

# 马尔可夫模型 (MM)

设系统有  $N$  个状态  $S = \{s_1, s_2, \dots, s_N\}$ , 随着时间的推移, 系统将  
从某一状态转移到另一状态, 设  $Q = \{q_1, q_2, \dots, q_T\}$  为一随机变量  
序列,  $q_t \in S, t=1, 2, \dots, T$ , 则系统在时间  $t$  处于状态  $s_j$  的概率  
取决于其在时间  $1, 2, \dots, t-1$  的状态

$$P(q_t = s_j | q_{t-1} = s_i, q_{t-2} = s_k, \dots)$$

即现在状态只与历史状态有关

一阶马尔可夫过程: 系统在时间  $t$  的状态只与其在  $t-1$  时间状态有关

$$P(q_t = s_j | q_{t-1} = s_i)$$

二元语法模型便是如此

马尔可夫模型: 初始状态  $\pi = (\pi_1, \pi_2, \dots, \pi_N)$   
 $\pi_i$  表示系统在初始时处于  $s_i$  状态的概率

以当下状态为横轴, 下一时刻状态为纵轴, 有:

	$s_i$	$s_j$	$s_k$	...
$s_i$	$a_{ii}$	$a_{ji}$	$a_{ki}$	
$s_j$	$a_{ij}$	$a_{jj}$	$a_{kj}$	...
$s_k$	$a_{ik}$	$a_{jk}$	$a_{kk}$	
...				

... 状态转移概率矩阵

$$P(q_{t+1} = s_j | q_t = s_i) = a_{ij}$$

$a_{ij}$  表示从  $s_i$  变成  $s_j$  的状态转移概率

$$a_{ij} \geq 0 \quad \sum_{j=1}^N a_{ij} = 1 \quad (\text{列和必为 } 1, \text{ 行和无限制})$$

