

4.3 汉语自动分词与 词性标注



汉语自动分词及词性标注

4.3.1 汉语自动分词概要

□ 汉语自动分词的重要性

- 自动分词是汉语句子分析的基础
- 词语的分析具有广泛的应用（词频统计，词典编纂，文章风格研究等）
- 文献处理以词语为文本特征
- “以词定字、以词定音”，用于文本校对、同音字识别、多音字辨识、简繁体转换

4.3.1 汉语自动分词概要

□ 汉语自动分词中的主要问题

◆ 汉语分词规范问题（《信息处理用限定汉语分词规范（GB13715）》）

一 汉语中什么是词？两个不清的界限：

（1）单字词与词素，如：新华社25日讯

（2）词与短语，如：花草，湖边，房顶，鸭蛋，小鸟，担水，一层，翻过？

4.3.1 汉语自动分词概要

◆ 歧义切分字段处理

1、中国人为了实现自己的梦想 (交集型歧义)

中国/ 人为/ 了/ 实现/ 自己/ 的/ 梦想

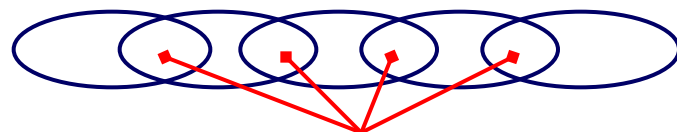
中国人/ 为了/ 实现/ 自己/ 的/ 梦想

中/ 国人/ 为了/ 实现/ 自己/ 的/ 梦想

例如：“大学生”、“研究生物”、“从小学起”、“为人民工作”、“中国产品质量”、“部分居民生活水平”等等

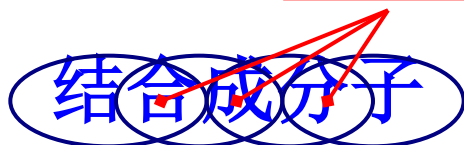
4.3.1 汉语自动分词概要

- ◆ **定义：链长** 一个交集型切分歧义所拥有的交集串的集合称为交集串链，它的个数称为链长。



交集串

例如，



“结合”、“合成”、“成分”和“分子”均构成词，交集串的集合为 {合，成，分}，因此，链长为3。



4.3.1 汉语自动分词概要

类似地,

(1) “为人民工作”

{人, 民, 工}, 歧义字段的链长为3;

(2) “中国产品质量”

{国, 产, 品, 质}, 歧义字段的链长为4;

(3) “部分居民生活水平”

{分, 居, 民, 生, 活, 水}, 歧义字段的链长为6。



4.3.1 汉语自动分词概要

2、门把手弄坏了。 (组合型歧义)

门/ 把/ 手/ 弄/ 坏/ 了/ 。

门/ 把手/ 弄/ 坏/ 了/ 。

例如，“将来”、“现在”、“才能”、“学生会”等，都是组合型歧义字段。

4.3.1 汉语自动分词概要

梁南元（1987）曾经对一个含有48,092字的自然科学、社会科学样本进行了统计，结果交集型切分歧义有518个，多义组合型切分歧义有42个。据此推断，中文文本中切分歧义的出现频度约为1.2次/100字，交集型切分歧义与多义组合型切分歧义的出现比例约为12:1。

4.3.1 汉语自动分词概要

◆ 未登录词的识别

1、人名、地名、组织机构名等，例如：

盛中国，令计划，令狐路线，张建国，蔡国庆，
党政法，蔡英文，水皮，雷地球，彭太发生，
平川三太郎，约翰·斯特朗，詹姆斯·埃尔德

2、新出现的词汇、术语、个别俗语等，例如：

博客，非典，禽流感，恶搞，裸退

4.3.1 汉语自动分词概要

例如：

- (1) 他还兼任何应钦在福州办的东路军军官学校的政治教官。
- (2) 大不列颠及北爱尔兰联合王国外交和英联邦事务大臣、议会议员杰克·斯特劳阁下在联合国安理会就伊拉克问题发言。
- (3) 坐落于江苏省南京市玄武湖公园内的夏璞墩是晋代著名的文学家、科学家夏璞的衣冠冢(zhong)。

4.3.1 汉语自动分词概要

□ 汉语自动分词的基本原则

1、语义上无法由组合成分直接相加而得到的字串应该合并为一个分词单位。(合并原则)

如：不管三七二十一（成语），或多或少（副词片语），十三点（定量结构），六月（定名结构），谈谈（重叠结构，表示尝试），辛辛苦苦（重叠结构，加强程度），进出口（合并结构）

4.3.1 汉语自动分词概要

2、语类无法由组合成分直接得到的字串应该合并为一个分词单位。 (合并原则)

(1)字串的语法功能不符合组合规律，如：好吃，好喝，好听，好看等

(2)字串的内部结构不符合语法规律，如：游水等

4.3.1 汉语自动分词概要

□ 汉语自动分词的辅助原则

操作性原则，富于弹性，不是绝对的。

1. 有明显分隔符标记的应该切分之 (切分原则)

分隔标记指标点符号或一个词。如：

上、下课 → 上/ 下课

洗了个澡 → 洗/ 了/ 个/ 澡

4.3.1 汉语自动分词概要

2. 附着性语(词)素和前后词合并为一个分词单位 (合并原则)

如：“吝”是一个附着语素，“不吝”、“吝于”等合并成一个词；

“员”：检查员、邮递员、技术员等；

“化”：现代化、合理化、多变化、民营化等。

4.3.1 汉语自动分词概要

3. 使用频率高或共现率高的字串尽量合并为一个分词单位 (合并原则)

如：“进出”、“收放”（动词并列）；

“大笑”、“改称”（动词偏正）；

“关门”、“洗衣”、“卸货”（动宾）；

“春夏秋冬”、“轻重缓急”、“男女”（并列）；

“象牙”（名词偏正）；“暂不”、“毫不”、“不再”、“早已”（副词并列）等

4.3.1 汉语自动分词概要

4. 双音节加单音节的偏正式名词尽量合并为一个分词单位 ([合并原则](#))

如：“线、权、车、点”等所构成的偏正式名词：“国际线、分数线、贫困线”、“领导权、发言权、知情权”、“垃圾车、交通车、午餐车”、“立足点、共同点、着眼点”等。



4.3.1 汉语自动分词概要

5. 双音节结构的偏正式动词应尽量合并为一个分词单位 (合并原则)

本原则只适合少数偏正式动词，如：“紧追其后”、“组建完成”等，不适合动宾及主谓式复合动词。

4.3.1 汉语自动分词概要

6. 内部结构复杂、合并起来过于冗长的词尽量切分 (切分原则)

(1) 词组带接尾词

太空/ 计划/ 室、塑料/ 制品/ 业

(2) 动词带双音节结果补语

看/ 清楚、讨论/ 完毕

(3) 复杂结构: 自来水/ 公司、中文/ 分词/ 规范/ 研究/ 计划

(4) 正反问句: 喜欢/ 不/ 喜欢、参加/ 不/ 参加

4.3.1 汉语自动分词概要

(5) 动宾结构、述补结构的动词带词缀时

写信/ 给、取出/ 给、穿衣/ 去

(6) 词组或句子的专名，多见于书面语，戏剧名、歌曲名等

鲸鱼/ 的/ 生/ 与/ 死、那/ 一/ 年/ 我们/ 都/ 很/ 酷

(7) 专名带普通名词

胡/ 先生、京沪/ 铁路



4.3.2 汉语自动分词基本算法

- 有词典切分/ 无词典切分
- 基于规则的方法/ 基于统计的方法

4.3.2 汉语自动分词基本算法



1. 最大匹配法 (Maximum Matching, MM)

