



智能科学技术著作丛书

# 智能检索技术

陆建江 张亚非 徐伟光 苗壮 编著



科学出版社  
[www.sciencep.com](http://www.sciencep.com)

智能科学技术著作丛书

# 智能检索技术

陆建江 张亚非 徐伟光 苗壮 编著

科学出版社(CP)



科学出版社

1990年

0003 Q228

图书 二 摘 单 签 本

新编辞书学研究会 中国辞书出版社

科学出版社

北京  
(英汉双语) 国际学术会议

## 内 容 简 介

面对海量信息,信息的精确检索就像大海捞针一样困难。智能检索技术吸取多个学科的研究成果,力图通过对文本、图像和视频信息的智能处理,实现信息的精确检索。本书系统地阐述了文本、图像和视频检索的理论方法和实现技术,并重点突出了本领域的最新研究成果。

本书可作为高等院校计算机科学与技术、模式识别与智能系统等学科方向高年级本科生和研究生的教材,也可作为相关领域学生的参考书。

著者：陆建江 张亚非 徐伟光 苗壮 编著

### 图书在版编目(CIP)数据

智能检索技术/陆建江,张亚非,徐伟光,苗壮编著.一北京:科学出版社,  
2009

(智能科学技术著作丛书)

ISBN 978-7-03-025328-6

I. 智… II. ①陆… ②张… ③徐… ④苗… III. 计算机网络-情报检索  
IV. G354.4

中国版本图书馆 CIP 数据核字(2009)第 148432 号

责任编辑:张海娜 / 责任校对:钟 洋  
责任印制:赵 博 / 封面设计:陈 敏

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

\*

2009 年 8 月第 一 版 开本:B5(720×1000)

2009 年 8 月第一次印刷 印张:16

印数:1—3 000 字数:302 000

定价: 48.00 元

(如有印装质量问题, 我社负责调换)

## 《智能科学技术著作丛书》编委会

名誉主编：吴文俊

主 编：涂序彦

副 主 编：钟义信 史忠植 何华灿 蔡自兴 孙增圻 谭 民

秘 书 长：韩力群

编 委：(按姓氏汉语拼音排序)

蔡庆生(中国科学技术大学)

杜军平(北京邮电大学)

何华灿(西北工业大学)

黄河燕(中国科学院计算语言研究所)

焦李成(西安电子科技大学)

刘 宏(北京大学)

秦世引(北京航空航天大学)

阮秋琦(北京交通大学)

孙增圻(清华大学)

涂序彦(北京科技大学)

王家钦(清华大学)

吴文俊(中国科学院系统科学研究所)

尹怡欣(北京科技大学)

张琴珠(华东师范大学)

庄越挺(浙江大学)

蔡自兴(中南大学)

韩力群(北京工商大学)

何 清(中国科学院计算技术研究所)

黄心汉(华中科技大学)

李祖枢(重庆大学)

刘 清(南昌大学)

邱玉辉(西南师范大学)

史忠植(中国科学院计算技术研究所)

谭 民(中国科学院自动化研究所)

王国胤(重庆邮电学院)

王万森(首都师范大学)

杨义先(北京邮电大学)

于洪珍(中国矿业大学)

钟义信(北京邮电大学)

## 《智能科学技术著作丛书》序

“智能”是“信息”的精彩结晶，“智能科学技术”是“信息科学技术”的辉煌篇章，“智能化”是“信息化”发展的新动向、新阶段。

“智能科学技术”(intelligence science&technology, IST)是关于“广义智能”的理论方法和应用技术的综合性科学技术领域，其研究对象包括：

- “自然智能”(natural intelligence, NI)，包括：“人的智能”(human intelligence, HI)及其他“生物智能”(biological intelligence, BI)。
- “人工智能”(artificial intelligence, AI)，包括：“机器智能”(machine intelligence, MI)与“智能机器”(intelligent machine, IM)。
- “集成智能”(integrated intelligence, II)，即：“人的智能”与“机器智能”人机互补的集成智能。
- “协同智能”(cooperative intelligence, CI)，指：“个体智能”相互协调共生的群体协同智能。
- “分布智能”(distributed intelligence, DI)，如：广域信息网，分散大系统的分布式智能。

