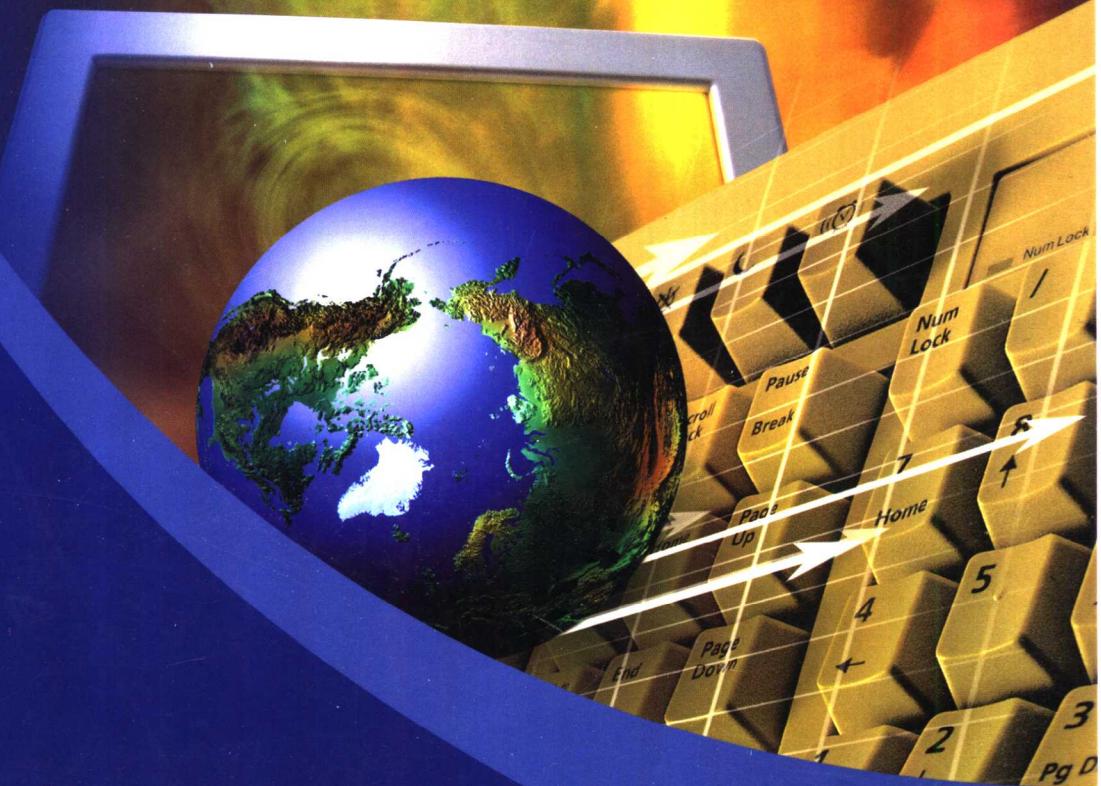


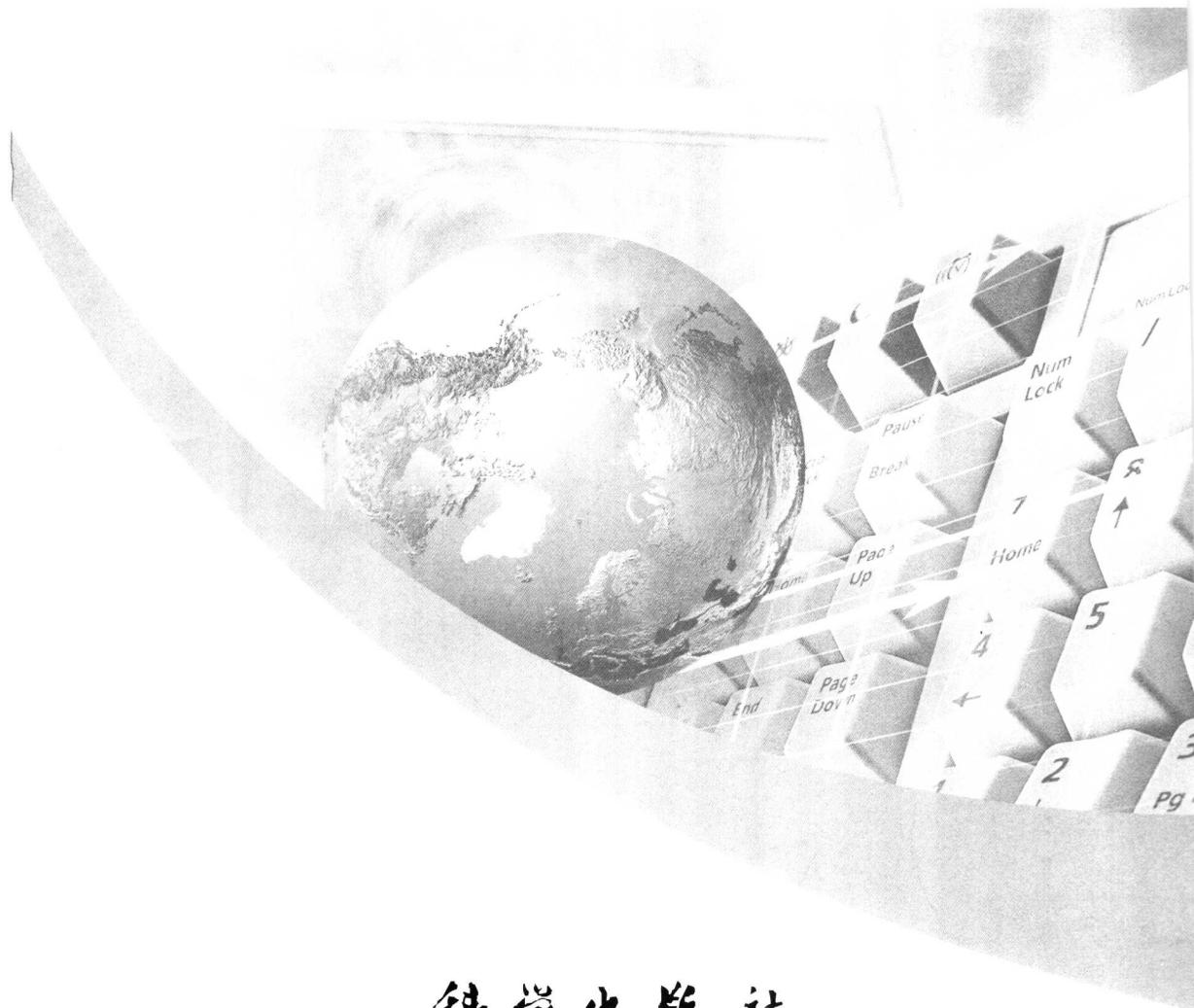
● 孙建军 成 颖 等 编著

# 信息检索技术



孙建军 成颖 丁芹  
李君君 宋玲丽 柯青 编著

# 信息检索技术



科学出版社

北京

## 内 容 简 介

本书系统地介绍了信息检索的原理与技术。讨论的中心问题是如何能迅速地检索到相关信息。具体内容包括：信息检索的布尔模型、向量空间模型、概率模型，以及逻辑模型；文献自动处理技术：自动分类、自动聚类、自动文摘；查询的扩展与精化、相关性、Z39.50，以及搜索引擎等。

本书可作为高等院校信息管理与信息系统专业本科生和研究生教材，也可作为信息机构有关信息服务人员、咨询人员、管理人员的参考用书。

### 图书在版编目(CIP)数据

信息检索技术/孙建军,成颖等编著.一北京:科学出版社,2004

ISBN 7-03-014244-6

I . 信 … II . ①孙 … ②成 … III . 情报检索 IV . G252.7

中国版本图书馆 CIP 数据核字(2004)第 087155 号

责任编辑:李 敏 / 责任校对:钟 洋

责任印制:钱玉芬 / 封面设计:东方上林

科学出版社 出版

北京东黄城根北街16号

邮政编码:100717

<http://www.sciencep.com>

新蕾印刷厂 印刷

科学出版社发行 各地新华书店经销

\*

2004年10月第一版 开本:B5 (720×1000)

2004年10月第一次印刷 印张:30 1/2 插页 1

印数:1—3 000 字数:615 000

定价:49.00 元

(如有印装质量问题,我社负责调换〈环伟〉)

