

Fundamentals of Hidden Markov Model (HMM)

Ren-yuan Lyu
Dept of Computer Science & Information Engineering,
Chang Gung University,
Guei-shan, Taoyuan, Taiwan
rylyu@mail.cgu.edu.tw

Reference:

1. X. Huang, “Spoken Language Processing”, Chap 8
2. L. Rabiner, “Fundamentals of Speech Recognition”, Chap 6
3. HTK Book, <http://htk.eng.cam.ac.uk/>

Hidden Markov Model (HMM)

- A *Hidden Markov Model* is a (powerful) statistical model to describe very complex random processes (sequences) with time-varying characteristics by
 - Building parametric models,
 - Incorporating Dynamic Programming,
 - Providing a unified scheme for pattern segmentation and pattern classification (recognition).

Applications of HMM

- Automatic Speech Recognition
- Statistical Language Modeling
- Machine Translation
- Bioinformatics
-
-

Definition of HMM

- The “Hidden” State Sequence $S(t)$
 - A HMM can be viewed as a *doubly embedded random process* with a “*hidden*” random process, usually called the state process $S(t)$, which is not directly observable.
 - $S(t)$ is assumed a (1st order) *Markov* chain
- The “Observable” Data Sequence $O(t)$
 - The observable random process $O(t)$ in a HMM is probabilistically associated with the hidden random process $S(t)$.
 - $O(t)$ is assumed dependent only on the value of $S(t)$ at the same time index t .

Specifications of HMM

- The State sequence (process), $S(t)$,
 - with the state space Ω_S

$t \in Z$: a set of the time index,

for a finite duration random process, $Z = [1...T]$

$S \in \Omega_S$: a set of state value, also called a statespace,

usually being a finite set with N distinct values,

A typical example is $\Omega_S = \{1,2,3,...,N\}$

- The Observation sequence, $O(t)$,
 - with the observation space Ω_o

$t \in Z$: a set of the time index,

for a finite duration random process, $Z = [1...T]$

$O \in \Omega_o$: a set of observation value, also called a observation space,

It may be a finite set with M distinct values,

A typical example is $\Omega_o = \{1,2,3,...,M\}$

It could also be a continuous set of values.

e.g., $\Omega_o = R$: the real number

Or even a continuous vector space

e.g., $\Omega_o = R^D$

- The state transitional probability distribution $\mathbf{A}=[a_{ij}]$, and initial state probability distribution $\boldsymbol{\pi}=[\pi_i]$

By assumming $S(t)$ be a 1st - order Markov chain,
the state transitional probability, $P(S(t) = j | S(t-1) = i)$, is time - invariant,
and can statistically completely describe the random process $S(t)$,
with the initial state probability $P(S(1) = i)$.

$$a_{ij} \equiv P(S(t) = j | S(t-1) = i), \quad i, j \in \Omega_S = \{1, 2, 3, \dots, N\}$$

$$\pi_i \equiv P(S(1) = i)$$

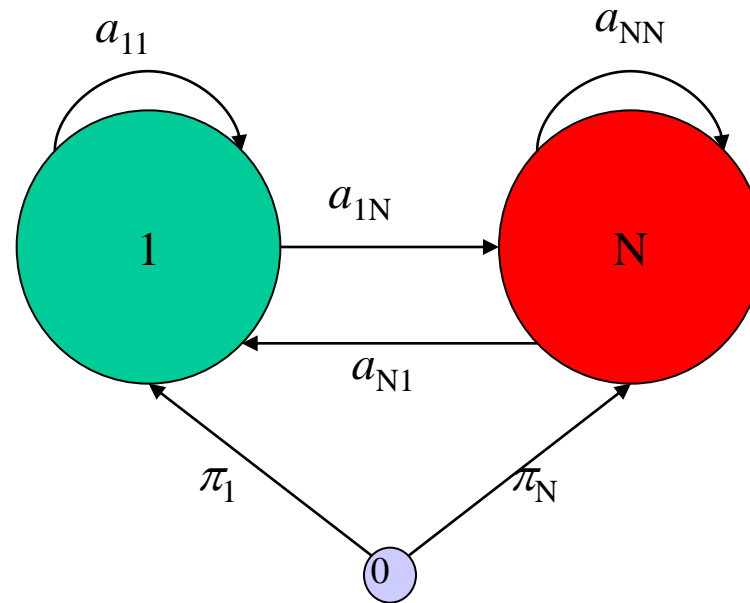
$$\underline{\underline{\mathbf{A}}} = [\underline{\underline{a}}_{ij}] = \begin{bmatrix} a_{11} & .. & a_{1N} \\ \vdots & & \vdots \\ a_{N1} & .. & a_{NN} \end{bmatrix},$$

