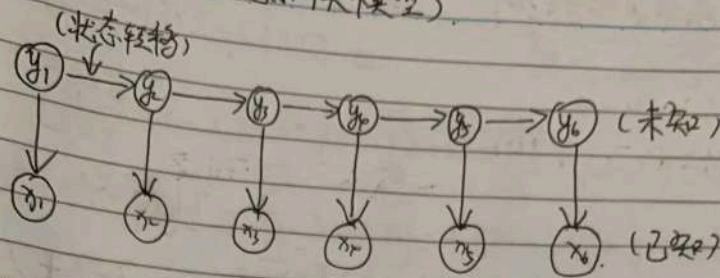


HMM (隐马尔可夫模型)



HMM:

- ① 用于标注问题
- ② 生成模型

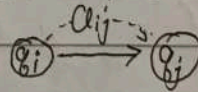
1) 定义:

HMM的三要素:

① 状态集合 Q \Rightarrow 状态转移矩阵 A : 种类

$$A = [a_{ij}]_{N \times N} \quad N \text{ 表示状态的数量}$$

* a_{ij} 表示 t 时刻处于状态 q_i 的条件下,
在 $t+1$ 时刻转移到状态 q_j 的概率

② 观测集合 V \Rightarrow 观测概率矩阵 B .

$$B = [b_{ik}]_{N \times m} \quad m \text{ 表示观测值的数量}$$

表示在时刻 t 处于 q_t 的条件下, 生成观测 v_k 的概率; 即 $b_{ik} = P(O_t = v_k | i = q_t)$

③

初始状态概率向量 π .

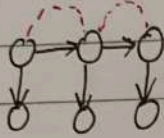
$$\pi = [\pi_i]$$

$\pi_i = P(i = q_1)$ 表示在时刻 t 处于 q_i 的概率.
(只在对初始状态起作用).

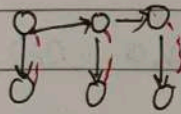
2) 两个假设:

① 齐次马尔可夫性假设。

任意时刻 t 的状态只依赖于其前一时刻的状态。



② 观测独立性假设。任意时刻的观测, 只依赖于该时刻的状态。



* 需写出以下内容 (对 HMM):

① 状态集合。

⑤ 状态转移的概率表

② 观测集合。

⑥ 观测概率表。

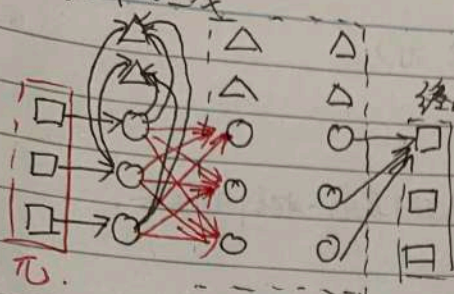
③ 状态序列和观测序列的长度。

④ 初始概率分布 (向量)

3) 观测序列的生成:

初始状态

3) 观测序列的生成

 Δ : 观测值 \bigcirc : 状态值

终止状态

4) HMM的3个基本问题:

1) 概率计算问题:

给定 $\lambda(A, B, \pi)$ 和观测序列 $O = (o_1, o_2, \dots, o_T)$, 计算条件概率:

$$P(O|\lambda)$$

2) 学习问题:

已知观测序列 $O = (o_1, o_2, \dots, o_T)$ 求 $\lambda(A, B, \pi)$ 的参数.

即: 用 MLE 求参数.

3) 预测问题(解码):

已知模型 $\lambda = (A, B, \pi)$ 和观测序列 $O = (o_1, o_2, \dots, o_T)$, 求条件概率:

$$P(I|O, \lambda)$$

5) 概率计算问题, 即求 $P(O|\lambda)$.① 直接计算法: (O : 观测序列, I : 状态序列)

$$P(O, I|\lambda) = P(O|I, \lambda) P(I|\lambda) \quad \text{即: 对固定 } I, \text{ 求 } P(O, I|\lambda)$$

$$P(O|\lambda) = \sum_I P(O, I|\lambda) \quad \text{即: 对所有可能的 } I \text{ 求和.}$$

时间复杂度: $O(TN^T)$ (N 表示状态种类, T 表示 O 的长度,此时 O 确定, 故和 m (观测值种类)

无关).