—有词典切分，机械切分

- 正向最大匹配算法 (Forward MM, FMM)
- 逆向最大匹配算法 (Backward MM, BMM)
- 双向最大匹配算法 (Bi-directional MM)

假设句子: $S = c_1c_2 \cdots c_n$, 某一词:

$w_i = c_1c_2 \cdots c_m$, m 为词典中最长词的字数。

4.3.2 汉语自动分词基本算法

◆ FMM 算法描述

- (1) 令 $i=0$ ，当前指针 p_i 指向输入字串的初始位置，执行下面的操作：
- (2) 计算当前指针 p_i 到字串末端的字数（即未被切分字串的长度） n ，如果 $n=1$ ，转(4)，结束算法。否则，令 m =词典中最长单词的字数，如果 $n < m$ ，令 $m=n$ ；

4.3.2 汉语自动分词基本算法

(3) 从当前 p_i 起取 m 个汉字作为词 w_i ，判断：

(a) 如果 w_i 确实是词典中的词，则在 w_i 后添加一个切分标志，转(c)；

(b) 如果 w_i 不是词典中的词且 w_i 的长度大于1，将 w_i 从右端去掉一个字，转(a)步；否则 (w_i 的长度等于1)，则在 w_i 后添加一个切分标志，将 w_i 作为单字词添加到词典中，执行 (c) 步；

(c) 根据 w_i 的长度修改指针 p_i 的位置，如果 p_i 指向字串末端，转(4)，否则， $i=i+1$ ，返回 (2)；

(4) 输出切分结果，结束分词程序。

4.3.2 汉语自动分词基本算法

例：假设词典中最长单词的字数为 7。

输入字串：他是研究生物化学的。

切分过程：他是研究生物化学的。

$p \uparrow$ |

... ..

他/ 是研究生物化学的。

$p \uparrow$ |

FMM 切分结果：他/ 是/ 研究生/ 物化/ 学/ 的/。

BMM 切分结果：他/ 是/ 研究/ 生物/ 化学/ 的/。

4.3.2 汉语自动分词基本算法

➤ 优点：

- 程序简单易行，开发周期短；
- 仅需要很少的语言资源（词表），不需要任何词法、句法、语义资源；

➤ 弱点：

- 歧义消解的能力差；
- 切分正确率不高，一般在95%左右。

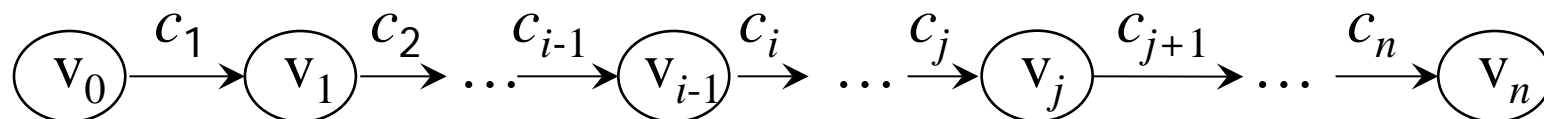
4.3.2 汉语自动分词基本算法



2. 最少分词法（最短路径法）

◆ 基本思想

设待分字串 $S=c_1 c_2 \dots c_n$ ，其中 $c_i (i=1,2,\dots,n)$ 为单个的字， n 为串的长度， $n \geq 1$ 。建立一个节点数为 $n+1$ 的切分有向无环图 G ，各节点编号依次为 $V_0, V_1, V_2, \dots, V_n$ 。

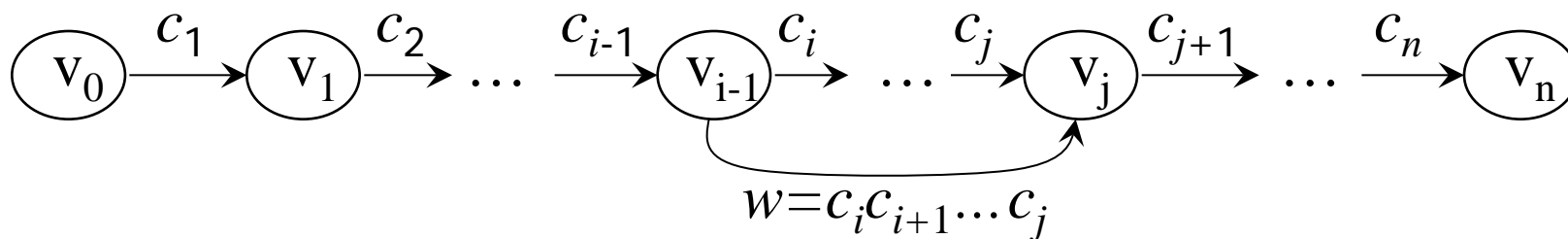


求最短路径：贪心法或简单扩展法。

4.3.2 汉语自动分词基本算法

◆ 算法描述:

- (1) 相邻节点 v_{k-1}, v_k 之间建立有向边 $\langle v_{k-1}, v_k \rangle$, 边对应的词默认为 c_k ($k=1, 2, \dots, n$)。
- (2) 如果 $w=c_i c_{i+1} \dots c_j$ ($0 < i < j \leq n$) 是一个词, 则节点 v_{i-1}, v_j 之间建立有向边 $\langle v_{i-1}, v_j \rangle$, 边对应的词为 w 。



- (3) 重复步骤(2), 直到没有新路径(词序列)产生。
- (4) 从产生的所有路径中, 选择路径最短的(词数最少的)作为最终分词结果。

4.3.2 汉语自动分词基本算法

例：(1) 输入字串：他只会诊断一般的疾病。

可能输出：他/ 只会/ 诊断/ 一般/ 的/ 疾病/。(7)

他/ 只/ 会诊/ 断/ 一般/ 的/ 疾病/。(8)

... ..

最终结果：他/ 只会/ 诊断/ 一般/ 的/ 疾病/ 。

(2) 输入字串：他说的确实在理。

可能输出：他/ 说/ 的/ 确实/ 在理/ 。（6）

他/ 说/ 的确/ 实在/ 理/ 。（6）

... ..

4.3.2 汉语自动分词基本算法



➤ 优点:

- 切分原则符合汉语自身规律
- 需要的语言资源（词表）也不多

➤ 弱点:

- 对许多歧义字段难以区分，最短路径有多条时，选择最终的输出结果缺乏应有的标准。
- 字串长度较大和选取的最短路径数增大时，长度相同的路径数急剧增加，选择最终正确的结果困难越来越大。

4.3.2 汉语自动分词基本算法

3. 基于统计语言模型的分词方法

◆ 方法描述:

设对于待切分的句子 S , $W = w_1 w_2 \dots w_k$
($1 \leq k \leq n$) 是一种可能的切分。

$$\begin{aligned} W^* &= \arg \max_W P(W | S) \\ &= \arg \max_W P(W)P(S | W) \end{aligned}$$

4.3.2 汉语自动分词基本算法

微软研究院把一个可能的词序列 W 转换成一个可能的词类序列 $C = c_1 c_2 \cdots c_N$ ，即：

- 专有名词的人名PN、地名LN、机构名ON分别作为一类；
- 实体名词中的日期 dat、时间tim、百分数per、货币mon 等作为一类；
- 对词法派生词MW和词表词LW，每个词单独作为一类。

