

词性标注

给定单词 w , 找到对应的词性 z

1. 问题概述

已知 (S)

Sentence = $w_1 w_2 w_3 w_4 w_5$ 单词
 ↓
 $z_1 z_2 z_3 z_4 z_5$ 词性

未知 $S' = w'_1 w'_2 w'_3 w'_4 w'_5$ 求标注词性

Noise channel model:

$$p(z|S) = \underset{\text{argmax}}{p(S|z)} \underset{\text{translation model}}{p(z)} \underset{\text{language model}}{p(z)}$$

$$= p(w_1 w_2 w_3 \dots w_n | z_1 z_2 \dots z_n) p(z_1 z_2 \dots z_n)$$

$$= \prod_{i=1}^n p(w_i | z_i) \cdot p(z_1) p(z_2 | z_1) \dots p(z_n | z_{n-1})$$

emission probability

给定一个词性, 在该 z_i 下单词出现概率是多少

$p(z_1)$ 该词性出现在句首的概率

$p(z_n | z_{n-1})$ bigram

$$\hat{z} = \underset{\text{argmax}}{p(z|S)}$$

$$= \underset{\text{argmax}}{\prod_{i=1}^n p(w_i | z_i) p(z_1) \cdot \prod_{i=1}^n p(z_n | z_{n-1})}$$

$$= \underset{\text{argmax}}{\log (\quad \quad \quad)}$$

$$= \underset{\text{argmax}}{\sum \log p(w_i | z_i) + \log p(z_1) + \sum \log p(z_n | z_{n-1})}$$

↓ 也就是求三个概率

$$\theta = \{A, B, \lambda\}$$

step 1. 求 A, B, λ

step 2. 求 z , (Viterbi Arg)

一般情况下 $p(w_1 w_2 w_3 | z_1 z_2 z_3) \neq p(w_1 | z_1) p(w_2 | z_2) p(w_3 | z_3)$
 我们在此简化了 $p(w_1 | z_1) p(w_2 | z_2) p(w_3 | z_3)$
 (也就是假设当前词 w_i 只依赖于 z_i)

2. 三个所要求的概率

$$\begin{aligned}\hat{z} &= \operatorname{argmax} p(z|s) \\ &= \operatorname{argmax} \prod_{i=1}^n p(w_i | z_i) p(z_1) \cdot \prod_{i=1}^n p(z_n | z_{n-1}) \\ &= \operatorname{argmax} \log (\quad \quad \quad) \\ &= \operatorname{argmax} \sum \log \underbrace{p(w_i | z_i)}_A + \log \underbrace{p(z_1)}_B + \sum \log \underbrace{p(z_n | z_{n-1})}_B\end{aligned}$$

1. $A \quad p(w_i | z_i)$

$$A = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}_N$$

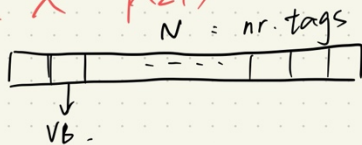
N : tags nr.
 M : words nr.

给定 i 假设 tag 是 vb.

[apple, study, ...]
 词库

在给定动词情况下每个单词出现概率.

2. $\sum p(z_i)$



一个句子以某个词性开始的 p

3. $B \quad p(z_n | z_{n-1})$ tags

$$B = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}_N$$

见 `jupyter9 pos tagger`

$\hat{z} = \operatorname{argmax}_{\vec{z}} \left[\sum_{i=1}^T \log P(w_i | z_i) + \log(z_1) + \sum_{t=2}^T P(z_t | z_{t-1}) \right]$

If Newsweek said it will introduce the

(NKP VBD PRP MD VB PT)