1956年，“人工智能”学科诞生，50年来，在起伏、曲折的科学征途上不断前进、发展，从狭义人工智能走向广义人工智能，从个体人工智能到群体人工智能，从集中式人工智能到分布式人工智能，在理论方法研究和应用技术开发方面都取得了重大进展。如果说，当年“人工智能”学科的诞生是生物科学技术与信息科学技术、系统科学技术的一次成功的结合，那么，可以认为，现在“智能科学技术”领域的兴起是在信息化、网络化时代又一次新的多学科交融。

1981年，“中国人工智能学会”(Chinese Association for Artificial Intelligence, CAAI)正式成立，25年来，从艰苦创业到成长壮大，从学习跟踪到自主研发，团结我国广大学者，在“人工智能”的研究开发及应用方面取得了显著的进展，促进了“智能科学技术”的发展。在华夏文化与东方哲学影响下，我国智能科学技术的研究、开发及应用，在学术思想与科学方法上，具有综合性、整体性、协调性的特色，在理论方法研究与应用技术开发方面，取得了具有创新性、开拓性的成果。“智能化”已成为当前新技术、新产品的发展方向和显著标志。

为了适时总结、交流、宣传我国学者在“智能科学技术”领域的研究开发及应用成果，中国人工智能学会与科学出版社合作编辑出版《智能科学技术著作丛书》。需要强调的是，这套丛书将优先出版那些有助于将科学技术转化为生产力以及对社会和国民经济建设有重大作用和应用前景的著作。

我们相信，有广大智能科学技术工作者的积极参与和大力支持，以及编委们的

共同努力,《智能科学技术著作丛书》将为繁荣我国智能科学技术事业、增强自主创新能力、建设创新型国家做出应有的贡献。

祝《智能科学技术著作丛书》出版,特赋贺诗一首:

**智能科技领域广**

**人机集成智能强**

**群体智能协同好**

**智能创新更辉煌**

涂序彦

中国人工智能学会荣誉理事长

2005年12月18日

## 前　　言

面对海量信息,信息的精确检索就像大海捞针一样困难。智能检索技术吸取多个学科的研究成果,力图通过对文本、图像和视频信息的智能处理,实现信息的精确检索。本书系统地阐述了文本、图像和视频检索的理论方法和实现技术,并重点突出了本领域的最新研究成果。

本书涵盖智能检索技术的主要内容,全书共分 14 章:第 1~4 章介绍文本的智能检索技术,包括文本检索技术、文本自动分词、概念语义空间、基于本体的文本检索技术等;第 5~10 章介绍图像的智能检索技术,包括 MPEG-7 标准中图像的视觉特征、图像的局部特征、基于视觉特征的图像检索技术、基于语义的图像检索技术、Web 图像的检索技术等;第 11~14 章介绍视频的智能检索技术,包括视频的结构化技术、语音识别技术、视频的标注技术等。

本书的成果是集体智慧的结晶,由陆建江、张亚非、徐伟光和苗壮负责撰稿。另外,感谢赵天忠、李阳、肖琪、谢正辉、周波、李冉、李言辉、康达周、王进鹏、王家宝、田豫龙等同学为本书撰写工作付出的辛勤工作,这些同学参与了全书的校对工作,在此深表感谢。全书每一章的内容组织和细节都经过多次讨论和修改才定稿,力求深入浅出,让读者轻松掌握相关的知识。尽管每一节、每一句、每篇参考文献,甚至每个标点我们都精心检查,但难免还存在一些缺点和遗漏,殷切希望广大读者批评指正。希望本书的出版能够对智能检索技术相关领域的研究人员有所裨益,并希望通过阅读本书,读者能够很快进行相关领域的研究工作。