4.3.2 汉语自动分词基本算法

那么,

$$C^* = \arg \max_C P(C) P(S | C) \quad (7-1)$$

语言模型 \rightarrow $P(C)$ \leftarrow $P(S | C)$ 生成模型

$P(C)$ 可采用3元语法:

$$P(C) = P(c_1)P(c_2 | c_1) \prod_{i=3}^N P(c_i | c_{i-2}c_{i-1}) \quad (7-2)$$

4.3.2 汉语自动分词基本算法

生成模型在满足独立性假设的条件下，可近似为：

$$P(S | C) \approx \prod_{i=1}^N P(s_i | c_i) \quad (7-3)$$

该公式的含意是，任意一个词类生成汉字串的概率只与自身有关，而与其上下文无关。例如，如果“教授”是词表里的词，那么 $P(s_i=\text{教授} | c_i=\text{LW})=1$ 。

4.3.2 汉语自动分词基本算法

词 类	生成模型 $P(S C)$	语言知识
词表词 (LW)	若 S 是词表词, $P(S LW)=1$, 否则为0;	分词词表
词法派生词 (MW)	若 S 是派生词, $P(S MW)=1$, 否则为0;	派生词词表
人名 (PN)	基于字的二元模型	姓氏表, 中文人名模板
地名 (LN)	基于字的二元模型	地名表、地名关键词表、地名简称表
机构名 (ON)	基于词类的二元模型	机关名关键词表, 机构名简称表
实体名 (FT)	若 S 可用实体名词规则集 G 识别, $P(S G)=1$, 否则为0。	实体名词规则集

4.3.2 汉语自动分词基本算法

模型的训练由以下三步组成：

- (1) 在词表和派生词表的基础上，用正向最大匹配法切分训练语料，专有名词通过一个专门模块标注，实体名词通过相应的规则和有限状态自动机标注，由此产生一个带词类别标记的初始语料；
- (2) 用带词类别标记的初始语料，采用最大似然估计方法估计语言模型的概率参数；
- (3) 用语言模型（公式(7-1)、(7-2)、(7-3)），对训练语料重新切分和标注，得到新的训练语料；
- (4) 重复(2)(3)步，直到系统的性能不再有明显的变化。

4.3.2 汉语自动分词基本算法

➤ 优点:

- 减少了很多手工标注的工作
- 在训练语料规模足够大和覆盖领域足够多时，可以获得较高的切分正确率

➤ 弱点:

- 训练语料的规模和覆盖领域不好把握
- 计算量较大

4.3.2 汉语自动分词基本算法

4. 基于HMM的分词方法

基本思想：

把输入字串(句子) S 作为HMM的输入；
(切分后的)单词串 S_w 为状态的输出，即观察序列 $S_w = w_1 w_2 \cdots w_n$ ($n \geq 1$)；词性序列 S_c 为状态序列，每个词性标记对应HMM中的一个状态 q_i ， $S_c = c_1 c_2 \cdots c_n$ 。

4.3.2 汉语自动分词基本算法

➤ 优点:

- 可以减少很多手工标注的工作
- 在训练语料规模足够大和覆盖领域足够多时，可以获得较高的切分正确率

➤ 弱点:

- 训练语料的规模和覆盖领域不好把握
- 模型实现复杂、计算量较大

4.3.2 汉语自动分词基本算法

5. 基于统计模型的分词与词性标注一体化方法

基本思想： 设句子 S 由词串组成 $W=w_1w_2\cdots w_n$ ($n\geq 1$), 单词 w_i 的词性标注为 t_i , 即句子 S 相应的词性标注符号序列可表达为 $T=t_1t_2\cdots t_n$ 。那么, 分词与词性标注的任务就是要在 S 所对应的各种切分和标注形式中, 寻找 T 和 W 的联合概率 $P(W, T)$ 为最优的词切分和标注组合。

4.3.2 汉语自动分词基本算法

如果把词性符号序列作为HMM的中间状态，词序列作为输出，那么， $P(W, T)$ 可以由HMM近似地表示为：

$$P(W, T) = P(W | T)P(T) \approx \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1} t_{i-2}) \quad (7-4)$$

生成模型

基于词性的
语言模型

4.3.2 汉语自动分词基本算法

反之，如果把词序列作为HMM的中间状态，词性符号作为输出，那么， $P(W, T)$ 的另一种形式为：

$$P(W, T) = P(T | W)P(W) \approx \prod_{i=1}^n P(t_i | w_i) P(w_i | w_{i-1} w_{i-2}) \quad (7-5)$$

生成模型

基于词的
语言模型

4.3.2 汉语自动分词基本算法

将上述(7-4)和(7-5)综合：

$$P^*(W, T) = \alpha \prod_{i=3}^n P(w_i | t_i) P(t_i | t_{i-1}, t_{i-2}) + \beta \prod_{i=1}^n P(t_i | w_i) P(w_i | w_{i-1}, w_{i-2}) \quad (7-6)$$

显然，这种综合模型的指导思想是希望通过调整参数 α 和 β 的值来确定两个子模型在整个分词与词性标注过程中所发挥作用的比重，从而获得分词与词性标注的整体最优。

4.3.2 汉语自动分词基本算法

从公式 (7-5) 得到的结果分析可知, $P(t_i | w_i)$ 对分词无帮助, 且在分词确定后对词性标注又会增添偏差。因此, 在实现这一模型时, 可仅取公式 (7-5) 中的语言模型部分, 而舍弃词性标注部分, 并令 $\alpha = 1$, 仅保留加权系统 β , 于是,

$$\begin{aligned} P^{\wedge}(W, T) = & \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1}, t_{i-2}) + \\ & \beta \prod_{i=1}^n P(w_i | w_{i-1}, w_{i-2}) \end{aligned} \quad (7-7)$$

4.3.2 汉语自动分词基本算法

在确定 β 系数值时，可根据词典中词汇 w 的个数和词性 t 的种类数目，取二者之比，即 $\beta = \text{词典中词 } w \text{ 的个数} / \text{词性 } t \text{ 的种类数}$ 。

4.3.2 汉语自动分词基本算法

➤ 优点:

- 可以减少很多手工标注的工作
- 在训练语料规模足够大和覆盖领域足够多，各类参数设定适当时，可以获得较高的切分正确率

➤ 弱点:

- 训练语料的规模和覆盖领域不好把握
- β 系数值难以把握

4.3.2 汉语自动分词基本算法

6. 由字构词的(基于字标注)分词方法 (Character-based tagging)

第一篇由字构词的汉语分词方法的论文[Xue, 2002]发表在2002年的第一届ACL汉语特别兴趣小组SIGHAN (<http://www.sighan.org/>) 组织的研讨会上, 在2005年和2006年的两次Bakeoff 评测中取得好成绩。

4.3.2 汉语自动分词基本算法

◆ **基本思想**：将分词过程看作是字的分类问题。该方法认为，每个字在构造一个特定的词语时都占据着一个确定的构词位置(即词位)。假定每个字只有4个词位：词首(B)、词中(M)、词尾(E)和单独成词(S)，那么，每个字归属一特定的词位。

这里所说的“字”不仅限于汉字，也可以指标点符号、外文字母、注音符号和阿拉伯数字等任何可能出现在汉语文本中的文字符号，所有这些字符都是由字构词的基本单元。

4.3.2 汉语自动分词基本算法

例如：

(1) 上海/ 计划/ 到/ 本/ 世纪/ 末/ 实现/ 人均/
国内/ 生产/ 总值/ 五千美元/ 。 /

(2) 上/B 海/E 计/B 划/E 到/S 本/S 世/B 纪/E
末/S 实/B 现/E 人/B 均/E 国/B 内/E 生/B 产/E 总
/B 值/E 五/B 千/M 美/M 元/E 。 /S

4.3.2 汉语自动分词基本算法

在字标注过程中，对所有的字根据预定义的特征进行词位特征学习，获得一个概率模型，然后在待切分字符串上，根据字与字之间的结合紧密程度，得到一个词位的分类结果，最后根据词位定义直接获得最终的分词结果。

◆ 工具：

- 支持向量机（SVM）
- 条件随机场（CRF）

最常用的两类特征是字本身和词位(状态)的转移概率

4.3.2 汉语自动分词基本算法

◆ 评价:

该方法的重要优势在于，它能够平衡地看待词表词和未登录词的识别问题，文本中的词表词和未登录词都是用统一的字标注过程来实现的。在学习构架上，既可以不必专门强调词表词信息，也不用专门设计特定的未登录词识别模块，因此，大大地简化了分词系统的设计[黄昌宁，2006]

4.3.2 汉语自动分词基本算法

□ 其他方法

- ◆ 全切分方法
- ◆ 串频统计和词形匹配相结合的分词方法
- ◆ 规则方法与统计方法相结合
- ◆ 多重扫描法

.....

