

情报检索语言与智能信息处理丛书

丛书主编 / 侯汉清

# 文本自动标引与

## 自动分类研究

章成志  
白振田 / 著



东南大学出版社  
SOUTHEAST UNIVERSITY PRESS

情报检索语言与智能信息处理丛书(丛书主编 侯汉清)

# 文本自动标引与 自动分类研究

章成志 白振田 著

东南大学出版社  
·南京·

## 图书在版编目(CIP)数据

文本自动标引与自动分类研究 / 章成志, 白振田著.  
南京: 东南大学出版社, 2009. 12

(情报检索语言与智能信息处理丛书/侯汉清主编)

ISBN 978 - 7 - 5641 - 1913 - 3

I. 文… II. ①章… ②白… III. 计算机应用—  
分类标引—研究 IV. G254 - 39

中国版本图书馆 CIP 数据核字(2009)第 200911 号

---

## 情报检索语言与智能信息处理丛书(侯汉清主编) 文本自动标引与自动分类研究

---

出版发行 东南大学出版社

出版人 江 汉

社 址 南京市四牌楼 2 号(邮编:210096)

印 刷 南京玉河印刷厂

责任编辑 李 正

(电话:025-83790887; E-mail:leezheng1978@sina.com)

经 销 新华书店

开 本 880 mm×1 230 mm 1/32

总印张 50.625(本册 6.625 印张)

总字数 1 310 千字(本册 173 千字)

版 次 2009 年 12 月第 1 版 2009 年 12 月第 1 次印刷

总 定 价 200.00 元(共 8 本)

---

\* 东大版图书若有印装质量问题, 请与读者服务部联系, 电话: 025 - 83792328

# 丛书总序

这部丛书包括下列八本专著：

- (1) 薛春香著《网络环境中知识组织系统构建与应用研究》；
- (2) 陆勇著《面向信息检索的汉语同义词自动识别》；
- (3) 杜慧平、仲云云著《自然语言叙词表自动构建研究》；
- (4) 章成志、白振田著《文本自动标引与自动分类研究》；
- (5) 张雪英著《情报检索语言的兼容转换》；
- (6) 刘华梅、戴剑波著《受控词表的互操作研究》；
- (7) 何琳著《领域本体的半自动构建及检索研究》；
- (8) 李运景著《基于引文分析可视化的知识图谱构建研究》。

这八本专著是侯汉清教授多年来指导博士生、硕士生们进行科学研究(有些是同他们合作研究)的具体成果的一部分。这些著作的主题内容，可以归结为“情报检索语言的自动化”和“自然语言检索”两个相关的问题，或者更概括地说，就是“信息检索自动化的升级问题”，属于当前信息检索学术研究的前沿课题。

这些专著，如果将其分散来看，或许不觉得分量之重；但如果把八本专著放到一起，就可以看出其成果之丰硕。侯汉清教授在带研究生中看准一个方向不断开拓、持之以恒的精神，可以出大成果，值得我们效法。南京农业大学在侯汉清教授领导下进行的有益的研究工作，我想一定会成为我国信息检索自动化发展史册之中浓浓的一笔。

这一类项目，本质上都是情报语言学的研究课题。所以，在研究中必须遵循情报语言学的理论，吸取情报语言学的已有成果，其结论应切合情报语言学的要求。它们只是利用计算机技术作为方法手段来达到研究目的而已，不能过分强调网络环境的特殊性而置情报语言学关于检索效率的基本要求于不顾。计算机技术应当与情报语言学密切结合。侯汉清教授和他的弟子们同时具备这两方面的知识，是顺利地较好地完成这些研究项目的关键。

这八个研究项目，大多采取实验研究法，故其成果具有较大的可信度和易理解性。其中有些项目，难度较大，甚至极难，专著只是作了认真、有益的探索；有些项目，虽然尚有一些不足，但作为中间成果，可在当前信息检索工作中推广应用，在应用中进一步完善。