$$\text{where } \sum_{j=1}^N a_{ij} = 1, \quad \forall i \in \Omega_S$$

$$\bar{\boldsymbol{\pi}} = [\pi_i] = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_N \end{bmatrix},$$

$$\text{where } \sum_{i=1}^N \pi_i = 1$$

- The state diagram



- The (state-dependent) observation probability distribution, $\underline{\mathbf{B}} = \{b_i(o)\}$

By assumming $O(t)$ be a state- dependent random process,
it is enough to specify $P(O(t) = o | S(t) = i)$ to completely describe $O(t)$,
as long as $S(t)$ is given.

Let

$$b_i(o) = P(O(t) = o | S(t) = i), \quad i \in \Omega_s = \{1, 2, 3, \dots, N\}, \\ o \in \Omega_o = \{1, 2, 3, \dots, M\},$$

$$\underline{\mathbf{B}} = \{b_i(o)\} = \left\{ \begin{array}{ccc} b_1(1) & \dots & b_1(M) \\ \vdots & & \vdots \\ b_N(1) & \dots & b_N(M) \end{array} \right\}, \quad \text{where } \sum_{o=1}^M b_i(o) = 1, \quad \forall i \in \Omega_s$$

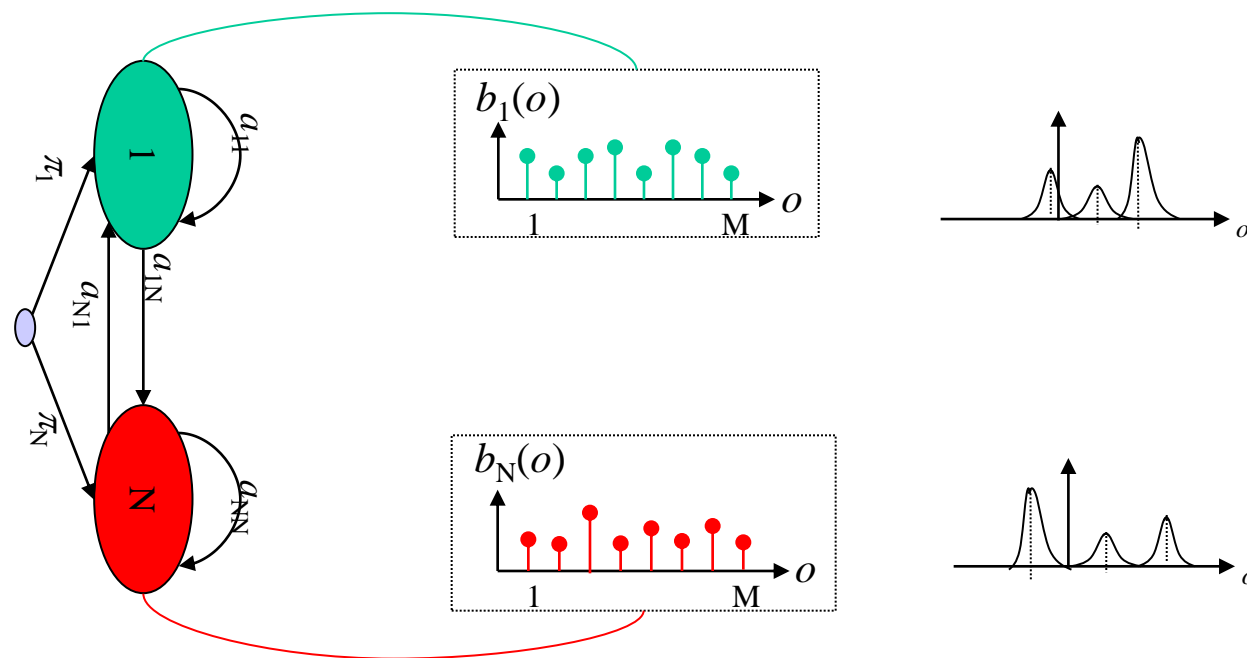
If $o \in \Omega_o = R$ (the set of real number),