w_1 w_2 w_3 w_4 w_5 w_6

{ NKP
VBD
:
:
:
}

54 {

$O(54^{6.5}) \approx O(P^n)$

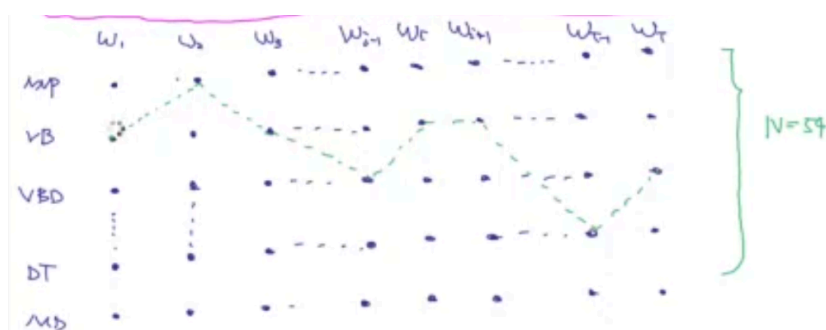
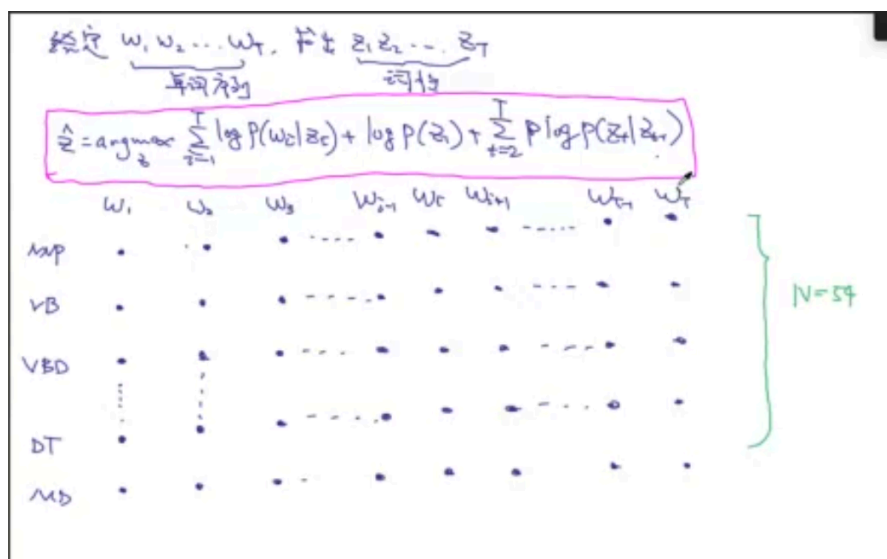
\Rightarrow (NKP, VB, VBD, DT, NKP, NKP)

$\log P(\text{Newsweek} | \text{NKP}) + \log P(\text{said} | \text{VB}) - \dots - \log P(\text{the} | \text{NKP})$

+ $\log P(\text{NKP}) + \log P(\text{VB} | \text{NKP}) + \log P(\text{VBD} | \text{VB})$

+ ... + $\log P(\text{NKP} | \text{NKP}) = \checkmark$

如果每个单词 w ，都去遍历一遍所有的 tag ，计算量太大。
Viterbi 算法找出的是全局最优解。



找路径，使得概率最小。

$$\text{Score}(n) = \log P(VB) + \log P(w_1 | VB) + \dots + w_1$$

$$+ \log P(NP | VB) + \log P(w_2 | NP) + \dots + w_2$$

$$+ \log P(VB | NP) + \log P(w_3 | VB) + \dots + w_3$$

$$+ \dots$$

~~find~~ dp

$dp[i, j]$: the best path such that ending at j at time i

(w_1, \dots, w_i) , $w_i \leftarrow$ 第 j 个 tag

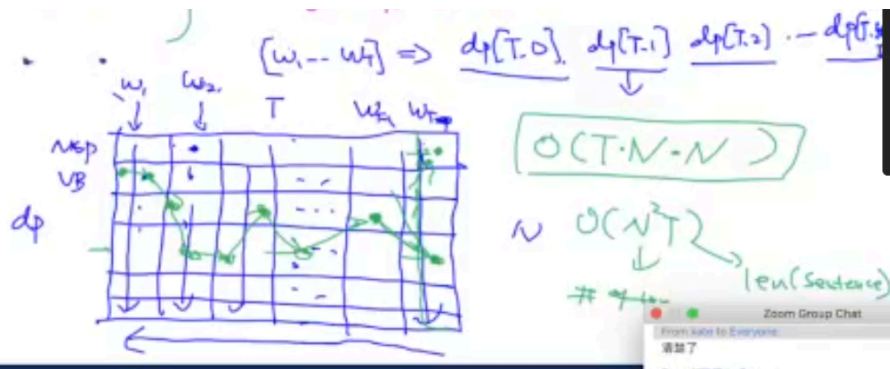
w_1, \dots, w_i v_1, \dots, v_i

$$dp[i, j] = \begin{cases} dp[i-1, 0] + \log P(DT | NP) + \log P(w_i | DT) \\ dp[i-1, 1] + \log P(DT | VB) + \log P(w_i | DT) \\ dp[i-1, 2] + \log P(DT | VBD) + \log P(w_i | DT) \\ \dots \end{cases}$$

$N=59$

$i = 1 \dots T$, $j = 0 \dots 54$

$(w_1, \dots, w_T) \Rightarrow dp[T, 0], dp[T, 1], dp[T, 2], \dots, dp[T, 54]$



参考 jupyter9