② 前向向后算法: (求 $P(O|\lambda)$) (根据观测, 只记符号频率)

a) 前向概率:

给定隐马尔可夫模型 λ , 定义到时刻 t 全部观测序列为 O_1, O_2, \dots, O_t 且状态为 q_t 的前向概率为前向概率, 记作:

$$\alpha_t(i) = P(O_1, O_2, \dots, O_t, i_t = q_i | \lambda).$$

表示 q_i , 可求

表示状态, 可求所有的状态值.

b) 算法流程:

输入: 隐马尔可夫模型 λ , 观测序列 O .

输出: 观测序列概率 $P(O|\lambda)$.

(1) 初值 $\alpha_1(i) = \pi_i b_i(O_1)$ $i = 1, 2, \dots, N$
($t=1$) (N 表示状态的种类).

(2) 递推, 对 $t = 1, 2, \dots, T-1$.

$$\alpha_{t+1}(i) = \left[\sum_{j=1}^N \alpha_t(j) a_{ji} \right] \underbrace{b_i(O_{t+1})}_{\text{观测概率}}, i = 1, 2, \dots, N.$$

(3) 终止:

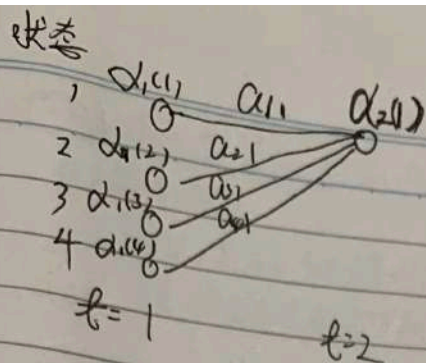
$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i).$$

(T 时刻所有状态的概率和).

$\alpha_t(j) a_{ji}$: 从 t 时刻, 状态 j 转移到状态 i 的概率 (此时, O_1, \dots, O_t 的概率由 $\alpha_t(j)$ 表示)

$\sum_{j=1}^N \alpha_t(j) a_{ji}$ 所有到 $t+1$ 时刻状态为 i 的概率和.
(当前状态, 可求之前所有状态转移过来).

$b_i(O_{t+1})$ 在状态 i , 观测 O_{t+1} 的概率.



MEMO NO. _____
DATE _____

* 时间复杂度 $O(N^2T)$,
每次引用前一次计算, 减少计算量.

③ 后向算法:

(1) 后向概率:

给定隐马尔可夫模型 λ 及在时刻 t 状态 q_i 的条件下, 从 $t+1$ 到 T 的部分观测序列 $O_{t+1}, O_{t+2}, \dots, O_T$ 的后向概率:

$$\beta_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T | i_t = q_i, \lambda).$$

(2) 后向算法:

输入: 隐马尔可夫模型 λ , 观测序列 O .

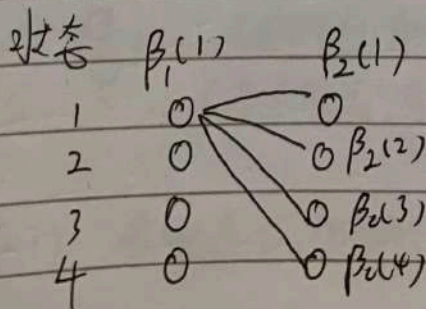
输出: 观测序列后向概率 $\beta(O|\lambda)$.