then $b_i(o)$ should be a continuous probability density function,

e.g., the Gaussian distribution function

or the others, like Mixture Gaussian distribution function

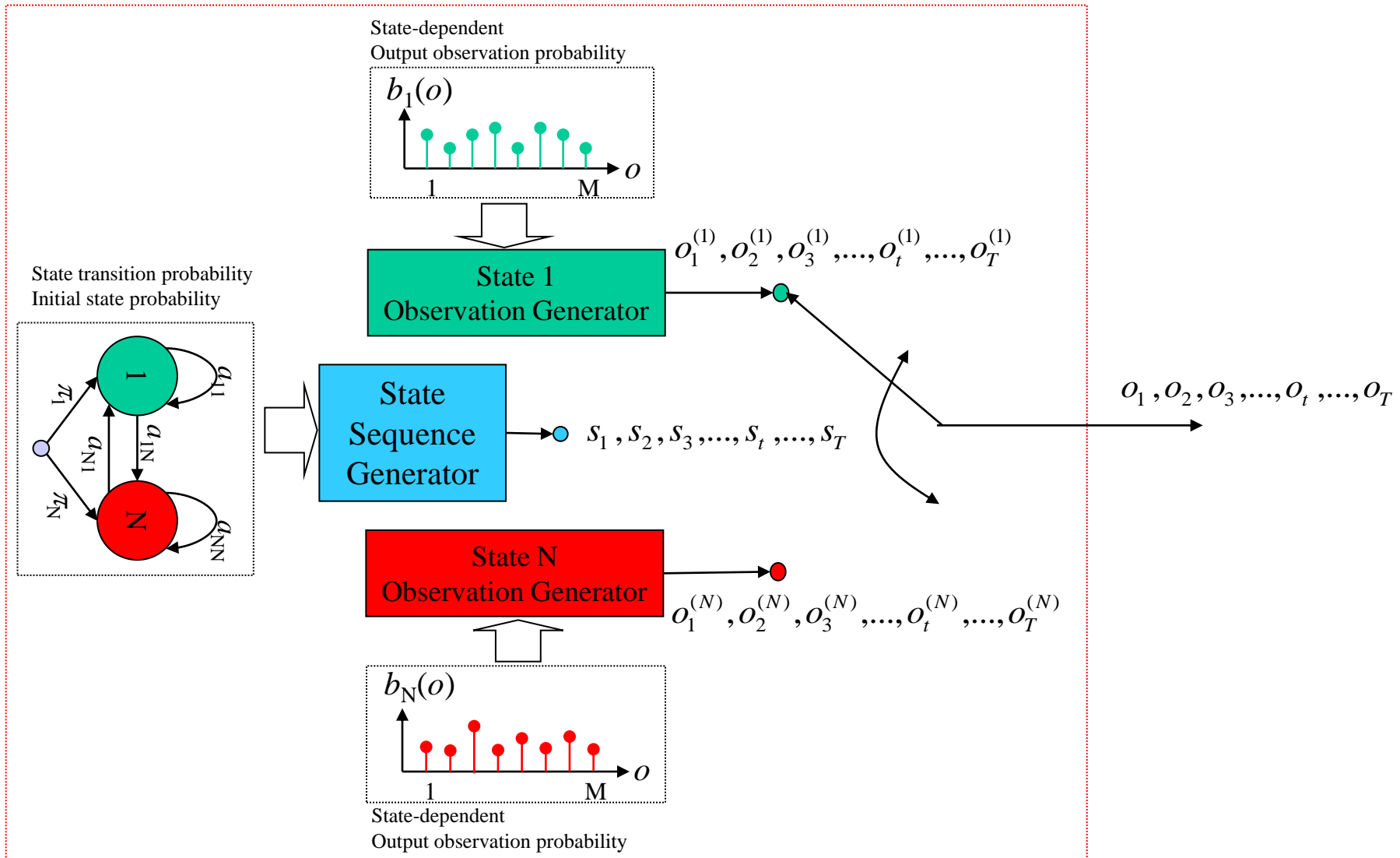
- The state-dependent observation distributions



