

自然语言处理

Natural Language Processing

权小军 教授

中山大学数据科学与计算机学院

quanxj3@mail.sysu.edu.cn

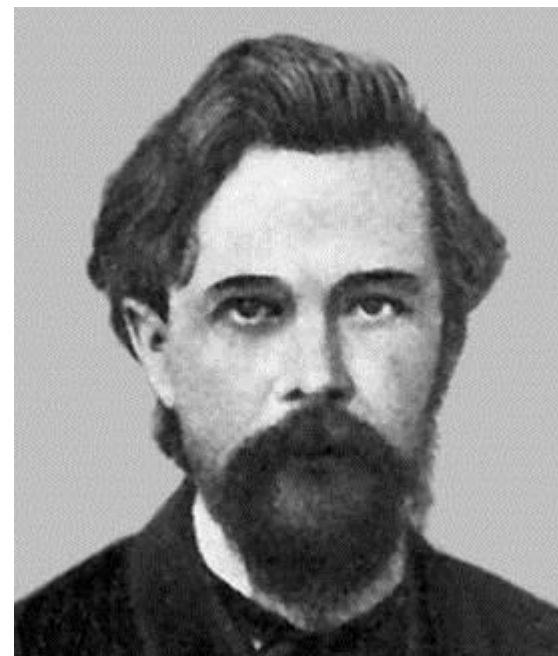
Lecture 5: 马尔可夫模型

Lecture 5.1 概述

马尔可夫模型

◆ 马尔可夫(Markov) (1856. 6. 14 ~ 1922. 7. 20)

前苏联数学家。在概率论、数论、函数逼近论和微分方程等方面卓有成就。他提出了用数学分析方法研究自然过程的一般图式——马尔可夫链，并开创了随机过程(马尔可夫过程)的研究。



马尔可夫模型

◆马尔可夫模型描述

存在一类重要的随机过程：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，该概率为：

$$p(q_t = S_j \mid q_{t-1} = S_i, q_{t-2} = S_k, \dots)$$

马尔可夫模型

●假设1:

如果在特定情况下，系统在时间 t 的状态只与其在时间 $t-1$ 的状态相关，则该系统构成一个离散的一阶马尔可夫链：

$$p(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots) = p(q_t = S_j | q_{t-1} = S_i) \\ \dots (1)$$

马尔可夫模型

●假设2:

如果只考虑公式(1)独立于时间 t 的随机过程, 即所谓的不动性假设, 状态与时间无关, 那么:

$$p(q_t = S_j \mid q_{t-1} = S_i) = a_{ij}, \quad 1 \leq i, j \leq N \quad \dots (2)$$

该随机过程称为**马尔可夫模型(Markov Model)**。

马尔可夫模型

在马尔可夫模型中，状态转移概率 a_{ij} 必须满足下列条件：

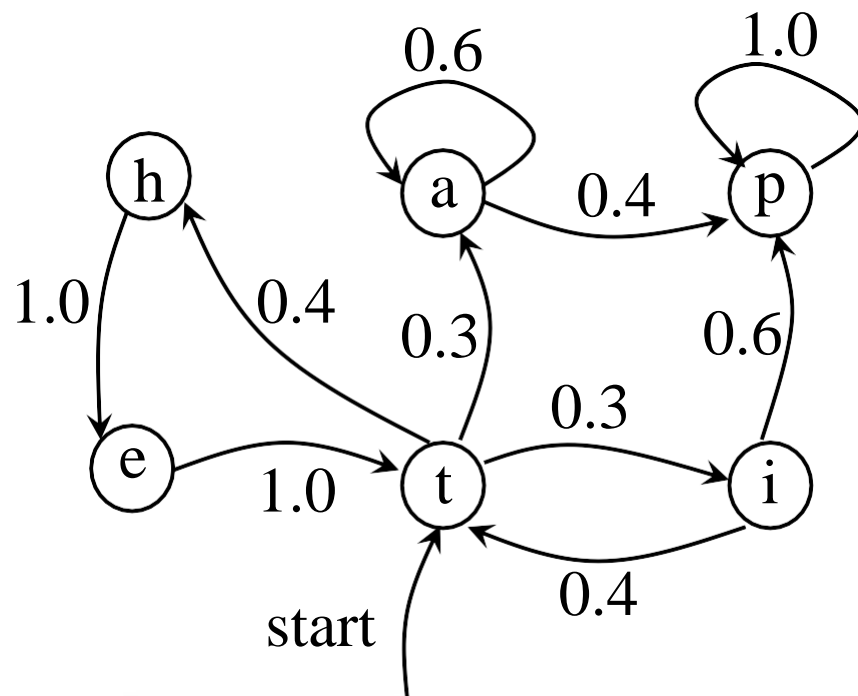
$$a_{ij} \geq 0 \quad \dots (3)$$

$$\sum_{j=1}^N a_{ij} = 1 \quad \dots (4)$$

马尔可夫模型

◆ 马尔可夫链可以表示成状态图（转移弧上有概率的非确定的有限状态自动机）

- 零概率的转移弧省略。
- 每个节点上所有发出弧的概率之和等于1。



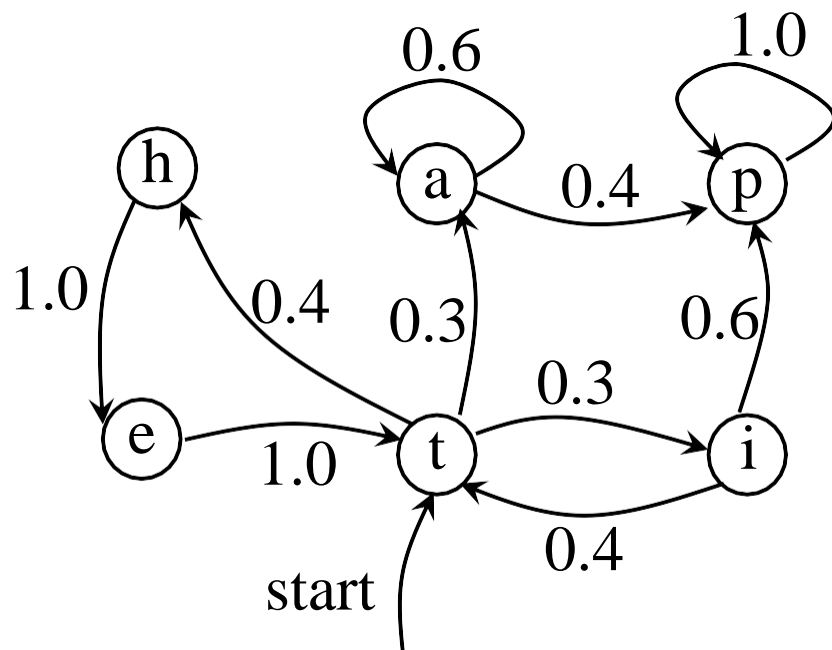
马尔可夫模型

状态序列 S_1, \dots, S_T 的概率:

$$\begin{aligned} p(S_1, \dots, S_T) &= p(S_1) \times p(S_2 | S_1) \times p(S_3 | S_1, S_2) \times \dots \times p(S_T | S_1, \dots, S_{T-1}) \\ &= p(S_1) \times p(S_2 | S_1) \times p(S_3 | S_2) \times \dots \times p(S_T | S_{T-1}) \\ &= \pi_{S_1} \prod_{t=1}^{T-1} a_{S_t S_{t+1}} \quad \dots (5) \end{aligned}$$

其中, $\pi_i = p(q_1 = S_i)$ 为初始状态的概率。

马尔可夫模型



$$\begin{aligned} p(t, i, p) &= p(S_1=t) \times p(S_2=i | S_1=t) \times p(S_3=p | S_2=i) \\ &= 1.0 \times 0.3 \times 0.6 \\ &= 0.18 \end{aligned}$$

Lecture 5.2 隐马尔可夫模型

隐马尔可夫模型

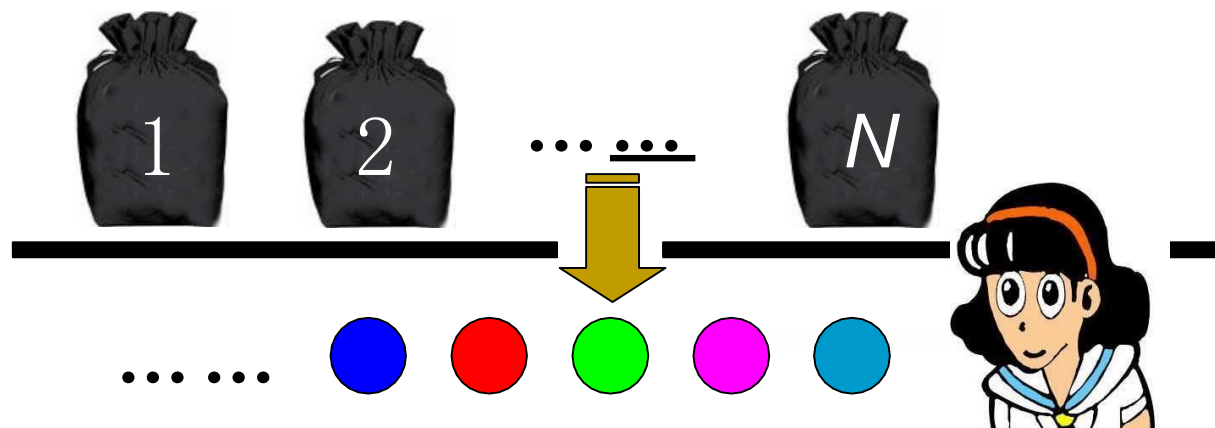
◆ 隐马尔可夫模型 (Hidden Markov Model, HMM)

描写: 该模型是一个双重随机过程，我们不知道具体的状态序列，只知道状态转移的概率，即模型的状态转换过程是不可观察的（隐蔽的），而可观察事件的随机过程是隐蔽状态转换过程的随机函数。

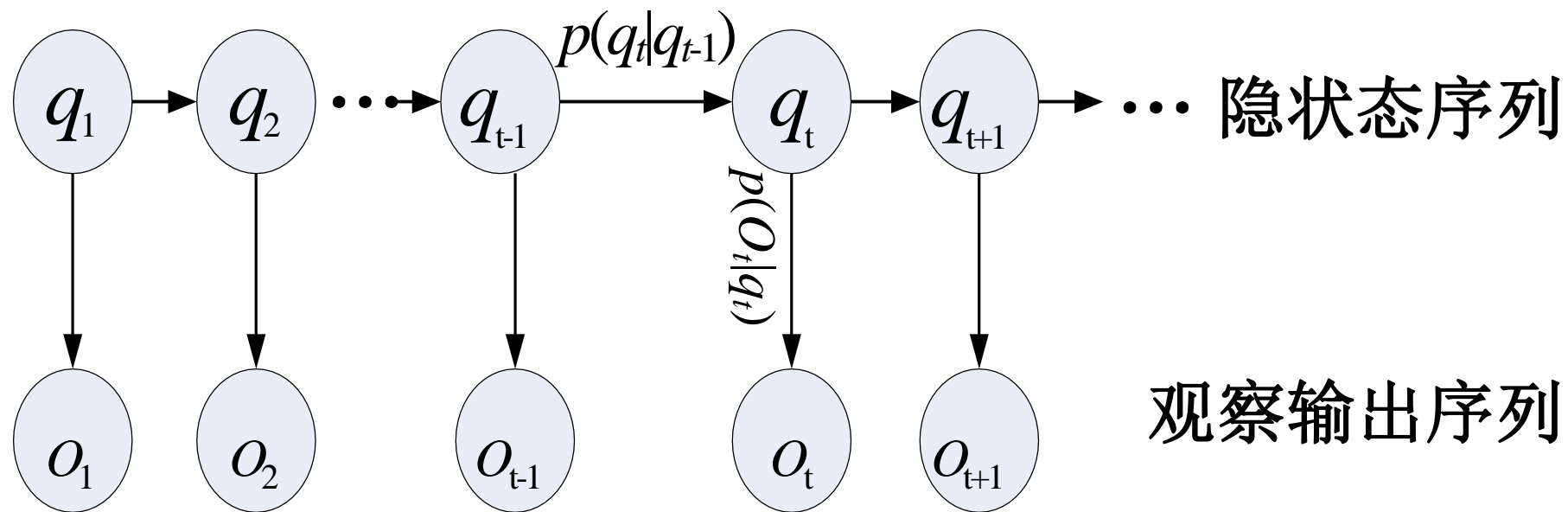
隐马尔可夫模型

例如： N 个袋子，每个袋子中有 M 种不同颜色的球。一实验员根据某一概率分布选择一个袋子，然后根据袋子中不同颜色球的概率分布随机取出一个球，并报告该球的颜色。

对局外人：可观察的过程是不同颜色球的序列，而袋子的序列是不可观察的。每只袋子对应HMM中的一个状态；球的颜色对应于HMM中状态的输出。



隐马尔可夫模型



HMM 图解

隐马尔可夫模型

◆HMM 的组成

1. 模型中的状态数为 N (袋子的数量)
2. 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)

