

Lev 10 语句标注

→ 标注任务

△ Part of speech tagging (POS) → 词性标注

2类词性标注：Open: 独立词
Closed: 语法功能词

→ 其他 Tasks 挑战：

{ parsing (syntactic)

translation

Sentiment/affective tasks → 语义素 adj.

text to speech

△ Chinese word segmentation (B/E/S)

△ Named entity recognition (B/E/S/O)

(PER, -LOC, -ORG)

△ Semantic role labeling (语义角色标注)

识别谓词及其语义角色 (施动、受动等)

(PRED, -ARG0, -ARG1) (B/E/S/O)

△ For POS Easiest Way:

逐个词预测固定标签

→ 未考虑上下文 / 和邻接词间关系

⇒ 可采用其他方法进行 POS Tagging!

⇒ HMM, CRF, Neural

* HMM Model (隐马尔可夫)

△ Model

$$P(x_1 \dots x_n, \gamma_0 \dots \gamma_{n+1}) = g(\gamma_{n+1} | \gamma_n) \prod_{t=1}^n g(\gamma_t | \gamma_{t-1}) e(x_t | \gamma_t)$$

转移 \downarrow $\gamma_0: \text{Start}$ $\gamma_{n+1}: \text{Stop}$
 观察 \downarrow $x_1 \dots x_n$
 标注 \downarrow $\gamma_1 \dots \gamma_n$

0 高级 > 2 阶关系：(High order HMM)

已知： $P(\gamma_t | \gamma_{t-1} \dots \gamma_{t-n+1}) \Rightarrow$ Similar to n-gram

△ Decoding

轨迹找到 $\hat{\gamma} = (\gamma_1, \gamma_2 \dots, \gamma_{n+1})$, 达到 Stop

$$\hat{\gamma} = \arg \max_{(\gamma_1 \dots \gamma_{n+1})} P(\gamma_1 \dots \gamma_{n+1} | x_1 \dots x_n) = \arg \max_{(\gamma_1 \dots \gamma_{n+1})} P(\gamma_1 \dots \gamma_{n+1}, x_1 \dots x_n)$$

$$\hookrightarrow = \prod_{t=1}^n P(\gamma_t | \gamma_{t-1}) P(x_t | \gamma_t) \quad (\text{转移 + 观察模型})$$

0 (暴力枚举) \rightarrow exp 复杂度

Viterbi Algorithm:

$\pi(i, \gamma_i)$: 前*i*位置、末尾为 γ_i 的最大概率

$$\pi(i, \gamma_i) = \max_{\gamma_{i-1}} P(x_i | \gamma_i) \pi(i-1, \gamma_{i-1})$$

\Rightarrow 只看最后一个 γ 建模

$$= \max_{\gamma_{i-1}} [e(x_i | \gamma_i) q(\gamma_i | \gamma_{i-1})] \pi(i-1, \gamma_{i-1}) \quad \text{DP转移方程}$$

$$= [e(x_i | \gamma_i) \cdot \max_{\gamma_{i-1}} q(\gamma_i | \gamma_{i-1}) \pi(i-1, \gamma_{i-1})]$$

① Initialize: $\pi(0, \text{Start}) = 1$, for other γ_0 : $\pi(0, \gamma_0) = 0$

② 第 i 步时：计算所有标签 $\gamma_i \rightarrow \pi(i, \gamma_i)$

$$③ P(\hat{\gamma}) = \max_{\gamma_n} p(\text{Stop} | \gamma_n) \pi(n, \gamma_n) = \pi(n+1, \text{Stop})$$

\hookrightarrow 从后往前推

复杂度：序列长 n, 标签数 L \Rightarrow $nL \times \pi(n, \gamma_n)$

\Rightarrow 高每个 $\pi(i, \gamma_i)$ 复杂度 $O(L)$ \Rightarrow 总复杂度 $O(nL^2)$

△ Marginal Inference

$$\text{计算 } P(x_1 \dots x_n) = \sum_{y_1 \dots y_n} P(x_1 \dots x_n, y_1 \dots y_{n+1})$$