Hidden Markov Machine

- To put all the elements of HMM together, we design a Hidden Markov Machine as follows:
 - There is a *state sequence generator*, which generates random sequence $S(t)$ under the control of HMM parameters $\{\pi, \mathbf{A}\}$.
 - There are N state-dependent *observation generators*, each of which generates observation sequence under the control of HMM parameter $\underline{\mathbf{B}} = \{b_i(o)\}$
 - The state value s_t , generated by the *state sequence generator* at time t , determines which observation generator's output, $o^{(i)}_t$, will be used as the output observation o_t at time t .

Hidden Markov Machine



3 Basic Problems of HMM

- The Evaluation Problem

$$P(o_1, o_2, o_3, \dots, o_t, \dots, o_T) = ?$$

- The Estimation (Learning) Problem

$$(\bar{\pi}^*, \underline{\underline{A}}^*, \underline{\underline{B}}^*) = \underset{\forall (\bar{\pi}, \underline{\underline{A}}, \underline{\underline{B}})}{\mathit{Argmax}} P(o_1, o_2, o_3, \dots, o_t, \dots, o_T \mid \bar{\pi}, \underline{\underline{A}}, \underline{\underline{B}}) = ?$$

- The Decoding Problem

$$(s_1^*, s_2^*, s_3^*, \dots, s_t^*, \dots, s_T^*) = \underset{\forall (s_1, s_2, s_3, \dots, s_t, \dots, s_T)}{\mathit{Argmax}} P(s_1, s_2, s_3, \dots, s_t, \dots, s_T \mid o_1, o_2, o_3, \dots, o_t, \dots, o_T) = ?$$

The Evaluation Problem of HMM

- Given a HMM, with parameters $\{\pi, \mathbf{A}, \underline{\mathbf{B}}\}$,

$$P(o_1, o_2, o_3, \dots, o_t, \dots, o_T) = ?$$

Markov Assumption

$$\underline{O} \equiv o_1, o_2, \dots, o_t, \dots, o_T$$

$$\underline{S} \equiv s_1, s_2, \dots, s_t, \dots, s_T$$

$$P(\underline{O}) = \sum_{\underline{S}} P(\underline{O}, \underline{S}) = \sum_{\underline{S}} P(\underline{O} | \underline{S}) \cdot P(\underline{S})$$

$$P(\underline{S}) = P(s_1, s_2, \dots, s_t, \dots, s_T)$$

$$= P(s_1)$$

$$\cdot P(s_2 | s_1)$$

$$\cdot P(s_3 | s_2 \text{ } \boxed{s_1})$$

....

$$\cdot P(s_t | s_{t-1} \text{ } \boxed{s_1, s_2, \dots, s_{t-2}, s_{t-1}})$$

....

$$\cdot P(s_T | s_{T-1} \text{ } \boxed{s_1, s_2, \dots, s_{T-2}, s_{T-1}})$$

$$P(\underline{O} | \underline{S}) = P(o_1, o_2, \dots, o_t, \dots, o_T | s_1, s_2, \dots, s_t, \dots, s_T)$$

$$= P(o_1 | s_1 \text{ } \boxed{s_2, \dots, s_t, \dots, s_T})$$

$$\cdot P(o_2 | \boxed{o_1, s_1} \text{ } s_2 \text{ } \boxed{s_3, \dots, s_t, \dots, s_T})$$

...

$$\cdot P(o_t | \boxed{o_1, \dots, o_{t-2}, o_{t-1}, s_1, s_2, \dots, s_t} \text{ } \boxed{s_{t+1}, \dots, s_T})$$