(1) $\beta_T(i) = 1 \quad i=1, 2, \dots, N.$

(2) $t = T-1, T-2, \dots, 1$

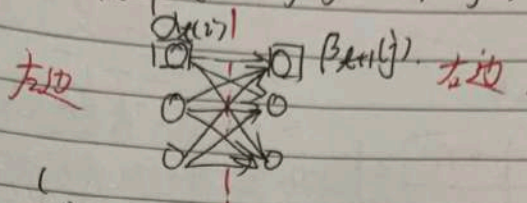
$$\beta_t(i) = \sum_{j=1}^N \alpha_{ij} b_j(O_{t+1}) \beta_{t+1}(j), \quad i=1, 2, \dots, N.$$

(3) $\beta(O|\lambda) = \sum_{i=1}^N \pi_i b_i(O_1) \beta_1(i).$



6) HMM 观测序列概率 $P(O|A)$ 公式:

$$P(O|A) = \sum_{i=1}^N \sum_{j=1}^N \alpha_i(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j), t=1, 3, \dots, T-1.$$



(其父节点一个联合概率)

* 前向后算法的图:

给出入, 计算 O 的概率
(观测值)

7) 学习算法: (求入(A, 元/的个数)

监督: MLE

非监督: Baum-Welch 算法, EM 算法.

8) 预测算法: (求 $P(I|O, A)$, 即状态序列的概率).

近似算法: 每个时刻最有可能的状态, 即贪心,

维特比算法: 全局最优.

9) 关键点:

① 初始时, π 表示某一状态的概率值.

② 假设求某一时刻 t 的状态概率, 前一时刻 $t-1$ 到某一状态的概率, 再求该状态下对应观测值的概率.

③ 不同状态之间可能相互转移, 即转移矩阵 A .

④ 某一观测值可由不同状态获取, 即观测概率 b_j .

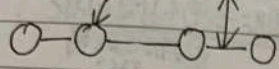
⑤ 先验的再观测 (第一个状态由 π 取得).

CRF (条件随机场)

又称为 '马尔可夫随机场'.

① 圈图表示概率分布

② 结点表示随机变量, 边表示随机变量之间的依赖关系.



③ 联合概率分布 $P(\mathbf{Y})$ 和表示它的无向图 G 中, 随机变量存在以下性质:

成对马尔可夫性
局部马尔可夫性
全局马尔可夫性

三者等价.

HMM 的两个假设:

齐次马尔可夫性假设 (当前状态仅依赖于前一个状态, 有向图)

观测独立性假设.

* 成对马尔可夫性:

两个变量 Y_i, Y_j , 没有边连接, 则两个变量在其余所有结点的条件下, 条件独立:

$$P(Y_i, Y_j | Y_o) = P(Y_i | Y_o) P(Y_j | Y_o)$$

(简单理解为, 没有边连接的两个变量, 条件独立).

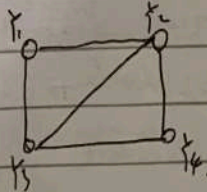
特点: (设 $P(\mathbf{Y})$ 表示概率无向图的联合概率分布).

① $P(\mathbf{Y})$ 满足 **马尔可夫性

② 易于因子分解.

2) 无向图模型与最大团

① 团: 任何两个结点均有边连接子集



② 最大团: 结点数量最多的团.

最大团 $\{Y1, Y2\}$ 和 $\{Y2, Y3, Y4\}$

3) 无向图模型的因式分解:

定义: 联合概率由最大团上的随机变量和互斥的乘积形式表示.

表示: (由 Hammersley-clifford 定理保证).

$$P(Y) = \frac{1}{Z} \prod_C \psi_C(Y_C)$$

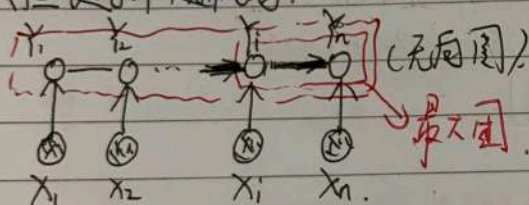
C : 为无向图上的最大团.

Y_C : C 上的随机变量.

Z : 规范化因子, $Z = \sum_Y \prod_C \psi_C(Y_C)$ (所有可能的 Y 取值和).

$\psi_C(Y_C)$: 势函数, 严格为正, $\psi_C(Y_C) = \exp\{-E(Y_C)\}$.

4) 线性链条件随机场.