4.3.2 汉语自动分词基本算法

□ 方法比较

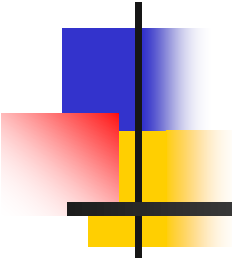
- (1) 最大匹配分词算法是一种简单的基于词表的分词方法，有着非常广泛的应用。这种方法只需要最少的语言资源（仅需要一个词表，不需要任何词法、句法、语义知识），程序实现简单，开发周期短，是一个简单实用的方法，但对歧义字段的处理能力不够强大。

4.3.2 汉语自动分词基本算法

- (2) 全切分方法首先切分出与词表匹配的所有可能的词，然后运用统计语言模型和决策算法决定最优的切分结果。这种切分方法的优点是发现所有的切分歧义，但解决歧义的方法很大程度上取决于统计语言模型的精度和决策算法，需要大量的标注语料，并且分词速度也因搜索空间的增大而有所缓慢。

4.3.2 汉语自动分词基本算法

- (3) 最短路径分词方法的切分原则是使切分出来的词数最少。这种切分原则多数情况下符合汉语的语言规律，但无法处理例外的情况，而且如果最短路径不止一条时，系统往往不能确定最优解。
- (4) 统计方法具有较强的歧义区分能力，但需要大规模标注 (或预处理) 语料库的支持，需要的系统开销也较大。



4.3.3 未登录词识别

4.3.3 未登录词识别

□ 命名实体(Named Entity, NE) (专有名词)

人名（中国人名和外国译名）、地名、组织机构名、数字、日期、货币数量

□ 其他新词

专业术语、新的普通词汇等。

4.3.3 未登录词识别

□ 关于中文姓名

- 台湾出版的《中国姓氏集》收集姓氏 5544 个，其中，单姓 3410 个，复姓 1990 个，3字姓 144 个
- 中国目前仍使用的姓氏共 737 个，其中，单姓 729 个，复姓 8 个
- 根据我们收集的 300 万个人名统计，姓氏有974 个，其中，单姓 952个，复姓 23 个，300万人名中出现汉字4064个。（曹文洁，2002a, 2002b）

4.3.3 未登录词识别

□ 中文姓名识别的难点

- 名字用字范围广，分布松散，规律不很明显。
- 姓氏和名字都可以单独使用用于特指某一人。
- 许多姓氏用字和名字用字（词）可以作为普通用字或词被使用，例如，姓氏：于（介词），张（量词），江（名词）等；名字：建国，国庆，胜利，文革等，全名本身也是普通词汇，如：万里，温馨，高山，高升，高飞，周密，江山等。

4.3.3 未登录词识别

➤ 缺乏可利用的启发标记。

例如: (1) 祝贺老总百战百胜。

(2) 林徽因此时已经离开了那里。

(3) 赵微笑着走了。

(4) 南京市长江大桥。

4.3.3 未登录词识别

□ 中文姓名识别方法

- ◆ 姓名库匹配，以姓氏作为触发信息，寻找潜在的名字
- ◆ 计算潜在姓名的概率估值及相应姓氏的姓名阈值(threshold value)，根据姓名概率评价函数和修饰规则对潜在的姓名进行筛选。

4.3.3 未登录词识别

□ 计算概率估计值

设姓名 $Cname = Xm_1m_2$ ，其中 X 表示姓， m_1m_2 分别表示名字首字和名字尾字。分别用下列公式计算姓氏和名字的使用频率：

$$F(X) = \frac{X \text{ 用作姓氏}}{X \text{ 出现的总次数}}$$

$$F(m_1) = \frac{m_1 \text{ 作为名字首字出现的次数}}{m_1 \text{ 出现的总次数}}$$

$$F(m_2) = \frac{m_2 \text{ 作为名字尾字出现的次数}}{m_2 \text{ 出现的总次数}}$$

4.3.3 未登录词识别

字串 $Cname$ 可能为姓名的概率估值:

$$P(Cname) = \begin{cases} F(X) \times F(m_1) \times F(m_2) & \text{复名情况} \\ F(X) \times F(m_2) & \text{单名情况} \end{cases}$$

□ 确定阈值

姓氏 X 构成姓名的最小阈值:

$$T_{\min}(X) = \begin{cases} F(X) \times \text{Min}(F(m_1) \times F(m_2)) & \text{复名情况} \\ F(X) \times \text{Min}(F(m_2)) & \text{单名情况} \end{cases}$$

4.3.3 未登录词识别

□ 设计评估函数

姓名的评价函数：

$$f = \ln P(Cname)$$

对于特定的姓氏 X 通过训练语料得到一
阈值 β_X ，当 f 大于 β_X 时，该识别的汉字串确
定为中文姓名。

4.3.3 未登录词识别

□ 使用修饰规则：

如果姓名前是一个数字，或者与“.”字符的距离小于 2 个字节，则否定此姓名。

◆ 确定潜在的姓名边界

➤ 左界规则：若潜在姓名前面是一称谓，或一标点符号，或者潜在姓名在句首，或者潜在的姓名的姓氏使用频率为100%，则姓名的左界确定。

4.3.3 未登录词识别

➤ 右界规则：若姓名后面是一称谓，或者是一指界动词(如，说，是，指出，认为等)或标点符号，或者潜在的姓名在句尾，或者潜在姓名的尾字使用频率为100%，则姓名的右界确定。

◆ 校正潜在的姓名

依据：含重合部分的潜在姓名不可能同时成立。利用各种规则消除冲突的潜在姓名。

4.3.3 未登录词识别

□ 中文地名识别方法

◆ 困难

- 地名数量大，缺乏明确、规范的定义。《中华人民共和国地名录》(1994)收集88026个，不包括相当一部分街道、胡同、村庄等小地方的名称。
- 真实语料中地名出现情况复杂。如地名简称、地名用词与其他普通词冲突、地名是其他专用名词的一部分，地名长度不一等。

4.3.3 未登录词识别

◆ 基本资源

- 建立地名资源知识库
 - 地名库、地名用字库、地名用词库
- 建立识别规则库
 - 筛选规则、确认规则、否定规则

4.3.3 未登录词识别

◆ 基本方法

- 统计模型
- 通过训练语料选取阈值
- 地名初筛选
- 寻找可以利用的上下文信息
- 利用规则进一步确定地名

4.3.3 未登录词识别

□ 中文机构名称的识别

◆ 中文机构名称的构成

- 词法角度：偏正式(修饰格式)的复合词
{名词|形容词|数量词|动词} + 名词
- 句法角度：“定语 + 名词性中心语”型的名词短语(定名型短语)
- 中心语：机构称呼词，如：大学，学院，研究所，学会，公司等。

