

HMM hidden markov model

Definition

Markov Chain

$Q = \{q_1, q_2, \dots, q_N\}$ a set of N states

$A = [a_{ij}]$ $a_{ij} \dots a_{ii}, a_{jj}$ transition probability matrix

$\pi = [\pi_1, \pi_2, \dots, \pi_N]$ initial probability distribution

Markov Assumption: $P(q_i=a | q_1, q_2, \dots, q_{i-1}) = P(q_i=a | q_{i-1})$

hidden Markov model

$Q = \{q_1, q_2, \dots, q_N\}$ a set of N states

$A = [a_{ij}]$ $a_{ij} \dots a_{ii}, a_{jj}$ transition probability matrix

$\pi = [\pi_1, \pi_2, \dots, \pi_N]$ initial probability distribution

$O = o_1 o_2 \dots o_T$ a sequence of T observations

$B = b_i(o_t)$ emission probability, each expressing the prob of
an observation o_t being generated from a state i

Markov Assumption: $P(q_i=a | q_1, q_2, \dots, q_{i-1}) = P(q_i=a | q_{i-1})$

HMM 在任何时刻 t 的状态只依赖于其前一时刻状态, 与其它时刻的状态和观测无关!

Output independence: $P(O_i | q_1, \dots, q_i, \dots, q_T, o_1, \dots, o_i, \dots, o_T) = P(o_i | q_i)$

HMM 在任何时刻 t 的输出只依赖于其前一时刻状态, 与其它时刻的状态和观测无关

HMM 的三个主要问题

1. likelihood: Given an HMM $\lambda = (A, B, \pi)$ and an observation sequence O ,
determine the like likelihood $P(O | \lambda)$

2. Decoding: Given an observation sequence O and an HMM $\lambda = (A, B, \pi)$,
discover the best hidden state sequence Q

3. Learning: Given an observation sequence O and the set of states
in the HMM - learn the HMM parameters A and B (MLE)

Question 1. likelihood \rightarrow forward & backward algorithm

问题: 给定 HMM: $\lambda = (A, B, \pi)$, 求 $P(O|\lambda)$

定义

答:

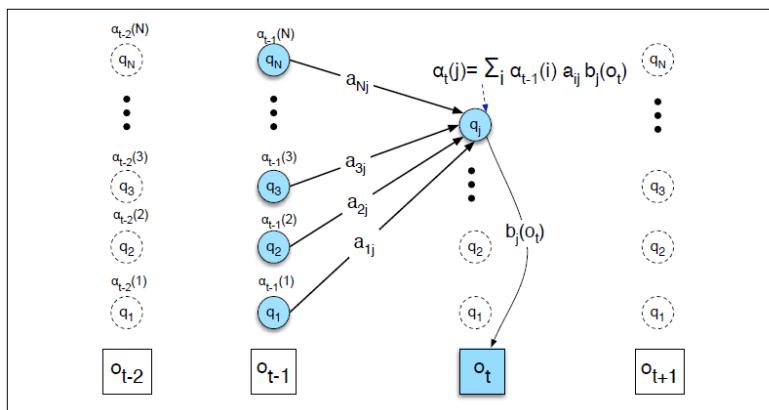
$\alpha_t(j)$: The prob of being in the State j after seeing the first t observations. The value of each cell $\alpha_t(j)$ is computed by summing over the prob. of each path that could lead to this cell.

$$\alpha_t(j) = P(O_1, O_2, \dots, O_t, q_t=j | \lambda)$$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(O_t)$$

a_{ij} : the transition prob from previous state q_i to current state q_j

$b_j(O_t)$: the state observation likelihood of the observation symbol O_t given the current state j



前向算法:

算法 10.2 (观测序列概率的前向算法)