定义: 设 $X = (X_1, X_2, \dots, X_n)$, $Y = (Y_1, Y_2, \dots, Y_n)$ 均为线性链条件随机变量序列, 若给定 X , 且 Y 满足马尔可夫性,

$$\text{即 } P(Y_i | X, Y_1, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n) = P(Y_i | X, Y_{i-1}, Y_{i+1})$$

$i = 1, 2, \dots, n$ (在 $i=1$ 和 n 时只考虑单边).

相邻节点

则称 $P(Y|X)$ 线性链条件随机场, X : 观测序列, Y : 状态序列.

5) 条件随机场的参数化形式:

设 $P(Y|x)$ 为线性链条件随机场, 则在 x 取值为 x 的条件下, 随机变量 Y 取值为 y 的条件概率具有如下形式:

$$P(y|x) = \frac{1}{Z(x)} \exp \left\{ \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,c} \mu_c s_c(y_i, x, i) \right\}$$

λ_k 权重 t_k 特征函数 μ_c 权重 s_c 特征函数
 $\sum_{i,k}$ 在置上的所有特征函数

(PR) = $\frac{1}{2} \log(P(x))$ ($Z(x)$ 的负对数参考例 11.2).

归一化因子: $Z(x) = \sum_y \exp \left\{ \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,c} \mu_c s_c(y_i, x, i) \right\}$

其中: t_k, s_c 是特征函数.

λ_k, μ_c 对应的权重.

* ① $t_k(y_{i-1}, y_i, x, i)$ 依赖于前位置和前-1位置, 称为转移特征.

② $s_c(y_i, x, i)$ 定义为当前结点, 称为状态特征.

③ t_k, s_c 一般取 0 或 1.

④ CRF 完全依赖于特征函数.

(例子参见例 11.1).

6) CRF 的简化形式:

①

$$f_k(y_{i-1}, y_i, x, i) = \begin{cases} t_k(y_{i-1}, y_i, x, i) & k=1, 2, \dots, K \\ s_c(y_i, x, i), & k=K+l, l=1, 2, \dots, K_c \end{cases}$$

② 对转移特征与状态特征在各位置求和, 记作:

$$f_k(y, x) = \sum_{i=1}^n f_k(y_{i-1}, y_i, x, i), \quad k=1, 2, \dots, K.$$

MEMO NO. _____
DATE . . .

③ 权值:

$$w_k = \begin{cases} w_k & k=1, 2, \dots, k_1 \\ w_\ell & k=k_1+\ell; \ell=1, 2, \dots, k_2 \end{cases}$$

$$④ \quad p(y|x) = \frac{1}{Z(x)} \exp \sum_{k=1}^K w_k f_k(y, x)$$

$$Z(x) = \sum_y \exp \sum_{k=1}^K w_k f_k(y, x)$$

⑤ 再简化:

$$w = (w_1, w_2, \dots, w_K)^T$$

$$F(y, x) = (f_1(w, x), f_2(w, x), \dots, f_K(w, x))^T$$

$$\text{or } p_w(y|x) = \frac{\exp(w \cdot F(y, x))}{Z_w(x)} \quad \text{⑥ 表示内积}$$

$$Z_w(x) = \sum_y \exp(w \cdot F(y, x))$$

7) CRF的矩阵形式:

引入特殊起始点和终止点标记 $y_0 = \text{start}$, $y_{n+1} = \text{stop}$.
对观测序列 x 的每一个位置 $i=1, 2, \dots, n+1$, 定义 m 阶矩阵,
(m 标记 y_i 取值的个数, 即状态的种类).

$$M_i(x) = [M_i(y_{i-1}, y_i | x)]$$

$$M_i(y_{i-1}, y_i | x) = \exp(w_i | y_{i-1}, y_i | x)$$

$$w_i(y_{i-1}, y_i | x) = \sum_{k=1}^K w_k f_k(y_{i-1}, y_i, x, i)$$

当前位置相邻状态权重和.

(CRF完全依赖于相邻函数).