隐马尔可夫模型

状态转移概率矩阵 $A = a_{ij}$, a_{ij} 为实验员从一只袋子 (状态 S_i) 转向另一只袋子 (状态 S_j) 取球的概率。其中,

$$\left\{ \begin{array}{l} a_{ij} = p(q_{t+1} = S_j | q_t = S_i), \quad 1 \leq i, j \leq N \\ a_{ij} \geq 0 \\ \sum_{j=1}^N a_{ij} = 1 \end{array} \right. \quad \dots (6)$$

隐马尔可夫模型

从状态 S_j 观察到某一特定符号 v_k 的概率分布矩阵为：

$B=b_j(k)$ ；其中， $b_j(k)$ 为实验员从第 j 个袋子中取出第 k 种颜色的球的概率。那么，

$$\left\{ \begin{array}{l} b_j(k) = p(O_t = v_k | q_t = S_j), \quad 1 \leq j \leq N, \quad 1 \leq k \leq M \\ b_j(k) \geq 0 \\ \sum_{k=1}^M b_j(k) = 1 \end{array} \right. \dots (7)$$

隐马尔可夫模型

初始状态的概率分布为： $\pi = \pi_i$ ，其中，

$$\left\{ \begin{array}{l} \pi_i = p(q_1 = S_i), \quad 1 \leq i \leq N \\ \pi_i \geq 0 \\ \sum_{i=1}^N \pi_i = 1 \end{array} \right. \quad \dots (8)$$

为了方便，一般将 HMM 记为： $\mu = (A, B, \pi)$ 或者 $\mu = (S, O, A, B, \pi)$ 用以指出模型的参数集合。

隐马尔可夫模型

◆ 给定HMM求观察序列

给定模型 $\mu = (A, B, \pi)$, 产生观察序列 $O = O_1 O_2 \dots O_T$:

隐马尔可夫模型

◆ 给定HMM求观察序列

给定模型 $\mu = (A, B, \pi)$, 产生观察序列 $O = O_1 O_2 \dots O_T$:

- (1) 令 $t=1$;
- (2) 根据初始状态分布 $\pi=\pi_i$ 选择初始状态 $q_1=S_i$;
- (3) 根据状态 S_i 的输出概率分布 $b_i(k)$, 输出 $O_t=v_k$;
- (4) 根据状态转移概率 a_{ij} , 转移到新状态 $q_{t+1}=S_j$;
- (5) $t=t+1$, 如果 $t < T$, 重复步骤 (3) (4), 否则结束。

隐马尔可夫模型

◆三个问题:

- (1) 在给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ 的情况下, 怎样快速计算概率 $p(O|\mu)$?

隐马尔可夫模型

◆三个问题:

- (1) 在给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ 的情况下, 怎样快速计算概率 $p(O|\mu)$?
- (2) 在给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ 的情况下, 如何选择在一定意义下“最优”的状态序列 $Q=q_1q_2 \dots q_T$, 使得该状态序列“最好地解释”观察序列?

隐马尔可夫模型

◆三个问题:

- (1) 在给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ 的情况下, 怎样快速计算概率 $p(O|\mu)$?
- (2) 在给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ 的情况下, 如何选择在一定意义下“最优”的状态序列 $Q=q_1q_2 \dots q_T$, 使得该状态序列“最好地解释”观察序列?
- (3) 给定一个观察序列 $O=O_1O_2 \dots O_T$, 如何根据极大似然估计来求模型的参数值? 即如何调节模型的参数, 使得 $p(O|\mu)$ 最大?

Lecture 5.3 前向算法

前向算法

◆ 问题1：快速计算观察序列概率 $p(O|\mu)$

给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O = O_1 O_2 \dots O_T$,
快速计算 $p(O|\mu)$:

前向算法

◆ 问题1：快速计算观察序列概率 $p(O|\mu)$

给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O = O_1 O_2 \dots O_T$,
快速计算 $p(O|\mu)$:

$$p(O|\mu) = \sum_Q p(O, Q|\mu) = \sum_Q p(Q|\mu) \times p(O|Q, \mu) \quad \dots (9)$$

前向算法

◆ 问题1：快速计算观察序列概率 $p(O|\mu)$

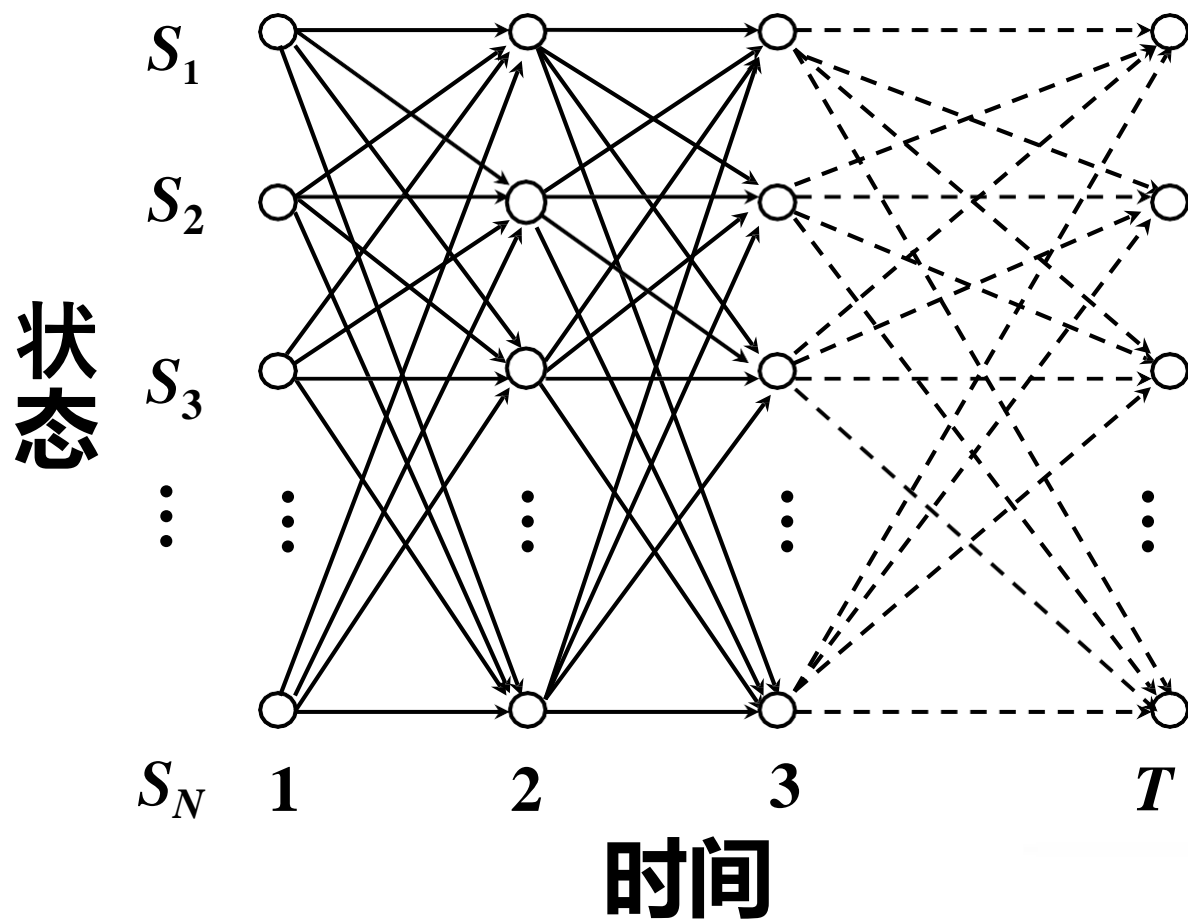
给定模型 $\mu=(A, B, \pi)$ 和观察序列 $O=O_1O_2 \dots O_T$ ，
快速计算 $p(O|\mu)$ ：

$$p(O|\mu) = \sum_Q p(O, Q|\mu) = \sum_Q \boxed{p(Q|\mu)} \times \boxed{p(O|Q, \mu)} \quad \dots (9)$$

$$p(Q|\mu) = \pi_{q_1} \times a_{q_1q_2} \times a_{q_2q_3} \times \dots \times a_{q_{t-1}q_T} \quad \dots (10)$$

$$p(O|Q, \mu) = b_{q_1}(O_1) \times b_{q_2}(O_2) \times \dots \times b_{q_T}(O_T) \quad \dots (11)$$

前向算法



● 困难：

如果模型 μ 有 N 个不同的状态，
时间长度为 T ，
那么有 N^T 个可能的状态序列，
搜索路径成指数级组合爆炸。

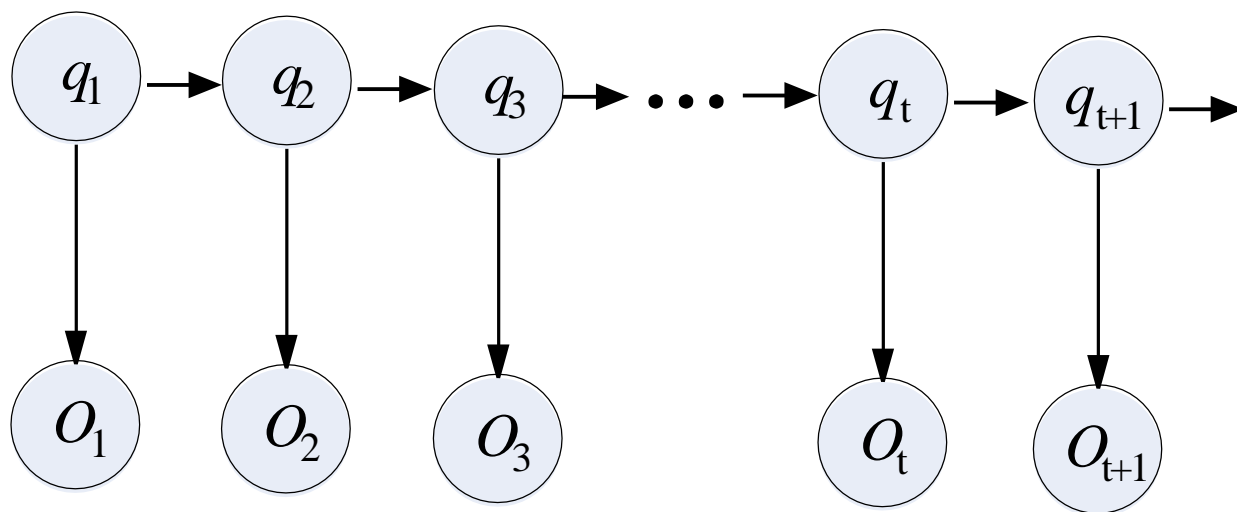
前向算法

- 解决办法：动态规划
前向算法(The forward procedure)
- 基本思想：定义前向变量 $\alpha_t(i)$ ：

$$\alpha_t(i) = p(O_1 O_2 \cdots O_t, \underline{q_t} = S_i | \mu) \quad \dots(12)$$