(输入序列概率)

放样标注 \Rightarrow 指数级复杂度

$\Rightarrow \alpha(i, y_i)$: 前*i*位置、末尾为*y_i*边的概

$$= P(x_1 \dots x_i, y_i)$$

$$= \sum_{y_{i+1}} e(x_i | y_i) q(y_i | y_{i-1}) P(x_1 \dots x_{i-1}, y_i \dots y_{i-1})$$

$$= e(x_i | y_i) \sum_{y_{i+1}} q(y_i | y_{i-1}) \alpha(i-1, y_{i-1}) \quad \text{递推方程}$$

① 初始化: $\alpha(0, \text{Start}) = 1$, 其余 $\alpha(0, y_0) = 0$

② 第*i*步算所有 $\alpha(i, y_i) \rightarrow$ 下一步

$$③ P(x_1 \dots x_n) = \sum_{y_n} \alpha(\text{Stop} | y_n) \alpha(n, y_n) = \alpha(n+1, \text{STOP}) \Rightarrow \xrightarrow{\alpha} \xleftarrow{\alpha} \text{Forward-Backward Algorithm}$$

复杂度: $nL \cdot L \Rightarrow O(nL^2)$

△ Supervised Learning

$$\text{MLE: } P(x, y) = \prod_{i=1}^{n+1} e(x_i | y_i) q(y_i | y_{i-1})$$

$$= \left(\prod_{i,j} q(j|i)^{c(i,j)} \right) \left(\prod_{j \in X} \prod_{i \in Y} e(j|i)^{c(i,j)} \right)$$

→ 转移部分 → 边界部分

(期望概率最大)

→ 计数函数

$$\text{计数: } c(u|v) = \frac{C(u,v)}{\sum_{v' \in X} C(u,v')}$$

$$c(r|o) = \frac{C(o,r)}{\sum_{r' \in Y} C(o,r')}$$

(梯度下降法: N-grams)

△ Unsupervised Learning

(\Rightarrow 只有 x 无 $y \Rightarrow$ 需要语料 y-words)

\rightarrow EM (Baum-Welch Algorithm) For HMM

① E 步: $\text{算出} \alpha \text{ 然后} \beta \rightarrow \text{分布} \rightarrow \exp$

⇒ 期望计数:

$$c(S) = E_{(y_1 \dots y_{n+1} | x_1 \dots x_n)} [\text{count}(S | x_1 \dots x_n, y_1 \dots y_{n+1})]$$

\Rightarrow 描述该序列表达式 S 为期望生成次数

$$c(NN) = \sum_j P(y_i = NN | x_1 \dots x_n)$$

$$c(NN \rightarrow VB) = \sum_j P(y_i = NN, y_{i+1} = VB | x_1 \dots x_n)$$

$$c(NN \rightarrow apple) = \sum_{i: y_i = apple} P(y_i = NN | x_1 \dots x_n)$$

⇒ $\xrightarrow{\alpha} \xleftarrow{\alpha}$ Forward-Backward Algorithm

② M 步: 计算新参数 (e.g.)

$$\text{期望: } E_{(y_1 \dots y_{n+1})} [\log P(x_1 \dots x_n, y_1 \dots y_{n+1})]$$

→ 待更新参数

$$(直接式解) \rightarrow \ell_{ML} \cdot \bar{q}_{ML}$$

* 除 EM 外: 也可用梯度下降求解

$$P(x_1 \dots x_n) :$$

计算边缘概率 $\alpha(n+1, \text{Stop})$ (前向法)

再在计算图上反向传播 \rightarrow e.g. 梯度下降

→ 与前-后向几乎一致

△ MEMM (最大熵 Markov)

Background: HMM 可解决标注(label)关系问题,

解决以下问题:

变化过程:

$$\text{HMM: } P(x_{1:n}, y_{1:n}) = \prod_t P(y_t | y_{t-1}) P(x_t | y_t)$$

→ generative model

$$\text{MEMM: } P(y_{1:n} | x_{1:n}) = \prod_t \underbrace{P(y_t | y_{t-1}, x_t)}_{\text{discriminative model}}$$

