	2.12		69
3	Retr	eval Evaluation	73
	3.1		73
	3.2		74
			75
			82
	3.3		84
		^	84
			91
		3.3.3 The Cystic Fibrosis Collection	94
	3.4		96
	3.5	Tall 11 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	96
4	Quer	Languages	99
	4.1		99
	4.2	Keyword-Based Querying	