28	· · · · ·	· · · · ·	· · · · ·	· · · · ·
72	· · · · ·	· · · · ·	· · · · ·	· · · · ·
08	· · · · ·	· · · · ·	· · · · ·	· · · · ·
08	· · · · ·	· · · · ·	· · · · ·	· · · · ·
58	· · · · ·	· · · · ·	· · · · ·	· · · · ·
<b>《智能科学技术著作丛书》序</b>				
<b>前言</b>				
<b>第1章 文本检索技术</b>				
11	1.1	基于索引的检索技术	· · · · ·	1
12	1.2	文本提取	· · · · ·	2
28	1.3	文本预处理	· · · · ·	3
38	1.3.1	停用词删除	· · · · ·	4
48	1.3.2	词干提取	· · · · ·	4
58	1.3.3	索引词选择	· · · · ·	4
68	1.3.4	建立词典	· · · · ·	5
78	1.4	索引	· · · · ·	5
88	1.5	文本检索模型	· · · · ·	7
98	1.5.1	布尔模型	· · · · ·	8
108	1.5.2	向量空间模型	· · · · ·	8
118	1.5.3	概率论模型	· · · · ·	9
128	1.5.4	PageRank 模型	· · · · ·	10
138	1.6	分布式搜索引擎	· · · · ·	12
148	1.6.1	分布式元搜索引擎	· · · · ·	13
158	1.6.2	散列式分布搜索引擎	· · · · ·	13
168	1.6.3	局部遍历型搜索引擎	· · · · ·	14
178	1.6.4	P2P 分布式搜索引擎	· · · · ·	15
188	参考文献	· · · · ·	· · · · ·	16
<b>第2章 文本自动分词</b>				
198	2.1	基于字符串匹配的正向最大匹配算法	· · · · ·	18
208	2.2	基于简码匹配的 Hash 分词算法	· · · · ·	20
218	2.2.1	简码匹配方式	· · · · ·	20
228	2.2.2	Hash 分词算法	· · · · ·	21
238	2.2.3	消歧融入切分过程	· · · · ·	22
248	2.2.4	基于简码的 Hash 算法	· · · · ·	23

2.2.5 平均匹配次数的理论分析	25
2.2.6 分词测试及结果	27
2.3 基于统计的分词方法	30
参考文献	30
<b>第3章 概念语义空间</b>	<b>32</b>
3.1 基于奇异值分解的潜在语义索引方法	32
3.2 基于非负矩阵分解的潜在语义索引方法	33
3.2.1 NMF 问题的提出	33
3.2.2 目标函数	34
3.2.3 NMF 方法的迭代规则	34
3.2.4 NMF 的非唯一性	35
3.2.5 基于 NMF 的概念语义生成	35
3.2.6 其他 NMF 方法	37
3.3 NMF 方法与 SVD 方法的比较	38
3.3.1 问题本质	38
3.3.2 概念语义向量的特点	38
3.3.3 概念语义向量的解释	39
3.3.4 NMF 方法与 SVD 方法敏感性的比较	39
3.3.5 NMF 方法与 SVD 方法检索性能的比较	40
参考文献	41
<b>第4章 基于本体的文本检索技术</b>	<b>42</b>
4.1 本体定义	42
4.2 描述逻辑	44
4.2.1 描述逻辑 ALC	44
4.2.2 描述逻辑 ALC 的构造子扩展	46
4.3 本体语言	49
4.3.1 可扩展标记语言 XML	50
4.3.2 资源描述框架 RDF	54
4.3.3 本体语言 OWL	59
4.4 基于本体的文本检索技术	64
4.4.1 本体构建	64
4.4.2 语义标注	70
4.4.3 语义查询	72
参考文献	77