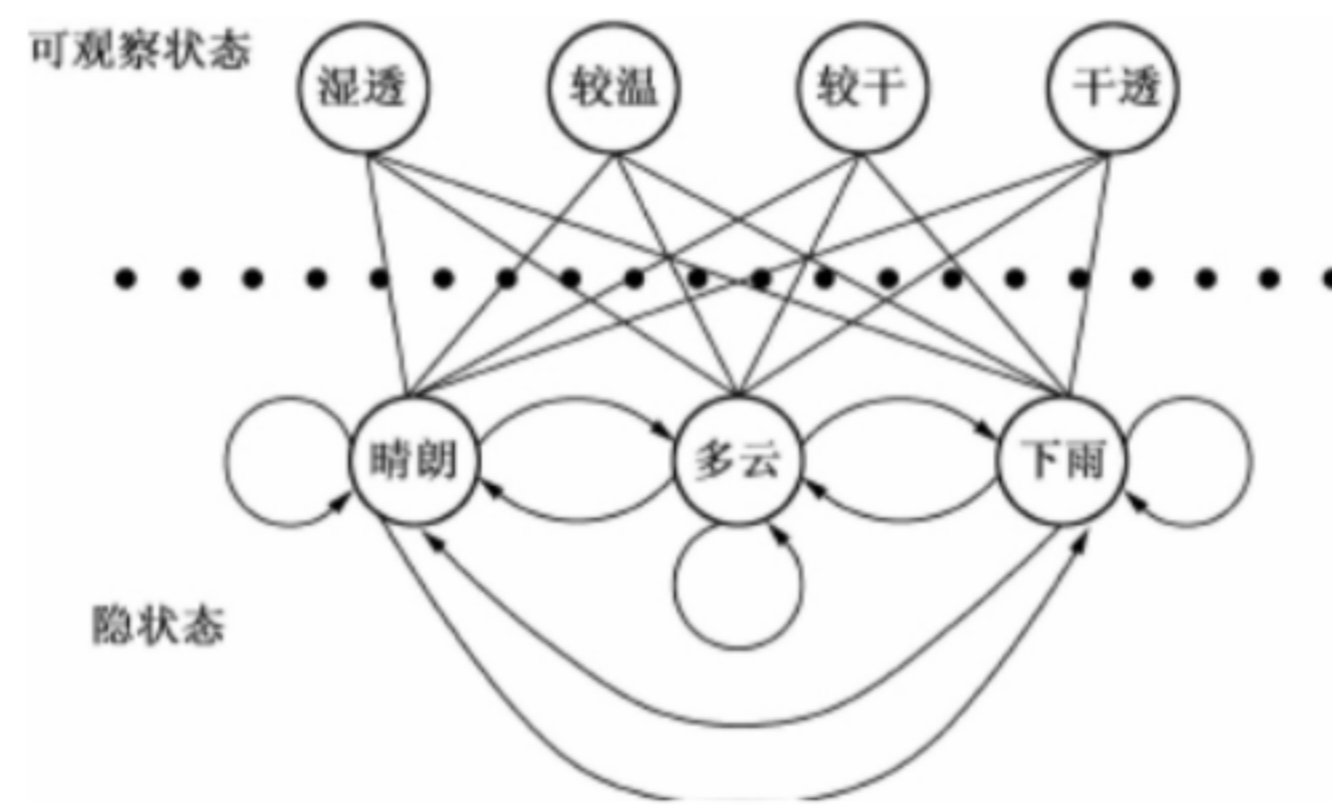
马尔可夫模型形式化定义:  $M = (S, A, \pi)$

$S = \{s_1, s_2, \dots, s_N\}$  状态集合

$A = \{a_{ij}\}$  状态转移概率矩阵  
 $\pi = \{\pi_1, \pi_2, \dots\}$  初始概率矩阵

## 隐马尔可夫模型 (HMM)

这是一种具有隐藏状态的马尔可夫链，其核心特征是：状态不可直接观测，只能通过观测序列推断。



形式定义： $\mu = (S, V, A, B, \pi)$  简记为  $\mu = (A, B, \pi)$

$S$ : 状态集合  $\{s_1, \dots, s_N\}$

$V$ : 观察值集合  $\{o_1, \dots, o_M\}$

$A$ : 转移概率矩阵  $[a_{ij}]$

$B$ : 发射概率矩阵  $[b_j(k)]$

$\pi$ : 初始状态矩阵  $\pi = [\pi_i]$

其中  $B$  表示每个隐藏状态生成每个观测值的概率， $b_j(k) = P(o_t = o_k | q_t = s_j)$  即状态为  $s_j$  时，生成观测  $o_k$  的概率

HMM 使用场景：语音识别 词性标注 拼音输入法  
 文字识别 ...

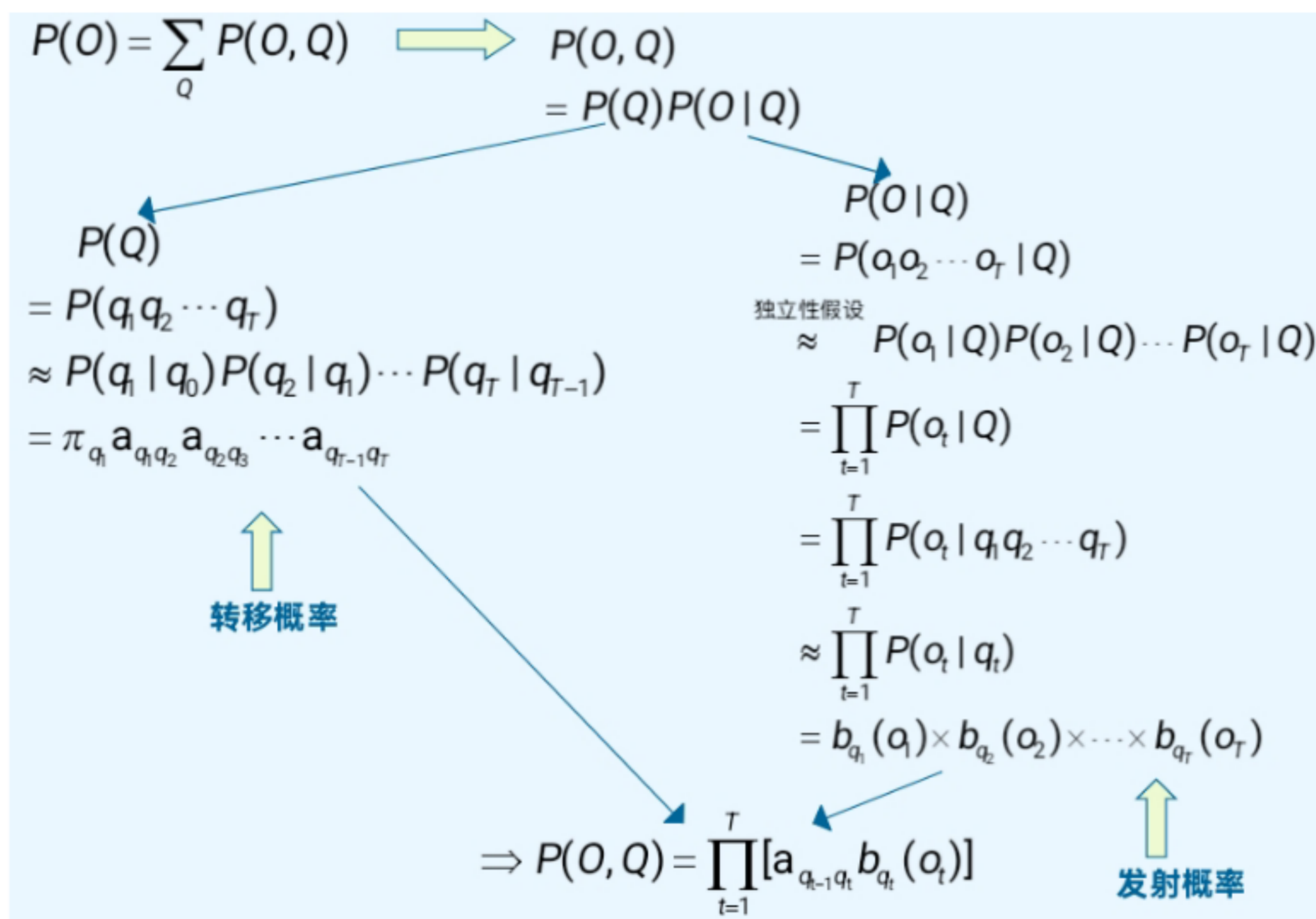
# HMM 三个基本问题:

评估: 给定模型  $u = (A, B, \pi)$ , 计算某个观测序列  $O = o_1 o_2 \dots o_T$  的概率  $P(O|u)$

预测: 给定观测序列  $O = o_1 o_2 \dots o_T$ , 如何有效地确定一个状态序列  $Q = q_1 q_2 \dots q_T$  以便最好地解释观测序列

训练: 给定观测序列  $O = o_1 o_2 \dots o_T$ , 如何找到一个能够为最好地解释这个观察序列模型, 即如何调节参数  $u = (A, B, \pi)$  使  $P(O|u)$  最大化

## 评估



时间复杂度  $O(N^T)$   
 $N$  是状态数  
 $T$  是观测状态数

注意到一阶马尔可夫过程与动态规划思想一致,

我们也可以使用动态规划来优化算法

$t$  时刻生成观测序列  $O = o_1 o_2 \dots o_t$  且达到状态  $S_i$

(记为  $\alpha_t(i)$ ) 的概率可由  $t-1$  时刻所有状态  $\alpha_{t-1}(j)$

结合状态转移概率  $a_{ji}$  和观测概率  $b_i(o_t)$  得到

$$\alpha_t(i) = \left( \sum_{j=1}^N \alpha_{t-1}(j) \cdot a_{ji} \right) \cdot b_i(O_t)$$

再累加最后一步所有状态便可知全局解

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

其中  $\alpha_t(i) = P(O_1, O_2, \dots, O_t, q_t = S_i)$ , 表示以时刻  $t$  的第  $i$  个状态  $S_i$  结束的所有路径概率之和, 称之为前向变量

• 前向算法:  $O(N^2T)$

① 初始化:  $\alpha_1(i) = \pi_i b_i(O_1)$

② 递归计算:  $\alpha_t(i) = \left( \sum_{j=1}^N \alpha_{t-1}(j) a_{ji} \right) b_i(O_t)$

③ 求和:  $P(O|u) = \sum_{i=1}^N \alpha_T(i)$

• 后向算法:

引入后向变量  $\beta_t(i)$ , 表示在时间  $t$  状态为  $S_i$  的条件下, 输出后续观测序列  $O_{t+1} O_{t+2} \dots O_T$  的概率

$$\beta_t(i) = P(O_{t+1} O_{t+2} \dots O_T | q_t = S_i)$$

即以时刻  $t$  的第  $i$  个状态开始的所有路径概率之和

$$\text{递归关系 } \beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $v_j$  转移      状态  $j$  被翻译为  $o_{t+1}$       后向变量