于:

$$P(y|x) = \frac{1}{Z_W(x)} \prod_{i=1}^{n+1} M_i(y_{i-1}, y_i | x)$$

$$Z_W(x) = (M_1(x) M_2(x) \dots M_{n+1}(x))_{\text{start}, \text{stop}}$$

$$\text{由 } M_i(y_{i-1}, y_i | x) = \exp\left(\sum_{k=1}^K w_k f_k(y_{i-1}, y_i, x, i)\right).$$

$$P(y|x) = \prod_{i=1}^{n+1} M_i(y_{i-1}, y_i | x)$$

$$= \exp\left(\sum_{i=1}^{n+1} \sum_{k=1}^K w_k f_k(y_{i-1}, y_i, x, i)\right).$$

$$= \exp\left(\sum_{i,k} w_k f_k(y_{i-1}, y_i, x, i)\right).$$

(参见例 11.2).

8) CRF 前向算法.

* CRF, HMM 前向算法的目标不同:

当前 CRF: 求状态序列的概率.**HMM:** 求观测序列的概率.对每个指标 $i=0, 1, \dots, n+1$, 定义前向向量 $\alpha_i(x)$:

$$\alpha_i(y|x) = \begin{cases} 1, & y = \text{start} \quad (\text{start 位置}) \\ 0, & \text{否则} \end{cases}$$

递推公式:

$$\alpha_i^*(y_i|x) = \alpha_{i-1}^*(y_{i-1}|x) M_i(y_{i-1}, y_i | x), i=1, 2, \dots, n+1.$$

$$\text{即: } \begin{matrix} (1 \times m) & (1 \times m) & (m \times m) \\ \vec{\alpha}_i^T(x) = \vec{\alpha}_{i-1}^T(x) \vec{M}_i(x) \end{matrix}$$

 $(1 \times m)$.

$\alpha_i(y_i|x)$ 表示在位置 i 的标记是 y_i 并回到位置 i 的前置标记序列的概率. y_i 为可取值 m .

9) CRF后向算法:

对每个指标 $i=0, 1, \dots, n+1$, 定义后向向量 $\beta_i(x)$:

$$\beta_{n+1}(y_{n+1}|x) = \begin{cases} 1 & y_{n+1} = \text{stop} \\ 0 & \text{否则} \end{cases}$$

$$\Rightarrow \beta_i(y_i|x) = M_i(y_i, y_{i+1}|x) \beta_{i+1}(y_{i+1}|x)$$

$$\Rightarrow \vec{\beta}_i(x) = M_{i+1}(x) \vec{\beta}_{i+1}(x)$$

$\beta_i(y_i|x)$ 表示在位置 i 的标记为 y_i 并且从 $i+1$ 到 n 后部分标记序列的非规范化概率。

or

$$Z(x) = \vec{\alpha}_n(x) \cdot \vec{1} = \vec{1}^T \cdot \vec{\beta}_1(x)$$

$\vec{1}$ 为列向量。

10) 概率计算:

位置 i 的 y_i 的条件概率:

$$P(Y_i = y_i | x) = \frac{\vec{\alpha}_i^T(y_i|x) \vec{\beta}_i(y_i|x)}{Z(x)}$$

在位置 $i-1$ 与 i 的标记 y_{i-1} 和 y_i 的条件概率:

$$P(Y_{i-1} = y_{i-1}, Y_i = y_i | x) = \frac{\vec{\alpha}_{i-1}^T(y_{i-1}|x) M_i(y_{i-1}, y_i|x) \beta_i(y_i|x)}{Z(x)}$$

* CRF前向后向算法的作用:

① 算 $P(y|x)$ ② 算 $P(y_i|x)$ ③ 算 $P(y_{i-1}, y_i|x)$ ④ 还有... ⑤

10) CRF 训练 —— 维归比算法。

① 目的：
给定条件随机场 $P(y|x)$ 和输入序列 (观测) 序列 x ，
求条件概率最大的输出序列 (标记序列) y^* ，即标注。

② 推导：

$$\text{由之前 } Z_W(x) = \sum_y \exp(W \odot F(y, x))$$

$$\text{即: } y^* = \arg \max_y P_W(y|x) \quad (\times) \text{ 即求}$$

$$= \arg \max_y \frac{\exp(W \odot F(y, x))}{Z_W(x)}$$

$$= \arg \max_y \exp(W \odot F(y, x))$$

$$= \arg \max_y [W \odot F(y, x)] \quad (\text{非标准范数})$$

(即对所有位置特征函数求和)