## 前　　言

目前,信息检索研究可以分为实务与原理两部分,图书馆、情报以及档案学界主要侧重于信息检索实务方面的研究与应用,成果颇丰,据不完全统计,国内已经有好几百种信息检索实务类的教材与教程。信息检索原理方面的研究则主要由计算机、数学及情报学界完成,20世纪90年代以后,该领域再次吸引了众多学者的目光,许多原先从事人工智能、自然语言处理等领域的学者也加入到信息检索的研究队伍中,近些年成果也非常丰富,但是对该领域系统总结的专著或者教材尚不多见,本书旨在完成该任务。

本书共12章,分别从信息检索模型、技术、标准以及应用4个方面系统地介绍了信息检索原理方面的最新进展。

引言(第1章)部分从宏观的角度鸟瞰了信息检索原理领域的全貌,构建了全书的框架,预测了原理研究的方向。

模型部分(第2、3、4、5章)系统总结了布尔模型、向量空间模型、概率模型以及逻辑模型,考虑部分兄弟院校可能会以本书作为教材使用,因此,在介绍这些模型时,我们兼顾了历史与进展。介绍这部分内容的目的,是让读者对信息检索原理有一个初步的、同时也是较为全面的理解与掌握。

技术部分(第6、7、8、9、10章)主要介绍信息检索领域的自动分类、自动聚类、自动文摘、查询扩展和精化以及相关性等研究,其研究比较深入,已经处于即将从实验室进入实用阶段或者已经部分实用的技术。通过这些章节的学习后,希望读者能对信息检索领域的主流技术有一个概要的认识,并能通过这些技术进行信息检索的相关研究以及应用工作。

标准部分(第11章)介绍了目前已经实际应用的标准——Z39.50,通过标准的学习,希望读者能够认识到标准对于信息检索研究与应用的意义。

应用部分(第12章)通过搜索引擎系统地阐述了前面章节中各种原理、技术以及标准的综合运用,旨在使读者深化前面内容的学习,从而做

## 前言

到原理、技术、标准与应用的有机结合。

本书由南京大学孙建军、成颖、丁芹、李君君、宋玲丽和柯青合作编写。编写工作的具体分工如下：孙建军组稿，提出大纲，第1章由孙建军、柯青撰写，第2、3、4章由丁芹撰写，第5、7、12章由李君君撰写，第6、8章由成颖撰写，第9、11章由宋玲丽撰写，第10章由孙建军、成颖撰写，最后由孙建军审订统稿。柯青协助完成了大量的校对工作。

本书得到了南京大学“985”工程的资助，特此致谢。

在本书写作过程中得到了科学出版社李敏编辑的大力支持与帮助，在此表示衷心的感谢。

国内外众多的信息检索研究为本书提供了良好的基础，本书的顺利完成也得益于参阅了大量的相关作者的成果，在此我们向这些文献的作者表示诚挚的谢意。

由于我们的学识、水平有限，书中难免存在一些缺陷以及不足，恳请专家与读者批评指正。

孙建军

2004年8月



# 目 录

## 前言

<b>1 引言</b>	<b>1</b>
<b>1.1 信息检索的起源与发展</b>	<b>1</b>
1.1.1 手工检索	1
1.1.2 脱机批处理检索	2
1.1.3 联机检索	2
1.1.4 光盘检索	2
1.1.5 网络化联机检索	3
<b>1.2 信息检索技术的研究内容</b>	<b>3</b>
1.2.1 检索模型研究	4
1.2.2 信息处理技术研究	5
1.2.3 技术应用研究	7
<b>1.3 信息检索技术的未来</b>	<b>8</b>
1.3.1 以人工智能为代表的信息检索自动化趋势	8
1.3.2 人工参与检索工具的信息组织是检索工具的发展趋势	8
1.3.3 多媒体信息检索技术的成熟与发展	9
1.3.4 多语种检索的支持	9
1.3.5 个人化的检索工具和专业化的检索工具	10
<b>2 布尔检索模型</b>	<b>11</b>
<b>2.1 传统布尔检索</b>	<b>11</b>
2.1.1 布尔运算	11
2.1.2 传统布尔检索模型	11
2.1.3 布尔查询的自动生成	14
2.1.4 传统布尔查询的评价	16
<b>2.2 扩展布尔检索方法</b>	<b>18</b>
2.2.1 研究背景	18
2.2.2 扩展布尔检索的思想基础	18
2.2.3 P-范式模型	21
2.2.4 P-范式模型的特点	24
2.2.5 P-范式模型的实现	25