① 初始化  $\beta_T(i) = 1$

② 递归  $\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$

③ 求和  $P(O|u) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i)$

• 前后结合:  $\begin{cases} \text{前: } P(O) = \sum \alpha_t(i) \\ \text{后: } P(O) = \sum \pi_i b_i(o_1) \beta_1(i) \end{cases}$

得  $P(O) = \sum \alpha_t(i) \beta_t(i)$

$$\begin{aligned}
 P(O) &= P(o_1 \cdots o_T, q_t = s_1) + P(o_1 \cdots o_T, q_t = s_2) + \cdots + P(o_1 \cdots o_T, q_t = s_N) \\
 &= \sum_{i=1}^N P(o_1 \cdots o_T, q_t = s_i) \\
 P(o_1 \cdots o_T, q_t = s_i) &= P(o_1 \cdots o_t, q_t = s_i, o_{t+1} \cdots o_T) \\
 &= P(o_1 \cdots o_t, q_t = s_i) \times P(o_{t+1} \cdots o_T | o_1 \cdots o_t, q_t = s_i) \\
 &\approx P(o_1 \cdots o_t, q_t = s_i) \times P(o_{t+1} \cdots o_T | q_t = s_i) \\
 &= \alpha_t(i) \beta_t(i)
 \end{aligned}$$

$\alpha_t(i) = P(o_1 o_2 \cdots o_t, q_t = s_i)$   
 $\beta_t(i) = P(o_{t+1} o_{t+2} \cdots o_T | q_t = s_i)$

$$P(O) = \sum_{i=1}^N \alpha_t(i) \beta_t(i), \quad 1 \leq t \leq T$$

预测: 给定模型  $(A, B, \pi)$  和观测序列  $O$ , 求概率最大的隐藏状态序列  $Q$ .

• 维特比变量:  $\delta_t(i)$  表示到  $t$  时刻, 第  $i$  个状态  $s_i$  的节点的所有可能路径里概率最大的那

# 某路径的概率值

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2, \dots, q_{t-1}, q_t = S_i, o_1 o_2 \dots o_t)$$

max: 所有P里的最大值

$q_t = S_i$ : t时刻状态为  $S_i$     $o_1 \dots o_t$ : 前t个观测序列

状态转移

- ① 初始化:  $\delta_1(i) = \pi_i b_i(o_1)$
- ② 递推:  $\delta_t(i) = [\max_j (\delta_{t-1}(j) \times a_{ji})] \times b_i(o_t)$   
到状态i的所有路径里概率最大的一条  
乘状态i被观测为  $o_t$  的概率

• 反向指针:  $\psi_t(i)$  表示到  $S_i$  的最优路径里  $S_i$  的前驱节点

$$\psi_t(i) = \operatorname{argmax}_j [(\delta_{t-1}(j) \times a_{ji}) \times b_i(o_t)]$$

arg: 取索引, 也就是  $S_i$  的前驱  $S_j$  的j

注: 其实  $\psi_t(i)$  值与  $b_i(o_t)$  无关,  $\psi_t(i) = \operatorname{argmax}_j (\delta_{t-1}(j) \times a_{ji})$

• Viterbi算法

① 初始化

$$\delta_1(i) = \pi_i b_i(o_1)$$

$$\psi_1(i) = 0$$

② 递归

$$\delta_t(i) = [\max_j (\delta_{t-1}(j) a_{ji})] b_i(o_t)$$

$$\psi_t(i) = \operatorname{argmax}_j [(\delta_{t-1}(j) a_{ji}) b_i(o_t)]$$

③ 路径回溯

$$q_T = \operatorname{argmax}_i \delta_T(i)$$

$$q_{t-1} = \psi_t(q_t)$$

训练：给一组观测序列  $O = o_1, o_2, \dots, o_T$ ，调节模型参数  $u = (A, B, \pi)$  使  $P(O|u)$  最大化，即  $\arg \max_u P(O|u)$

情况1：状态序列  $Q = q_1, \dots, q_T$  已知

方法：最大似然估计，直接用  $O, Q$  得到  $A, B, \pi$

$$\pi_i = \frac{\text{状态 } i \text{ 作为初始状态次数}}{\text{总序列数}} = \delta(q_1, s_i)$$

$$\begin{aligned} a_{ij} &= \frac{\text{状态 } i \text{ 转移到 } j \text{ 的次数}}{\text{状态 } i \text{ 出现总次数}} = \frac{C(s_i s_j)}{C(s_i)} \\ &= \frac{\sum_{t=1}^T \delta(q_t, s_i) \times \delta(q_{t+1}, s_j)}{\sum_{t=1}^T \delta(q_t, s_i)} \end{aligned}$$

$$\begin{aligned} b_j(k) &= \frac{\text{状态 } j \text{ 生成观察状态 } v_k \text{ 的次数}}{\text{状态 } j \text{ 出现总次数}} = \frac{C(s_j v_k)}{C(s_j)} \\ &= \frac{\sum_{t=1}^T \delta(q_t, s_j) \times \delta(o_t, v_k)}{\sum_{t=1}^T \delta(q_t, s_j)} \end{aligned}$$

