

COMP9414：人工智能第5b讲。语言模型

韦恩·沃布克

电邮：w. wobcke@unsw.edu.au

概率语言模型

- 基于从大型文本/语音语料库中得出的统计数据
 - ▲ 布朗语料库（1960年代）--100万字
 - ▲ 佩恩树库（1980年代）--700万字
 - ▲ 北美新闻（1990年代）--3.5亿字
 - ▲ IBM - 10亿字

本讲座

- 语篇标签化
 - ▲ n-gram模型
 - ▲ 隐马尔科夫模型
 - ▲ 维特比算法
- 词义歧义
 - ▲ 互惠信息
 - ▲ 基于类的模型

- ▲ 谷歌、脸书和微软 - 数以万亿计的文字

- 与认为语言能力基于（先天）知识的观点相反的是
- 我的想法是语言能力是可以学习的。有足够的数据...

潘恩树库标签集

Tag	Description	Example	Tag	Description	Example
CC	coord. conjunction	<i>and, or</i>	RB	adverb	<i>extremely</i>
CD	cardinal number	<i>one, two</i>	RBR	adverb, comparative	<i>never</i>
DT	determiner	<i>a, the</i>	RBS	adverb, superlative	<i>fastest</i>
EX	existential there	<i>there</i>	RP	particle	<i>up, off</i>
FW	foreign word	<i>noire</i>	SYM	symbol	<i>+, %</i>
IN	preposition or sub-conjunction	<i>of, in</i>	TO	"to"	<i>to</i>
JJ	adjective	<i>small</i>	UH	interjection	<i>oops, oh</i>
JJR	adject., comparative	<i>smaller</i>	VB	verb, base form	<i>fly</i>
JJS	adject., superlative	<i>smallest</i>	VBD	verb, past tense	<i>flew</i>
LS	list item marker	<i>1, one</i>	VBG	verb, gerund	<i>flying</i>
MD	modal	<i>can, could</i>	VBN	verb, past participle	<i>flown</i>
NN	noun, singular or mass	<i>dog</i>	VBP	verb, non-3sg pres	<i>fly</i>
NNS	noun, plural	<i>dogs</i>	VBZ	verb, 3sg pres	<i>flies</i>
NNP	proper noun, sing.	<i>London</i>	WDT	wh-determiner	<i>which, that</i>
NNPS	proper noun, plural	<i>Azores</i>	WP	wh-pronoun	<i>who, what</i>
PDT	predeterminer	<i>both, lot of</i>	WP\$	possessive wh-	<i>whose</i>
POS	possessive ending	<i>'s</i>	WRB	wh-adverb	<i>where, how</i>
PRP	personal pronoun	<i>he, she</i>			

语篇标签化

- DT大/JJ陪审团/NN评论/VBD在/IN一个/DT数字/NN的/IN其他/J主题/NNS ./。
- 那里/EX是/VBP 70/CD儿童/NNS那里/RB
- 初步/JJ发现/NNS是/VBD报告/VBN在/IN今天/NN's/POS新/NNP英格兰/NNP杂志/NNP的/IN医学/NNP ./。

为什么这很难？

含糊不清，比如说背部

- 盈利增长处于次要地位/JJ
- 后面的一个小建筑/NNN
- 绝大多数参议员支持/VBP该法案
- 戴夫开始向后退/VB走向门口
- 使国家能够回购/偿还债务
- 那时候我21岁了/RB

概率论的表述

- 事件。词 w 的出现，带有标签 t 的词的出现
- 给定单词序列 w_1, \dots, w_n ，选择 t_1, \dots, t_n ，以便 $P(t_1, \dots, t_n | w_1, \dots, w_n)$ 是最大的。
- 应用贝叶斯规则
 - ▲ $P(t_1, \dots, t_n | w_1, \dots, w_n) = \frac{P(w_1, \dots, w_n | t_1, \dots, t_n) \cdot P(t_1, \dots, t_n)}{P(w_1, \dots, w_n)}$
 - ▲ 因此最大化 $P(w_1, \dots, w_n | t_1, \dots, t_n) \cdot P(t_1, \dots, t_n)$