→ discriminative model

$$= \frac{1}{Z(y_{1:n}, x_{1:n})} \exp(s(y_{t-1}, y_t, x_t))$$

$$= W^T f(y_{t-1}, y_t, x_t)$$

(label bias problem)

(最终模型)

$$\text{MEMMM: } P(y_{1:n} | x_{1:n}) = \prod_t \underbrace{P(y_t | y_{t-1}, x_{1:n})}_{\text{解: 不考虑标记串的多样性}}$$

$$\rightarrow \sum_{y_{t-1}} \frac{1}{Z(y_{1:n}, x_{1:n})} \exp(s(y_{t-1}, y_t, x_{1:n})) \rightarrow \text{label bias}$$

(但解决了以下问题)

Sol: 不考虑标记串的多样性

△ CRF

$$P(y_{1:n} | x_{1:n}) = \frac{1}{Z(x_{1:n})} \prod_t \exp(s(y_{t-1}, y_t, x_{1:n}))$$

① Decoding

$$\text{找到 } y^* = \arg \max_{y_{1:n}} \sum_{t=1}^{n+1} \exp(s(y_{t-1}, y_t, x_{1:n}))$$

$$= \arg \max_{y_{1:n}} \sum_{t=1}^{n+1} s(y_{t-1}, y_t, x_{1:n})$$

\Rightarrow Viterbi 算法或 max

$$\pi(i, y_i) = \max_{y_{i-1}} [s(y_{i-1}, y_i, x_{1:n}) + \pi(i-1, y_{i-1})]$$

$$(\pi(0, \text{Start}) = 0) \Rightarrow \text{初始得分}$$

② 监督学习

1. 通过 γ^* 和 x 来最小化

① Forward Algorithm 计算 $Z(x_{1:n}) \Rightarrow$ 导出 $L = -\log P(\gamma^* | x)$

② 梯度下降

梯度包含潮流计算 $c(S) \rightarrow$ 前-后向算法
auto-differentiation

2. 最小化方法推导 \rightarrow SSVM

$$L_{\text{SSVM}} = \max_{y_{1:n}} \left(S(y_{1:n}) + \underbrace{\Delta(y_{1:n}, \gamma_{1:n}^*) - S(\gamma_{1:n}^*)}_{\text{梯度}} \right)$$

\rightarrow 二分类, ≥ 0

\rightarrow 考虑 decision boundary 和分布整体分布

\hookrightarrow 即: 又注重状态一概率分布, 又注重标签分布

\rightarrow 如果 Δ 可操作性分布 $\rightarrow L_{\text{SSVM}} \leftarrow (\#) \text{Viterbi 算法}$

优化: $\begin{cases} \text{Subgradient} \\ \text{Quadratic programming} \end{cases}$

③ 无监督学习

\hookrightarrow Encoder \rightarrow CRF

直接通过 y_t 发射出 x_t (Decoder)

$$L = -\log P(x_{1:n} | x_{1:n})$$

$$= \sum_{y_{1:n}} P(y_{1:n} | x_{1:n}) P(x_{1:n} | y_{1:n})$$

$$= \frac{1}{Z} \prod_t \exp(s_t)$$

\Rightarrow 可用前向算法计算梯度下降优化

* Neural Method

△ 静态 embeddings 无法表示 Label \Rightarrow 2 problems

△ RNN \rightarrow 只能处理上一个

Bidirectional RNN \rightarrow 上下文 ✓

△ Transformer ✓

△ Neural CRF : 神经方法计算 CRF potentials

(如发射得分)

For all:

◦ Inference: $\left\{ \begin{array}{l} \text{只用 CRF: neural softmax: 通过量纲} \\ \text{且用 CNN} \end{array} \right.$

只用 CNN

◦ Learning: $\left\{ \begin{array}{l} \text{只用 CRF: Viterbi} \\ \text{且对数条件似然} \\ \text{动态规划} \end{array} \right.$