4.3.3 未登录词识别

◆ 中文机构名称的类型

- 地名，如：北京大学，武汉大学
- 人名，如：中山大学，哈佛大学
- 学科、专业 and 部门系统，如：公安部，教育委员会
- 研究、生产或经营等活动的对象，如：软件研究所，卫星制造厂
- 上述情况的综合，如：白求恩医科大学



4.3.3 未登录词识别

- 大机构、团体、组织和职业的名称，如：中国人民解放军洛阳外国语学院，中国发明家学会等
- 专造的机构名，如：复旦大学，四通公司，微软研究院
- 创办、工作的方式，如：某某股份公司，中央电视大学

4.3.3 未登录词识别

◆ 机构名称识别方法

- 找到一机构称呼词
- 根据相应规则往前逐个检查名词作为修饰名词的合法性，直到发现非法词
- 如果所接受的修饰词同机构称呼词构成一个合法的机构名称，则记录该机构名称
- 统计模型



4.3.4 词性标注面临的问题

4.3.4 词性标注面临的问题

□ 概要

词性(part-of-speech, POS)标注(tagging)的主要任务是消除词性兼类歧义。

例如，在英语中：

1) **Time flies like an arrow.**

2) **I want you to web our annual report.**

对 **Brown** 语料库的统计，**55%**词次兼类。汉语中常用词兼类现象严重，《现代汉语八百词》兼类占 **22.5%**。

4.3.4 词性标注面临的问题

◆ 汉语中的词性兼类现象

(1) 形同音不同，如：“好（hao3，形容词）、好（hao4，动词）”

这个人什么都好，就是好酗酒。

(2) 同形、同音，但意义毫不相干，如：“会（会议，名词）、会（能够、动词）”

每次他都会在会上制造点新闻。

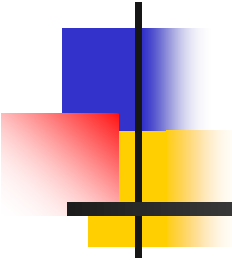
4.3.4 词性标注面临的问题

(3) 具有典型意义的兼类词，如：“典型(名词或形容词)”、“教育(名词或动词)”

让孩子接受那样的教育简直是对教育事业的侮辱。

(4) 上述情况的组合，如：“行(xing2, 动词/形容词; hang2, 名词/量词)”

每当他走过那行白杨树时，他都感觉好像每一棵树都在向他行注目礼。



4.3.5 词性标注集

4.3.5 词性标注集

□ 标注集的确定原则：

不同语言中，词性划分基本上已经约定俗成。
自然语言处理中对词性标记要求相对细致。

◆ 一般原则：

- 标准性：普遍使用和认可的分类标准和符号集；
- 兼容性：与已有资源标记尽量一致，或可转换；
- 可扩展性：扩充或修改。

4.3.5 词性标注集

◆ UPenn Treebank 的词性标注集确定原则：

- 可恢复性(recoverability)：从标注语料能恢复原词汇或借助于句法信息能区分不同词类；
- 一致性(consistency)：功能相同的词应该属于同一类；
- 不明确性(indeterminacy)：为了避免标注者在不明确的情况下任意决定标注类型，允许标注者给出多个标记（限于一些特殊情况）。

—[Marcus et al., 1993]

4.3.5 词性标注集

◆ UPenn Treebank 的词性标注集

□ 33 类

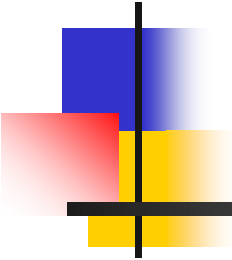
□ **NN** 名词、**NR** 专业名词、**NT** 时间名词、**VA** 可做谓语的形容词、**VC** “是”、**VE** “有”作为主要动词、**VV** 其他动词、**AD** 副词、**M** 量词等。

4.3.5 词性标注集

◆ 北大计算语言研究所的词性标注集

□ 26个基本词类代码，74个扩充代码，标记集中共有106个代码。

名词(n)、时间词(t)、处所词(s)、方位词(f)、数词(m)、量词(q)、区别词(b)、代词(r)、动词(v)、形容词(a)、状态词(z)、副词(d)、介词(p)、连词(c)、助词(u)、语气词(y)、叹词(e)、拟声词(o)、成语(i)、习用语(l)、简称(j)、前接成分(h)、后接成分(k)、语素(g)、非语素字(x)、标点符号(w)。



4.3.6 词性标注方法



4.3.6 词性标注方法

- 基于规则的词性标注方法
- 基于统计模型的词性标注方法
- 规则和统计方法相结合的词性标注方法
- 基于有限状态变换机的词性标注方法
- 基于神经网络的词性标注方法



4.3.6 词性标注方法

□ 基于规则的词性标注方法

◆ TAGGIT 词性标注系统(Bwon University)

- 86 种词性, 3300 规则
- 手工编写词性歧义消除规则
- 机器自动学习规则

4.3.6 词性标注方法

□ 山西大学的词性标注系统 [刘开瑛, 2000]

◆ 手工编写消歧规则

➤ 建立非兼类词典

➤ 建立兼类词典

- 词性可能出现的概率高低排列

➤ 构造兼类词识别规则

4.3.6 词性标注方法

(1) 并列鉴别规则

如：体现了人民的要求(N/V ?)和愿望(N, 非兼类)。

(2) 同境鉴别规则

如：一个优秀的企业必须具备一流的产品(名词, 非兼类)、一流的管理(N/V ?)和一流的服务(N/V ?)。

4.3.6 词性标注方法

(3) 区别词鉴别规则(区别词只能直接修饰名词)

如：他们搞的这次大型(鉴别词，非兼类) 调查(V/N ?)历时半年。

(4) 唯名形容词鉴别规则(有些形容词只能直接修饰名词)

如：重大（唯名形容词）损失（N/V ?）

巨大（唯名形容词）影响（N/V ?）

4.3.6 词性标注方法

➤ 根据词语的结构建立词性标注规则

(1) 词缀（前缀、后缀）规则

- 形容词：蓝茵茵，绿油油，金灿灿，...
- 数量词：一片片，一次次，一回回，...
- 人名简称：李总，张工，刘老，...
- 其他：年轻化，知识化，...{化}
 篮球赛，足球赛，...{赛}



4.3.6 词性标注方法

(2) 重叠词规则

一 看看，瞧瞧，高高兴兴，热热闹闹，...