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

$$\alpha_1(i) = \pi_i b_i(o_1), \quad i = 1, 2, \dots, N \quad (10.15)$$

(2) 递推 对 $t = 1, 2, \dots, T - 1$,

$$\alpha_{t+1}(i) = \left[\sum_{j=1}^N \alpha_t(j) a_{ji} \right] b_i(o_{t+1}), \quad i = 1, 2, \dots, N \quad (10.16)$$

(3) 终止

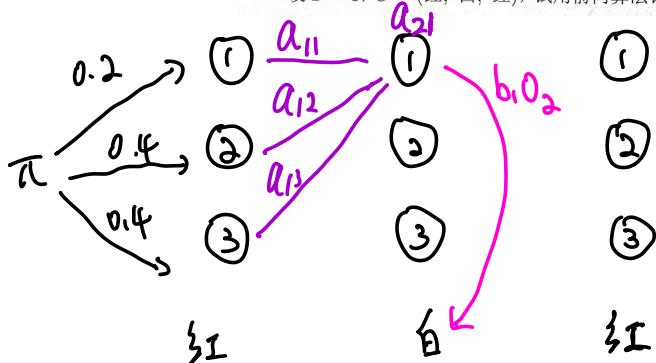
$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i) \quad (10.17)$$

例题：

例 10.2 考虑盒子和球模型 $\lambda = (A, B, \pi)$, 状态集合 $Q = \{1, 2, 3\}$, 观测集合 $V = \{\text{红, 白}\}$,

$$A = \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.5 & 0.2 \\ 0.2 & 0.3 & 0.5 \end{bmatrix}, \quad B = \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \\ 0.7 & 0.3 \end{bmatrix}, \quad \pi = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.4 \end{bmatrix}$$

设 $T = 3$, $O = (\text{红, 白, 红})$, 试用前向算法计算 $P(O|\lambda)$.



起始

$$a_1(1) = \pi_1 b_1(O_1) = 0.2 \times 0.5 = 0.1$$

$$a_1(2) = \pi_2 b_2(O_1) = 0.4 \times 0.4 = 0.16$$

$$a_1(3) = \pi_3 b_3(O_1) = 0.4 \times 0.7 = 0.28$$

$$\underline{a_2(1)} = \left(\sum_{j=1}^3 a_1(j) a_{j1} \right) b_1(O_2)$$

$$= (0.1 \times 0.5 + 0.16 \times 0.3 + 0.28 \times 0.2) \times 0.5$$

$$= (0.05 + 0.048 + 0.056) \times 0.5$$

$$= 0.077$$

$$a_2(2) = 0.1104$$

$$a_2(3) = 0.0606$$

$$a_3(1) = 0.04187$$

$$a_3(2) = 0.03551$$

$$a_3(3) = 0.05284$$

$$\underline{\underline{P(O|\lambda)}}: P(O|\lambda) = \sum_{i=1}^3 a_3(i) = 0.04187 + 0.03551 + 0.05284 = 0.13022$$

后向算法：

后向概率 (backward prob):

定义 10.3 (后向概率) 给定隐马尔可夫模型 λ , 定义在时刻 t 状态为 q_i 的条件下, 从 $t+1$ 到 T 的部分观测序列为 $o_{t+1}, o_{t+2}, \dots, o_T$ 的概率为后向概率, 记作

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | i_t = q_i, \lambda) \quad (10.18)$$

可以用递推的方法求得后向概率 $\beta_t(i)$ 及观测序列概率 $P(O|\lambda)$ 。

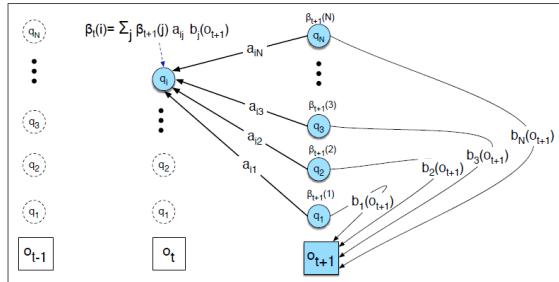


Figure A.11 The computation of $\beta_t(i)$ by summing all the successive values $\beta_{t+1}(j)$ weighted by their transition probabilities a_{ij} and their observation probabilities $b_j(o_{t+1})$.

