

<<智能信息处理>>

图书基本信息

书名：<<智能信息处理>>

13位ISBN编号：9787030291356

10位ISBN编号：7030291352

出版时间：2010-10

出版时间：科学出版社

作者：郑家恒 等著

页数：318

版权说明：本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问：<http://www.tushu007.com>

## 前言

从20世纪90年代开始,国际自然语言处理领域发生了一些重大变化,重要特征之一就是转向对大规模真实文本的研究和处理。

以大规模真实文本为基础的语料库研究和知识自动获取受到高度重视。

显然,大规模真实文本的处理是计算语言学今后一个时期的战略目标,建设高质量的大规模语料库是中文信息处理领域的基础性工程。

基于语料库的语言研究是计算语言学的一个重要领域,语料库的建立为语言学的研究提供了丰富的语言现象,为计算语言学学者从加工的语料库中获取语言知识、建立语言模型、研究语言信息处理技术提供了翔实的语言信息数据。

作为研究资源的语料库的价值是通过对语料的加工来体现的,对语料库加工的层次越高,语料库的应用价值就越高。

希望本书的出版能促进语料库加工方法和技术的发展,为基于语料库的相关研究和应用提供支撑。

作者及其课题组从事语言信息处理的教学与研究已有二十多年。

近年来,作者有幸承担了若干国家863计划项目(中文文本自动切词和词性标注软件及其评测技术研究(863-306-03-09-4)、大规模中文文本语料库深加工质量检验技术研究(2001AAlI4031))、国家自然科学基金项目(大规模中文文本语料库分词与词性标注一致性检验技术研究(60473139)、基于中文文本的计算机中介通信中欺骗检测研究(60775041))、省部级项目及横向合作项目等。

这些项目的研究成果为本书的编写提供了关键性支持。

多年来,刘开瑛、黄昌宁等诸位学术前辈都为作者的相关研究思路和方法提供了许多指导。

本书编写过程中,山西大学梁吉业、李德玉、李茹、王文剑、王素格等教授为作者提供了多方面的支持。

魏善德、任玉、魏莉、魏丽霞、樊勇、王振宇、刘博、张剑锋、何苑、温艳霞、毋菲等同学也为本书的出版做了许多文字校对方面的工作,谨在此一并表示深深的感谢。

## 内容概要

本书以作者主持的国家项目、省部级项目及合作项目等为依托，以课题组近年来的研究成果为基础，重点介绍语料库深加工中的若干技术和方法，涉及分词、词性标注、句法分析、语义标注以及相关加工中的自动校对和一致性检验技术。