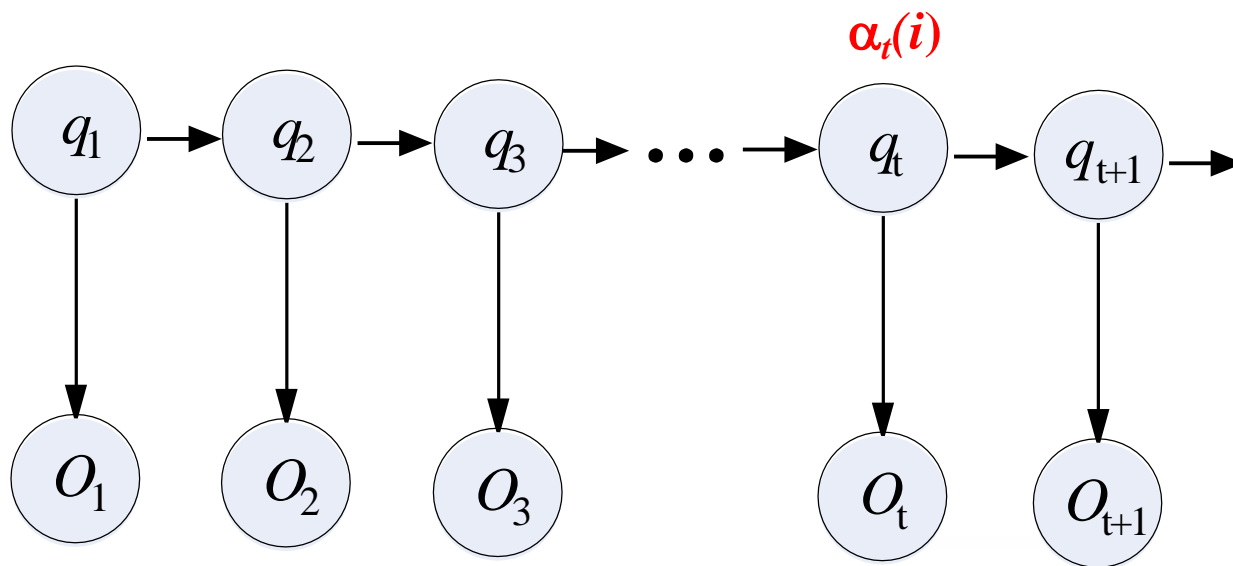
如果可以高效地计算 $\alpha_t(i)$ ，就可以高效地求得 $p(O|\mu)$ 。

