

6.1 马尔可夫模型

如果一个系统有 n 个状态 s_1, s_2, \dots, s_n , 随着时间的推移, 该系统从某一状态转移到另一状态, 如果用 q_t 表示系统在时间 t 的状态变量, 那么, t 时刻的状态取值为 s_j ($1 \leq j \leq n$) 的概率取决于前 $t-1$ 个时刻 ($1, 2, \dots, t-1$) 的状态.

$$\text{概率为 } P(q_t = s_j | q_{t-1} = s_1, q_{t-2} = s_2, \dots) \quad (6.1)$$

假设1: 在特定情况下, 系统在时间 t 的状态只与前一个时刻的 $t-1$ 状态相关, 离散的- n 马尔可夫链:

$$P(q_t = s_j | q_{t-1} = s_1, q_{t-2} = s_2, \dots) = P(q_t = s_j | q_{t-1} = s_1) \quad (6.1')$$

假设2: 如果只考虑 (6.1) 独立于时间 t 的随机过程, 即不动性假设, 状态与时间

无关, 那么 $P(q_t = s_j | q_{t-1} = s_i) = a_{ij}, 1 \leq i, j \leq n$ (6.2)

(与时间无关)

a_{ij} 为状态转移概率

$$\sum_{j=1}^n a_{ij} = 1$$

状态图表示: (非确定性的非确定性的有限状态自动机)

0 概率省略

每个节点所有发出边的概率之和等于 1

状态序列 s_1, \dots, s_T 的概率

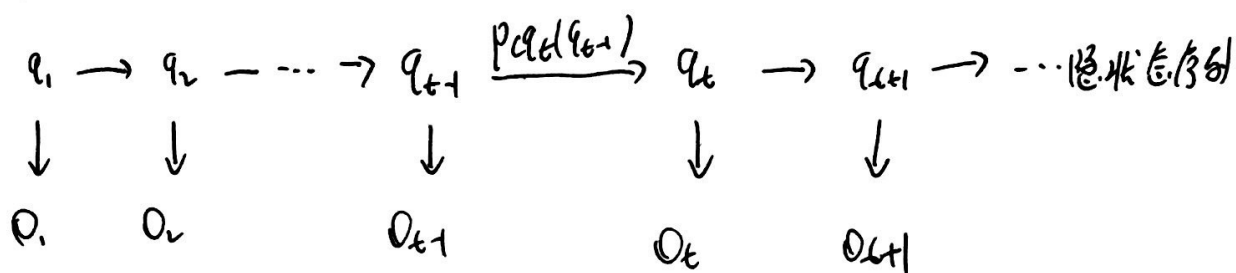
$$\begin{aligned} P(s_1, \dots, s_T) &= P(s_1) \times P(s_2 | s_1) \times P(s_3 | s_1, s_2) \times \dots \times P(s_T | s_1, \dots, s_{T-1}) \\ &= P(s_1) \times P(s_2 | s_1) \times P(s_3 | s_2) \times \dots \times P(s_T | s_{T-1}) \\ &= \pi_{s_1} \prod_{t=1}^{T-1} a_{s_t s_{t+1}} \end{aligned}$$

其中, $\pi_{s_1} = P(q_1 = s_1)$, 为初始状态的概率.



6.2 隐马尔可夫模型

双重随机过程. 不知道具体的状态序列, 只知道状态转移的概率, 即模型的状态转移过程是不可见的(隐藏的), 而可见事件的随机过程是隐藏状态转移过程的随机函数.