算法 10.3 (观测序列概率的后向算法)

输入: 隐马尔可夫模型 λ , 观测序列 O ;

输出: 观测序列概率 $P(O|\lambda)$ 。

(1)

$$\beta_T(i) = 1, \quad i = 1, 2, \dots, N \quad (10.19)$$

(2) 对 $t = T-1, T-2, \dots, 1$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \underline{\beta_{t+1}(j)}, \quad i = 1, 2, \dots, N \quad (10.20)$$

先看π:

(3)

$$P(O|\lambda) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i) \quad (10.21)$$

Question 2: Decoding - Viterbi alg (最优化问题)

Decoding: Given an observation sequence O and an HMM $\lambda = (A, B, \pi)$.

discover the best hidden state sequence $Q = q_1, q_2, \dots, q_T$

(给定观测序列, 找最有可能的存在的状态序列) $P(Q|O)$

定义:

1. 定义在时刻 t 状态为 i 的所有单个路径中概率最大值为

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_t = i, q_{t-1}, \dots, q_1, O_t, \dots, O_1 | \lambda) \quad i=1, 2, \dots, N$$

2. 递推公式

$$\begin{aligned} \delta_{t+1}(i) &= \max_{q_1, \dots, q_t} P(q_{t+1} = i, \dots, q_1, O_{t+1}, \dots, O_1 | \lambda) \\ &= \max_{1 \leq j \leq N} [\delta_t(j) a_{ji}] b_i(O_{t+1}) \quad i=1, 2, \dots, N \\ &\quad t=1, 2, \dots, T-1 \end{aligned}$$

2. 定义在时刻 t 状态为 i 的所有单个路径中概率最大路径的第 $t+1$ 个结点为:

$$\psi_t(i) = \arg \max_{1 \leq j \leq N} [\delta_{t-1}(j) a_{ji}] \quad i=1, 2, \dots, N$$

Viterbi Alg:

输入: 模型 $\lambda = (A, B, \pi)$ 和观测 $O = (o_1, o_2, \dots, o_T)$;

输出: 最优路径 $I^* = (i_1^*, i_2^*, \dots, i_T^*)$.

(1) 初始化

$$\delta_1(i) = \pi_i b_i(o_1), \quad i = 1, 2, \dots, N$$

$$\psi_1(i) = 0, \quad i = 1, 2, \dots, N$$

forward alg:

(2) 递推。对 $t = 2, 3, \dots, T$

$$\delta_t(i) = \max_{1 \leq j \leq N} [\delta_{t-1}(j) a_{ji}] b_i(o_t), \quad i = 1, 2, \dots, N$$

$$\psi_t(i) = \arg \max_{1 \leq j \leq N} [\delta_{t-1}(j) a_{ji}], \quad i = 1, 2, \dots, N$$

2) 递推 对 $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(i) = \left[\sum_{j=1}^N \alpha_t(j) a_{ji} \right] b_i(o_{t+1}), \quad i = 1, 2, \dots, N \quad (10.16)$$

(3) 终止

$$P^* = \max_{1 \leq i \leq N} \delta_T(i)$$

$i_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)]$ 由于 $\delta_T(i)$ 是 $\delta_t(i)$ 的子集

所以 i_T^* 也是 $\delta_t(i)$ 的子集, 即 i_T^* 是 $\delta_t(i)$ 的最大值对应的结点, 也就是说 i_T^* 就是结点 i 。

(4) 最优路径回溯。对 $t = T-1, T-2, \dots, 1$

$$i_t^* = \psi_{t+1}(i_{t+1}^*)$$

求得最优路径 $I^* = (i_1^*, i_2^*, \dots, i_T^*)$.