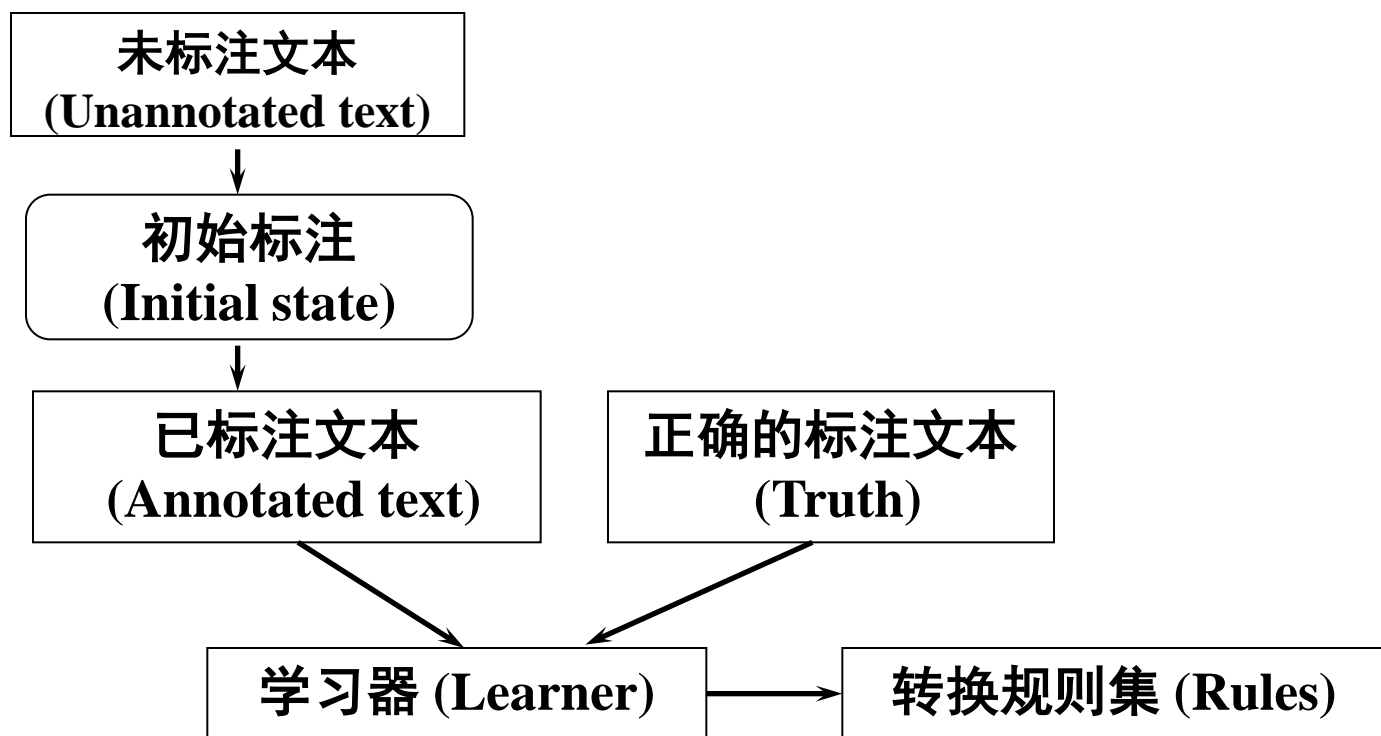
4.3.6 词性标注方法

□ 基于错误驱动的机器学习方法

- 初始词性赋值
- 对比正确标注的句子，自动学习结构转换规则
- 利用转换规则调整初始赋值

— [E. Brill, 1992]

4.3.6 词性标注方法



基于转换规则的错误驱动的机器学习方法

4.3.6 词性标注方法

□ 基于统计模型的词性标注方法

◆ 基于 n -gram 的语言模型

应用系统: (1) 1983年 Mashall 提出的 LOB 语料库的标注系统: CLAWS (Constituent-Likelihood Automatic Word-tagging System)
(2) DeRose 对 CLAWS 改进后 VOLSUNGA 系统 (bi-gram)。

4.3.6 词性标注方法

◆ 基于 HMM 的词性标注方法

- 状态集 (词性序列, 状态数: 词类符号数)
- 输出符号 (单词序列, 词汇量)
- 初始状态概率
- 状态转移概率
- 符号输出概率

— [Manning, 2001] pp. 357-359:

. **Jelink's Method**

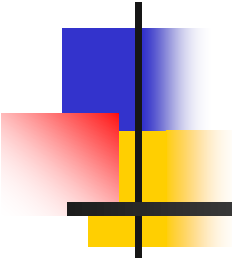
. **Kupier's Method**

4.3.6 词性标注方法

□ 规则和统计相结合的词性标注方法

- ◆ 规则消歧，统计概率引导
- ◆ 或者统计方法赋初值，规则消歧

—[周强，1995；张民，1998]



4.3.7 分词与词性标注 结果评测

4.3.7 分词与词性标注结果评测

□ 两种测试

- 封闭测试 / 开放测试
- 专项测试 / 总体测试

4.3.7 分词与词性标注结果评测

□ 评测指标

◆ **正确率**(Correct ratio/Precision, C): 测试结果中正确结果的个数占系统所有输出结果的比例。假设系统输出 N 个, 其中, 正确的结果为 n 个, 那么,

$$C = \frac{n}{N} \times 100\%$$

4.3.7 分词与词性标注结果评测

◆ 召回率(找回率) (Recall ratio, R): 测试结果中正确结果的个数占标准答案总数的比例。假设系统输出 N 个结果, 其中, 正确的结果为 n 个, 而标准答案的个数为 M 个, 那么,

$$R = \frac{n}{M} \times 100\%$$

4.3.7 分词与词性标注结果评测

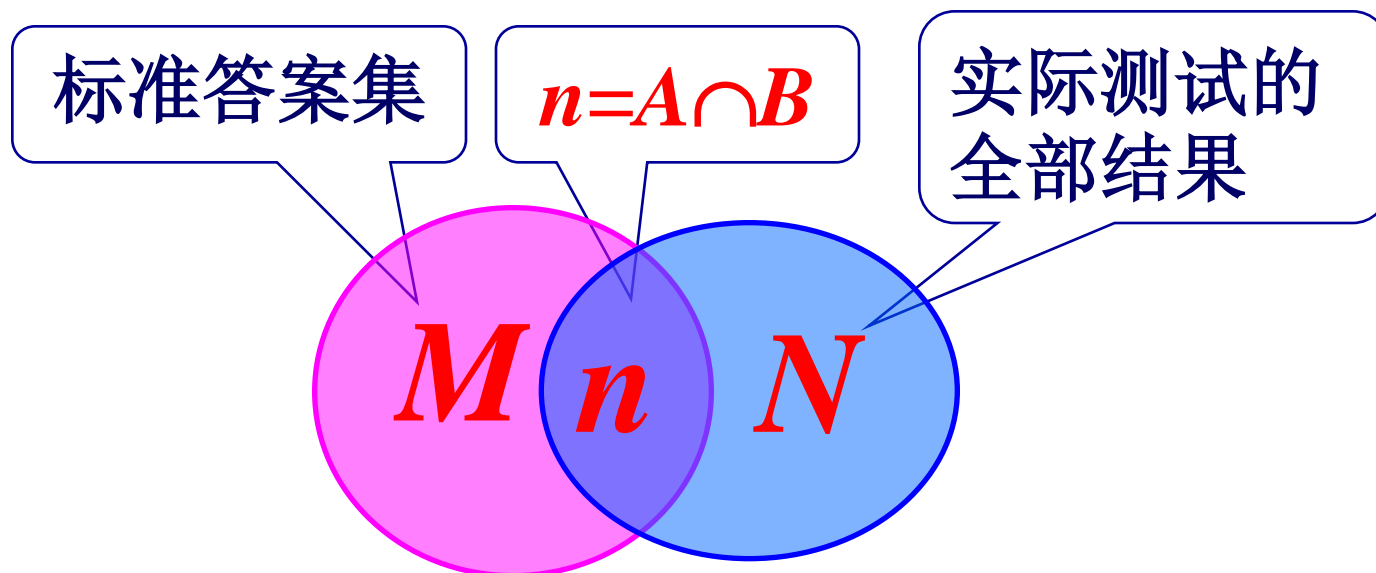
◆ F-测度值(F-Measure): 正确率与召回率的综合值。计算公式为:

$$F - measure = \frac{(\beta^2 + 1) \times C \times R}{\beta^2 \times C + R} \times 100\%$$

一般地, 取 $\beta = 1$, 即

$$F1 = \frac{2 \times C \times R}{C + R} \times 100\%$$

4.3.7 分词与词性标注结果评测



$$C = \frac{n}{N} \times 100\%$$

$$R = \frac{n}{M} \times 100\%$$



4.3.7 分词与词性标注结果评测

□ 2003年国家863评测部分结果

◆ 分词

● 最好成绩: $C=93.44\%$, $R=93.69\%$,
 $F1=93.46\%$

● 最差成绩: $C=91.42\%$, $R=89.27\%$,
 $F1=90.33\%$

4.3.7 分词与词性标注结果评测

◆ 词性标注

- 最好成绩: $C=87.47\%$, $R=87.52\%$,
 $F1=87.50\%$
- 最差成绩: $C=68.65\%$, $R=68.99\%$,
 $F1=68.82\%$