HMM 图解

组成 { 模型中的状态数为 N (符号的数量)

从每一个状态可能输出的不同的符号数 M (不同颜色球的数目)

状态转移概率矩阵 $A = a_{ij}$ $\left\{ \begin{array}{l} a_{ij} = P(q_{t+1} = j | q_t = i), 1 \leq i, j \leq N \\ a_{ij} \geq 0 \\ \sum_{j=1}^N a_{ij} = 1 \end{array} \right.$

从状态 s_j 观察到某一特定符号 v_k 的概率分布矩阵 $B = b_j(v_k)$

$\left\{ \begin{array}{l} b_j(v_k) = P(o_t = v_k | q_t = s_j), 1 \leq j \leq N, 1 \leq k \leq M \\ b_j(v_k) \geq 0 \\ \sum_{k=1}^M b_j(v_k) = 1 \end{array} \right.$

初始状态的概率分布 $\pi = \pi_i$ $\left\{ \begin{array}{l} \pi_i = P(q_1 = s_i), 1 \leq i \leq N \\ \pi_i \geq 0 \\ \sum_{i=1}^N \pi_i = 1 \end{array} \right.$

将 HMM 记为: $\lambda = (A, B, \pi)$ 或 $\lambda = (s, o, A, B, \pi)$

给定 HMM 求观察序列:

给定 $\lambda = (A, B, \pi)$, 求 $O = o_1 o_2 \dots o_T$



1. $s_t = i$

2. 根据初始状态分布 $\pi = \pi_i$ 选择初始状态 $q_1 = i$

3. 根据状态 s_t 的输出概率分布 $b_i(o_t)$, 输出 $o_t = v_t$

4. 根据状态转移概率 a_{ij} , 转移到新状态 $q_{t+1} = j$

5. $t = t+1$, 如果 $t < T$, 重复 3, 4, 5 直到结束

问题: 1. 快速计算概率 $P(o|\mu)$ (给定 μ 和 o)

2. 如何在有限-定意义下“最佳”的状态序列, 使用此状态序列 = “最佳地解释”

3. 给定观察序列 o , 如何根据最大似然估计来求模型的参数值? 即如何求得使 $P(o|\mu)$ 最大的参数?

6.3 前向算法

解决问题 1. 对于给定的状态序列 $Q = q_1 q_2 \dots q_T$, $P(o|\mu) = ?$

$$P(o|\mu) = \sum_Q P(o, Q|\mu) = \sum_Q P(Q|\mu) \times P(o|Q, \mu)$$

$$\frac{P(o, Q, \mu)}{P(\mu)} = \frac{P(Q, \mu) P(o|Q, \mu)}{P(\mu)} \quad \text{4.1.1.1.1.1}$$

$$P(Q|\mu) = \pi_{q_1} \times a_{q_1 q_2} \dots a_{q_{T-1} q_T}$$

$$P(o|Q, \mu) = b_{q_1(o_1)} \times b_{q_2(o_2)} \dots b_{q_T(o_T)}$$

困难: 搜索路径或指数级爆炸

解决方法: 动态规划

定义前向变量 $\alpha_t(i)$:

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = i | \mu)$$

(当前时刻) 序列概率, 每走一步从 0, 算到

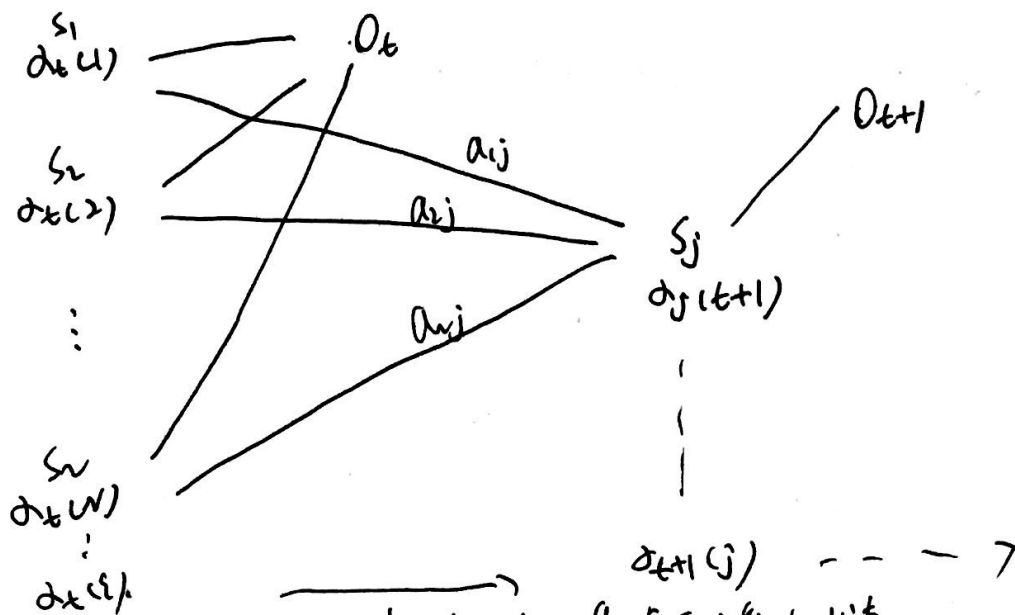
$$P(o|\mu) = \sum_i P(o_1, o_2, \dots, o_T, q_T = i | \mu) = \sum_i \alpha_T(i)$$