Figure A.9 shows pseudocode for the Viterbi algorithm. Note that the Viterbi algorithm is identical to the forward algorithm except that it takes the **max** over the previous path probabilities whereas the forward algorithm takes the **sum**. Note also that the Viterbi algorithm has one component that the forward algorithm doesn't

have: **backpointers**. The reason is that while the forward algorithm needs to produce an observation likelihood, the Viterbi algorithm must produce a probability and also the most likely state sequence. We compute this best state sequence by keeping track of the path of hidden states that led to each state, as suggested in Fig. A.10, and then at the end backtracing the best path to the beginning (the Viterbi **backtrace**).

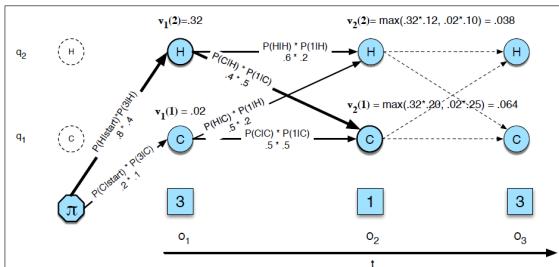


Figure A.8 The Viterbi trellis for computing the best path through the hidden state space for the ice-cream eating events 3 / 3. Hidden states are in circles, observations in squares. White (unfilled) circles indicate illegal transitions. The figure shows the computation of $v_i(j)$ for two states at two time steps. The computation in each cell follows Eq. A.14: $v_i(j) = \max_{1 \leq i \leq N-1} v_i(i) \cdot a_{ij} b_j(o_i)$. The resulting probability expressed in each cell is Eq. A.13: $v_i(j) = P(q_0, q_1, \dots, q_{i-1}, o_1, o_2, \dots, o_i, q_i = j | \lambda)$.

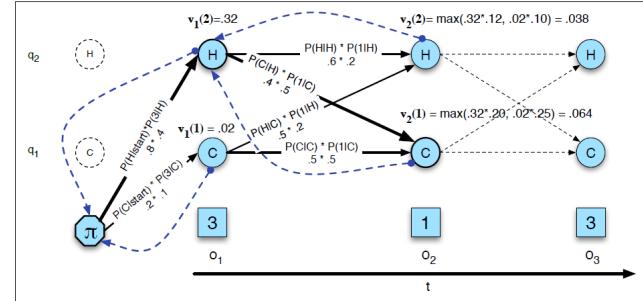


Figure A.10 The Viterbi backtrace. As we extend each path to a new state account for the next observation, we keep a backpointer (shown with broken lines) to the best path that led us to this state.

Question 3: learning (Forward-backward alg and Baum-Welch alg)

问题：给定观测序列 O 和隐状态序列 Q，推断参数 A, B, π

$$\lambda_{MLE} = \arg \max_{\lambda} p(O|\lambda)$$

$$\lambda = (A, B, \pi)$$

Baum-Welch algorithm

review.

1. EM algorithm

算法 9.1 (EM 算法)

输入: 观测变量数据 Y, 隐变量数据 Z, 联合分布 $P(Y, Z|\theta)$, 条件分布 $P(Z|Y, \theta)$;
输出: 模型参数 θ .

(1) 选择参数的初值 $\theta^{(0)}$, 开始迭代;

(2) E 步: 记 $\theta^{(i)}$ 为第 i 次迭代参数 θ 的估计值, 在第 $i+1$ 次迭代的 E 步, 计算

$$Q(\theta, \theta^{(i)}) = E_Z [\log P(Y, Z|\theta) | Y, \theta^{(i)}] \\ = \sum_Z \log P(Y, Z|\theta) P(Z|Y, \theta^{(i)}) \quad (9.9)$$

这里, $P(Z|Y, \theta^{(i)})$ 是在给定观测数据 Y 和当前的参数估计 $\theta^{(i)}$ 下隐变量数据 Z 的条件概率分布;