信息检索自动化的初级阶段已在我国普遍实现。但要晋升一级，扩大自动化过程的范围和提高自动化的水平，当前的研究还属起步，发表的科研成果尚少见，学术研究有待扩大和深入。这部丛书起了很好的开拓作用，为继续研究打下了基础，是研究者很好的学习和参考用书，希望对此感兴趣的读者能从中受益。

张琪玉

2009年7月

# 序 言

随着因特网及其相关技术的不断普及深入,信息检索和文本挖掘成为人们获取知识的不可或缺的手段。但是面对异构的、动态的、海量的网络信息,如何快速找到用户感兴趣的信息并有效地加以利用是摆在我们面前亟待解决的重大问题。

近十年来信息检索技术得到了迅猛的发展,成为解决海量信息查找的有利工具。然而由于网络信息资源的更新速度和文本的非规范性,信息检索后的结果也往往数量巨大,“噪音”很多。所以需要快速有效的信息组织方式将这些信息分门别类加以分析处理,使用户能够更加有效地利用信息。文本的自动标引和自动分类是信息资源加工与组织中的关键技术。

《文本自动标引与自动分类研究》一书对文本挖掘的研究和应用给出了一个全新的视角。他们从情报检索的角度,以厚实的专业基础和丰富的领域经验,对于文本自动标引和自动分类技术进行了有益的探索和周密的实验,取得了丰硕的研究成果。

他们研究的特点是以数字图书馆、网络信息处理等为主要的应用背景,结合目前中文信息处理和文本处理的最新发展状况,在自动标引和自动分类方面有所继承,更有创新,主要研究特色在于:



一是自动标引技术的研究与应用，自动标引包括抽词标引与赋词标引两种类型。从信息组织的角度，自动标引相当于词汇级的摘要，而文本摘要则是句子级的摘要。而前者由于短小精悍、代表性强、使用效率高，成为信息检索、信息抽取、问答系统、主题识别和跟踪的有利工具。作者对自动标引中抽词词典构造、基于多特征选择及权值计算、标引源权重设置等问题进行了深入的研究，给出了高效的自动标引算法。

二是采用基于知识库的文本自动分类方法，以大规模的人工标引关键词词串与分类号对应记录为基础，生成关键词词串-分类号对应的分类知识库，并将其用于新文本的自动分类；对于文本分类难点问题——多层次分类问题，进行了深入的研究，给出了高效的层次分类算法。分类体系为《中国图书馆分类法》和行业分类表，分类体系庞大，总的分类类目数量多达到成千上万的规模。

很荣幸，能够先睹为快，拜读大作，掩卷之余，感受颇多。

目前在情报检索领域，专门讨论自动标引和自动分类的专著很少，这本书将情报检索技术和海量文本处理技术相结合，将自然语言处理和机器学习等引入情报检索领域，系统全面地介绍了文本自动标引和自动分类的研究工作，为情报检索的研究探索出有价值的研究方向，为文本挖掘的研究提供了很好的应用平台。

林鸿飞

2009年7月25日于大连理工大学创新园大厦

# 目 次

## 第一部分

|                                |    |
|--------------------------------|----|
| <b>第1章 引言</b> .....            | 3  |
| 1.1 研究背景 .....                 | 3  |
| 1.2 自动标引与自动分类的作用 .....         | 5  |
| 1.3 本书的内容与章节安排 .....           | 7  |
| <b>第2章 文本自动标引与分类研究进展</b> ..... | 10 |
| 2.1 自动标引研究综述.....              | 10 |
| 2.2 文本分类研究综述.....              | 23 |
| 2.3 本章小结.....                  | 33 |