前向算法



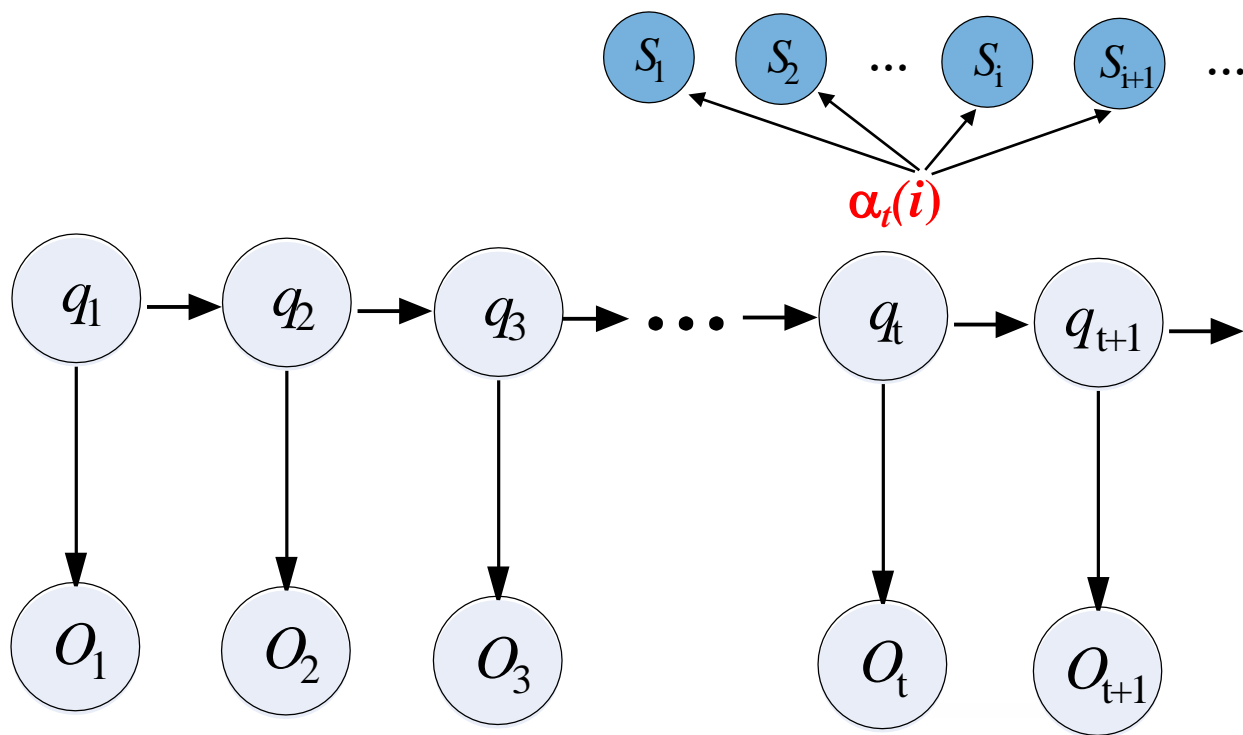
前向算法

$$\alpha_t(i) = p(O_1 O_2 \cdots O_t, q_t = S_i | \mu)$$



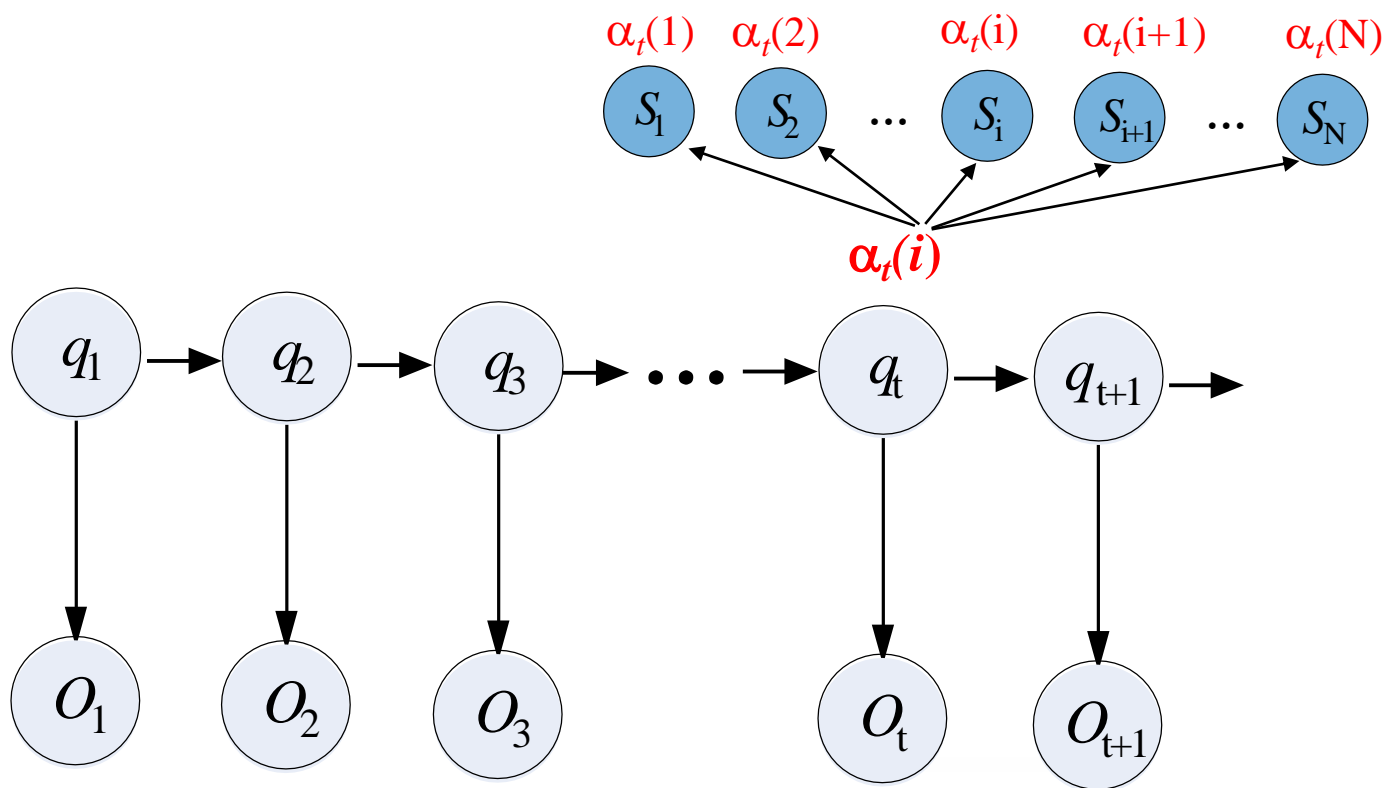
前向算法

$$\alpha_t(i) = p(O_1 O_2 \cdots O_t, q_t = S_i | \mu)$$



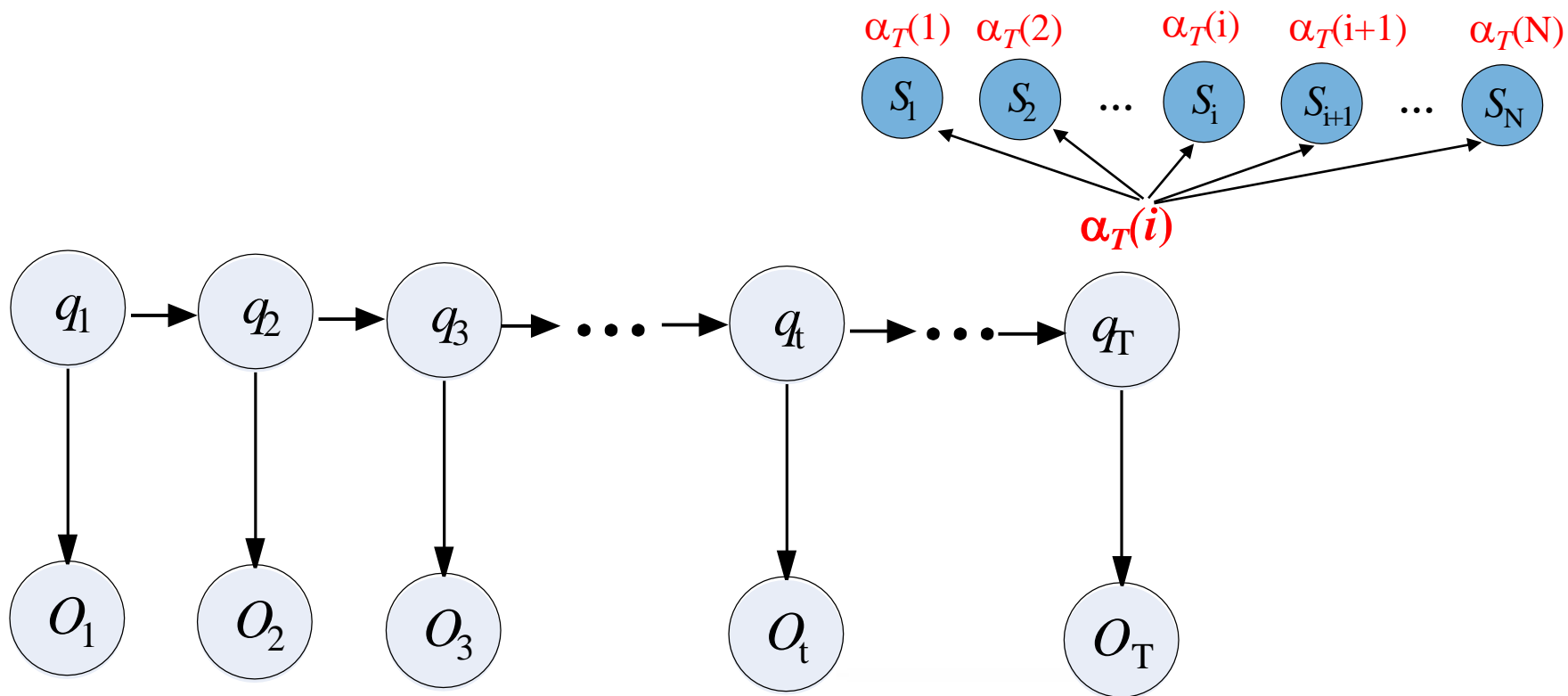
前向算法

$$\alpha_t(i) = p(O_1 O_2 \cdots O_t, q_t = S_i | \mu)$$



前向算法

$$\alpha_t(i) = p(O_1 O_2 \cdots O_t, q_t = S_i | \mu)$$



前向算法

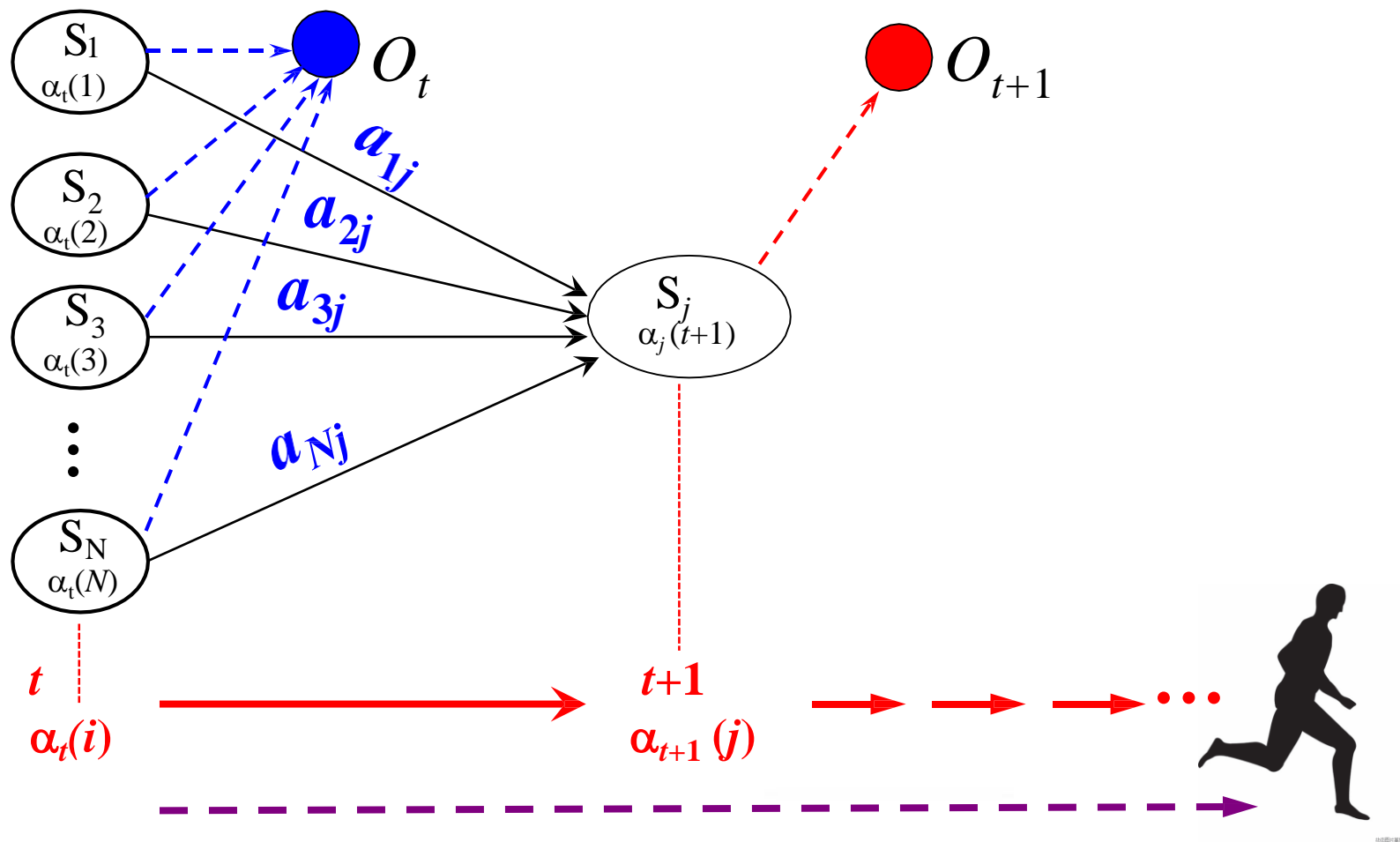
因为 $p(O|\mu)$ 是在到达状态 q_T 时观察到序列 $O = O_1 O_2 \dots O_T$ 的概率(所有可能的概率之和)：

$$p(O|\mu) = \sum_{S_i} p(O_1 O_2 \dots O_T, q_T = S_i | \mu) = \sum_{i=1}^N \alpha_T(i) \quad \dots (13)$$

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：

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] \times b_j(O_{t+1})$$

前向算法



前向算法

● 算法1：前向算法描述

(1) 初始化： $\alpha_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N$

(2) 循环计算：

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] \times b_j(O_{t+1}), \quad 1 \leq t \leq T-1$$

(3) 结束，输出：

$$p(O | \mu) = \sum_{i=1}^N \alpha_T(i)$$

前向算法

$\alpha_I(1)$

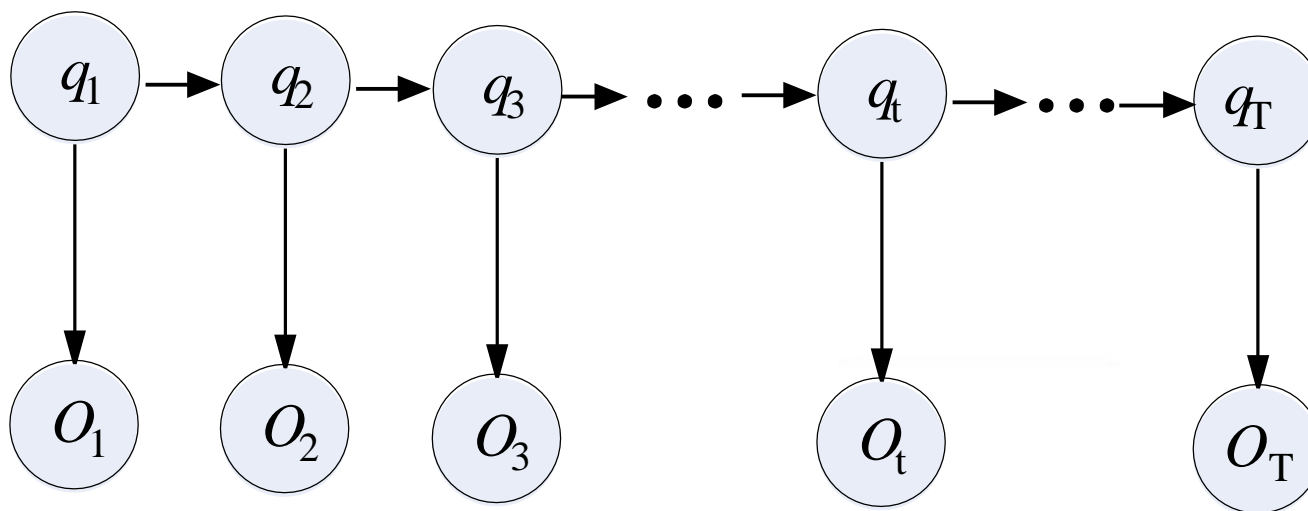
$\alpha_I(2)$

\vdots

$\alpha_I(i)$

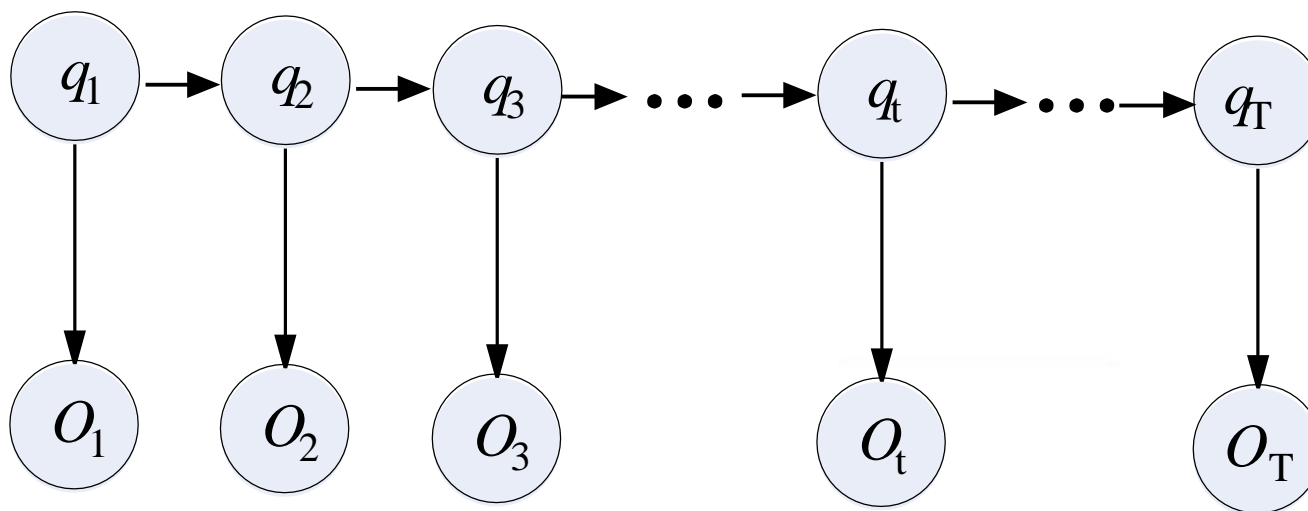
\vdots

$\alpha_I(N)$



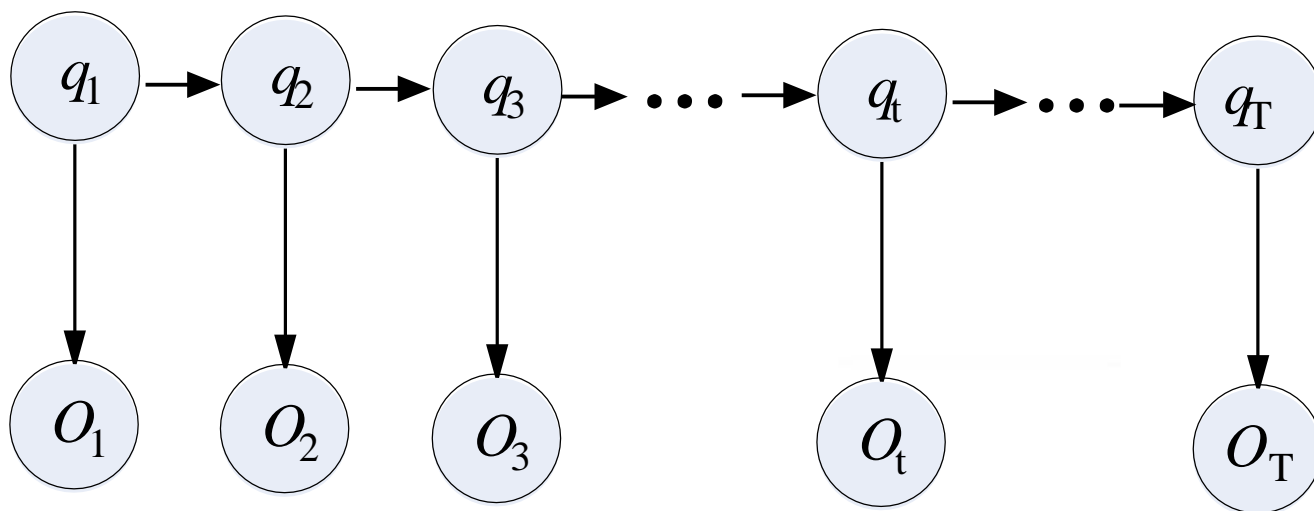
前向算法

$\alpha_1(1)$	$\alpha_2(1)$
$\alpha_1(2)$	$\alpha_2(2)$
\vdots	\vdots
$\alpha_1(i)$	$\alpha_2(i)$
\vdots	\vdots
$\alpha_1(N)$	$\alpha_2(N)$



前向算法

$$\left(\begin{array}{ccc} \alpha_1(1) & \alpha_2(1) & \dots & \alpha_T(1) \\ \alpha_1(2) & \alpha_2(2) & & \alpha_T(2) \\ \vdots & \vdots & & \vdots \\ \alpha_1(i) & \alpha_2(i) & & \alpha_T(i) \\ \vdots & \vdots & & \vdots \\ \alpha_1(N) & \alpha_2(N) & \dots & \alpha_T(N) \end{array} \right)$$



