



- 什么是语料库
- 语料库的类型
- 语料库的基本加工过程与规范
- 语料库的应用
- NLTK



■ 语料库

存放语言材料的数据库(文本集合),其复数形式为 Corpora。库中的文本通常经过整理,具有既定的格式和标记 ,特指计算机存储的数字化语料库

- 语料库的基本认识
 - 存放的是在语言实际使用中真实出现过的语言材料
 - 以计算机为载体承载语言知识的基础资源
 - 真实语料需要经过加工(分析与处理)



2.1什么是语料库语言学

- 语料库语言学
 - 根据篇章材料对语言的研究称为语料库语言学。[Aijmer, 1991]
 - 基于现实生活中语言运用的实例进行的语言研究称为语料库语言学。一[McEnery, 1996]
 - 以语料为语言描写的起点或以语料为验证有关语言的假说的 方法称为语料库语言学。[Crystal, 1991]

语料库语言学主要研究机器可读自然语言文本的采集、存储、检索、统计、语法标注、句法语义分析,以及具有上述功能的语料库在语言定量分析、词典编纂、作品风格分析、自然语言理解和机器翻译等领域中的应用。



语料库语言学的历史

- 第一代(70-80)
 - 百万词级
 - 以语言研究为导向
- 第二代 (80-90)
 - 千万词级
 - 词典编纂-应用为导向
- 第三代(1990-)
 - 超大规模(上亿词级)
 - -标准编码体系
 - -深度标注/多语种
 - –NLP应用
- 第四代(?)
 - 互联网作为语料库



2.1语料库语言学

- 语料库语言学研究的内容:
 - 语料库的建设与编纂
 - 语料库的加工和管理技术
 - 语料库的使用



第一代语料库(1)

- Brown语料库
 - 始建于1960年代初
 - W.N.Francis和H.Kucera发起
 - 美国Brown大学建立
 - 世界上第一个根据系统性原则采集样本的标准语料库
 - 主要代表当代美国英语
 - 规模100万词次



第一代语料库(2)

- LOB语料库
 - 始建于1970年代初
 - 由英国Lancaster大学著名语言学家Geoffrey Leech倡议
 - 挪威Oslo大学StigJohansson主持完成
 - 安装在挪威Bergen大学挪威人文科学计算中 心
 - 规模于Brown语料库相当
 - 主要代表当代英国英语



第二代语料库(1)

- COBUILD语料库
 - 建于1980年代
 - ■以词典编撰为应用背景
 - 有英国Birminghan大学与Collins出版社合作完成
 - 规模达2000万词次
 - 基于该语料库出版的Collins Cobuild词典(1987)受到了广泛的好评



第一代语料库(3)

- London-Lund语料库
 - 1960年代初,由Randolph Quirk主持
 - 收集**2000**小时的谈话和广播等口语素材并整 理成书面材料
 - 由瑞典Lund大学J. Svartvik主持全部录入计算机
 - 1975年建成



第二代语料库(2)

- Longman语料库
 - 建于1980年代
 - ■包括三个语料库
 - LLELC语料库(Longman/Lancaster英语语料库)
 - LSC语料库(Longman口语语料库)
 - •LCLE(Longman英语学习语料库)
 - 目标是编撰英语学习词典,为外国人学习英语服务
 - 规模达5000万词次



第三代语料库(1)

- ACL/DCI语料库
 - 美国ACL倡议发起
 - 收集语料范围广泛
 - 华尔街日报
 - Collins英语词典
 - •Brown语料库
 - PennTreeBank
 - 既有己标注的语料,也有未标注语料
 - 制定了语料库文件的格式标注
 - 采用统一的SGML标注语言
 - 语料标注依照TEI(Text Encoding Initiative)标准



中文语料库

■ 中文语料库建设状况

自**1979**年以来,中国就开始进行机器可读语料库的建设,早期在中国建立的主要机器可读语料库有:

汉语现代文学作品语料库(1979年),527万字,武汉大学现代汉语语料库(1983年),2000万字,北京航空航天大学中学语文教材语料库(1983年),106万8千字,北京师范大学现代汉语词频统计语料库(1983年),182万字,北京语言学院

早期语料库有如下特点:

- ◆ 多数采用手工输入方式建立,耗时耗力,缺乏规范,规模较小,重用性差。
- ◆ 发现了汉语文本切分歧义的两种类型。即交集型歧义和多义组合型歧义。
- ◆ 建立了初步的分词规范。即GB13715《信息处理用现代汉语分词规范》。



第三代语料库(2)