## 第二部分

|                                |    |
|--------------------------------|----|
| <b>第3章 文本分词技术及抽词词典构造</b> ..... | 45 |
| 3.1 文本分词技术概述.....              | 45 |
| 3.2 分词模式设计及其原理.....            | 48 |
| 3.3 原始抽词词典的构造.....             | 49 |
| 3.4 词典约简算法实验.....              | 51 |
| <b>第4章 基于多特征选择及权值计算</b> .....  | 57 |
| 4.1 特征选择方法概述.....              | 57 |
| 4.2 算法设计原理.....                | 62 |
| 4.3 结果分析.....                  | 66 |



|                           |     |
|---------------------------|-----|
| <b>第5章 自动标引中标源权重方案确定</b>  | 68  |
| 5.1 标引源权重研究综述             | 68  |
| 5.2 标引源权重方案的确定            | 69  |
| 5.3 本章小结                  | 78  |
| <br><b>第三部分</b>           |     |
| <b>第6章 分类知识库的制作</b>       | 83  |
| 6.1 概述                    | 83  |
| 6.2 关键词(串)一分类号关联研究综述      | 85  |
| 6.3 关键词(串)一分类号关联方法        | 90  |
| 6.4 分类知识库的制作              | 93  |
| 6.5 分类知识库的性能测评            | 101 |
| 6.6 篇名知识库的制作              | 105 |
| 6.7 本章小结                  | 109 |
| <b>第7章 基于语义体系的词语相似度计算</b> | 111 |
| 7.1 概述                    | 111 |
| 7.2 词语相似度研究综述             | 112 |
| 7.3 基于语义体系的词语相似度算法        | 116 |
| 7.4 基于语义相似度的同义词挖掘         | 128 |
| 7.5 本章小结                  | 137 |
| <b>第8章 基于知识库的文本自动分类</b>   | 141 |
| 8.1 文本自动系统总体设计            | 141 |
| 8.2 文本自动分类系统的测评           | 143 |
| 8.3 《全国报刊索引》自动标引与自动分类系统介绍 | 146 |
| 8.4 本章小结                  | 147 |
| <br><b>第四部分</b>           |     |
| <b>第9章 统计与决策规则双重分类算法</b>  | 151 |
| 9.1 分类器概述                 | 151 |

---

|                                          |            |
|------------------------------------------|------------|
| 9.2 双重分类原理 .....                         | 158        |
| 9.3 分类规则提取 .....                         | 160        |
| 9.4 双重分类过程 .....                         | 161        |
| 9.5 实验结果及分析 .....                        | 163        |
| <b>第 10 章 层次分类算法实验 .....</b>             | <b>166</b> |
| 10.1 层次分类原理 .....                        | 166        |
| 10.2 层次分类算法设计 .....                      | 168        |
| 10.3 实验结果及分析 .....                       | 171        |
| <b>第 11 章 基于统计与规则相结合的文本分类系统的实现 .....</b> | <b>174</b> |
| 11.1 系统实验用语料选择及分析 .....                  | 174        |
| 11.2 系统总体框架与模块介绍 .....                   | 178        |
| 11.3 系统测试分析 .....                        | 181        |
| 11.4 本章小结 .....                          | 182        |
| <b>名称索引 .....</b>                        | <b>185</b> |
| <b>主题索引 .....</b>                        | <b>190</b> |
| <b>后记 .....</b>                          | <b>197</b> |

## 图表目次

|                                         |     |
|-----------------------------------------|-----|
| 图 1-1 文本挖掘任务框架 .....                    | 5   |
| 图 1-2 本书章节安排示意简图 .....                  | 7   |
| 图 2-1 术语、主题词、标引词包含关系图 .....             | 12  |
| 图 2-2 信息描述的颗粒度 .....                    | 12  |
| 图 2-3 自动标引研究路线图 .....                   | 16  |
| 图 2-4 基于机器学习的自动抽词方法逻辑视图 .....           | 18  |
| 图 4-1 不同权值计算方法的实验结果 .....               | 66  |
| 图 5-1 统计工作流程图 .....                     | 72  |
| 图 6-1 《中图法》分类知识库构建流程图 .....             | 94  |
| 图 6-2 分类知识库样例(规模:8万余条) .....            | 99  |
| 图 7-1 《词林》语义空间 .....                    | 117 |
| 图 7-2 语义距离的计算 .....                     | 118 |
| 图 7-3 最短路径计算原型示意图 .....                 | 119 |
| 图 7-4 词汇间的语义相似度计算流程图 .....              | 121 |
| 图 7-5 同义词挖掘系统流程图 .....                  | 132 |
| 图 7-6 同义词挖掘系统结构图 .....                  | 132 |
| 图 7-7 同义词挖掘系统界面 .....                   | 133 |
| 图 7-8 同义词测试界面 .....                     | 133 |
| 图 7-9 同义词挖掘界面 .....                     | 133 |
| 图 7-10 同义词挖掘用数据库维护界面 .....              | 133 |
| 图 8-1 中文文本自动标引和分类系统结构图 .....            | 142 |
| 图 8-2 《全国报刊索引》自动标引与自动分类系统<br>主界面 .....  | 147 |
| 图 9-1 基于最短距离法与规则匹配法的双重分类过程示<br>意图 ..... | 162 |
| 图 10-1 类别体系的层次结构示意图 .....               | 167 |