前向算法

● 算法的时间复杂性:

每计算一个 $\alpha_t(i)$ 必须考虑从 $t-1$ 时的所有 N 个状态转移到状态 S_i 的可能性, 时间复杂性为 $O(N)$, 对应每个时刻 t , 要计算 N 个前向变量: $\alpha_t(1), \alpha_t(2), \dots, \alpha_t(N)$, 所以, 时间复杂性为: $O(N) \times N = O(N^2)$ 。又因 $t = 1, 2, \dots, T$, 所以前向算法总的复杂性为: $O(N^2T)$ 。

Lecture 5.4 后向算法

后向算法

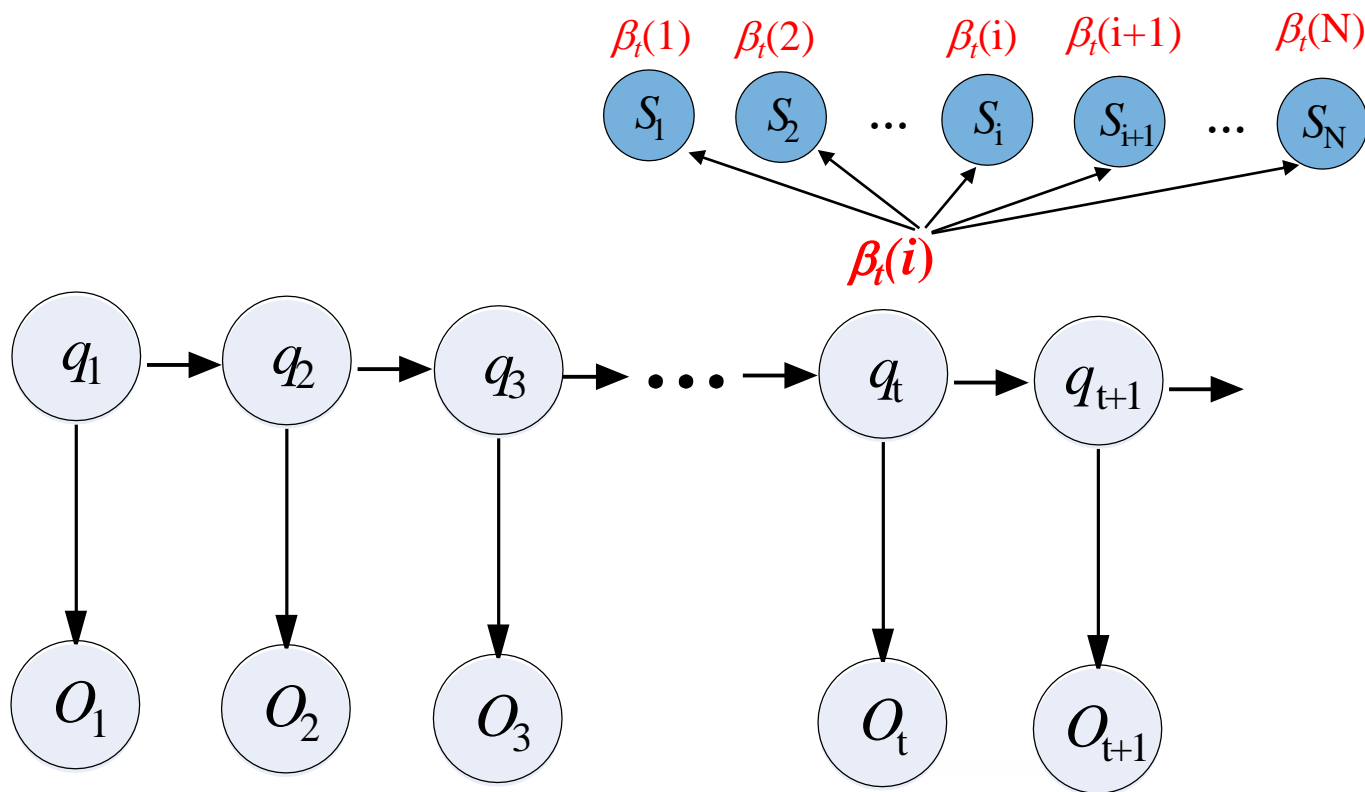
- 后向算法 (The backward procedure)

定义后向变量 $\beta_t(i)$ 是在给定了模型 $\mu = (A, B, \pi)$ 和假定在时间 t 状态为 S_i 的条件下, 模型输出观察序列 $O_{t+1}O_{t+2}\cdots O_T$ 的概率:

$$\beta_t(i) = p(O_{t+1}O_{t+2}\cdots O_T | q_t = S_i, \mu) \quad \dots (15)$$

后向算法

$$\beta_t(i) = p(O_{t+1}O_{t+2}\cdots O_T | q_t = S_i, \mu)$$



后向算法

与前向变量一样，运用动态规划计算后向量：

- (1) 从时刻 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，模型输出观察序列 $O_{t+2}O_{t+3}\cdots O_T$ 。

后向算法


第一步的概率： $a_{ij} \times b_j(O_{t+1})$

第二步的概率按后向变量的定义为： $\beta_{t+1}(j)$

于是，有归纳关系：

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \times \beta_{t+1}(j) \quad \dots (16)$$

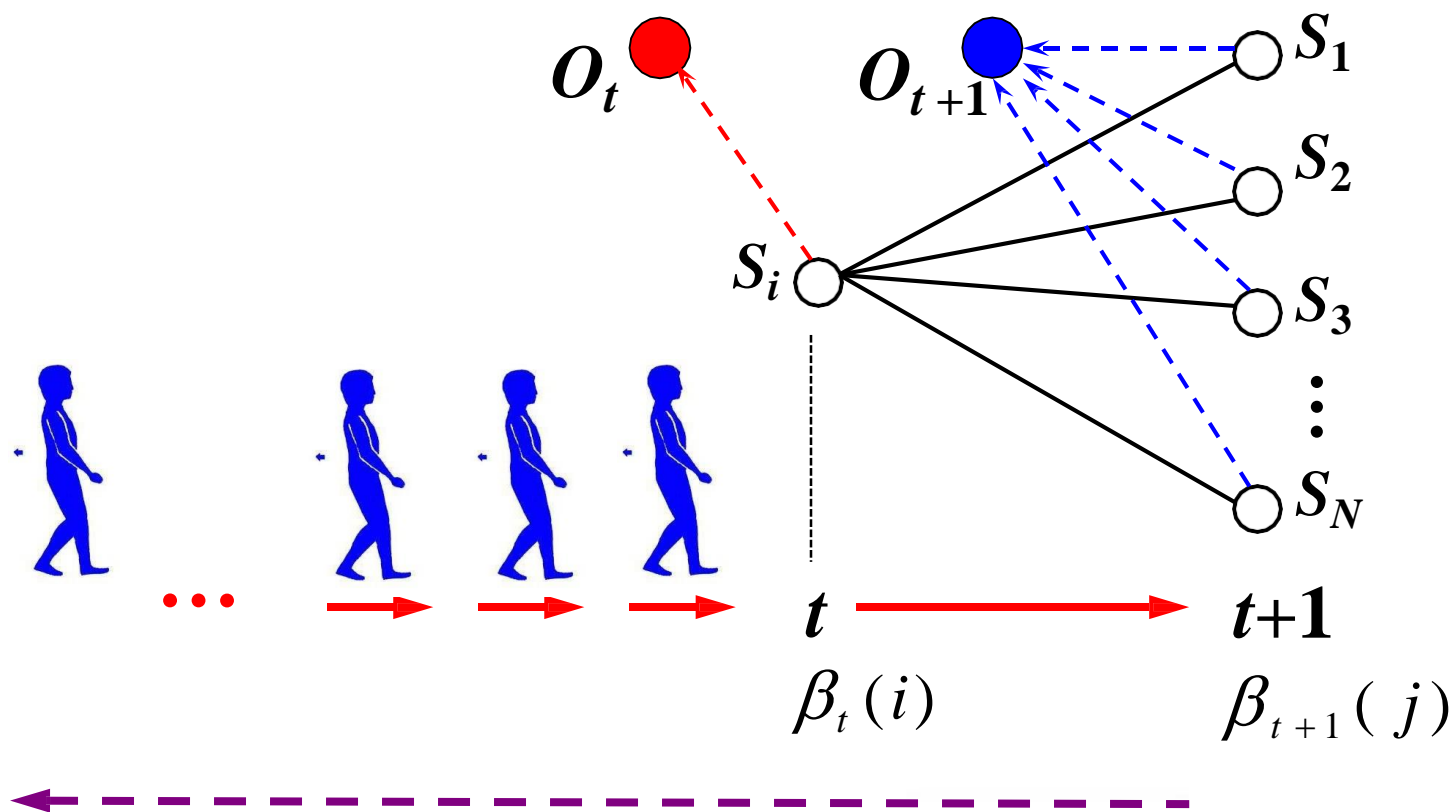
归纳顺序： $\beta_T(x), \beta_{T-1}(x), \dots, \beta_1(x)$



(x 为模型的状态)

后向算法

算法图解：



后向算法

● 算法2：后向算法描述

(1) 初始化： $\beta_T(i) = 1, \quad 1 \leq i \leq N$

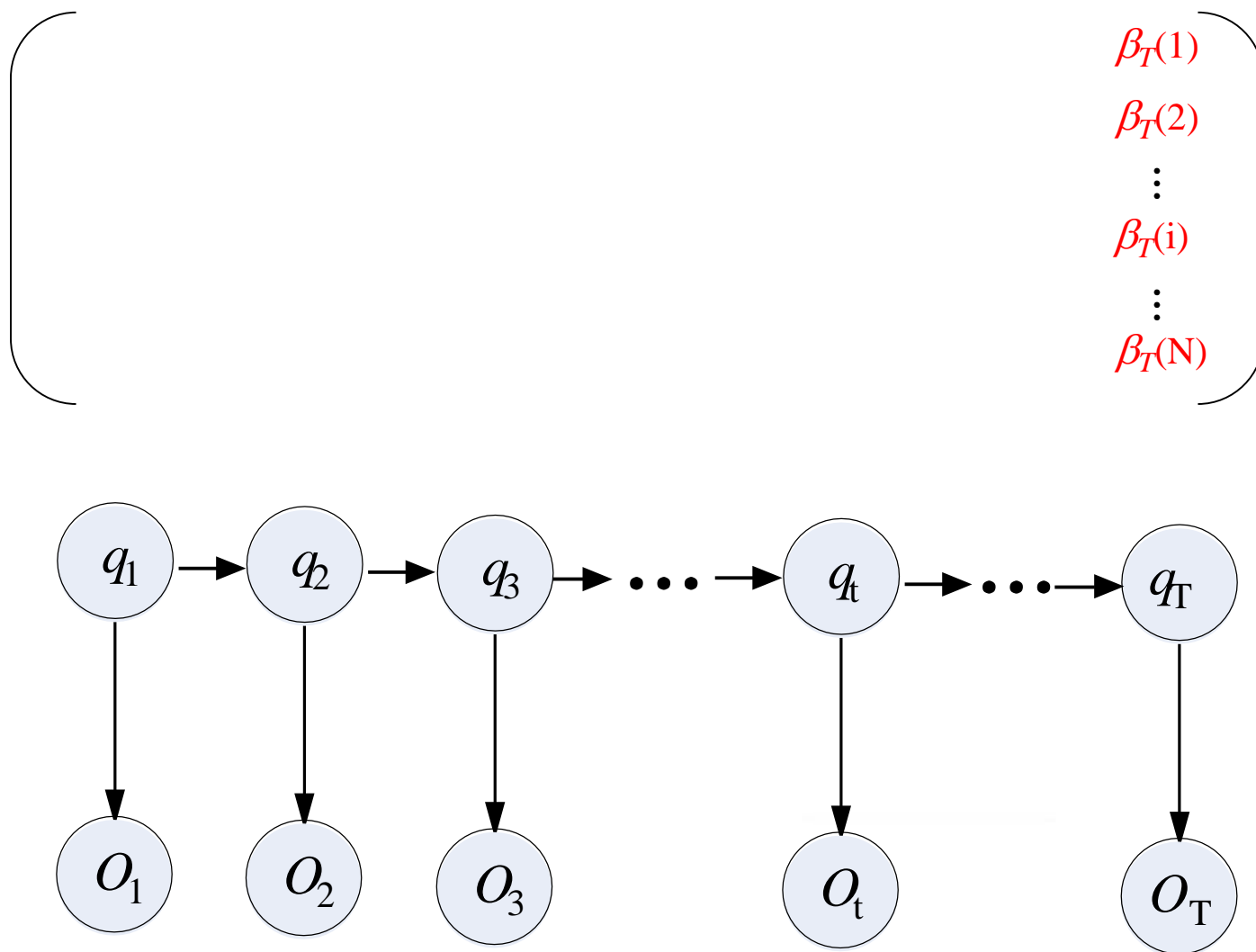
(2) 循环计算：

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \times \beta_{t+1}(j), \quad T-1 \geq t \geq 1, \quad 1 \leq i \leq N$$