注：这里的  $\delta$  是克罗内克函数  $\delta(a, b) = \begin{cases} 1 & a=b \\ 0 & a \neq b \end{cases}$   
不是维特比变量。

情况2：状态序列  $Q$  未知。

目前无任何方法来求取全局最优解，但可以用 Baum-Welch 算法（前向后向算法）来求局部最优解，这种算法是 EM 算法的特例，下面介绍 EM 算法。

## EM算法

一种迭代优化算法，用于在存在隐变量的情况下进行参数估计。它通过交替执行两个步骤来优化参数：

E步：在当前参数下，计算隐变量（状态序列）的期望  
即 Expectation

M步：通过E步的期望更新参数，以求最大化似然函数  
即 Maximization.

Baum-Welch算法：① 初始化：随机给参数赋值，但要满足以下约束：

$$\sum_{i=1}^N \pi_i = 1 \quad \sum_{j=1}^N a_{ij} = 1 \quad \sum_{k=1}^M b_i(k) = 1$$

由此得到模型  $u_0$ ，令  $i=0$ ，执行②

② EM计算：

E：由  $u_i$  参数计算  $\gamma_t(i)$ （时刻  $t$  处于状态  $i$  的概率）

$$\begin{aligned} \gamma_t(i) &= \tilde{P}(q_t = s_i | O, u) \\ &= \frac{\tilde{P}(q_t = s_i, O | u)}{\tilde{P}(O | u)} \\ &= \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)} \end{aligned}$$

$\tilde{P}(q_t = s_i | O, u)$ ：给定观测序列  $O$ ，模型  $u$ ，求时刻  $t$  处于状态  $s_i$  的概率。

$\tilde{P}(q_t = s_i, O | u)$ ：在模型  $u$  下，时刻  $t$  处于状态  $s_i$  且观测序列为  $O$  的概率。

$\tilde{P}(O | u)$ ：在模型  $u$  下观测到  $O$  的边缘概率。

$\alpha_t(i)$ ：前向概率       $\beta_t(i)$ ：后向概率。

M：由  $\gamma_t(i)$  估计  $\pi_i, a_{ij}, b_i(k)$ ，得到  $u_{i+1}$

$$\pi_i = \tilde{P}(q_1 = s_i | O, u) = \gamma_1(i)$$

$$a_{ij} = \frac{\sum_{t=2}^{T-1} \tilde{P}(q_t = s_i, q_{t+1} = s_j | O, u)}{\sum_{t=2}^{T-1} \tilde{P}(q_t = s_i | O, u)}$$

$$z = \frac{\sum_{t=2}^{T-1} \frac{\tilde{P}(q_t = s_i, q_{t+1} = s_j, O | u)}{\tilde{P}(O | u)}}{\sum_{t=2}^{T-1} \gamma_t(i)}$$

$$= \frac{\sum_{t=2}^{T-1} \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{j=2}^N \alpha_t(j) \beta_t(j)}}$$

$$b'_i(k) = \frac{\sum_{t=2}^{T-1} \tilde{P}(q_t = s_i, O_t = v_k | O, u)}{\sum_{t=2}^{T-1} \tilde{P}(q_t = s_i | O, u)}$$

$$= \frac{\sum_{t=2}^T \gamma_t(i) \times \delta(O_t, v_k)}{\sum_{t=2}^T \gamma_t(i)}$$

③ 循环:  $i++$ , 执行②, 直到  $T$ ,  $A, B$  收敛

A、B、 $\pi$  的初始化: 刚提到, 算法开始时我们会对 A、B、 $\pi$  随机赋值, 而算法只能找到这个起点附近的局部最优解. 要想得到全局最优解, 起点就必须在其附近. 在实践中我们发现, 好的初始估计对 B 十分重要, 而 A、 $\pi$  随机赋值就够了.