- PennTreeBank (宾州大学树库)
 - 美国Pennsylvania大学1980年代末开始发起
 - 由该校计算机系M.Marcus主持
 - **1993**年,完成了对近**300**万英语词的句子语 法结构标注
 - 2000年完成了中文树库(第一版): 10万 词次,4185个句子



中文语料库

1991年,国家语言文字工作委员会开始建立国家级的大型汉语语料库,计划期规模达到7000万字。

现在该课题已经结项,国建语委语言文字应用研究所成立了"汉语语料库深加工"的课题组,准备对国家级语料库的**2000**万字的核预料进行深加工,逐步把这个生语料库变为熟语料库。

1992年以来,大量的语料库在中国研究中文信息处理的单位建立起来,建设了大规模真实文本语料库的单位有:

《人民日报》光盘数据库:

北京大学计算语言学研究所;

北京语言文化大学

清华大学:

山西大学 等。



2.2语料库的类型

- 按照语料库的研究目的和用途
 - 通用的(systematic)
 - 充分考虑语料的动态和静态问题、代表性和平衡问题以及语料库的规模等问题。
 - 专用的(specialized)
 - 如:北美的人文科学语料库



2.2语料库的类型

- 平衡语料库
 - ■平衡语料库着重考虑语料的代表性与平衡性。
 - 语料采集的七项原则:语料的真实性、可靠性、科学性、代表性、权威性、分布性和流通性。其中,语料的分布性还要考虑语料的科学领域分布、地域分布、时间分布和语体分布等。 [张普,2003]



2.2语料库的类型

□ 按语言种类划分

- ◆ 单语的
- ◆ 双语的或多语的 篇章对齐 / 句子对齐 / 结构对齐
- □ 是否标注?

两个术语:

- 具有词性标注
- > 生语料
- 句法结构信息标注(树库)
- ▶ 熟语料
- 语义信息标注



2.2语料库的类型

- 平行语料库
 - 两种含义:一种是指在同一种语言的语料.上的平行,例如正在建立的"国际英语语料库",共有20个平行的子语料库,分别来自以英语为母语或官方语言和主要语言的国家,如英国、美国、加拿大、澳大利亚、新西兰等。其平行性表现为语料选取的时间、对象、比例、文本数、文本长度等几乎是一致的。建库的目的是对不同国家的英语进行对比研究。



2.2语料库的类型

》另一种平行语料库是指在两种或多种语言 之间的平行采样和加工,例如,机器翻译中 的双语对齐语料库

C: 早晨好!

C: 早晨1好2!3

E: Good, morning, .3

E: Good morning.

C: 您能给我一杯咖啡吗?

E: Could you give me a cup of coffee?

... ..



2.2语料库的类型

- + 共时语料库与历时语料库.
 - 所谓共时语料库是为了对语言进行共时(同一时段)研究而建立的语料库。研究大树的横断面所见的细胞和细胞关系,即研究一个共时平面中的元素与元素的关系。

4

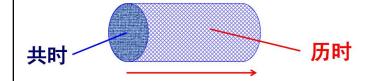
已有的双语资源

- 加拿大议会会议录 录 (Canadian Hansards) http://www.isi.edu/naturallanguage/download/hansard/
- 克姆尼茨英- 德翻译库 语料库 (Chemnitz E-G Translation Corpus) http://www.tu-chemnitz.de/phil/english/chairs/linguist/real /independent/transcorpus/index.htm
- 英语- 挪威语平行库 语料库 (ENPC)
 https://www.hf.uio.no/ilos/english/services/omc/enpc/
- 葡- 英双向平行库 语料库 (Compara)
 http://www.linguateca.pt/COMPARA/Welcome.html
- 香港立法委员会录 会议记录 (Hong Kong Hansards) http://catalog.ldc.upenn.edu/LDC2000T50
- 新闻 香港新闻 (Hong Kong News)
- 法律 香港法律 (Hong Kong Laws)



2.2语料库的类型

历时语料库是为了对语言进行历时研究而建立的语料库。研究大树的纵剖面所见的每个细胞和细胞关系的演变,即研究一个历时切面中元素与元素关系的演化





2.2语料库的类型

- 判断历时语料库的4 条原则 [张普, 2003]
 - 是否 动态: 语料库必须是开放的、动态的
 - 文本是否具有量化的流通度属性: 所有的语料都应来源于大众传媒, 具有与传媒特色相应的流通度属性。 其量化的属性值也是动态的。
 - 深加工是否基于动态的加工方法: 随语料的动态变化采集, 并进行动态地加工。
 - 是否取得动态的加工结果: 语料的加工结果也应是 动态的和历时的。