...

$$\cdot P(o_T | \boxed{o_1, \dots, o_{T-2}, o_{T-1}, s_1, s_2, \dots, s_t, \dots, s_T})$$

1. The current state depends only on the previous state
2. The observation depends only on the current state

$$P(\underline{S}) = P(s_1, s_2, \dots, s_t, \dots, s_T)$$

$$= P(s_1)$$

$$\cdot P(s_2 | s_1)$$

$$\cdot P(s_3 | s_2)$$

....

$$\cdot P(s_t | s_{t-1})$$

....

$$\cdot P(s_T | s_{T-1})$$

$$P(\underline{O} | \underline{S}) = P(o_1, o_2, \dots, o_t, \dots, o_T | s_1, s_2, \dots, s_t, \dots, s_T)$$

$$= P(o_1 | s_1)$$

$$\cdot P(o_2 | s_2)$$

...

$$\cdot P(o_t | s_t)$$

...

$$\cdot P(o_T | s_T)$$

Markov Assumption

1

2

Direct computation of $P(\underline{O})$

$$\underline{O} \equiv o_1, o_2, \dots, o_t, \dots, o_T$$

$$\underline{S} \equiv s_1, s_2, \dots, s_t, \dots, s_T$$

$$\begin{aligned} P(\underline{O}) &= \sum_{\underline{S}} P(\underline{O}, \underline{S}) = \sum_{\underline{S}} P(\underline{S}) \cdot P(\underline{O} | \underline{S}) \\ &= \sum_{\underline{S}} \left\{ P(s_1) \cdot P(s_2 | s_1) \cdots P(s_t | s_{t-1}) \cdots P(s_T | s_{T-1}) \right. \\ &\quad \left. \cdot P(o_1 | s_1) \cdot P(o_2 | s_2) \cdots P(o_t | s_t) \cdots P(o_T | s_T) \right\} \\ &\quad \begin{cases} P(s_1) = \pi_{s_1} = a_{s_0 s_1}, s_0 \equiv 0 = "B": \text{the beginning state at } t = 0 \\ P(s_t | s_{t-1}) = a_{s_{t-1} s_t} \\ P(o_t | s_t) = b_{s_t}(o_t) \end{cases} \\ &= \sum_{\underline{S}} \left\{ a_{s_0 s_1} \cdot a_{s_1 s_2} \cdots a_{s_{t-1} s_t} \cdots a_{s_{T-1} s_T} \right. \\ &\quad \left. \cdot b_{s_1}(o_1) \cdot b_{s_2}(o_2) \cdots b_{s_t}(o_t) \cdots b_{s_T}(o_T) \right\} \\ &= \sum_{s_T=1}^N \cdots \sum_{s_t=1}^N \cdots \sum_{s_2=1}^N \sum_{s_1=1}^N \left\{ a_{s_0 s_1} \cdot a_{s_1 s_2} \cdots a_{s_{t-1} s_t} \cdots a_{s_{T-1} s_T} \right. \\ &\quad \left. \cdot b_{s_1}(o_1) \cdot b_{s_2}(o_2) \cdots b_{s_t}(o_t) \cdots b_{s_T}(o_T) \right\} \end{aligned}$$

Totally, there are N^T possible sequences for \underline{S} ,
i.e., N^T terms of $\{ \dots \}$ need to be calculated.

Forward/Backward Algorithm

- The forward algorithm

$$P(\underline{O}) = \sum_{s_T=1}^N \cdots \sum_{s_t=1}^N \cdots \sum_{s_2=1}^N \sum_{s_1=1}^N a_{s_0 s_1} \cdot b_{s_1}(o_1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \cdots a_{s_{t-1} s_t} \cdot b_{s_t}(o_t) \cdots a_{s_{T-1} s_T} \cdot b_{s_T}(o_T)$$