(3) M 步: 求使 $Q(\theta, \theta^{(i)})$ 极大化的 θ , 确定第 $i+1$ 次迭代的参数的估计值 $\theta^{(i+1)}$

$$\theta^{(i+1)} = \arg \max_{\theta} Q(\theta, \theta^{(i)}) \quad (9.10)$$

(4) 重复第 (2) 步和第 (3) 步, 直到收敛。 ■

式 (9.9) 的函数 $Q(\theta, \theta^{(i)})$ 是 EM 算法的核心, 称为 Q 函数 (Q function)。

定义 9.1 (Q 函数) 完全数据的对数似然函数 $\log P(Y, Z|\theta)$ 关于在给定观测数据 Y 和当前参数 $\theta^{(i)}$ 下对未观测数据 Z 的条件概率分布 $P(Z|Y, \theta^{(i)})$ 的期望称为 Q 函数, 即

$$Q(\theta, \theta^{(i)}) = E_Z [\log P(Y, Z|\theta) | Y, \theta^{(i)}] \quad (9.11)$$

下面关于 EM 算法作几点说明:

步骤 (1) 参数的初值可以任意选择, 但需注意 EM 算法对初值是敏感的。

步骤 (2) E 步求 $Q(\theta, \theta^{(i)})$. Q 函数式中 Z 是未观测数据, Y 是观测数据。注意, $Q(\theta, \theta^{(i)})$ 的第 1 个变元表示要极大化的参数, 第 2 个变元表示参数的当前估计值。每次迭代实际在求 Q 函数及其极小值。

步骤 (3) M 步求 $Q(\theta, \theta^{(i)})$ 的极大化, 得到 $\theta^{(i+1)}$, 完成一次迭代 $\theta^{(i)} \rightarrow \theta^{(i+1)}$ 。后面将证明每次迭代使似然函数增大或达到局部极值。

步骤 (4) 给出停止迭代的条件, 一般是对较小的正数 ϵ_1, ϵ_2 , 若满足

$$\|\theta^{(i+1)} - \theta^{(i)}\| < \epsilon_1 \quad \text{或} \quad \|Q(\theta^{(i+1)}, \theta^{(i)}) - Q(\theta^{(i)}, \theta^{(i)})\| < \epsilon_2$$

则停止迭代。如果迭代过程中似然函数不再增加, 则认为已经达到局部极值。

$$\text{EM: } \theta^{(t+1)} = \arg \max_{\theta} \int_Z \log P(X, Z | \theta) P(Z | X, \theta^{(t)}) dZ$$

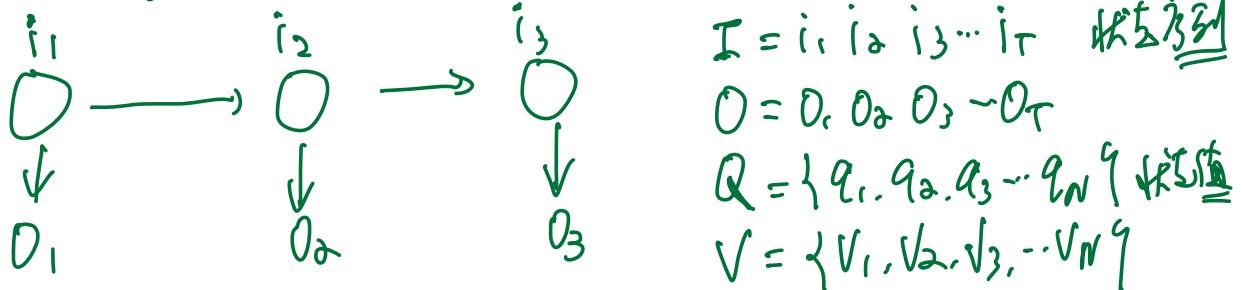
↑
 current Step
 ↓
 X : 观测值 Z : 隐变量 θ : 参数
 ↓ ↓ ↓