---

<b>第5章 基于内容的图像检索</b>	82
5.1 基于内容的图像检索的原因	82
5.2 基于内容的图像检索概述	82
5.2.1 基于视觉特征的图像检索	83
5.2.2 基于对象类型的图像检索	83
5.2.3 基于抽象属性的图像检索	83
5.3 Web 图像检索	83
参考文献	84
<b>第6章 MPEG-7 标准中图像的视觉特征</b>	86
6.1 图像的颜色特征	86
6.1.1 颜色空间	86
6.1.2 颜色量化	90
6.1.3 主颜色	91
6.1.4 可伸缩颜色	93
6.1.5 颜色布局	94
6.1.6 颜色结构	97
6.2 图像的纹理特征	100
6.2.1 同质纹理	100
6.2.2 纹理浏览	106
6.2.3 边缘直方图	110
6.3 图像的形状特征	113
6.3.1 基于区域的形状	113
6.3.2 基于轮廓的形状	116
参考文献	120
<b>第7章 图像的局部特征</b>	122
7.1 图像兴趣点和兴趣区域的发现器	122
7.1.1 Harris 兴趣点发现器	122
7.1.2 Harris-Laplace 兴趣区域发现器	124
7.1.3 Hessian-Laplace 兴趣区域发现器	125
7.1.4 高斯差分金字塔	125
7.2 尺度不变特征变换 SIFT	125
7.2.1 SIFT 特征的提取	125
7.2.2 SIFT 兴趣点的匹配	129
7.2.3 与 SIFT 有关的其他局部特征	130

7.3 方向可调滤波器 .....	130
7.4 形状上下文 .....	132
7.5 矩不变量 .....	133
参考文献.....	134
<b>第8章 基于视觉特征的图像检索技术.....</b>	<b>135</b>
8.1 图像分割技术 .....	136
8.1.1 图像分割概念 .....	136
8.1.2 图像分割算法 .....	137
8.1.3 分割方法存在的问题 .....	141
8.2 相似性度量 .....	142
8.2.1 几何模型 .....	142
8.2.2 相关计算模型 .....	143
8.2.3 关联系数模型 .....	144
8.3 索引 .....	144
8.3.1 高维索引方法 .....	144
8.3.2 降维方法 .....	146
8.3.3 近似最近邻方法 .....	147
8.3.4 单一维空间映射方法 .....	148
8.3.5 多重空间填充曲线方法 .....	148
8.3.6 基于过滤的方法 .....	148
8.4 相关反馈技术 .....	149
8.5 图像检索系统性能的评价准则 .....	151
8.6 基于视觉特征的图像检索系统 .....	151
参考文献.....	153
<b>第9章 基于语义的图像检索技术.....</b>	<b>157</b>
9.1 图像标注技术的概况 .....	157
9.2 图像标注系统的工作原理 .....	159
9.3 基于 MPEG-7 的图像标注技术 .....	160
9.3.1 SVM 分类器 .....	160
9.3.2 基于 MPEG-7 的图像标注技术 .....	166
9.4 基于特征选择的图像标注技术 .....	167
9.4.1 遗传算法的基本思想 .....	167
9.4.2 基于二进制编码遗传算法的最优特征子集选择方法 .....	168
9.4.3 基于双编码遗传算法的最优加权特征子集选择方法 .....	170

---

9.4.4 基于特征选择的图像标注技术 .....	172
9.5 基于 Adaboost 算法的图像标注技术 .....	172
9.5.1 Adaboost 算法 .....	173
9.5.2 $k$ -NN 分类器 .....	175
9.5.3 主从式并行遗传算法的实现 .....	176
9.5.4 图像标注技术 .....	179
9.6 基于类对特征选择的图像标注技术 .....	180
9.7 实验结果 .....	181
9.8 大规模图像的标注技术 .....	182
9.8.1 WordNet 简介 .....	183
9.8.2 基于 WordNet 的图像标注技术 .....	184
9.8.3 小结 .....	187
参考文献 .....	187
<b>第 10 章 Web 图像的检索技术 .....</b>	<b>190</b>
10.1 Web 图像搜索引擎的工作原理 .....	190
10.2 Web 图像的抓取 .....	191
10.3 网页文本信息的挖掘 .....	193
10.3.1 网页上的文本信息源 .....	193
10.3.2 标注精炼 .....	195
10.4 图像排序 .....	200
10.5 搜索结果重排 .....	200
10.5.1 基于相关反馈的结果重排 .....	200
10.5.2 基于 PageRank 的结果重排 .....	202
参考文献 .....	204
<b>第 11 章 基于内容的视频检索技术 .....</b>	<b>205</b>
11.1 基于内容的视频检索技术的基础 .....	205
11.2 当前的基于内容的视频检索技术 .....	207
11.3 存在的问题 .....	209
参考文献 .....	210
<b>第 12 章 视频的结构化技术 .....</b>	<b>212</b>
12.1 镜头的边界检测 .....	213
12.1.1 非压缩域内镜头边界检测算法 .....	213
12.1.2 压缩域内镜头边界检测算法 .....	216
12.2 镜头关键帧的提取 .....	217