$\alpha_{s_1}(1)$
 $\alpha_{s_2}(2)$
 $\alpha_{s_t}(t)$
 $\alpha_{s_T}(T)$

$$P(\underline{O}) = \sum_{s_T=1}^N \dots \sum_{s_t=1}^N \dots \sum_{s_2=1}^N \sum_{s_1=1}^N a_{s_0 s_1} \cdot b_{s_1}(o_1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \cdots a_{s_{t-1} s_t} \cdot b_{s_t}(o_t) \cdots a_{s_{T-1} s_T} \cdot b_{s_T}(o_T)$$

$$\alpha_{s_1}(1) = a_{s_0 s_1} \cdot b_{s_1}(o_1) \quad , s_1 \in [1..N], s_0 \equiv 0, a_{s_0 s_1} \equiv \pi_{s_1}$$

$$\alpha_{s_2}(2) = \sum_{s_1=1}^N \alpha_{s_1}(1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \quad , s_2 \in [1..N]$$

:

$$\alpha_{s_t}(t) = \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1} s_t} \cdot b_{s_t}(o_t) \quad , s_t \in [1..N]$$

:

$$\alpha_{s_T}(T) = \sum_{s_{T-1}=1}^N \alpha_{s_{T-1}}(T-1) \cdot a_{s_{T-1} s_T} \cdot b_{s_T}(o_T) \quad , s_T \in [1..N]$$

$$\sum_{s_T=1}^N \alpha_{s_T}(T) ===== P(\underline{O})$$

The computation complexity reduces to be **T·N²**

- The meaning of the forward probability

$$\begin{aligned}\alpha_{s_1}(1) &= a_{s_0 s_1} \cdot b_{s_1}(o_1) \\ &= P(s_1 | s_0) \cdot P(o_1 | s_1) \\ &= P(s_1) \cdot P(o_1 | s_1) \dots \dots \dots \{ \because s_0 \equiv 0 \text{ is a determined value} \} \\ &= P(o_1, s_1)\end{aligned}$$

$$\begin{aligned}\alpha_{s_2}(2) &= \sum_{s_1=1}^N \alpha_{s_1}(1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \\ &= \sum_{s_1=1}^N P(o_1, s_1) \cdot P(s_2 | s_1) \cdot P(o_2 | s_2) \\ &\quad \left\{ \begin{array}{l} \because \\ P(o_1, o_2, s_1, s_2) = P(o_1, s_1) \cdot P(o_2, s_2 | o_1, s_1) \\ P(o_2, s_2 | o_1, s_1) = P(s_2 | o_1, s_1) \cdot P(o_2 | o_1, s_1, s_2) \\ P(s_2 | o_1, s_1) = P(s_2 | s_1) \dots \text{this is Markov assumption} \\ P(o_2 | o_1, s_1, s_2) = P(o_2 | s_2) \dots \text{the assumption} \\ \quad \text{that observation depends only on current state} \end{array} \right. \\ &= \sum_{s_1=1}^N P(o_1, o_2, s_1, s_2) \\ &= P(o_1, o_2, s_2)\end{aligned}$$

$$\alpha_{s_t}(t) = \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) = P(o_1, o_2, \dots, o_t, s_t)$$

$$\alpha_{s_T}(T) = \sum_{s_{T-1}=1}^N \alpha_{s_{T-1}}(T-1) \cdot a_{s_{T-1}s_T} \cdot b_{s_T}(o_T) = P(o_1, o_2, \dots, o_T, s_T)$$

$$\sum_{s_T=1}^N \alpha_{s_T}(T) = \sum_{s_T=1}^N P(o_1, o_2, \dots, o_T, s_T) = P(o_1, o_2, \dots, o_T) = P(\underline{O})$$

- The backward algorithm

$$P(\underline{O}) = \sum_{s_T=1}^N \cdots \sum_{s_t=1}^N \cdots \sum_{s_2=1}^N \sum_{s_1=1}^N a_{s_0 s_1} \cdot b_{s_1}(o_1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \cdots a_{s_{t-1} s_t} \cdot b_{s_t}(o_t) \cdots a_{s_{T-1} s_T} \cdot b_{s_T}(o_T)$$