HMM Φ :

$$\begin{aligned} \lambda^{(t+1)} &= \arg \max_{\lambda} \sum_I \log P(D, I | \lambda) \frac{P(I | O, \lambda^{(t)})}{P(I, O | \lambda^{(t)})} \\ &= \arg \max_{\lambda} \sum_I \log P(D, I | \lambda) \frac{P(I, O | \lambda^{(t)})}{P(O, \lambda^{(t)})} \leftarrow C \\ &= \arg \max_{\lambda} \sum_I \log P(O, I | \lambda) P(I, O | \lambda^{(t)}) \end{aligned}$$

E_t

Q_t \propto $\sum_I \log P(O, I | \lambda) P(I, O | \lambda^{(t)})$



$$P(I | \lambda) = \pi_1 a_{12} a_{23} \dots a_T a_{T+1}$$

$$P(O | I, \lambda) = b_1 O_1 b_2 O_2 \dots b_T O_T$$

$$P(O, I | \lambda) = \pi_1 a_{12} b_1 O_1 a_{23} b_2 O_2 \dots a_{T+1} b_T O_T$$

∴ Q_t \propto

$$\sum_I \left[\log \pi_{i_t} + \sum_{t=1}^T \log a_{i_t i_{t+1}} + \sum_{t=1}^T \log (b_{i_t} O_t) \right] P(I, O | \lambda^{(t)})$$

Target : $\arg \max_{\lambda} P(D|\lambda)$, $\lambda = (A, B, \pi)$

M-Step:

$$\begin{aligned}
 \textcircled{1} \quad \pi^{(t+1)} &= \arg \max_{\pi} \sum_I \log \pi_i, p(I, D | x^{(t)}) \\
 &= \arg \max_{\pi} \sum_I \log \pi_i, p(i_1, i_2, i_3, \dots, i_n, D | \lambda^{(t)}) \\
 &= \arg \max_{\pi} \sum_{i_1} \sum_{i_2} \sum_{i_3} \log \pi_{i_1} p(i_1, i_2, i_3, \dots, i_n, D | x^{(t)}) \\
 &= \arg \max_{\pi} \sum_{i_1} \log \pi_{i_1} p(i_1, D | \lambda^{(t)}) \\
 &= \arg \max_{\pi} \sum_{i=1}^N \left[\log \pi_i p(D, i_1 = q_i | \lambda^{(t)}) \right]
 \end{aligned}$$

Besides,

$$\text{lagrange: } \sum_{i=1}^N \log \pi_i + \eta \left(\sum_{i=1}^N \pi_i - 1 \right)$$

$$\begin{aligned}
 \frac{\partial L}{\partial \pi_i} &= \sum_{i=1}^N \left[\frac{-p(D, i_1 = q_i | \lambda^{(t)})}{\pi_i} + \eta \right] = 0 \\
 &= \sum_{i=1}^N \left[p(D, i_1 = q_i | \lambda^{(t)}) + \underbrace{\pi_i \eta}_{\sum_{i=1}^N \pi_i = 1} \right] = 0 \\
 &= p(D | \lambda^{(t)}) + \eta = 0 \quad \sum_{x_1, x_2, \dots, x_N} p(x, y) = p(y)
 \end{aligned}$$

$$\Rightarrow \eta = -p(D | \lambda^{(t)})$$

$$\begin{aligned}
 & \left(p(O, i_1 = q_i | \lambda^t) + \pi_{i_1} \right) = 0 \\
 \Rightarrow \quad & \pi_{i_1} = -\frac{1}{q_i} p(O, i_1 = q_i | \lambda^t) \\
 \pi_i^{(t+1)} &= \frac{p(O, i_1 = q_i | \lambda^t)}{p(O | \lambda^t)} \\
 \Pi^{(t+1)} &= (\pi_1^{(t+1)}, \pi_2^{(t+1)} \dots \pi_n^{(t+1)})
 \end{aligned}$$

②