4.3.7 分词与词性标注结果评测

◆ 人名识别

- 最好成绩: $C=72.35\%$, $R=78.07\%$,
 $F1=68.33\%$
- 最差成绩: $C=27.27\%$, $R=43.29\%$,
 $F1=33.46\%$

4.3.7 分词与词性标注结果评测

◆ 机构名识别

● 最好成绩: $C=81.51\%$, $R=77.38\%$,
 $F1=68.56\%$

● 最差成绩: $C=4.65\%$, $R=10.60\%$,
 $F1=6.52\%$

4.3.7 分词与词性标注结果评测

□ 2005年SIGHAN 汉语分词评测结果(使用MSR语料)

评测方式	系统排名	性能指标				
		召回率	精确率	F-值	R_{ooV}	R_{iV}
封闭测试	最好	0.962	0.966	0.964	0.717	0.968
	最差	0.898	0.896	0.897	0.327	0.914
开放测试	最好	0.980	0.965	0.972	0.59	0.99
	最差	0.788	0.818	0.803	0.37	0.8

R_{ooV} 表示集外词的召回率, R_{iV} 表示集内词的召回率。

4.3.7 分词与词性标注结果评测

◆ 说明：

如果汉语自动分词与词性标注一体化进行，对于词性标注来说，可以用“召回率”衡量词性标注系统的性能，但是，如果不是分词与词性标注一体化进行，而是词性标注系统对已经切分好的汉语词汇进行词性标注，那么，一般不采用“召回率”指标衡量词性标注系统的性能。



4.4 词义消歧与标注技术

4.4.1 词义消歧

□ 词义消歧问题

(word sense disambiguation, WSD)

例如：

英文： bank: 银行/ 河岸

plant: 工厂/ 植物

汉语： 打： play/ take/ dial/ weave ...

包： package/ guarantee / ...



4.4.1 词义消歧

□基本方法

- ◆早期基于规则的消歧方法
- ◆统计机器学习消歧方法
 - 有监督学习方法
 - 无监督学习方法

基本思路：一个词的不同语义一般发生在不同的上下文中。

- ◆基于词典信息的消歧方法

4.4.1 词义消歧

□有监督的词义消歧方法

基本思路：通过建立分类器，利用划分多义词的上下文类别的方法来区分多义词的词义。

(1)基于互信息的消歧方法 (Brown *et al.*, 1991)

假设我们有一个双语对齐的平行语料库，以法语和英语为例，通过词语对齐模型每个法语单词可以找到对应的英语单词，一个多义的法语单词在不同的上下文中对应多种不同的英语翻译。

4.4.1 词义消歧

例子:

- *prendre une mesure* → **to take** a measure
- *prendre une décision* → **to make** a decision

也就是说，法语动词 *prendre* 可以被翻译成 to take，也可以被翻译成 to make，这取决于它所带的宾语是 *mesure* 还是 *décision*。



4.4.1 词义消歧

可以把一个多义法语单词被翻译成的英语单词看作是这个法语单词的语义解释，而决定法语多义词语义的条件看作是语义指示器(indicator)，如：前面例子中法语单词 *prendre* 所带的宾语。因此，只要我们知道了多义词的语义指示器，也就确定了该词在特定上下文中的语义。这样，多义词的词义消歧问题就变成了语义指示器的分类问题。

4.4.1 词义消歧

假设 T_1, T_2, \dots, T_m 是多义法语词的翻译(或语义), V_1, V_2, \dots, V_n 是指示器可能的取值, 利用 Flip-Flop 算法来解决指示器分类问题:

- (1) 随机地将 T_1, T_2, \dots, T_m 划分为两个集合 $P = \{P_1, P_2\}$
- (2) 执行如下循环:
 - (a) 找到 V_1, V_2, \dots, V_n 的一种划分 $Q = \{Q_1, Q_2\}$, 使 Q 与 P 之间的互信息最大;
 - (b) 找到的一种改进的划分 P' , 使 P' 与 Q 的互信息最大。

4.4.1 词义消歧

根据互信息的定义：

$$I(P; Q) = \sum_{x \in P} \sum_{y \in Q} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

从上面的 Flip-Flop 算法我们可以看出，每次迭代互信息 $I(P; Q)$ 都应该单调增加，因此，算法终止的条件自然是互信息 $I(P; Q)$ 不再增加或者增加甚少。

4.4.1 词义消歧

一旦指示器的取值划分确定了，词义消解就变成了如下简单的过程：

- (1) 对于出现的歧义词确定其指示器值 V_i ；
- (2) 如果 V_i 在 Q_1 中，指定该歧义词的语义为语义1，如果在 Q_2 中，指定其语义为语义2。

4.4.1 词义消歧

(2)基于贝叶斯分类器的消歧方法(Gale *et al.*, 1992)

在双语语料库中多义词的语义取决于该词所处的上下文语境 C ，如果某个多义词 w 有多个语义 $s_i (i \geq 2)$ ，那么，可以通过计算 $\arg \max_{s_i} P(s_i | C)$ 确定 w 的词义。

根据贝叶斯公式：

$$P(s_i | C) = \frac{P(C | s_i)P(s_i)}{P(C)}$$

4.4.1 词义消歧

考虑分母的归一化，并运用如下独立性假设：

$$P(C | s_i) = \prod_{w \in C} P(w | s_i)$$

因此，

$$\hat{s}_i = \arg \max_{s_i} [P(s_i) \prod_{w \in C} P(w | s_i)]$$

概率 $P(w | s_i)$ 和 $P(s_i)$ 都可以通过最大似然估计求得。

4.4.1 词义消歧

基于上述思想的词义消歧算法：

■ 训练过程：

(1) 对于多义词 a 的每个语义 s_i 执行如下循环：

对于词典中所有的词 w 计算 $P(w|s_i) = \frac{N(w, s_i)}{N(s_i)}$

(2) 对于多义词 a 的每个语义 s_i 计算：

$$P(s_i) = \frac{N(s_i)}{N(a)}$$

4.4.1 词义消歧

■ 消歧过程：

对于多义词 a 的每个语义 s_i 计算 $P(s_i)$ ，并根据上下文中的每个词 w 计算 $P(w|s_i)$ ，选择

$$\hat{s}_i = \arg \max_{s_i} [P(s_i) \prod_{w \in C} P(w|s_i)]$$

在实际算法实现中，通常将概率 $P(w|s_i)$ 和 $P(s_i)$ 的乘积运算转换为对数加法运算：

$$\hat{s}_i = \arg \max_{s_i} [\log P(s_i) + \sum_{w \in C} \log P(w|s_i)]$$

4.4.1 词义消歧

□ 基于词典的词义消歧方法

(1) 基于语义定义的消歧

词典中词条本身的定义作为判断其语义的条件。例如 `cone` 在词典中有两个定义：一个是指“松树的球果”，另一个是指“用于盛放其他东西的锥形物，比如，盛放冰激凌的锥形薄饼”。如果在文本中，“树(`tree`)”或者“冰(`ice`)”与 `cone` 出现在相同的上下文中，那么，`cone` 的语义就可以确定了，`tree` 对应 `cone` 的语义1，`ice` 对应 `cone` 的语义2。