$$= \sum_{s_1=1}^N \sum_{s_2=1}^N \cdots \sum_{s_t=1}^N \sum_{s_{t+1}=1}^N \cdots \sum_{s_{T-1}=1}^N \sum_{s_T=1}^N \underbrace{b_{s_T}(o_T) \cdot a_{s_{T-1} s_T}}_{\beta_{s_{T-1}}(T-1)} \cdot \underbrace{b_{s_{T-1}}(o_{T-1}) \cdot a_{s_{T-2} s_{T-1}} \cdots b_{s_{t+1}}(o_{t+1}) \cdot a_{s_t s_{t+1}}}_{\beta_{s_{T-2}}(T-2)} \cdot \underbrace{b_{s_t}(o_t) \cdot a_{s_{t-1} s_t} \cdots b_{s_2}(o_2) \cdot a_{s_1 s_2}}_{\beta_{s_t}(t)} \cdot \underbrace{b_{s_1}(o_1) \cdot a_{s_0 s_1}}_{\beta_{s_{t-1}}(t-1)} \cdot \underbrace{a_{s_0 s_1}}_{\beta_{s_1}(1)} \cdot \underbrace{a_{s_0 s_1}}_{\beta_{s_0}(0)}$$

$$\beta_{s_{T-1}}(T-1) = \sum_{s_T=1}^N b_{s_T}(o_T) \cdot a_{s_{T-1}s_T}$$

$$\beta_{s_{T-2}}(T-2) = \sum_{s_{T-1}=1}^N \beta_{s_{T-1}}(T-1) \cdot b_{s_{T-1}}(o_{T-1}) \cdot a_{s_{T-2}s_{T-1}}$$

:

$$\beta_{s_t}(t) = \sum_{s_{t+1}=1}^N \beta_{s_{t+1}}(t+1) \cdot b_{s_{t+1}}(o_{t+1}) \cdot a_{s_t s_{t+1}}$$

:

$$\beta_{s_1}(1) = \sum_{s_2=1}^N \beta_{s_2}(2) \cdot b_{s_2}(o_2) \cdot a_{s_1 s_2}$$

$$\beta_{s_0}(0) = \sum_{s_1=1}^N \beta_{s_1}(1) \cdot b_{s_1}(o_1) \cdot a_{s_0 s_1} ===== P(\underline{O})$$

- The meaning of the backward probability

$$\begin{aligned}\beta_{s_{T-1}}(T-1) &= \sum_{s_T=1}^N b_{s_T}(o_T) \cdot a_{s_{T-1}s_T} \\&= \sum_{s_T=1}^N P(o_T | s_T) \cdot P(s_T | s_{T-1}) \\&= \sum_{s_T=1}^N P(o_T | s_T, s_{T-1}) \cdot P(s_T | s_{T-1}) \\&= \sum_{s_T=1}^N P(o_T, s_T | s_{T-1}) \\&= P(o_T | s_{T-1}) \\ \beta_{s_{T-2}}(T-2) &= \sum_{s_{T-1}=1}^N \beta_{s_{T-1}}(T-1) \cdot b_{s_{T-1}}(o_{T-1}) \cdot a_{s_{T-2}s_{T-1}} = P(o_{T-1}, o_T | s_{T-2}) \\&\vdots\end{aligned}$$

$$\beta_{s_t}(t) = \sum_{s_{t+1}=1}^N \beta_{s_{t+1}}(t+1) \cdot b_{s_{t+1}}(o_{t+1}) \cdot a_{s_t s_{t+1}} = P(o_{t+1}, o_{t+2}, \dots, o_{T-1}, o_T \mid s_t)$$

:

$$\beta_{s_1}(1) = \sum_{s_2=1}^N \beta_{s_2}(2) \cdot b_{s_2}(o_2) \cdot a_{s_1 s_2} = P(o_2, o_3, \dots, o_{T-1}, o_T \mid s_1)$$