---

第 12 章	12.3 视频的特征提取	218
12.4 视频结构化中的关键技术	220	
参考文献	221	
<b>第 13 章</b>	<b>语音识别技术</b>	<b>222</b>
13.1 语音识别技术的发展历程	222	
13.2 语音识别系统的工作原理	223	
13.3 梅尔频率倒谱系数	224	
13.3.1 语音信号预处理	224	
13.3.2 离散 Fourier 变换	225	
13.3.3 取能量	226	
13.3.4 梅尔尺度滤波器组	226	
13.3.5 取对数	226	
13.3.6 离散余弦变换	226	
13.3.7 梅尔频率倒谱系数	226	
13.4 HMM 模型	227	
13.5 语言模型	228	
参考文献	229	
<b>第 14 章</b>	<b>视频的标注技术</b>	<b>230</b>
14.1 视频的标注技术概述	230	
14.2 特定领域内的视频标注技术	231	
14.2.1 视频场景分析	232	
14.2.2 视频精彩片段提取	232	
14.2.3 视频事件检测	232	
14.3 视频的多标签标注技术	233	
14.3.1 独立概念标注技术	234	
14.3.2 概念融合标注技术	234	
14.3.3 同时发现概念和概念间相互关系的标注技术	235	
14.4 主动学习方法在视频标注中的应用	236	
参考文献	238	

## 第 1 章 文本检索技术

20 世纪 90 年代,当基于 Web 的超文本应用兴起的时候,越来越多的人开始对文本检索技术感兴趣。随着时间的推移,基于 Web 的信息越来越多,如何在海量的信息中获取自己真正需要的信息成为一个巨大挑战。强大的搜索引擎 Google 能搜索上百亿之多的 Web 页面,主要原因是引擎采用了基于索引的检索技术,减少了搜索的响应时间。

### 1.1 基于索引的检索技术

在海量信息中获取自己真正需要的信息,顺序搜索的响应时间将变得不可忍受。解决搜索响应时间的办法是对文本文档库中的文本进行预处理,为文本文档库建立一种便于搜索的数据结构——索引。基于索引的检索技术(见图 1.1)非常适用于大规模、稳定的或周期性变化的文本文档库,如今绝大部分搜索引擎采用的都是基于索引的检索技术。

