

一、从HMM到CRF

1. 任务类型: Sequence Labelling 输入和输出都是一个序列

$$\begin{array}{cccc} x_1 & x_2 & \dots & x_L (\text{seq}) \\ \downarrow & \downarrow & \dots & \downarrow \\ y_1 & y_2 & \dots & y_L (\text{seq}) \end{array}$$

2. 案例: POS Tagging (词性标注)

3. HMM 步骤

① 首先构建词性(隐藏状态)的转移矩阵 & 初始概率,

$$P(PN \ V \ D \ N) = \frac{P(PN | \text{start}) P(V | PN) P(V | PN) P(D | V)}{P(N | D)} \quad \begin{array}{l} \text{计算最有可能的词性序列} \\ \text{(联合概率)} \end{array}$$

② 然后计算联合概率的前提下,由发射矩阵,确定取当前单词的条件概率

$$P(\text{John saw the saw} | PN \ V \ D \ N) = \frac{P(\text{John} | PN) * P(\text{saw} | V) * P(\text{the} | D) * P(\text{saw} | N)}{P(N | D)}$$

$$\textcircled{3} \text{ 令 } P(PN \ V \ D \ N) = P(y)$$

$$P(\text{John saw the saw} | PN \ V \ D \ N) = P(x|y)$$

$$P_{\text{joint}} = P(x, y) = P(y) P(x|y)$$

公式:

$$\text{step 1: } P(y) = P(y_1 | \text{start}) \times \prod_{l=1}^{L-1} P(y_{l+1} | y_l) \times P(\text{end} | y_L) \quad \text{转移概率}$$

$$\text{step 2: } P(x|y) = \prod_{l=1}^L P(x_l | y_l) \quad \text{发射概率}$$

$$4. \text{HMM 问题求解: } y = \arg \max P(y | x) = \arg \max \frac{P(x, y)}{P(x)}$$

$$\because P(x) \text{ 已知 } y = \arg \max P(x, y)$$

- ① 穷举所有的 y , 找到 $p(x, y)$ 最大的解
- ② 维特比算法

5. 问题:

HMM 的矩阵是基于语料库的, 对于未出现在语料库的词汇 HMM 会统计出一个值, 但这个值不一定很小



所以这一情况符合当数量级很小的情况, HMM 效果一般会更好



原因是转移矩阵和发射矩阵之间是独立的, 要分别 model



在使用上述矩阵的前提下, 解决 HMM 的问题

6. CRF:

① 模型假设 $p(x, y) \propto \exp(w \cdot \phi(x, y))$

• $\phi(x, y)$ 特征矩阵

• w 权重 (字训练字词的)

• $\exp(w \cdot \phi(x, y))$: 确保值大于 1

$$p(x, y) = \frac{\exp(w \cdot \phi(x, y))}{R}$$

$$p(y|x) = \frac{p(x, y)}{\sum_{y'} p(x, y')} = \frac{\exp(w \cdot \phi(x, y)) \cancel{X_R}}{\sum_{y'} \exp(w \cdot \phi(x, y')) \cancel{X_R}}$$

$$= \frac{\exp(w \cdot \phi(x, y))}{\sum_{y'} \exp(w \cdot \phi(x, y'))} \quad (\sum_{y'} \text{ 所以 } y \text{ 的可能})$$



归一化系数, 总和为 1

$$\therefore = \exp(w \cdot \phi(x, y)) / Z(x)$$

② 样本说明: 假设一个 Batch N 个样本, 一个序列 R 为 n
上标 i : 第 i 个样本 下标 t : 第 t 个样本

(l_1, l_2, \dots)

(s_1, s_2, \dots)

$$\begin{aligned} x_1^1, x_2^1, \dots, x_n^1 &\rightarrow y_1^1, y_2^1, \dots, y_n^1 \\ x_1^2, x_2^2, \dots, x_n^2 &\rightarrow y_1^2, y_2^2, \dots, y_n^2 \\ &\vdots \\ x_1^N, x_2^N, \dots, x_n^N &\rightarrow y_1^N, y_2^N, \dots, y_n^N \end{aligned}$$

③ 标签说明: $P(y|x) = \frac{e^{w \cdot \phi(x,y)}}{\sum_y e^{w \cdot \phi(x,y)}}$ 设一共 L 个单词, S 个标签

- 单词和标签之间组合: $|L| \times |S|$
- 标签之间组合: $|S| \times |S|$
- $\text{start} \rightarrow s$ 和 $s \rightarrow \text{end}$: $2 \times |S|$

所以 $\phi_k(x,y)$ 是一个计数结果 $N_{\text{start}(x,y)}$, 表示特点样本 (x,y) 中发生的次数,

已知 HMM $P(x,y) = P(y) P(x|y)$

$$P(y, | \text{start}) \times \prod_{l=1}^{L-1} P(y_{l+1} | y_l) \times P(\text{end} | y_L) \times \prod_{l=1}^L P(x_l | y_l)$$

$$\log P(x,y) = \log P(y, | \text{start}) + \sum_{l=1}^{L-1} \log P(y_{l+1} | y_l) + \log P(\text{end} | y_L) + \sum_{l=1}^L \log P(x_l | y_l)$$

比如: x The dog ate the homework
 y D N V D N

这里 (The, D) 出现 2 次, 所以 $\log P(x,y)$ 可以按次数进行统计:

$$\begin{aligned} \log P(x,y) = & \sum_s \log P(s | \text{start}) \times N_{\text{start}, s}(x,y) + \\ & \sum_{s,s'} \log P(s' | s) \times N_{s, s'}(x,y) + \\ & \sum_s \log P(\text{end} | s) \times N_{s, \text{end}}(x,y) + \\ & \sum_{s,t} \log P(t | s) \times N_{s, t}(x,y) \end{aligned}$$

\downarrow
word + tags

所以 $\log P(x, y)$ 可以看作是两列向量乘积

$$\begin{bmatrix} \vdots \\ \log P(H|S) \\ \vdots \\ \log P(S|\text{start}) \\ \vdots \\ \log P(S'|S) \\ \vdots \\ \log P(\text{end}|S) \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} \vdots \\ N_{s,+}(x, y) \\ \vdots \\ N_{\text{start},s}(x, y) \\ \vdots \\ N_{s,s'}(x, y) \\ \vdots \\ N_{s,\text{end}}(x, y) \\ \vdots \end{bmatrix} = w \cdot \phi(x, y)$$

不同样本 (x^i, y^i) 的 $\phi(x^i, y^i)$ 不同, 但 w 训练参数共享

④ CRF++ 区别:

1. 没有 start 和 end 标签 (更灵活)
2. 特征使用模板抽取, 可以跨词之间的组合

↓

相同的词的模板不同, 它的词缀也不同