动态规划/计算 $\alpha_t(i)$: 在时间 $t+1$ 的前向变量可以根据时间 t 的

前向变量 $\alpha_t(1) \dots \alpha_t(n)$ 的值递推计算

$$\alpha_{t+1}(j) = \left[\sum_i \alpha_t(i) a_{ij} \right] \times b_j(o_{t+1})$$





s_i 表示状态, q_t 表示时间 t 状态.

算法描述: (1). 初始化: $\alpha_i(1) = \pi_i b_i(O_1)$, $1 \leq i \leq N$

(2). 循环计算: $\alpha_{t+1}(j) = [\sum_i \alpha_t(i) a_{ij}] \times b_j(O_{t+1})$, $1 \leq t \leq T-1$

(3). 结束, 输出: $p(O|M) = \sum_{i=1}^N \alpha_T(i)$

时间复杂度分析: 从 $t+1$ 时所有 N 个状态转移到 s_i 的可能性, 每个时间 t 要计算 N^2 个向量的乘积, 因此时间复杂度为 $O(N^2T)$

6.4 后向算法

同样解决向题 (1).

后向变量 $\beta_t(i) = p(O_{t+1} O_{t+2} \dots O_T | q_t = s_i, M)$

运用动态规划计算后向量:

(1). 从时间 t 到 $t+1$, 模型由状态 s_i 转移到状态 s_j , 并从 s_j 输出 O_{t+1}

(2). 在时间 $t+1$, 状态为 s_j 的条件下, 序列 $O_{t+2} O_{t+3} \dots O_T$

第一步的转移: $a_{ij} \times b_j(O_{t+1})$

第二步转移: $\beta_{t+1}(j)$

归纳条件: $\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \times \beta_{t+1}(j)$

初始条件: $\beta_T(1), \beta_T(2), \dots, \beta_T(N)$



6.5 Viterbi 搜索算法

解决问题2: 如何发现最优状态序列, 能够最好的解释观察序列

一种解释: 状态序列中的每个状态, 都单独地具有概率, 对于每个时间 $t (1 \leq t \leq T)$, 寻找 q_t 使得 $\gamma_t(z) = P(q_t = z | O, \mu)$ 最大.

$$\gamma_t(z) = P(q_t = z | O, \mu) = \frac{P(q_t = z, O | \mu)}{P(O | \mu)}$$

分解过程: (1). 模型在时间 t 到达状态 z , 并且输出 $O = O_1 O_2 \dots O_T$
根据前向变量的定义, 实现这一步的概率为 $\alpha_t(z)$