|                                        |     |
|----------------------------------------|-----|
| 图 10-2 多层次分类过程流程示意图 .....              | 170 |
| 图 10-3 分类过程示例图 .....                   | 171 |
| 图 11-1 不同训练样本的分类结果对比 .....             | 176 |
| 图 11-2 取不同维数后的分类结果对比 .....             | 177 |
| 图 11-3 系统功能结构图 .....                   | 178 |
| 图 11-4 系统总体流程图 .....                   | 179 |
| 图 11-5 词典维护模块 .....                    | 180 |
|                                        |     |
| 表 2-1 近五十年比较有代表性的自动标引方法 .....          | 13  |
| 表 2-2 自动标引方法的分类 .....                  | 16  |
| 表 2-3 国外较有代表性的自动分类研究(包括相关<br>系统) ..... | 26  |
| 表 2-4 国内较有代表性的自动分类研究(包括相关<br>系统) ..... | 30  |
| 表 3-1 不同词典对分类结果的影响 .....               | 55  |
| 表 4-1 各权值计算方法结果示例 .....                | 65  |
| 表 5-1 主题表达能力的抽样统计数据来源一览表 .....         | 70  |
| 表 5-2 自动标引标引源统计表样例 .....               | 71  |
| 表 5-3 自动标引词频统计样例 .....                 | 71  |
| 表 5-4 文章字数分布情况统计 .....                 | 72  |
| 表 5-5 文章段落数分布情况统计 .....                | 73  |
| 表 5-6 自动标引词数统计 .....                   | 73  |
| 表 5-7 标引源人工打分结果统计 .....                | 74  |
| 表 5-8 标引源人工打分统计(300 篇经济类文本) .....      | 75  |
| 表 5-9 样本标引词数分布情况表 .....                | 77  |
| 表 6-1 标引经验知识库中的关键词—分类对应形式<br>举例 .....  | 89  |
| 表 6-2 事件 A、B 的可能出现频次表 .....            | 92  |
| 表 6-3 经去重处理后的记录样例 .....                | 95  |

|                                     |     |
|-------------------------------------|-----|
| 表 6-4 经过分类辅助用词过滤处理后的数据样例 .....      | 96  |
| 表 6-5 关键词权值处理结果样例 .....             | 97  |
| 表 6-6 排序后结果样例 .....                 | 97  |
| 表 6-7 相关度计算结果样例 .....               | 98  |
| 表 6-8 数据库规模与强规则对应表 .....            | 99  |
| 表 6-9 分类知识库实际标引测试对照表 .....          | 101 |
| 表 6-10 分类知识库数量比较 .....              | 101 |
| 表 6-11 分类知识库抽样统计结果表 .....           | 103 |
| 表 6-12 篇名知识库样例 .....                | 108 |
| 表 7-1 词汇语义相似度计算结果样例 .....           | 127 |
| 表 7-2 封闭实验结果对照表 .....               | 134 |
| 表 7-3 开放实验结果对照表 .....               | 135 |
| 表 7-4 同义词挖掘运行效率对照表(单位:秒) .....      | 136 |
| 表 8-1 文本自动标引和分类(全文)、自动标引结果比较表 ..... | 144 |
| 表 8-2 自动分类(全文)与人工分类结果比较表 .....      | 145 |
| 表 8-3 文本自动标引和分类(简化)、自动标引结果比较表 ..... | 145 |
| 表 8-4 自动分类(简化)与人工分类结果比较表 .....      | 146 |
| 表 8-5 分类知识库规模 .....                 | 147 |
| 表 9-1 知识表示例表 .....                  | 155 |
| 表 9-2 决策表 .....                     | 157 |
| 表 9-3 向量距离分类结果示意(片段) .....          | 159 |
| 表 9-4 决策信息表 .....                   | 161 |
| 表 9-5 修正正确的部分示例结果 .....             | 163 |
| 表 9-6 加入规则补充分类的测试结果 .....           | 164 |
| 表 10-1 非层次分类下的分类结果(片段) .....        | 171 |
| 表 11-1 训练样本的分类合理性分析 .....           | 177 |
| 表 11-2 系统总体实验结果评估 .....             | 182 |