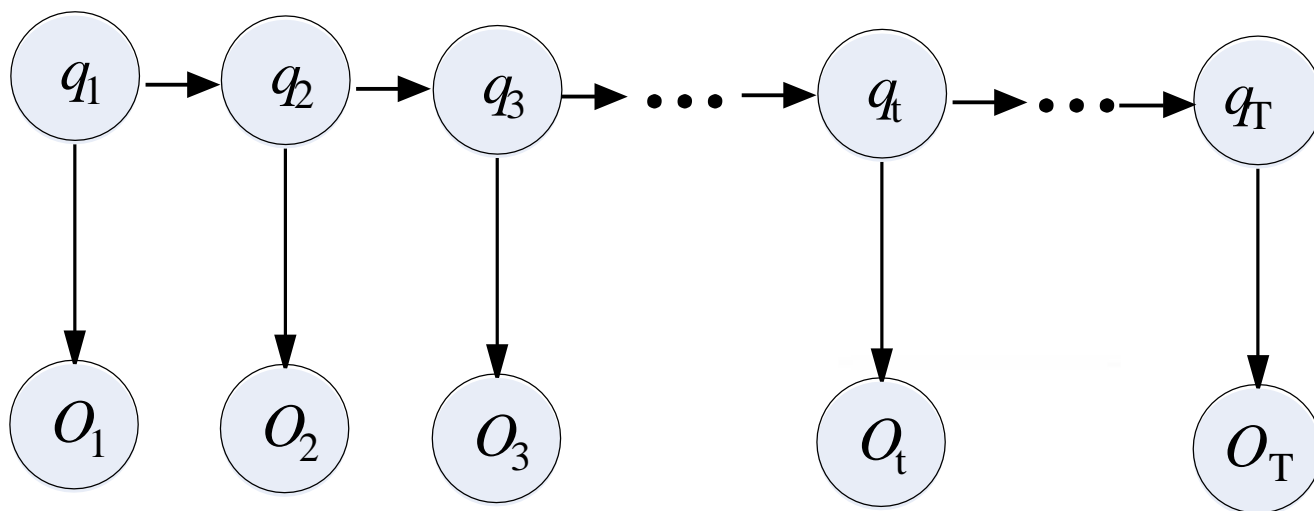
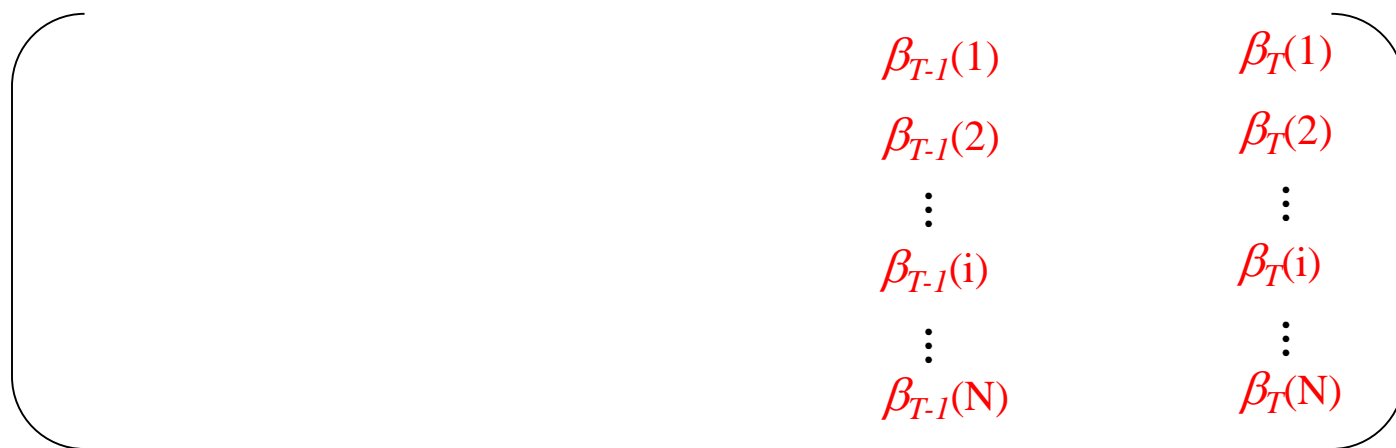
(3) 输出结果： $p(O|\mu) = \sum_{i=1}^N \beta_1(i) \times \pi_i \times b_i(O_1)$

算法的时间复杂性： $O(N^2T)$

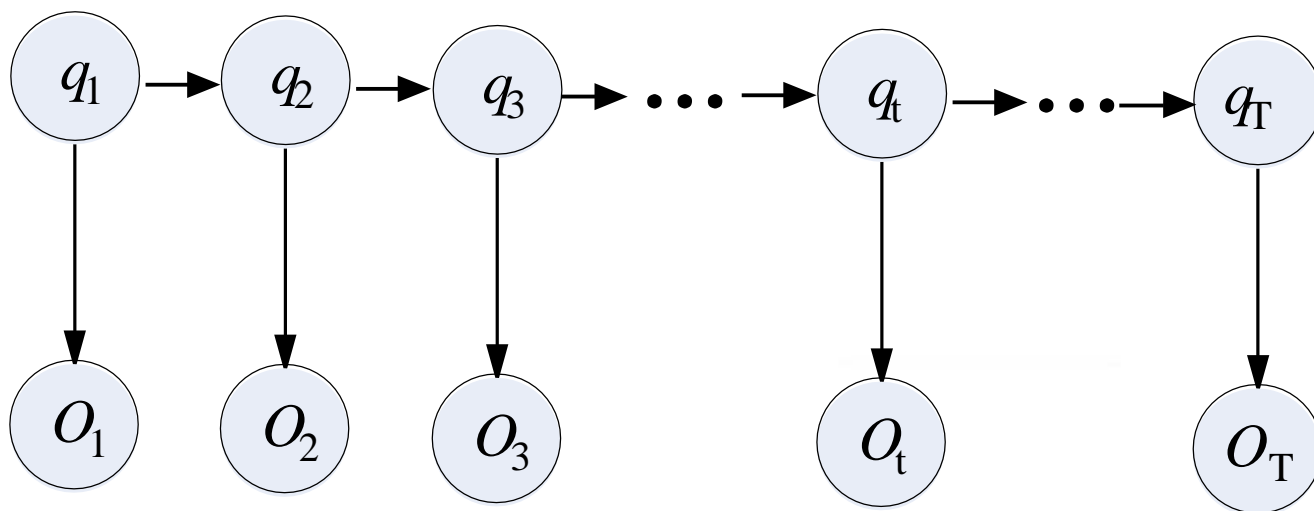
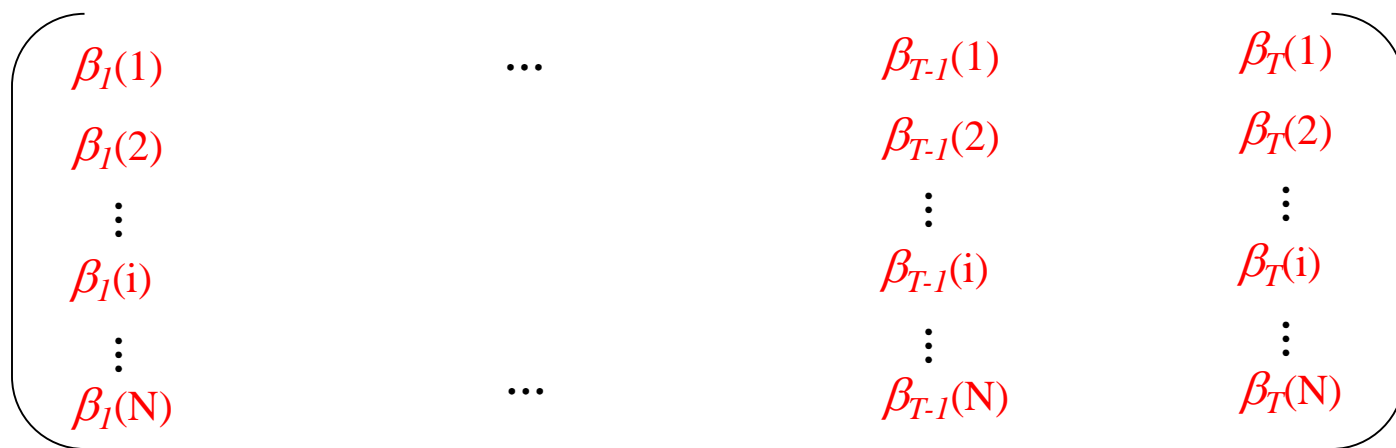
前向算法



前向算法



前向算法



Lecture 5.5 Viterbi搜索算法

Viterbi 搜索算法

◆问题2—如何发现“最优”状态序列能够“最好地解释”观察序列？

在给定模型 μ 和观察序列 O 的条件下求概率最大的状态序列：

$$\hat{Q} = \arg \max_Q p(Q | O, \mu) \quad \dots (21)$$

Viterbi 搜索算法

◆问题2—如何发现“最优”状态序列能够“最好地解释”观察序列？

在给定模型 μ 和观察序列 O 的条件下求概率最大的状态序列：

$$\begin{aligned}\hat{Q} &= \arg \max_Q p(Q | O, \mu) && \dots (21) \\ &= \arg \max_Q \frac{p(Q, O | \mu)}{p(O | \mu)}\end{aligned}$$

Viterbi 搜索算法

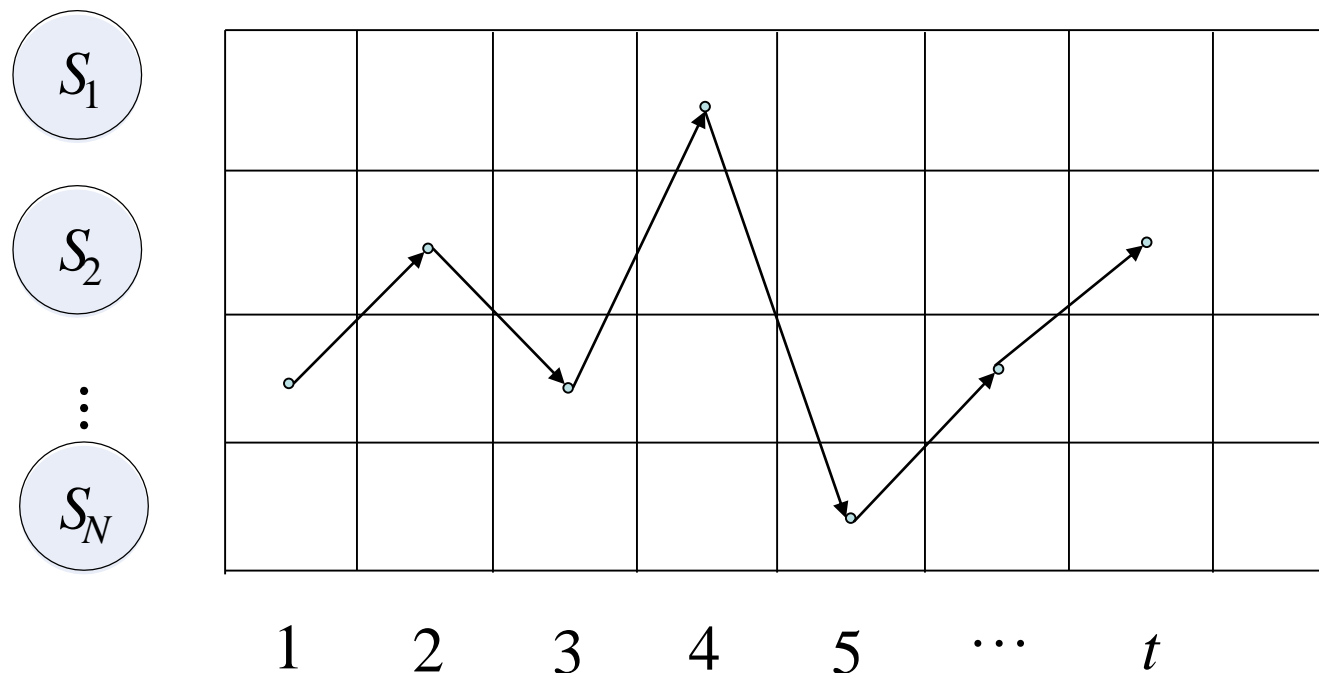
Viterbi 算法：动态搜索最优状态序列。

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

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} p(q_1, q_2, \dots, q_t = S_i, O_1 O_2 \dots O_t | \mu) \quad \dots (22)$$