(2). 从时间 t , 状态 z 出发, 模型输出 $O = O_1 O_2 \dots O_T$
根据后向变量定义, 实现这一步的概率为 $\beta_t(z)$

$$\text{于是 } P(q_t = z, O | \mu) = \alpha_t(z) \times \beta_t(z)$$

$$P(O | \mu) \text{ 与时间无关, 因此 } P(O | \mu) = \sum_{z=1}^N \alpha_t(z) \times \beta_t(z)$$

$$\text{即 } \gamma_t(z) = \frac{\alpha_t(z) \times \beta_t(z)}{\sum_{i=1}^N \alpha_t(i) \times \beta_t(i)}$$

$$\text{因此 } t \text{ 时刻的最优状态与 } q_t^* = \arg \max_{1 \leq i \leq N} (\gamma_t(i))$$

问题: 两个最优状态之间的转移概率可能为0, 因此每一个状态单独最优不能保证整个状态序列最优

另一种解释: 在给定 μ 和 O 的条件下求概率最大的状态序列

$$Q = \arg \max P(Q | O, \mu)$$

(防止不收敛也, 取对数)

Viterbi 算法: 动态搜索最优状态序列

定义: Viterbi 变量 $\delta_t(z)$ 是在时间 t 时, 模型沿某一条路径到达 z , 输出观察序列 $O = O_1 O_2 \dots O_t$ 的最大概率为:

$$\delta_t(z) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2, \dots, q_{t-1}, q_t = z, O_1 O_2 \dots O_t | \mu)$$

$$\text{递归计算 } \delta_{t+1}(i) = \max_j [\delta_t(j) \cdot a_{ji}] \cdot b_i(O_{t+1})$$



算法描述

1). 初始化 $\delta_i(i) = \pi_i b_i(O_i), 1 \leq i \leq n$

$$\psi_i(i) = 0$$

2). 递推计算:

$$\delta_t(j) = \max_{1 \leq l \leq n} [\delta_{t-1}(l) \cdot a_{lj}] \cdot b_j(O_t), 2 \leq t \leq T, 1 \leq j \leq n$$

$$\psi_t(j) = \arg \max_{1 \leq l \leq n} [\delta_{t-1}(l) \cdot a_{lj}] \cdot b_j(O_t), 2 \leq t \leq T, 1 \leq j \leq n$$

3). 结束, 回溯得到状态序列

$$q_t = \psi_{t+1}(q_{t+1})$$

时间复杂度: $O(n^2 T)$

6.6 参数学习

解决的问题: 如何调节模型 λ 的参数, 使得 $p(O|\lambda)$ 最大.

前向后向算法 (Baum-Welch or forward-backward procedure)
(如果状态序列已知)

$$\pi_i = \delta(q_i, s_i)$$

$$\bar{a}_{ij} = \frac{\text{从状态 } q_i \text{ 转移到 } q_j \text{ 的次数}}{\text{从状态 } q_i \text{ 转移到任一状态 (包括自身) 的次数}}$$

$$= \frac{\sum_{t=1}^{T-1} \delta(q_t, s_t) \times \delta(q_{t+1}, s_{t+1})}{\sum_{t=1}^{T-1} \delta(q_t, s_t)}$$

$$\delta(x, y) = \begin{cases} 1 & x=y \\ 0 & \text{else} \end{cases}$$

是 Kronecker 克罗内克 δ 函数

$$\text{类似地, } \bar{b}_j(i) = \frac{\sum_{t=1}^T \delta(q_t, s_t) \times \delta(O_t, V_i)}{\sum_{t=1}^T \delta(q_t, s_t)}$$

否则如果不存在大量标注的样本 (实际很少用, 效果较差):

① EM 期望值最大 (EM): 随机猜测值迭代计算.