4.4.1 词义消歧

(2) 基于义类辞典(thesaurus) 的消歧

多义词的不同义项在使用时往往具有不同的上下文
语义类，即通过上下文的语义范畴可以判断多义词的使用义项。

如 *crane* 的两个词义“鹤”和“起重机”分别属于语义类“ANIMAL”和“MACHINERY”。不同的语义类往往具有不同的上下文环境，如，经常表示

“ANIMAL”语义类的共现词语为“species、family、eat”等，而表示“MACHINE”语义类的共现词语则为“tool、engine、blade”等。因此，只要确定多义词的上下文词的义类范畴，也就确定了多义词的词义。



4.4.1 词义消歧

(3) 基于双语词典的消歧

需要消歧的语言称为第一语言，把需要借助的另一种语言称为第二语言。建立多义词 x 与相关词 y 之间的搭配关系，然后，在第二种语言的语料库中统计对应 x 不同词义的翻译与相关词 y 的翻译同现的次数，同现次数高的搭配对应的义项即为消歧后的词义。

4.4.1 词义消歧

例如：单词 *plant* 有两个含义，一个为“植物”，另一个为“工厂”。当对 *plant* 进行词义消歧时，需要首先识别出含有 *plant* 的短语，如：*manufacturing plant*，然后，在汉语语料库中搜索与这个短语对应的汉语短语实例，由于 *manufacturing* 的汉语翻译“制造”只和“工厂”共现，因此，可以确定在这个短语中 *plant* 的词义为“工厂”。而短语 *plant life* 在汉语翻译中，“生命(*life*)”与“植物”共现的机会更多，因此，可以确定在短语 *plant life* 中 *plant* 的词义为“植物”。

4.4.1 词义消歧

(4) Yarowsky 消歧算法

基于词典的词义消歧算法都是分别处理每个出现的歧义词，并且对歧义词有两个限制：

- 每篇文本只有一个意义：在任意给定的文本中，目标词的词义具有高度的一致性；
- 每个搭配只有一个意义：目标词和周围词之间的相对距离、词序和句法关系，为目标词的意义提供了很强的一致性的词义消歧线索。

4.4.1 词义消歧

在Yarowsky 消歧算法中的处理方法:

- 对于第一个约束, 如果一个给定的多义词第一次出现时使用某个义项, 那么, 它在后面出现时也很可能使用这个义项。
- 对于第二个约束, Yarowsky (1995) 采用基于自举 (bootstrapping) 的(半监督) 学习技术。搭配特征依据如下比率排序:

$$\frac{P(s_{k_1} | f)}{P(s_{k_2} | f)}$$

两个义项与特征同
现的次数之比。

其中, s_{k_i} 为词义, f 为搭配特征。

4.4.1 词义消歧

□ 无监督的词义消歧方法

H. Schütze (1998) 提出的上下文分组辨识 (context-group discrimination) 方法是无监督的词义消歧方法的典型代表。

与(Gale, 1992) 方法类似，对于一个具有 k 个义项的词 w ，估计使用义项 s_i ($k \geq i \geq 1$) 的上下文中出现词 v_j 的概率，即 $P(v_j | s_i)$ 。

4.4.1 词义消歧

但是，在该方法中参数 $P(v_j | s_i)$ 的估计不是根据有标注的训练语料，而是在无标注的语料上进行，开始时随机地初始化参数，然后根据EM算法重新估计该概率值。

主要问题在于，很多同义词的同一个意义出现的上下文往往有很大的差异，因此，很难保证同一个意义的上下文被划分到同一个等价类中。

4.4.1 词义消歧

为了解决这个问题, H. Schütze (1992) 对词汇集中的每一个词 w 定义了关联向量(associate vector), 该向量为 w 的平均上下文。

$$A(w) = \sum_{i=1}^n \delta(w_k, w) \langle c_k^1, c_k^2, \dots, c_k^w \rangle$$

上标表示词汇集中的词形(type), 如 w^j 表示词汇集中的第 j 个词; 下标表示一个词在语料库中的一次具体使用, 简称为“词用(token)”, w_k 表示语料库中的第 k 个词; n 为词的个数, 即语料库大小; c_k^j 为词形 w^j 出现在 w_k 的上下文中的次数; $\delta(x, y)$ 为Kronecker函数。

4.4.1 词义消歧

关于该工作的详细介绍请参阅：

[Schütze, 1992a] Schütze, Hinrich. 1992. Context Space. In *Working Notes of the AAAI Fall Symposium on Probabilistic Approaches to Natural Language*, Menlo Park, CA. AAAI Press. Pages 113-120.

[Schütze, 1992b] Schütze, Hinrich. 1992. Word Sense Disambiguation with Sublexical Representation. In *Proceedings of the 1992 AAAI Workshop on Statistically-based Natural Language Programming Techniques*. Pages 90-94.



4.4.1 词义消歧

严格地讲，利用完全无监督的消歧方法进行词义标注是不可能的，因为词义标注毕竟需要提供一些关于语义特征的描述信息。但是，词义辨识 (word sense discrimination) 却可以利用完全无监督的机器学习方法实现。

本章小结-句法分析部分

□ 短语结构语法与乔姆斯基语法体系

- ◆ 0型文法（短语结构） ◆ 1型文法（上下文有关）
- ◆ 2型文法（上下文无关） ◆ 3型文法（正则文法）

□ 基于短语结构语法的句法分析方法

- ◆ 自顶向下的并行法 ◆ 自底向上的回溯法

□ 递归转移网络与扩充转移网络

□ 词汇功能语法

□ 依存句法及分析方法

□ 格语法

本章小结-词法分析部分

- 词法分析的任务（英语汉语有所不同）
- 英语形态分析
 - ◆ 单词识别
 - ◆ 形态还原
- 汉语自动分词
 - ◆ 汉语分词中的主要问题
 - ◆ 基本原则和辅助原则
 - ◆ 几种基本方法

(MM、最少分词法、统计法等)



本章小结

- 未登录词识别
 - ◆ 人名、地名、机构名等
- 词性标注
 - 问题(兼类、标注集、规范)
 - 方法(规则方法、统计方法、综合方法)
- 分词与词性标注结果评测
 - 正确率、找回率、F-测度值



本章小结

- 语义分析的基本任务及其面临的困难
- 语义计算研究概括及常见的语义理论
- 格语法(定义、格框架约束分析)
- 语义网络(概念、关系、语义网络表示、事件的语义关系、基于语义网络的推理分析)
- CD 理论(三个层次：基本动作、剧本、计划)
- 汉语语义计算研究概括
- 词义消歧(规则方法、统计方法、词典法)



习题

1. 设计并实现算法用于还原英语动词。
2. 编写程序实现汉语逆向最大分词算法（可采用有限词表），并利用该程序对一段中文文本进行分词实验，校对切分结果，计算该程序分词的正确率、召回率及F-测度。
3. 设计并实现一个汉语未登录词（汉族人名）的识别算法(可限定条件)，并通过实验分析该算法的优缺点。

习题

4. 了解目前常见的几种汉语词性标注集，比较它们的差异，并阐述你个人的观点。
5. 掌握各种词性标注方法的要点，了解目前汉语词性标注的几种主要方法。
6. 试参考前人的工作，提出消除汉语自动分词中组合歧义的几点设想。
7. 阅读《信息处理用现代汉语分词规范》(中华人民共和国国家标准 GB13715)，了解规范的基本内容。



习题

1. 阅读有关 HowNet 和HNC 理论的文献，了解相关工作及其《同义词词林》在自然语言处理中的应用。
 2. 了解蒙塔格语法(Montague Grammar)。
 3. 阅读有关词义消歧的论文，了解词义消歧的相关工作。
-



Thanks

谢谢!