## 目 录

2.2.6 扩展布尔操作符 .....	31
<b>3 向量空间检索 .....</b>	<b>36</b>
3.1 传统向量空间检索 .....	36
3.1.1 向量空间模型介绍 .....	36
3.1.2 向量空间模型的评价 .....	40
3.2 广义向量空间检索 .....	41
3.2.1 布尔代数的向量表示 .....	41
3.2.2 无权重的标引词项的向量表示 .....	42
3.2.3 广义的向量空间模型 .....	43
3.2.4 用 GVSM 来处理布尔查询 .....	49
3.3 项的权重模式 .....	53
3.3.1 项向量的规范化 .....	54
3.3.2 项权重模式 .....	61
3.4 相似度的计算 .....	65
3.4.1 内积相似度运算 .....	66
3.4.2 余弦相似度 .....	66
3.4.3 “距离”相似度运算 .....	67
3.4.4 以项匹配的个数作为相似度计算的依据 .....	67
3.4.5 一种基于概率向量的相似度计算方法 .....	69
3.5 潜在语义标引 .....	71
3.5.1 模型的提出 .....	71
3.5.2 潜在语义标引模型 .....	72
3.5.3 空间中各种向量的匹配 .....	76
3.5.4 应用于布尔查询的潜在语义标引 .....	78
3.5.5 模型的评价 .....	79
<b>4 概率检索 .....</b>	<b>81</b>
4.1 概率信息检索的背景 .....	81
4.1.1 信息检索中概率模型的历史 .....	81
4.1.2 概率检索理论的背景知识 .....	82
4.2 基于相关性概率估计的检索模型 .....	86
4.2.1 作为一个决策策略的概率模型 .....	86
4.2.2 二元独立模型 .....	87
4.2.3 基于概率标引的检索模型 .....	99
4.2.4 逻辑回归模型 .....	102
4.2.5 2-泊松模型 .....	105

4.3 推理网络模型 .....	108
4.3.1 推理网络的总体介绍 .....	108
4.3.2 应用于文献检索的推理网络 .....	109
4.3.3 推理网络与其他模型的比较 .....	122
5 逻辑模型 .....	129
5.1 逻辑模型的建构 .....	129
5.1.1 逻辑模型的基本思想 .....	129
5.1.2 逻辑模型的建构方法 .....	130
5.2 古典逻辑与古典逻辑模型 .....	131
5.2.1 古典逻辑 .....	131
5.2.2 古典逻辑模型 .....	133
5.3 van Rijsbergen 的非古典逻辑模型 .....	134
5.3.1 逻辑蕴涵模型 .....	134
5.3.2 不确定性原理 .....	136
5.4 逻辑蕴涵程度的测算 .....	140
5.4.1 逻辑蕴涵程度测算方法 .....	140
5.4.2 向量空间模型的测算 .....	142
5.4.3 布尔模型的测算 .....	143
5.4.4 概率模型的测算 .....	144
5.5 信息检索逻辑模型 .....	145
5.5.1 基于可能世界的逻辑模型 .....	145
5.5.2 基于情景理论的信息检索模型 .....	148
5.5.3 基于术语逻辑的信息检索模型 .....	150
5.5.4 信息检索的元模型 .....	153
5.5.5 信息检索逻辑模型的特征 .....	156
6 自动分类 .....	160
6.1 引言 .....	160
6.2 基本概念 .....	161
6.2.1 定义 .....	161
6.2.2 分类 .....	162
6.2.3 应用 .....	163
6.2.4 训练集与测试集 .....	165
6.3 特征选取 .....	166
6.3.1 预处理 .....	166
6.3.2 标引 .....	167
6.4 降维技术 .....	170