6.7 HMM 应用举例

1. 汉语的自动分词与句性标注:

汉语分词结果: $O = O_1 O_2 \dots O_T$

即求解 $\theta = \arg \max_{\theta} p(O|\mu)$

句性标注: 求解 $Q = \arg \max_Q p(Q|O, \mu)$

1). 估计 HMM $\mu = (A, B, C)$ 的参数

2). 针对任意给定的一个输入句子及其可能的输出序列 O , 求解所有

可能的 O 中使 $p(O|\mu)$ 最大的群

3). 快速地选择最优状态序列

问题 1). 模型参数

① 观察序列: 单词序列

② 状态序列: 词类标注序列

③ 状态数目 N

④. 输出符号数 M

如果有标注语料, 需要有句性标注词类并用 EM 迭代 (无监督)

否则有标注子句: 用最大似然估计参数.

$$\bar{\pi}_{pos_i} = \frac{\text{pos}_i \text{ 出现在句首}}{\text{所有句首}}$$

$$\bar{a}_{ij} = \frac{\text{pos}_i \rightarrow \text{pos}_j}{\text{pos}_i \rightarrow \text{pos}}$$

$$\bar{b}_j(u) = \frac{\text{pos}_j \rightarrow u}{\text{pos}_j}$$

问题 2). 观察序列

能够切出所有可能切分.

~ (老师说非常简单?? 可我不会啊QAQ)



6.8 CRFs 及其应用 (条件随机场)

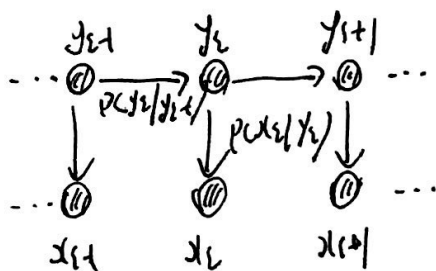
基本思路: 给定观察序列 X , 输出标识序列 Y , 通过计算 $P(Y|X)$ 求解最优标识序列。

设 $G=(V, E)$ 是一个无向图, $Y=\{Y_v | v \in V\}$, 则 V 中每个节点对应于一个随机变量 Y_v , 取值范围为可能的标识集合 \mathcal{Y}

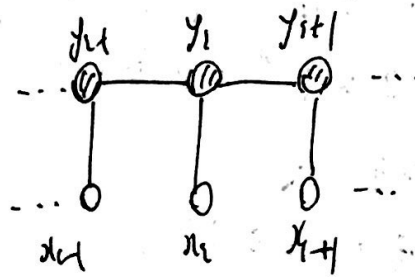
如果从观察序列 X 为条件, 每个随机变量 Y_v 都满足以下马尔可夫特性:

$$P(Y_v | X, Y_w, w \neq v) = P(Y_v | X, Y_w, w \sim v)$$

其中 $w \sim v$ 表示两个节点在图中是邻近节点, 那么 (X, Y) 为一个条件随机场 G 的图结构是任意的



MRFs



CRFs

CRFs 的节点并不是由模型生成的, 本身存在, 只是给打标签而已。

给定观察序列 X , 某个特定标识序列 Y 的概率可以定义为:

$$\exp(\sum_j \lambda_j \underbrace{f_j(y_{t-1}, y_t, x, \epsilon)}_{\text{转移概率}} + \sum_k \mu_k \underbrace{g_k(y_t, x, \epsilon)}_{\text{标识概率}})$$

转移概率

标识概率

λ_j, μ_k 由训练样本估计得出。

$$f_j(y_{t-1}, y_t, x, \epsilon)$$

无参, 与保距形式上一样。

$$b(x, \epsilon) = \begin{cases} 1 & x \text{ 的 } \epsilon \text{ 位置为某个特定标识} \\ 0 & \text{else} \end{cases}$$

$$f_j(y_{t-1}, y_t, x, \epsilon) = \begin{cases} b(x, \epsilon) & y_{t-1} \text{ 和 } y_t \text{ 满足某种条件} \\ 0 & \text{else} \end{cases}$$

$$F_j(Y, X) = \sum_{t=1}^n f_j(y_{t-1}, y_t, x, \epsilon)$$

↓ 表示 s 或 t 。

$$P(Y|X) = \frac{1}{Z(X)} \exp(\sum_j \lambda_j \cdot F_j(Y, X))$$

$Z(X)$ 为归一化因子