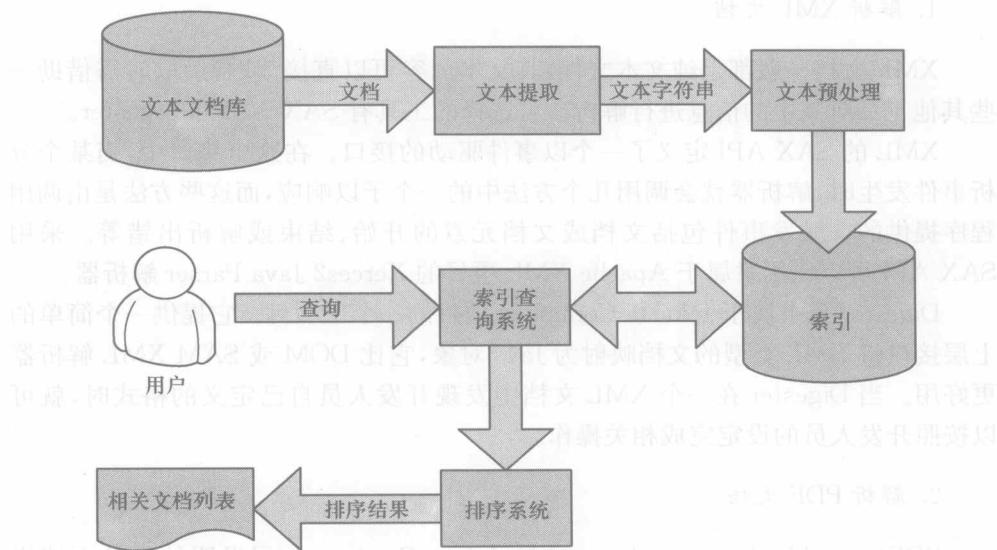


图 1.1 基于索引的检索技术

首先检索系统将所有的检索对象收集起来,构建集中的本地文本文档库。例如对于 Web 搜索引擎,其检索对象主要是 Web 网页,因此搜索引擎需要从互联网

上抓取尽可能多的网页保存到本地文本文档库中,一般这个过程由程序自动完成,本书不关注这个过程的细节,有兴趣的读者可以参考相关资料<sup>[1]</sup>。本地文本文档库构建完成之后,检索系统提取文本文档库中文档的文本字符串,并进行文本预处理。在有了文本预处理结果后,需要建立文档的索引。利用文档索引可以大大提高信息检索的速度。目前有多种建立文档索引的方法,但是对于大规模的文本文档库来说,用得最多的是倒排索引。在信息库的文档建立索引后,就可以对其进行检索。用户提交查询后,检索系统将直接访问索引。由于索引是一个可以便于搜索的数据结构,检索系统可以通过索引快速获得与查询相关的文档集合。在获取相关的文档后,由排序系统评价相关文档与查询的相关程度并对其排序,最后返回给用户。本章接下来的内容将对基于索引的检索技术的各个模块进行详细介绍。

## 1.2 文本提取

文本提取过程主要是提取各种格式文档中的字符串。文本检索系统不仅面向互联网的 Web 网页,还面向各种文档类型,例如 XML、PDF、Microsoft Word 或者 Excel 等类型的文档,下面以 XML、PDF 和 Microsoft Word 格式信息为例,介绍如何从多种常用格式的文档中提取文本内容。

### 1. 解析 XML 文档

XML 文档一般都是纯文本文件,其文本内容可以直接读取,读取时需借助一些其他工具对其中的信息进行解析。可选择的工具有 SAX API 和 Digester。

XML 的 SAX API 定义了一个以事件驱动的接口。在这个接口中,当某个分析事件发生时,解析器就会调用几个方法中的一个予以响应,而这些方法是由调用程序提供的。触发事件包括文档或文档元素的开始、结束或解析出错等。采用 SAX API 可以使用隶属于 Apache XML 项目的 Xerces2 Java Parser 解析器。

Digester 是隶属于 Jakarta Commons 项目的一个子项目。它提供一个简单的上层接口将 XML 类型的文档映射为 Java 对象,它比 DOM 或 SAM XML 解析器更好用。当 Digester 在一个 XML 文档中发现开发人员自己定义的格式时,就可以按照开发人员的设定完成相关操作。

### 2. 解析 PDF 文档

PDF(portable document format)是 Adobe Systems 公司发明的一种文档格式。这个格式打破了简单的文本内容的局限,它允许作者插入图片、超链接、颜色等。PDFBox 是由 Ben Litchfield 编写的一个免费开源库。使用该开源库,从 PDF 中获取文本内容变得非常简单,只需如下几条 Java 语句即可得到一篇 PDF 文档

中的文本内容：

```
PDFParser parser=new PDFParser(new InputStream(path));
parser.parse();
CosDocument cos=parser.getDocument();
PDFTextStripper stripper=new PDFTextStripper();
String text=stripper.getText(new PDDocument(cosDoc));
```

### 3. 解析 Microsoft Word 文档

与其他格式文档不同,Microsoft Word 文档的格式是保密的。Microsoft 公司对这个格式进行了保密处理,这使得其他人很难编写应用程序去读写 Microsoft Word 格式的文档。但是一些开源的项目解决了这个问题,例如 Jakarta POI 和 TextMining.org 项目。