## 目 录

6.4.1 特征选择 .....	170
6.4.2 特征重构 .....	175
6.5 分类方法 .....	178
6.5.1 Rocchio's 算法及改进 .....	178
6.5.2 朴素贝叶斯分类方法 .....	180
6.5.3 K 最近邻算法 .....	183
6.5.4 决策树方法 .....	185
6.5.5 支持向量机 .....	193
6.5.6 基于投票的方法 .....	196
6.6 文档分类的评估指标 .....	198
6.6.1 多重二元分类任务 .....	198
6.6.2 多重分类和多重标识分类 .....	200
7 聚类 .....	201
7.1 聚类检索 .....	201
7.1.1 聚类策略 .....	202
7.1.2 检索步骤 .....	203
7.2 文献相似度 .....	203
7.2.1 距离 .....	204
7.2.2 相似系数 .....	204
7.2.3 基于提问式的文献相似度 .....	205
7.3 层次聚类法 .....	206
7.3.1 合成聚类法 .....	207
7.3.2 分解聚类法 .....	218
7.4 启发式聚类法 .....	219
7.4.1 密度测试法 .....	219
7.4.2 线性时间法 .....	220
7.5 增量式聚类法 .....	222
7.5.1 单遍聚类法 .....	222
7.5.2 后缀树法 .....	222
7.6 聚类浏览 .....	229
7.6.1 聚类浏览概述 .....	229
7.6.2 聚类浏览算法 .....	230
8 自动文摘 .....	232
8.1 语料库 .....	232
8.1.1 语料库的分类 .....	233
8.1.2 语料库的设计与建设 .....	233



8.1.3 语料库的研究方法 .....	234
8.1.4 概率论基础知识 .....	236
8.1.5 Ngram 语法 .....	237
<b>8.2 词法分析 .....</b>	<b>238</b>
8.2.1 自动分词 .....	238
8.2.2 歧义切分 .....	242
8.2.3 未登录词 .....	244
8.2.4 词性标注 .....	244
<b>8.3 句法分析 .....</b>	<b>245</b>
8.3.1 句法分析中的知识表示 .....	245
8.3.2 句法分析算法 .....	258
<b>8.4 自动摘要 .....</b>	<b>268</b>
8.4.1 自动摘要的步骤 .....	268
8.4.2 自动摘要的不足 .....	271
<b>8.5 基于理解的自动文摘 .....</b>	<b>272</b>
8.5.1 基本步骤 .....	272
8.5.2 篇章意义的机内表示 .....	273
8.5.3 理解文摘的不足 .....	274
<b>8.6 信息抽取 .....</b>	<b>275</b>
8.6.1 信息抽取研究的发展历史 .....	275
8.6.2 信息抽取系统的体系结构 .....	276
8.6.3 命名实体识别 .....	276
<b>8.7 基于结构的自动文摘 .....</b>	<b>278</b>
8.7.1 关联网络 .....	278
8.7.2 修辞结构 .....	278
8.7.3 语用功能 .....	279
<b>8.8 文摘评估方法 .....</b>	<b>279</b>
8.8.1 直接评价方法 .....	279
8.8.2 基于任务的评价方法 .....	281
8.8.3 基于目标的评估方法 .....	282
<b>8.9 自动文摘研究所取得的成绩和面临的问题 .....</b>	<b>284</b>
<b>9 查询扩展和精化 .....</b>	<b>286</b>
9.1 查询扩展和精化概述 .....	286
9.1.1 查询扩展和精化的意义 .....	286
9.1.2 查询扩展的类型 .....	287
9.2 相关反馈技术 .....	288

## 目 录

9.2.1 相关反馈技术介绍 .....	288
9.2.2 向量空间模型中的相关反馈 .....	292
9.2.3 概率模型中的相关反馈 .....	297
9.2.4 布尔模型中的相关反馈 .....	302
9.2.5 相关反馈技术的改进 .....	305
9.3 查询检索词选择方案 .....	308
9.3.1 检索词选择概述 .....	308
9.3.2 自动查询扩展中的检索词选择 .....	309
9.3.3 交互式查询扩展中的检索词选择方法 .....	313
9.4 词表扩展技术 .....	315
9.4.1 人工词表 WordNet .....	316
9.4.2 自动构建词表 .....	317
9.5 整体分析技术和局部分析技术 .....	318
9.5.1 整体分析技术 .....	319
9.5.2 局部分析技术 .....	328
9.6 查询的重用 .....	335
9.6.1 steepest descent 算法 .....	336
9.6.2 查询相似度计算 .....	337
<b>10 相关性 .....</b>	<b>341</b>
10.1 相关性的研究历史 .....	341
10.1.1 第一阶段的研究 .....	342
10.1.2 第二阶段的研究 .....	343
10.1.3 第三阶段的研究 .....	347
10.2 相关性研究的学派 .....	352
10.2.1 面向系统的相关性 .....	353
10.2.2 面向用户的相关性 .....	356
10.2.3 结论 .....	359
10.3 相关性模型 .....	360
10.3.1 相关性模型:其他学科的视角 .....	360
10.3.2 相关性模型:信息科学的视角 .....	361
10.4 属性与类别 .....	368
10.4.1 基本属性 .....	368
10.4.2 类别 .....	369
10.4.3 属性与类别的关系 .....	371
10.4.4 属性与类别之间关系的修正 .....	376