生语料与熟语料

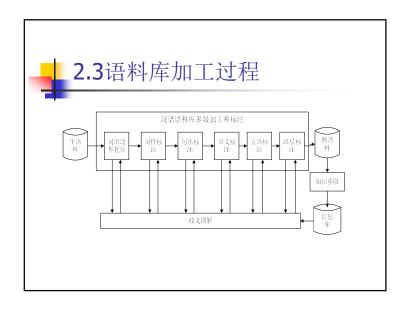
文本通常有两种形式:生文本和标注文本。"标注"这个术语是指把某种分类代码插入到一个计算机文件中,这种分类代码通常并不是文件的组成部分,但是通过这些分类代码我们可以了解文件的结构或格式信息。

由于语料库的来源复杂语料中可能存在无法处理 的各种各样的格式或者内容,如:文档页眉和分隔符、 排版代码、表和图标等。因此在进一步对语料处理前 需要一个过滤器来过滤掉这些无用内容。



2。**3**汉语语料库的基本加工规 范

- 生语料与熟语料
- 汉语语料库的加工思路
- 汉语语料库的加工规范
- 汉语文本词性标注标记集





生语料与熟语料

制定《现代汉语语料库加工手册——词语切分与词性标注》的基本思路如下:

- 1.词语切分的规范应尽可能的同已有的中国国家标准GB13715《信息处理用现代汉语分词规范》保持一致。
- 2.词性规范使用小标记集。除了使用《语法信息词典》中的26个词 类标记外,增加了专有名词分类标记、 语素g按其子类标注以及 动词和形容词的某些功能标记。
- 3.与已有资源的配合。



汉语语料库加工规范

标注规范:标注规范包括一般词性的标注和专有名词的标注。

- 一般词性的标注即切分单位的标记。包括:
- 1.标记集以26个词类标记为基准,名动词、副动词、名形词和专有名词的标记是在动词代码V、形容词代码a、名词代码n后增加一个小写字母,语素标记是在语素代码q前增加一个大写字母。
- 2.一个词语若在词法词典中已属于一个或若干个词类,标注时不要 轻易增加词性。
- 3.当语法词典给某个词确定的词性确实不对或不完备时,当然也要 修正或补充。



汉语语料库加工规范

北大计算语言所制定的《现代汉语预料库加工——词语切 分与词性标注规范》分为三个部分:**切分规范,标注规范, 切分和标注相结合的规范。**

切分规范:主要规定将汉字串形式的句子切分为词序列的原则,即 什么样的汉字组合可以作为一个切分单位。

进行切分通常要有一部"分词词典"。北大的《语法信息词典》收录的词条已超过7.3万,本规范规定词典的词条一般都是切分单位,这使得对"切分单位"的把握有了基本的参照。但处理大规模真实文本时总会遇到词典中没有的"未定义词"



汉语语料库加工规范

- 4. 即使语法词典中的简称实际上是指团体、机构、组织名称或地名, 标注时仍标以**j**,而不是改为**nt**或**nz**。
- 5. "唐朝"、"宋代"等历史朝代虽然也是专名,因语法词典已作为时间收入,标注时仍标以t,不改为nz。

这里的专有名词标注,对"专有名词"含义进行了扩展。短语型的地名、团体机构名称及其他专有名词在词的切分基础上用ASCII的方括号括起来,并在右方括号后标注以相应的ns,nt,nz,方括号不嵌套。



汉语语料库加工规范

切分和标注相结合的规范:规范中还给出了一些基于词性描述的 构词规则,规定了什么样的组合处理为一个切分单位,并给出了新 组合词的词性。

由语素构造合成词的方式有"复合"、"附加"、"重叠",但运用这三种方式将两个成分结合成一个较大的单位时,这个单位是否作为切分单位,不可一概而论。以"附加"方式的后接成分"者"为例:

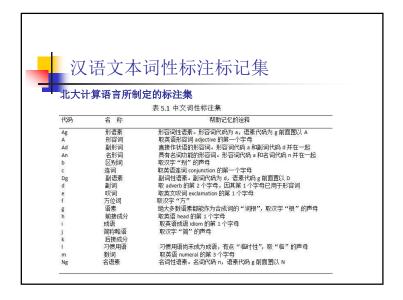
"者"接在语素或词的后面构成合成词,自然为切分单位,标记/n,如:死者/n, 笔者/n, 当局者/n, 旁观者/n, 屡教不改者/n