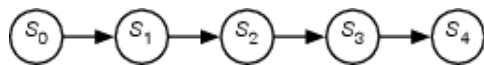
单元模式

最大化 $P(w_1, \dots, w_n | t_1, \dots, t_n) \cdot P(t_1, \dots, t_n)$

- 应用独立假设
 - $P(w_1, \dots, w_n | t_1, \dots, t_n) = P(w_1 | t_1) \cdot \dots \cdot P(w_n | t_n)$
 - 由 t 产生的词 w 的概率与上下文无关
 - $P(t_1, \dots, t_n) = P(t_1) \cdot \dots \cdot P(t_n)$
 - 标签序列的概率与顺序无关
- 估计概率
 - $P(w | t) = \#(w \text{ 与标签 } t \text{ 一起出现}) / \#(\text{有标签 } t \text{ 的词})$
 - $P(t) = \#(\text{带有标签 } t \text{ 的词}) / \# \text{ 词}$
 - 选择使 $\prod P(t_i | w_i)$ 最大化的标签序列。
 - 为每个词选择最常见的标签

- 准确率约为90%--
但每个句子中仍有 ≈ 1 个单词的错误。

马尔科夫链



■ 贝叶斯网络

- △ $P(S_0)$ 指定初始条件
- △ $P(S_{i+1} | S_i)$ 指定动态（如果每个 i 都相同，则为**静止**）。

■ 独立性假设

- △ $p(s_{i+1} | s_0, \dots, s_i) = p(s_{i+1} | s_i)$
- △ 过渡概率**仅**取决于当前状态 S_i - 与达到该状态的历史**无关** S_0, \dots, S_{i-1}
- △ 鉴于现在，未来是独立于过去的。

Bigram模型

最大化 $P(w_1, \dots, w_n | t_1, \dots, t_n)$ 。 $P(t_1, \dots, t_n)$

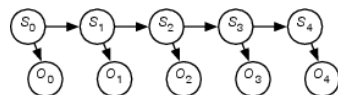
■ 应用独立假设（马尔可夫假设）。

- △ $P(w_1, \dots, w_n | t_1, \dots, t_n) = \prod P(w_i | t_i)$
- △ 观察（词）**只**取决于状态（标签）。
- △ $P(t_1, \dots, t_n) = P(t_n | t_{n-1}) \dots P(t_1 | \phi)$ ，其中 ϕ = 开始
- △ 大图模型：状态（标签）**只**取决于先前的状态（标签）。

■ 估计概率

- △ $P(t_i | t_j) = \#((t_j, t_i) \text{ 发生}) / \#(t_j \text{ 开始一个大数})$
- △ 选择使 $\prod P(w_i | t_i)$ 最大化的标签序列。 $P(t_i | t_{i-1})$
- △ 由有限状态机生成的语篇

隐马尔科夫模型



■ 贝叶斯网络

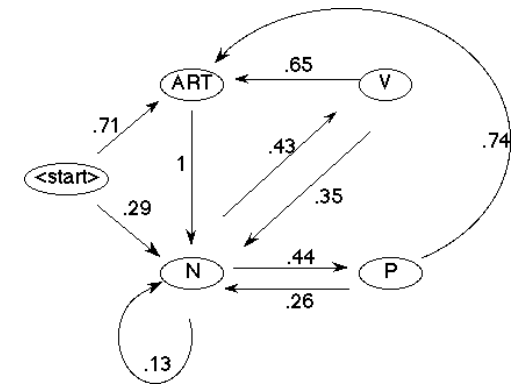
- △ $P(S_0)$ 指定初始条件
- △ $P(S_{i+1} | S_i)$ 指定动态
- △ $P(O_i | S_i)$ 指定 "观察"。

■ 独立性假设

- △ $P(S_{i+1} | S_0, \dots, S_i) = P(S_{i+1} | S_i)$ （马尔可夫链）。
- △ $p(o_i | s_0, \dots, s_{i-1}, s_i, o_0, \dots, o_{i-1}) = p(o_i | s_i)$