Viterbi 搜索算法

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} p(q_1, q_2, \dots, q_t = S_i, O_1 O_2 \dots O_t | \mu) \quad \dots (22)$$



递归计算: $\delta_{t+1}(i) = \max_j [\delta_t(j) \cdot a_{ji}] \cdot b_i(O_{t+1}) \quad \dots (23)$

Viterbi 搜索算法

● 算法3：Viterbi 算法描述

(1) 初始化： $\delta_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N$

概率最大的路径变量： $\psi_1(i) = 0$

(2) 递推计算：

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) \cdot a_{ij}] \cdot b_j(O_t), \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) \cdot a_{ij}] \cdot b_j(O_t), \quad 2 \leq t \leq T, \quad 1 \leq i \leq N$$

Viterbi 搜索算法

(3) 结束:

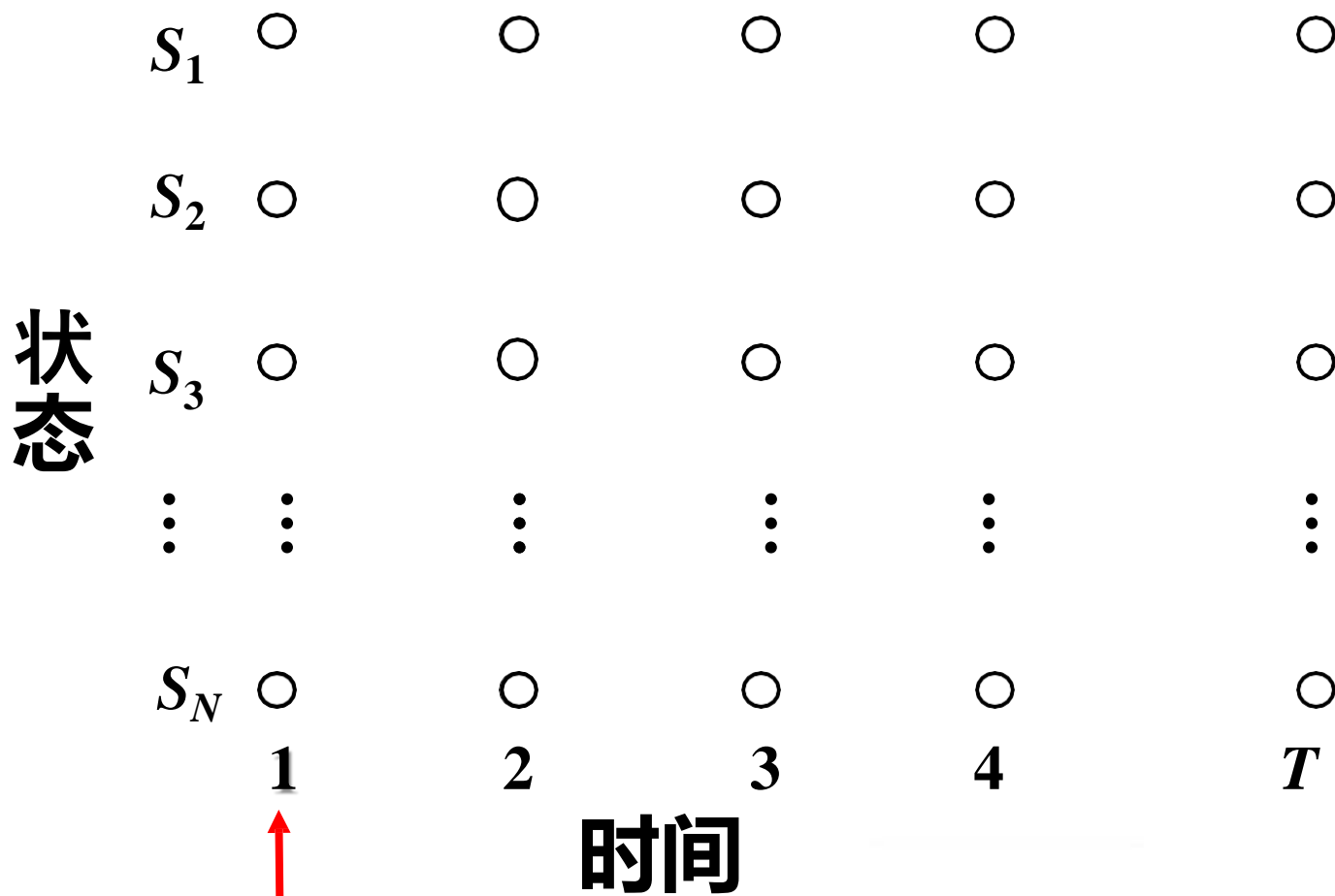
$$\hat{Q}_T = \operatorname{argmax}_{1 \leq i \leq N} [\delta_T(i)], \quad \hat{p}(\hat{Q}_T) = \max_{1 \leq i \leq N} \delta_T(i)$$

(4) 通过回溯得到路径（状态序列）：

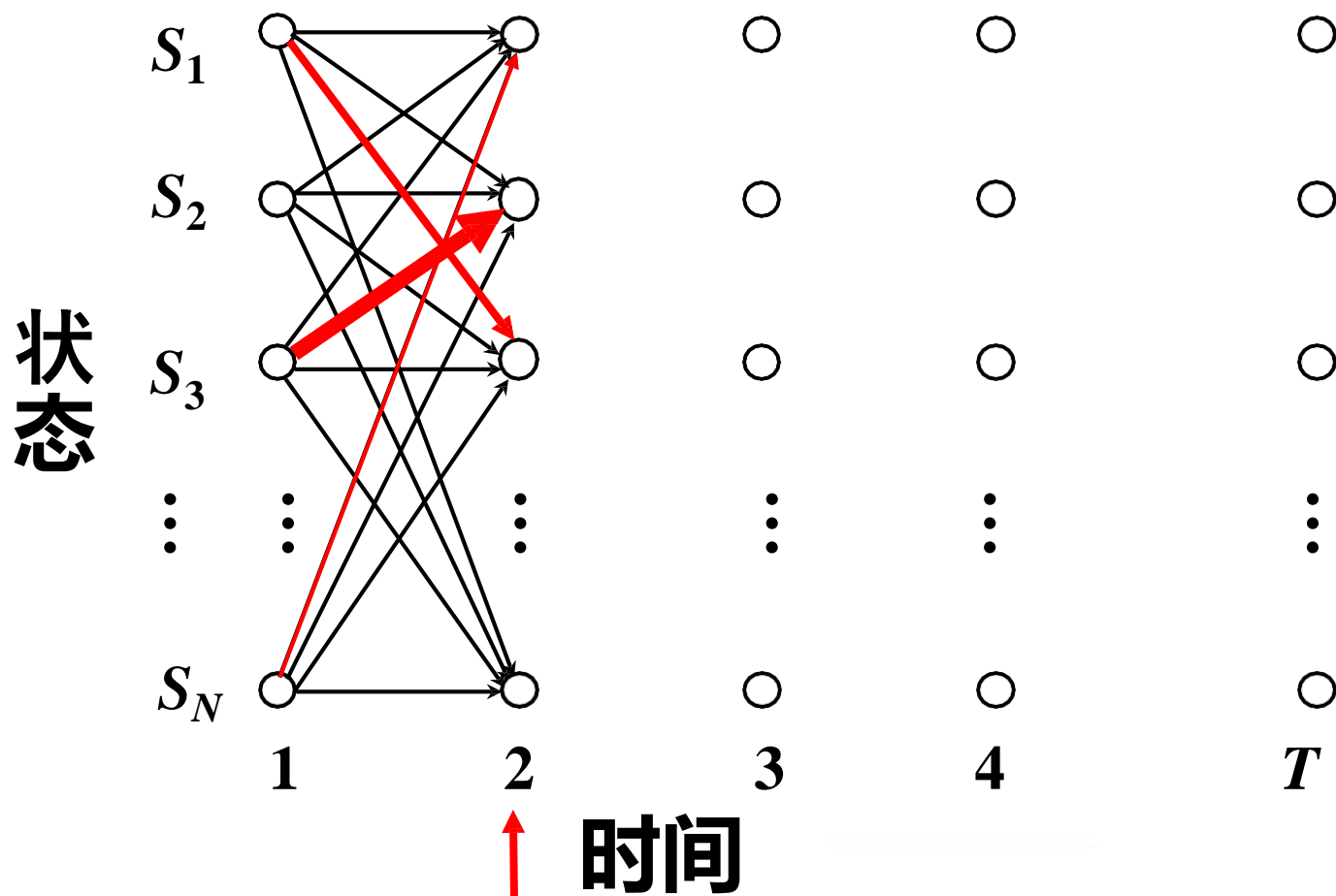
$$\hat{q}_t = \psi_{t+1}(\hat{q}_{t+1}), \quad t = T-1, T-2, \dots, 1$$