③ 化简：

$$\text{设 } W = (w_1, w_2, \dots, w_k)^T$$

$$F(y, x) = (f_1(y, x), f_2(y, x), \dots, f_k(y, x))^T$$

$$f_k(y, x) = \sum_{i=1}^n f_k(y_{i-1}, y_i, x, i), k=1, 2, \dots, K.$$

$$\text{则: } W \odot F(y, x) = (w_1, w_2, \dots, w_k)^T \odot (f_1(y, x), f_2(y, x), \dots, f_k(y, x))^T$$

$$= (w_1, w_2, \dots, w_k)^T \odot \left(\sum_{i=1}^n f_1(y_{i-1}, y_i, x, i), \dots, \sum_{i=1}^n f_k(y_{i-1}, y_i, x, i) \right)^T$$

$$= \sum_{k=1}^K w_k \sum_{i=1}^n f_k(y_{i-1}, y_i, x, i)$$

$$= \sum_{i=1}^n \sum_{k=1}^K w_k f_k(y_{i-1}, y_i, x, i)$$

$$= \sum_{i=1}^n w \cdot F_i(y_{i-1}, y_i, x, i).$$

$$\text{其中 } F_i(y_{i-1}, y_i, x) = (f_1(y_{i-1}, y_i, x, i), f_2(y_{i-1}, y_i, x, i), \dots, f_k(y_{i-1}, y_i, x, i))^T$$

④ 定义:

符号: $\delta_i(j)$ 表示在位置 i 时, label 为 j 的非累积概率 $\varphi_i(j)$ 表示在位置 i 时, label 为 j 时, 上一位置
能使 $\delta_i(j)$ 最大的 $i-1$ 位置.

求位置 1 时对于每种标记的非累积概率:

$$\delta_1(j) = W F(y_{i-1}, y_i, x, 1), j = 1, 2, \dots, m.$$

标记的种类.

(2) 递推求 $\delta_i(l)$ (i 表示位置, l 表示 label)

$$\delta_i(l) = \max_{1 \leq j \leq m} \{ \delta_{i-1}(j) + W F(y_{i-1}=j, y_i=l, x) \}$$

 $l = 1, 2, \dots, m.$ 同时记录 $\varphi_i(l) = \arg \max_{1 \leq j \leq m} \{ \delta_{i-1}(j) + W F(y_{i-1}=j, y_i=l, x) \} \quad l = 1, 2, \dots, m.$ (3) 直到位置 n .

$$\max_y (W \cdot F(y, x)) = \max_{1 \leq j \leq m} \delta_n(j).$$

最终位置 n 的 label $y_n^* = \arg \max_{1 \leq j \leq m} \delta_n(j)$

(4) 求经过点图:

$$y_i^* = (\varphi_i, y_{i+1}^*) \quad i = n-1, \dots, 1.$$

获得最终的最优路径:

$$y^* = (y_1^*, y_2^*, \dots, y_n^*)^T.$$

(见例 11.3)

MEMO NO. _____
DATE _____

1) CRF 的 学习 算法.

{ MLE

正则化的 MLE.

* 现在实现 算法:

1) 改进的 坐标 尺度 法 IIS

2) 梯度 下降 法

3) 拟 牛顿 法.

2) NN-CRF.

$$P(y|X) = \frac{e^{S(X,y)} e^{S(X,y)}}{\sum_{y \in K} e^{S(X,y)}}$$

$$S(X,y) = \sum_{i=0}^n A_{y_i, y_{i+1}} + \sum_{i=1}^n P_{y_i}$$

LSTM hidden state.

i : 位置

y_i : label

A : 转移矩阵.

A_{ij} : label i 转移到 label j 的 分数.