用于POS标签的马尔科夫模型

- △ 观察（词）**只**取决于当前状态（标签）。



词义歧义

例子

我应该换掉那把愚蠢的锁，让你留下钥匙，如果我知道你会回来打扰我，只要一秒钟。

- 锁定 = - - - -
- 离开 = - - - -
- 第二 = - - - -
- 返回 = - - - -

简单的（编造的）例子

词语	桥梁/结构	桥梁/牙齿	任何窗口
牙齿	1	10	300
悬挂	200	1	2000
的	5500	180	500 000
牙医	2	35	900
总数	5651	194	501 500

$P(\text{桥梁/结构}) = 5651 / 501\,500 = 0.0113$
 $P(\text{桥/牙}) = 194 / 501\,500 = 3.87 \times 10^{-4}$
 $P(\text{牙|桥/结构}) = 1 / 5651 = 1.77 \times 10^{-4}$
 $P(\text{牙|桥/牙}) = 10 / 194 = 0.052$

如果窗口包含 "牙齿", 则首选桥梁/牙齿

窗户

- 考虑在一个关于w的窗口中的共同发生情况

w_1					w					w_n
-------	--	--	--	--	-----	--	--	--	--	-------

- 词的感觉应该与 "相关 "的词类共同出现
- 选择w的意义s, 使 $P(w \text{ as } s | w_1, \dots, w_n)$ 最大化。
- 应用贝叶斯规则
 - 最大化 $\frac{P(w_1, \dots, w_n | w \text{ as } s) \cdot P(w \text{ as } s)}{P(w_1, \dots, w_n)}$
- 应用独立假设
 - $P(w_1, \dots, w_n | w \text{ as } s) = \prod P(w_i | w \text{ as } s)$

相互信息

$$MI(x, y) = \log \frac{P(x, y)}{P(x) \cdot P(y)}$$

$$\log \frac{MI(sense(w), w)}{N \cdot \#(sense(w_1), w_2)}$$

- 估计概率。 $P(w_i | w \text{ as } s)$
 - $\Delta \#(w_i \text{ 在 } w \text{ 为 } s \text{ 周围的 } n \text{ 个词窗口中}) / \#(w \text{ 为 } s \text{ 的窗口})$

$$\frac{1}{N} \sum_{w_2} \frac{1}{\#(w_2)}$$

$$2^{\#(sense(w_1))}.$$

$\#(w_2)$

其中 N 是语料库的大小

- $MI = 0$: $sense(w_1)$ 和 w_2 是有条件的独立。

- $MI < 0$: $sense(w_1)$ 和 w_2 一起出现的次数少于随机次数
- $MI > 0$: $sense(w_1)$ 和 w_2 一起出现的次数多于随机出现的次数
- 增加相互信息就相当于假设了独立性
- 选择 w 的感觉 $s = \operatorname{argmax}_{s \in \text{senses}(w)} \sum_{w_i \in \text{wzbn}(w)} MI(s, w_i)$

基于类的方法

- 使用预定义的 "意义类", 例如WordNet、Wikipedia
 - ▲ 锁 → 机械装置 ← 工具、曲柄、齿轮、 - - - -
 - ▲ 锁 → 身体部位 ← 头发、眼睛、手、 - - -
- 通过添加词的意义来计算词的意义计数
- 把所有的 $P(w|s)$ 仅仅看作是 $P(s)$, 即 $P(\text{任何具有}s\text{意义的词})$ 。
- 优势
 - ▲ 减少空间和时间的复杂性
 - ▲ 降低数据稀疏度
 - ▲ 允许无监督的学习

总结

- 统计学（和神经网络）模型在许多任务上表现良好
 - ▲ 语料部分标记
 - ▲ 词义歧义
 - ▲ 对传统分析器的控制
 - ▲ 概率分析法
- 问题
 - ▲ 不切实际的简化假设（似乎有效）。
 - ▲ 要求有非常多的（标示）文本
 - ▲ 在（甚至是大型）文本语料库中的单词出现的稀疏程度

▲ 随着时间的推移，词语使用的变化（例如，参议员奥巴马）。