同时，对语料库加工质量的评价技术和语料库的相关应用做了详细介绍。

各章节的顺序展示了语料库加工中由浅入深的发展过程。

本书可作为计算机、语言学等专业高年级本科生、研究生教材，也可作为自然语言处理和计算语言学研究人员的参考书。

## 书籍目录

《智能科学技术著作丛书》序前言第1章 绪论 1.1 语料库的定义和作用 1.1.1 什么是语料库 1.1.2 语料库的作用 1.2 语料库的建立 1.2.1 什么是语料库标注 1.2.2 语料库标注的原则 1.2.3 建立语料库需要考虑的几个问题 1.2.4 语料库标注和建立的方法 1.2.5 语料库的质量检验 1.3 本书的编排 参考文献第2章 自动分词 2.1 自动分词概述 2.1.1 自动分词的意义 2.1.2 自动分词的主要难点 2.1.3 自动分词方法简介 2.1.4 自动分词评测 2.2 分词规范 2.2.1 制定分词规范的目的和意义 2.2.2 几种典型的分词规范介绍 2.3 歧义字段的切分技术 2.3.1 歧义字段现象分析 2.3.2 基于统计的歧义字段排歧 2.4 未登录词识别 2.4.1 专有名词识别 2.4.2 新词语识别 2.5 缩略语识别 2.5.1 缩略语特征分析 2.5.2 缩略语资源库的建立 2.5.3 缩略语识别模型 2.5.4 缩略语的还原 2.6 分词一致性检验 2.6.1 分词不一致性现象分析 2.6.2 基于规则的分词一致性检验方法 2.6.3 基于统计的分词一致性检验方法 2.6.4 分词一致性检验系统 参考文献第3章 词性标注 3.1 词性标注概述 3.1.1 词性标注的意义 3.1.2 词性标注的难点 3.1.3 词性标注方法简介 3.1.4 常用语料库 3.2 词性标注规范 3.2.1 制定词性标注规范的目的和意义 3.2.2 几种典型的词性标注规范介绍 3.3 兼类词的标注 3.3.1 什么是兼类词 3.3.2 典型的兼类词标注方法 3.4 词性标注一致性检验 3.4.1 问题描述和分析 3.4.2 一致性检验模型的建立 3.4.3 实验结果和分析 3.4.4 方法评价 3.5 词性标注自动校对 3.5.1 基于分类的词性标注自动校对 3.5.2 基于决策表的词性标注自动校对 参考文献第4章 句法分析 4.1 完全句法分析 4.1.1 完全句法分析概述 4.1.2 形式语法体系 4.1.3 树库资源的建设 4.1.4 汉语句法分析的特点 4.1.5 句法分析方法 4.1.6 相关会议及评测 4.1.7 句法分析模型的评价方法 4.2 浅层句法分析 4.2.1 浅层句法分析概述 4.2.2 组块库的获取 4.2.3 组块的类型及其标注规范 4.2.4 组块分析方法 4.2.5 相关会议及评测 4.2.6 评价参数 4.3 句法树库的一致性检验 4.3.1 不一致现象分析 4.3.2 不一致的发现和消解 参考文献第5章 语义标注语料库 5.1 语义标注范围 5.1.1 词义标注 5.1.2 句义标注 5.1.3 篇章级的语义标注 5.2 语义标注语料库的建立方法 5.2.1 传统的以人工标注为主的方法 5.2.2 自动构建语义标注语料库 5.3 主要的语义标注语料库 5.3.1 词义标注语料库 5.3.2 句义标注语料库 5.3.3 语篇关系标注语料库 5.3.4 时间关系标注语料库 5.3.5 信息抽取方面的语料库 5.3.6 生物医药领域中的语义标注语料库 参考文献第6章 语料库评测 6.1 语料库评测的意义 6.2 语料库分词质量评价 6.2.1 评价样本的抽样 6.2.2 抽样样本的聚类及评价 6.2.3 实验及分析 6.3 语料库可用性评价 6.3.1 可用性评价体系 6.3.2 可用性评价计算 6.3.3 评价结果分析 参考文献第7章 基于语料库的应用研究 7.1 网页信息处理 7.1.1 重复网页分析 7.1.2 基于语义的网页去重 7.1.3 基于网页文本结构的网页去重 7.2 特殊领域的信息抽取 7.2.1 基于HMM的农业信息抽取 7.2.2 基于NLP的土壤污染数据抽取 7.2.3 基于BOotstrapping的交通工具名识别 7.3 基于大规模语料库的汉语韵律边界研究 7.3.1 基于统计语言模型建立二叉树结构 7.3.2 基于树结构的汉语韵律边界预测 7.4 基于大规模语料库的欺骗行为检测 7.4.1 欺骗性语料库的建设 7.4.2 欺骗检测的特征线索 7.4.3 文本特征抽取 7.4.4 欺骗行为检测方法 7.4.5 实验结果和分析 参考文献

## 章节摘录

插图：关于语料库（corpus）的定义主要有以下几种：（1）McEnery和Wilson指出：“总体来说，多篇文本的集合就是语料库，但在现代语言学中使用语料库这个术语时，更倾向于包含更多的内涵，主要有采样（sampling）收集、有代表性（representativeness）、规模有限（finite size）、机器可读（machine-readable）、标准参考数据（a standard reference）等内涵特征。”

（2）语料库就是某种语言在实际运用中的大量实例集合，这些例子可以是书面文本，也可以是语音形式的文本。

（3）语料库是根据外部原则选择的电子形式的文本或文本片段的集合。

该集合能够代表一种语言，或一种语言的分支，或一种语言的变体，并可作为语言学研究使用的数据来源[引]。

这里外部原则（external criteria）是指通过文本的交流功能来选择文本的原则。

与外部原则相对的一个概念就是内部原则（internal criteria），具体指按照文本反映的语言细节来选择文本。

在上述的几种定义中，定义（1）使用最多，认为语料库不是简单收集的文本集合，而是通过采样收集，具有代表性，规模大小可以确定，是机器可读的标准数据。

但是Kilgarriff和Grefenstette提出了异议，认为McEnery和Wilson混淆了“什么是语料库”和“什么是好的、适合于某项语言研究的语料库”这两个问题，他们认为语料库就是文本的集合。

然而在具体使用中，有些研究者认为有许多文本的集合并不一定是语料库。

最具有争议的莫过于万维网（WWW）了。

WWW刚出现时，人们因为不了解搜索引擎，也不清楚对WWW如何采样，觉得WWW相当神秘。

因此，文献指出：“WWW不是语料库，因为其维度未知且不断变化，而且WWW最初也不是从语言学角度来设计的。”

编辑推荐

《智能信息处理:汉语语料库加工技术及应用》：智能科学技术著作丛书

版权说明

本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问:<http://www.tushu007.com>