<b>11 Z39.50 检索标准</b>	378
11.1 Z39.50 标准概述	378
11.2 Z39.50 标准的起源	379
11.2.1 与 Z39.50 标准相关机构与标准制定过程	379
11.2.2 标准的沿革	380
11.2.3 版本间的关系	381
11.3 Z39.50 的功能	381
11.3.1 建立虚拟联合目录	382
11.3.2 联合编目	382
11.3.3 馆际互借	383
11.3.4 光盘检索	383
11.3.5 定题服务	383
11.3.6 万维网检索和信息过滤	384
11.4 Z39.50 的工作原理	384
11.4.1 Z39.50 的运行机制	384
11.4.2 Z39.50 的实现模型	385
11.4.3 Z39.50 源端和目标端的主要功能	387
11.5 Z39.50 协议简介	396
11.5.1 Z39.50 协议的信息检索服务	396
11.5.2 Z39.50 协议说明	417
11.6 下一代 Z39.50	424
<b>12 Web 信息检索工具——搜索引擎</b>	427
12.1 搜索引擎的工作原理与结构	427
12.1.1 信息采集	428
12.1.2 信息标引	429
12.1.3 索引数据库	430
12.1.4 信息检索	432
12.2 搜索引擎的分类	433
12.2.1 目录式搜索引擎	433
12.2.2 Robot 搜索引擎	435
12.2.3 元搜索引擎	435
12.2.4 智能搜索引擎	438
12.3 搜索引擎的检索功能	441
12.3.1 基本检索功能	442
12.3.2 高级检索功能	444
12.3.3 与检索相关的功能	445

## 目 录

12.4 信息采集 Robot 的实现 .....	447
12.4.1 Robot 的组成模块 .....	447
12.4.2 Robot 的搜索算法 .....	448
12.4.3 Robot 的遍历策略 .....	449
12.4.4 Robot 的专用协议 .....	451
12.4.5 Robot 优化策略 .....	452
12.5 搜索引擎的发展 .....	453
12.5.1 第一代搜索引擎——基于关键词的检索 .....	453
12.5.2 第二代搜索引擎——基于超链接的检索 .....	455
12.5.3 第三代搜索引擎——基于概念的检索 .....	458
参考文献 .....	463



# 1 引言

## 1.1 信息检索的起源与发展

信息检索是指信息用户为处理解决各种问题而查找、识别、获取相关的事实、数据、文献的活动及过程。作为人类社会活动不可分割的一部分，信息检索有着悠久的历史。信息检索研究则是伴随着科学技术的发展和信息数量的剧增而兴起的研究领域。英国科学家詹姆斯·马丁认为：人类的科学知识在 19 世纪是每 50 年增加 1 倍，20 世纪中叶是每 10 年增加 1 倍，在 20 世纪 70 年代就已经缩短到每 5 年增加 1 倍；同时，信息分散，交叉引用频繁，人类信息的生产能力超过了人类对信息的处理、组织和吸收能力，从而产生了信息爆炸的危机。人们越来越关注如何从浩如烟海的信息源中迅速而准确地查找到学习和研究所需要的资料，因而，信息检索的战略地位也就显得日益重要。其主要研究范围包括：信息检索理论、信息检索语言、信息检索系统的建构及评价、信息检索技术与方法等。

20 世纪中叶以前，信息存储和传播主要以纸质介质为载体，信息检索活动也围绕着文献的获取和控制展开。因此，信息检索研究关注的是如何检索、利用文献中记载的信息，从而导致文献检索成为信息检索的同义词。早期的文献中不使用“信息检索”这一概念。50 年代以后，社会信息传播与存储载体呈现多元化，人们不再拘泥于纸质载体研究信息检索，于是开始广泛使用情报检索一词。由于汉语中“信息”一词较“情报”一词的含义更为宽泛，加上英文 information 可以理解为“信息”和“情报”，近年来人们越来越倾向于将情报检索研究和文献检索研究归为信息检索研究这一更具兼容性的概念，以便将各种不同的检索综合起来，使该研究领域取得更多、更为实用的研究成果，对信息检索实践起到更全面的指导作用。