算法的时间复杂度： $O(N^2T)$

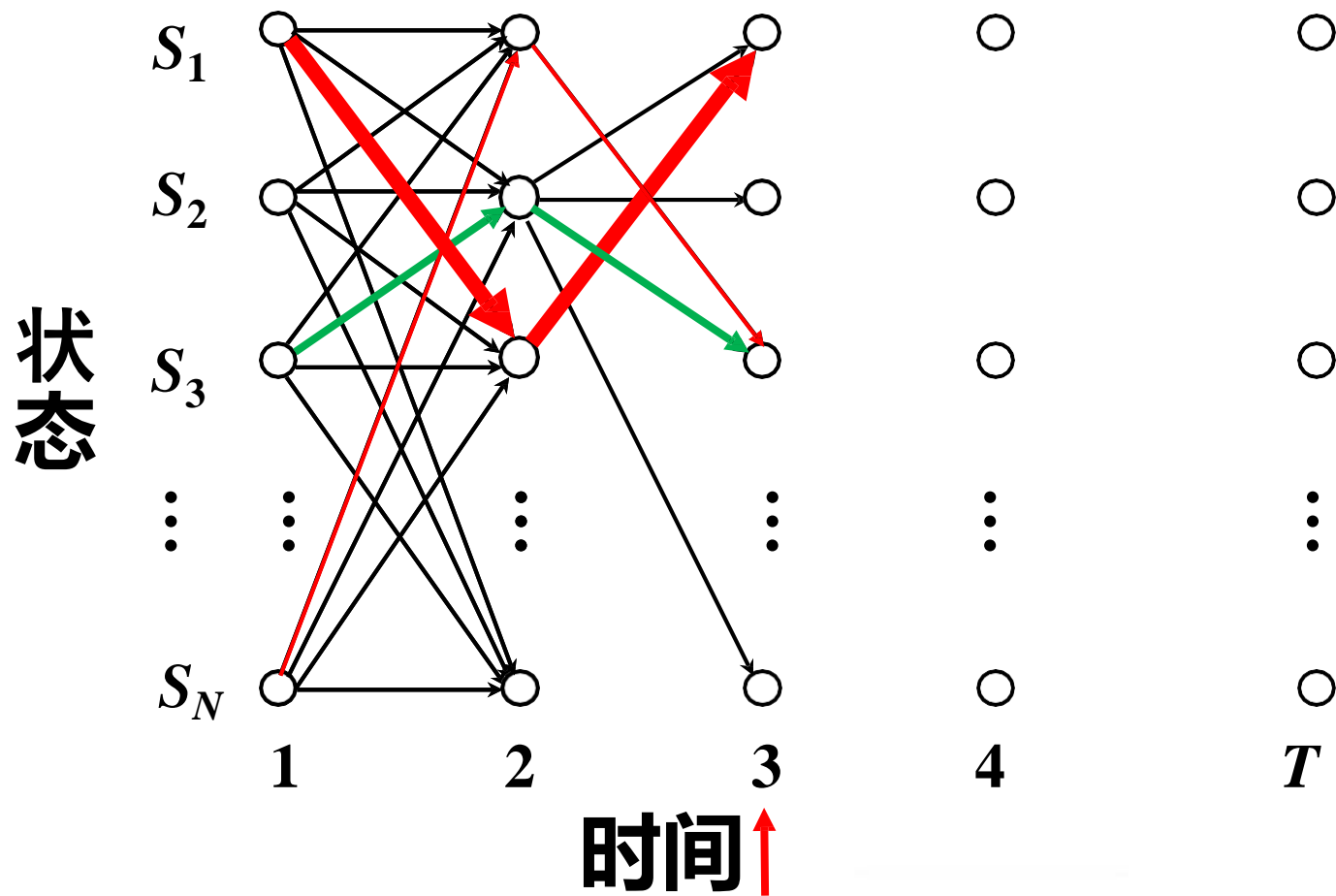
Viterbi 搜索算法



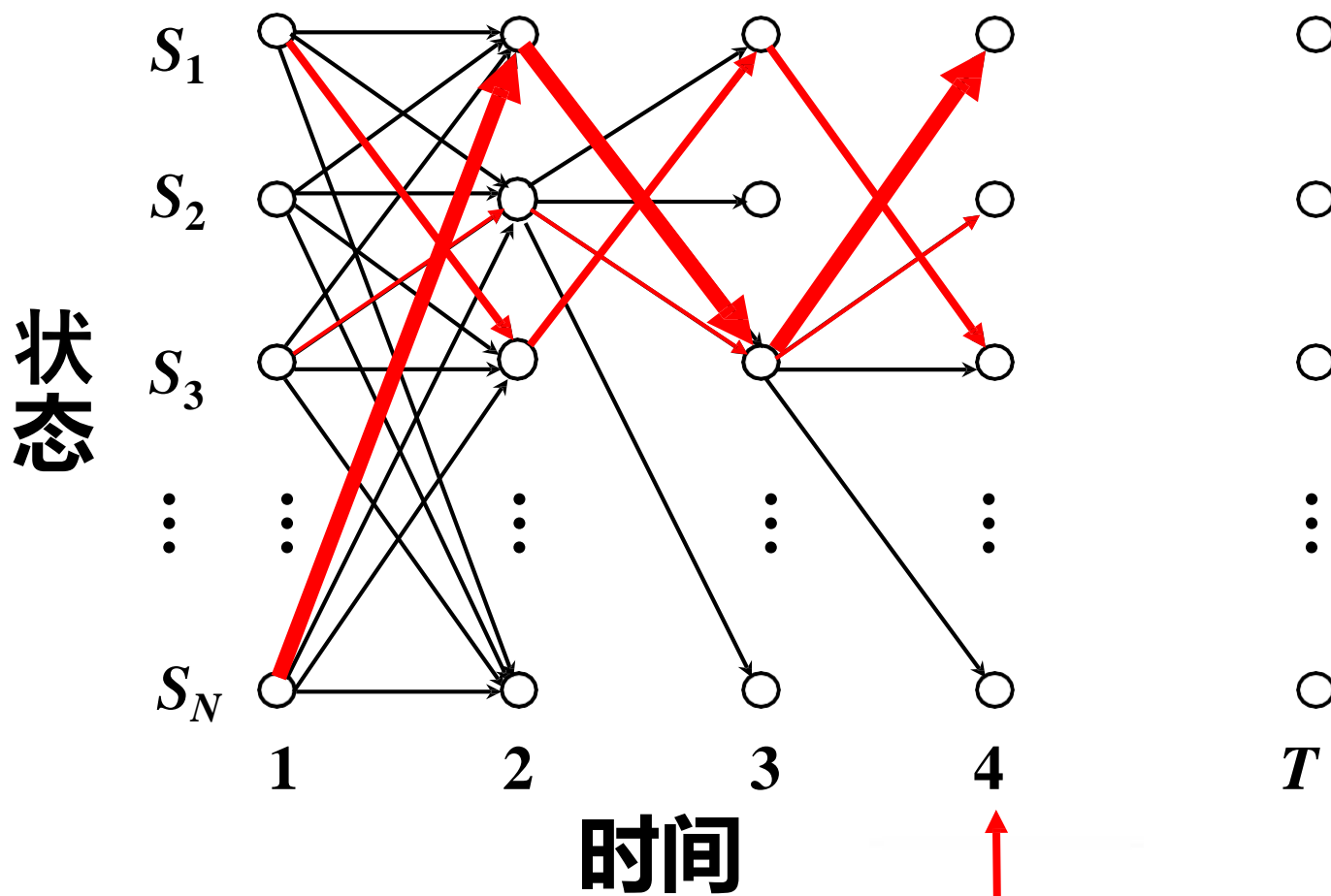
Viterbi 搜索算法



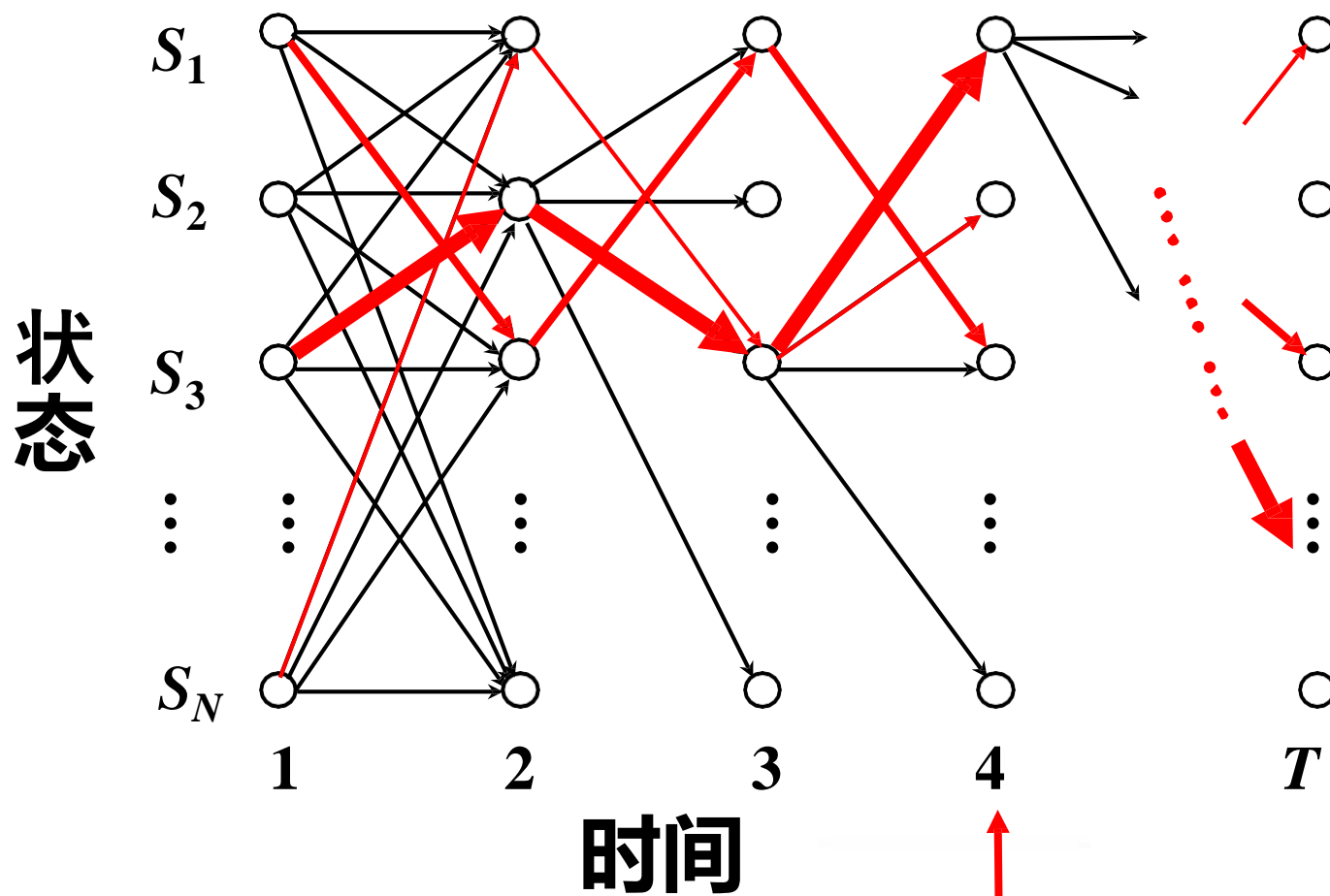
Viterbi 搜索算法



Viterbi 搜索算法



Viterbi 搜索算法



Lecture 5.6 参数学习

参数学习

◆问题3－模型参数学习

给定一个观察序列 $O = O_1 O_2 \dots O_T$ ，如何根据最大似然估计来求模型的参数值？或者说如何调节模型 μ 的参数，使得 $p(O|\mu)$ 最大？即估计模型中的 $\pi_i, a_{ij}, b_j(k)$ 使得观察序列 O 的概率 $p(O|\mu)$ 最大。

参数学习

如果产生观察序列 O 的状态 $Q = q_1q_2\dots q_T$ 已知(存在大量标注的样本), 可以用极大似然估计来计算 μ 的参数:

$$\pi_i = \delta(q_1, S_i)$$

$$\bar{a}_{ij} = \frac{Q \text{ 中从状态 } q_i \text{ 转移到 } q_j \text{ 的次数}}{Q \text{ 中所有从状态 } q_i \text{ 转移到另一状态(包括 } q_j \text{ 自身)的总数}}$$

... (24)

参数学习

类似地，

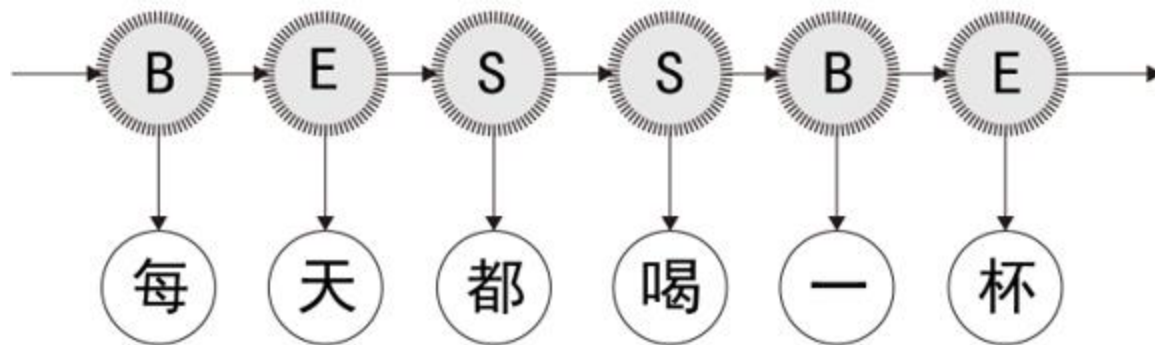
$$\bar{b}_j(k) = \frac{Q \text{ 中从状态 } q_j \text{ 输出符号 } v_k \text{ 的次数}}{Q \text{ 到达 } q_j \text{ 的总次数}}$$

... (6.25)

Lecture 5.7 HMM应用举例

HMM应用举例

- 将状态值集合 Q 置为 $\{B, E, M, S\}$ ，分别表示词的开始、结束、中间（begin、end、middle）及字符独立成词（single）；观测序列即为中文句子。
- 例如，“每天都喝一杯”通过HMM(Viterbi 算法)求解得到状态序列“B E S S B E”，则分词结果为“每天/都/喝/一杯”。



Thank you!

权小军 中山大学数据科学与计算机学院