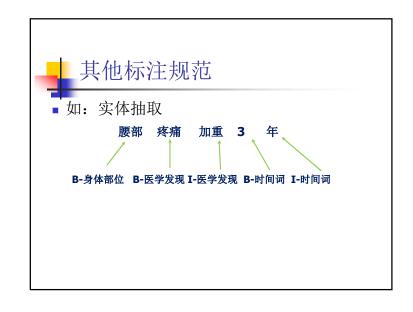
"者"接在较长短语或句子后,却应分开,将"者"单独标记为/k,如:经过/d 苦苦/d 追求/v 而/c 获得/v 幸福/a 者/k 。



汉语文本词性标注标记集

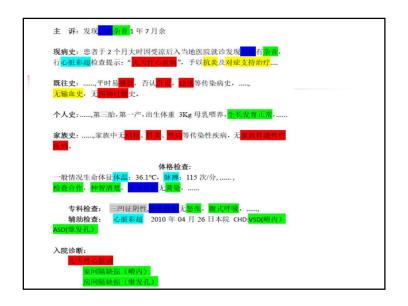
The state of the s		
代码	名 称	帮助记忆的诠释
n	名词	取英语名词 noun 的第 1 个字母
Nr	人名	名词代码 n 和 "人 (ren)"的声母并在一起
Ns	地名	名词代码 n 和处所词代码 s 并在一起
Nt	机构团体	"团"的声母 t,名词代码 n 和 t 并到一起
Nz	其他专名	"专"的声母的第一个字母 w 为 z,名词代码 n 和 z 并在一起
0	拟声词	取英语拟声词 onomatopoeia 的第 2 个字母
p	介词	取英语介词 prepositional 的第 1 个字母
q	量词	取英语 quanity 的第 1 个字母
r	代词	取英语代词 pronoun 的第 2 个字母,因 p 已用于介词
S	处所词	取英语 space 的第 1 个字母
Tg	时语素	时间词性语素。时间词代码为 t, 在语素的代码 g 前面置以 T
t	时间词	取英语 time 的第 1 个字母
u	助词	取英语组词 auxiliary 的第 2 个字母,因 a 已用于形容词
Vg	动语素	动词性语素。动词代码为 v。在语素的代码 g 前置以 v
v	动词	取英语动词 verb 的第 1 个字母
Vd	副动词	直接做状语的动词。动词和副词的代码并在一起
Vn	名动词	指具有名词功能的动词。动词和名词的代码并在一起
w	标点符号	
x	非语素词	非语素字只有一个符号,字母x通常用于代表未知数、符号
у	语气词	取汉字"语"的声母
Z	状态词	取汉字"状"的声母的前一个字母





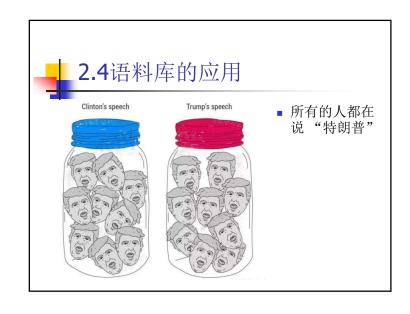




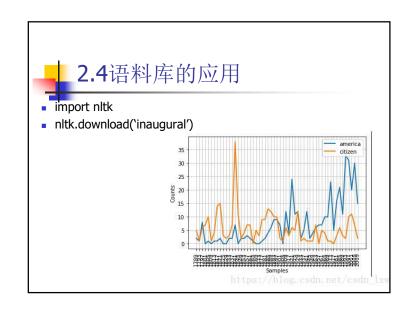




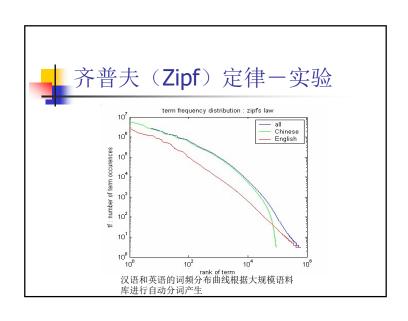














齐普夫 (Zipf) 定律

■ 齐普夫定律是美国学者G.K.齐普夫于本世纪40年代提出的词频分布定律。它可以表述为:

如果把一篇较长文章中每个词出现的频次统计起来,按照高频词在前、低频词在后的递减顺序排列,并用自然数给这些词编上等级序号,即频次最高的词等级为1,频次次之的等级为2,,频次最小的词等级为D。

若用f表示频次,r表示等级序号,则有:







题旨

■ 到http://www.icl.pku.edu.cn下载北京大学《人民日报》切分标注语料库(1个月),研究汉语"动词+名词"可能构成哪些歧义结构?