随着科学技术的发展，尤其是计算机的应用，信息检索经历了从手工检索到机械检索再到计算机化检索的过程。

### 1.1.1 手工检索

手工检索是指仅用手工的方式来处理和查找文献工具，如文摘、索引、目录、参考工具书等。它是一种传统而又基础的检索手段。

手工检索因其不需要特殊设备，查找简单、灵活，而且用户可以随时修改检索策略，检索费用较低等优点而在某些部门领域仍然使用，但是，利用手工检索往往

费时较多、效率低下、查全率也较低。

### 1.1.2 脱机批处理检索

20世纪40年代中期世界上第一台电子计算机问世后,50年代初就有人开始研究其在信息检索领域的应用。50年代中期至60年代中后期是信息检索的脱机批处理阶段。当时,计算机硬件发展很快,但还没有连接通信网,也没有远程终端装置,不能提供问答服务(Q-A)的检索方式,只能进行现刊文献的定题检索(SDI)和过期文献的追溯检索(RS),同时利用计算机编辑出版检索性刊物。所谓脱机批处理方式,是指定期由专职检索人员把许多用户课题汇总,批量处理提问要求并把结果提供给用户。在美国,这个时期出现了3个重要系统:1959~1963年美国武装部队技术情报局的ASTIA(即后来的国防文献中心DDC)系统;1962年美国国家航空和航天局的NASA系统;1964年美国国家医学图书馆创建的医学文献分析与检索系统MEDLARS。MEDLARS不仅可以进行逻辑“或”、“与”、“非”多种运算,而且可以从多种途径检索文献。

脱机批处理能同时进行多项检索,对复杂的检索词也具有处理能力,因此,在生产普通印刷索引、专题书目、回溯检索和定题检索服务等方面得到广泛使用。但是,它也有许多不足之处,如缺乏与用户的交互过程,检索结果获得不及时以及信息需求和检索结果之间存在一定误差等,这些缺点限制了脱机批处理的发展。

### 1.1.3 联机检索

20世纪70年代计算机分时系统的出现,通信技术的改进,使得多终端、远距离两地检索信息的技术得以推广,计算机检索技术从脱机阶段进入联机信息检索时期。所谓联机检索,就是用户使用终端设备,通过通信线路与中央计算机连接,直接与计算机对话进行检索,结果由终端输出。第一个大规模的联机检索系统是美国NASA的RECON系统,1969年全面投入运行。随后许多著名的联机检索系统相继出现,如1970年洛克希德火箭公司的DIALOG系统,美国的MEDLARS于1970年发展的MEDLINE联机系统等。

联机检索无需委托,直接面向最终用户,在检索过程中是“人机对话”方式,具有很强的交互功能,而且能及时取得检索结果,但是检索指令复杂,需要依赖专业检索人员。

### 1.1.4 光盘检索

利用国际联机检索系统检索到的电子文献信息具有较高的使用价值,但国际联机检索费用昂贵,一般用户难以承受。人们开始努力寻求一种低廉的存储、检索电子信息的方式,光盘存储技术则适应了这一要求。CD-ROM光盘是20世纪80年

代在计算机技术、激光技术等现代新科技成果的基础上发展起来的新型电子出版物。它具有信息存储密度高、容量大、读取速度快、存储的信息类型多等优点,备受人们的青睐。光盘技术与光盘产品的发展相当迅速,品种和数量激增,而且更新换代快,功能日益完善。光盘塔和光盘网络的出现和广泛应用提高了单张光盘的利用率,使光盘的多用户检索和共享成为现实。