# **第一部分**

---

第1章 引言

第2章 文本自动标引与分类研究进展



# 第1章 引言

## 1.1 研究背景

随着计算机及网络的普及,数字化载体逐渐融入人们的生产、生活中,成为人们获取信息资源不可或缺的途径、方法和手段。根据第 23 次中国互联网络发展状况统计报告显示,目前我国上网网民已经突破 1 亿,网民用于上网的时间每周平均在 14 个小时以上<sup>[1]</sup>。我国网民人数的增加、上网时间的增长,从一个侧面说明,数字化载体十分具有吸引力,能够方便、快捷地为人们提供所需要的信息资源。在过去的 20 年中,万维网的迅速发展使其成为世界上规模最大的公共数据源。万维网数据量巨大且不断增长、数据类型丰富、信息异构、信息包含噪音等特点,使得挖掘有用的信息和知识的任务变得十分有趣,并富有挑战<sup>[2]</sup>。

我们正处于“信息爆炸”的时代,但为什么当各类信息像洪水

一样向我们涌来时,我们仍然缺乏所需要的信息呢?这是因为在信息社会之中,“没有控制和没有组织的信息不再是一种资源。它倒反而成为信息工作者的敌人”<sup>[3]</sup>。

然而,在这些纷繁复杂的信息资源中,最主要的还是非结构化或半结构化的文本信息资源。人们上网获取信息资源的要求是快捷方便,而要快捷方便,通常的做法是对文本信息资源进行预先处理,运用某种方式组织和存贮起来。如何对异构、动态的海量信息资源进行快速的加工与组织,以智能化、个性化的方式为用户提供高效的信息服务,是信息组织部门、信息组织研究者等共同面临的难题。

一方面,数字化信息资源数量高速增长,其中包含着对人们极有潜在价值的知识和信息;另一方面,人们运用网络获取信息资源的数量也在高速增长,而人们获得的有效信息资源的比例却在下降。其原因除在于互联网的政策——任何人可以发布任何未经加工的信息,这些未经加工的信息难以获取效率,而经过加工的信息能够提高人们的获取效率。未经加工的信息越多,人们的信息资源检索效率就越低。

然而,要解决信息资源无限增长和检索效率低下的矛盾,究其原因是多方面的,非某一种技术所能解决,存在的困难也是多方面的。但主要的原因还是信息资源的多样化和海量化,且没有经过规范的加工处理。在这些杂乱无序的信息海洋中,用户要迅速准确地找到自己所需要的信息,困难重重。

关于信息资源的加工与组织方法比较多,其中文本的自动标引和自动分类是比较关键的技术,并且有广泛的应用。文本自动标引(本书是指狭义上的自动标引,即文本的关键词抽取或主题词获取)是对文本根据其表达的内容或主题,进行关键词或主题词自动获取的过程。分类是人们浏览和查找信息的主要手段之一。文本自动分类是根据某一特定的分类体系,将文本资源分门别类地进行自动组织的方法。

本书从文本的自动标引和自动分类角度出发,调研文本自动标引和自动分类技术的相关理论研究,结合实际应用,以中文文本