POI 是 Jakarta 的一个子项目,它为操作基于微软的 OLE2 Compound Document 格式的各种文件格式提供了一个 Java API 包,使用 POI 可以从 Microsoft Word 文档中提取文本内容,同样也可以对 Excel 或其他 OLE2 Compound Document 格式的文档进行操作。从 Microsoft Word 中提取文本内容,只需如下几条 Java 语句即可:

```
WordDocument wd=new WordDocument(new InputStream(path));
StringWriter docTextWriter=new StringWriter();
wd.writeAllText(new PrintWriter(docTextWriter));
docTextWriter.close();
String text=docTextWriter.toString();
```

TextMining 包提供了类似于 Jakarta POI API 的接口,并对 Microsoft Word 文档的接口进行了优化。用 TextMining 只需一条语句即可获得 Word 的文本内容:

```
String t= new WordExtractor().extractText(new InputStream(path));
```

### 1.3 文本预处理

提取出文本字符串后,还需对文本字符串进行预处理以选择合适的词来建立索引。文本预处理首先将文本中包含的词分析出来,即分词。英文单词天然地被空格隔开,很容易切分,而中文的词是由连续的单个字符组成,因此中文分词相对困难一些,本书将在第 2 章重点介绍中文分词的相关技术。在语义表达方面并不是所有词的表达能力都是同等的,因此除分词之外,文本预处理还包括停用词删除、词干提取、索引词选择和建立词典等操作。

### 1.3.1 停用词删除

我们知道如果一个词在某个文本中多次出现,那么这个词就很有可能与文本的主题密切相关。然而如果一个词在多个文本中出现,而且频率过高,那么它对文本的区别能力就非常低。一般地,在文档库的文本中出现频率超过 80% 的词对检索过程根本起不到作用。这部分词被称为停用词(stopword)。在选择构建索引的词时,停用词需要被过滤,以提高索引效率。一般地,冠词、介词、连词等都是停用词,实际检索系统都会设置一个停用词表。

尽管删除停用词可以大大缩小索引空间的大小,一般可以缩小 40% 左右。删除停用词的缺点是可能会影响检索系统的查准率,有的文本检索系统为了克服这一缺点采用全文索引,并不剔除停用词,对所有的词都建立索引。

### 1.3.2 词干提取

词干提取是为了解决英文检索中存在的问题而采取的操作。词干是指将词的词缀(前缀和后缀)删除后剩下的部分。例如单词“compete”是它的变形“competes”、“competitor”、“competition”、“competing”和“competed”的词干。在英文检索中,如果用户输入的词是信息库中某个相关文本中词的一种变形,词的变形可以是该词的复数、动名词或者过去分词形式等,那么这些相关文本将被视作与查询无关的文本,这将大大影响召回率。为解决这个问题,在构建索引时,用词干来代替词干的所有变形。

词干提取不仅在很大程度上提高召回率,改善信息检索的性能,同时由于词干的众多变形都由词干代替,用于构建索引的词数量也大大减少,索引空间也进一步缩小。

目前,词干提取技术可以分为:词缀删除、表格查询、后续变形、N-连字。词缀删除技术比较直观、简单、有效。在词缀删除中,最重要的就是对词中后缀的删除,因为大多数词的变形是通过后缀来实现的。目前已经有多种关于词缀删除的算法,其中,Porter 算法以其简单性和有效性而得到广泛应用。表格查询技术通过在表格中查找某个词的词干来实现,表格中的信息依赖于整个语言中词的词干,因此通常需要相当大的存储空间来存放表格,这就制约了表格查询技术的应用。后续变形技术主要是通过结构化语言的知识来确定词素的边界,这种技术比词缀删除技术复杂。N-连字技术判断单词中的字母是否连在一起,这一过程实际上是词条聚类的过程。

### 1.3.3 索引词选择

如果采用全文索引,那么文档库中的所有词都要建立索引,而对有些语义表达