利用光盘检索系统费用大大低于联机检索,利用 CD-ROM 存储信息方便、易于携带,除可提供追溯检索、定题服务外,还可用于“自建库”和做联机检索前预处理。

### 1.1.5 网络化联机检索

国际联机检索和光盘检索为我们提供了大量的信息资源,但各自又都有着或多或少的缺点,例如联机检索费用昂贵,指令复杂,而光盘检索得到的信息又不十分及时等。因此,极有必要产生一种新型的信息检索方式。1993 年,美国政府提出国家基础设施建设(NII)计划,兴建以 Internet 为雏形的信息高速公路,网络资源如潮水般涌来。在信息爆炸的当今社会,单个计算机所能完成的工作和所存储的信息都极为有限,而把单机连起来的计算机网络则能在局部或更大范围内实现通信和信息共享。由于电话网、电传网、公共数据通信网都可能为信息检索传输数据,世界各大检索系统纷纷进入各种通信网络,每个系统的计算机成为网络上的节点,每个节点连接多个检索终端,各节点之间以通信线路彼此相连。网络上的任何一个终端都可联机检索所有的数据库的数据。这就是网络化联机信息检索。网络联机信息检索是联机信息检索的高级阶段,它的实现使人们可以在很短的时间里查遍全球的信息资料,使人类实现信息资源共享成为可能。除了传统的文献信息,网络信息源还包括电子论坛、各种软件资料、图像文件、声音文件等。值得指出的是,网络信息环境的出现,使得信息检索研究的对象和范围不断扩大,研究队伍也突破了原有的以图书情报领域的专家学者为主的框架,众多的信息公司加入到研究开发信息检索系统的行列。可以说,网络使计算机信息检索技术进入一个崭新发展阶段,而网络信息检索又使得网上信息源利用率提高,信息组织更为有序和高效。

## 1.2 信息检索技术的研究内容

随着网络信息资源的日益丰富和复杂化,为满足不同用户能够检索到所需信息,检索系统朝着自然语言检索、用户界面友好的方向发展,这给信息检索技术提出了更高的要求。因此,当前信息检索技术的研究主要包括以下方面:

### 1.2.1 检索模型研究

信息组织是实现信息检索的基础,原始的文档中包括文本、图像、视频、音频等数据,不能直接进行检索,需要从这些原始数据中抽取逻辑视图,支持信息检索。用户用查询来表示他的信息需求。检索系统根据查询的表示,搜索文档集,获取与用户查询相关的文档。信息检索的匹配是相似性匹配,查询的结果按序返回。以上过程实际上涉及3个重要的处理:文档集的逻辑表示、查询的表示、相似匹配及其排序。对这些检索的因素和过程建模,就产生了各种不同的信息检索模型。

我们把信息检索模型定义如下:

一个信息检索模型是将文档表示、查询以及它们之间关系进行建模的框架,它由三元体

$$F[D, Q, R(q_i, d_j)]$$

表示。其中, $D$ 是文档集中的一组文档逻辑视图(或称为文档的表示); $Q$ 是一组用户信息需求的逻辑视图(表示),这种视图(表示)被称为查询; $R(q_i, d_j)$ 是一个排序函数,该函数输出一个与查询 $q_i \in Q$ 和文档表示 $d_j \in D$ 有关的实数。这样就在文档之间根据查询 $q_i$ 定义了一个顺序。

信息检索中4个传统模型是:布尔模型、向量空间模型、概率模型和逻辑模型。近些年来,研究人员对于每种传统的模型都提出了各种不同的改进模式,如在基于集合论的检索模型中,提出了模糊布尔模型和扩展布尔模型;在代数型模型中,衍生出广义矢量模型、隐含语义索引模型和神经网络模型等3种;在概率检索模型中,发展出推理网络模型和信念网络模型等。除了涉及文本的内容之外,模型还应该涉及文本的结构。在这种情况下,就应该还有表示文本结构的结构模型。对于文本的结构模型,主要有两种类型:非重叠链表模型和邻近节点模型。

对这些检索模型可以用图1-1来表示它们的层次关系<sup>①</sup>:

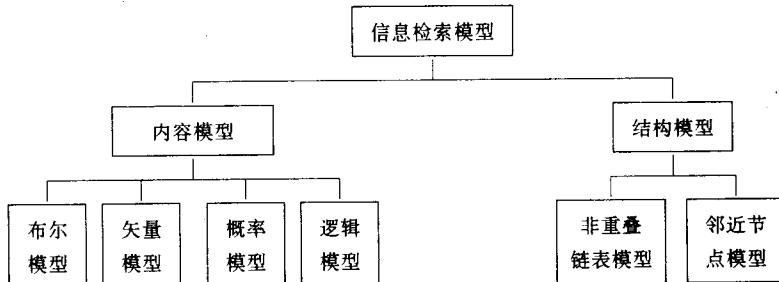


图1-1 信息检索模型的分类

<sup>①</sup> 李国辉,汤大权,武德峰. 信息组织与检索. 北京:科学出版社,2003. 91~105