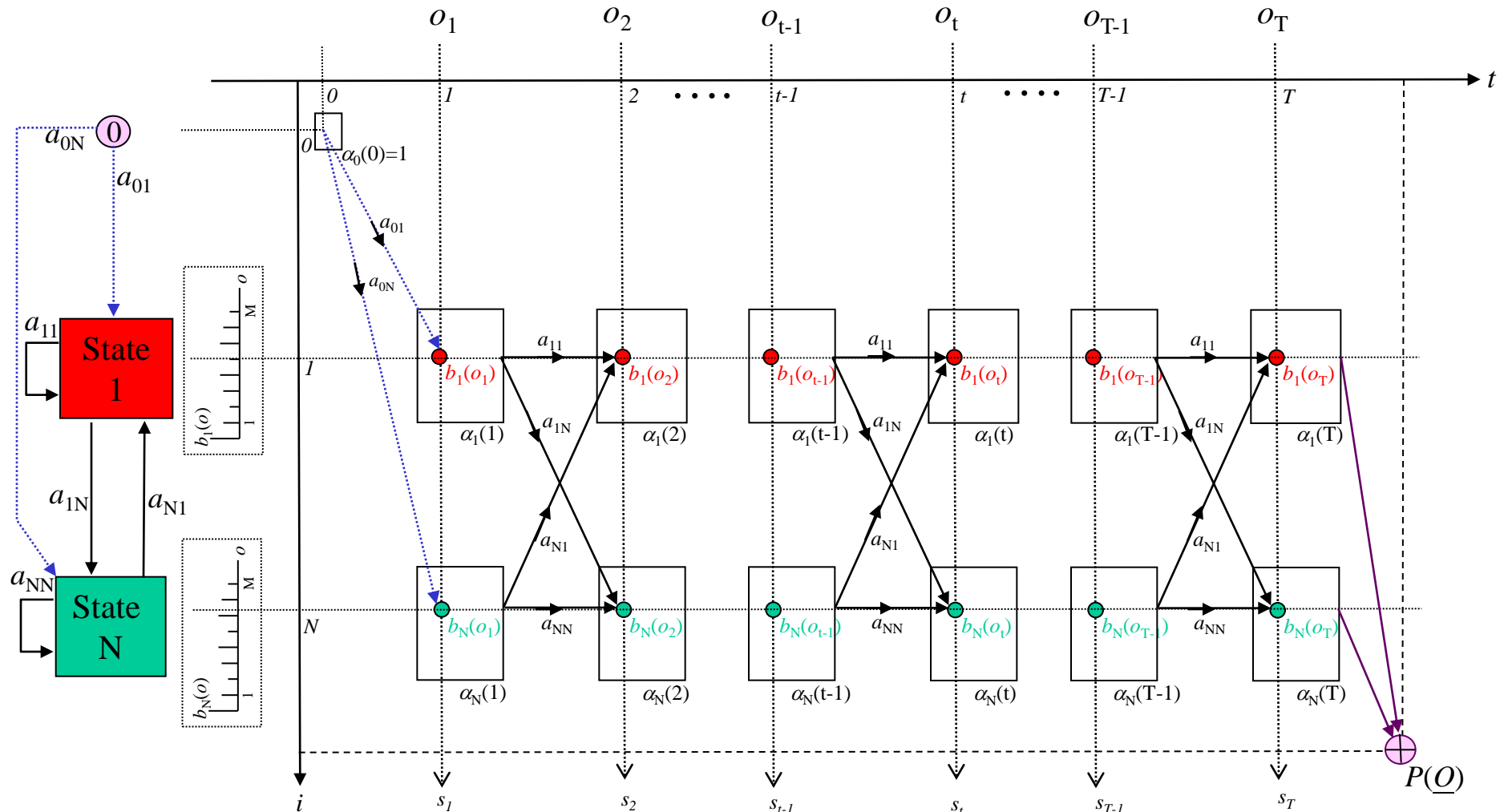
$$\beta_{s_0}(0) = \sum_{s_1=1}^N \beta_{s_1}(1) \cdot b_{s_1}(o_1) \cdot a_{s_0 s_1} = P(o_1, o_2, o_3, \dots, o_{T-1}, o_T \mid s_0)$$

$$\begin{cases} \because s_0 (\equiv 0), \text{ is determined,} \\ \text{and independent of } o_1, \dots, o_T \end{cases}$$

$$= P(o_1, o_2, o_3, \dots, o_{T-1}, o_T)$$

$$= P(\underline{O})$$

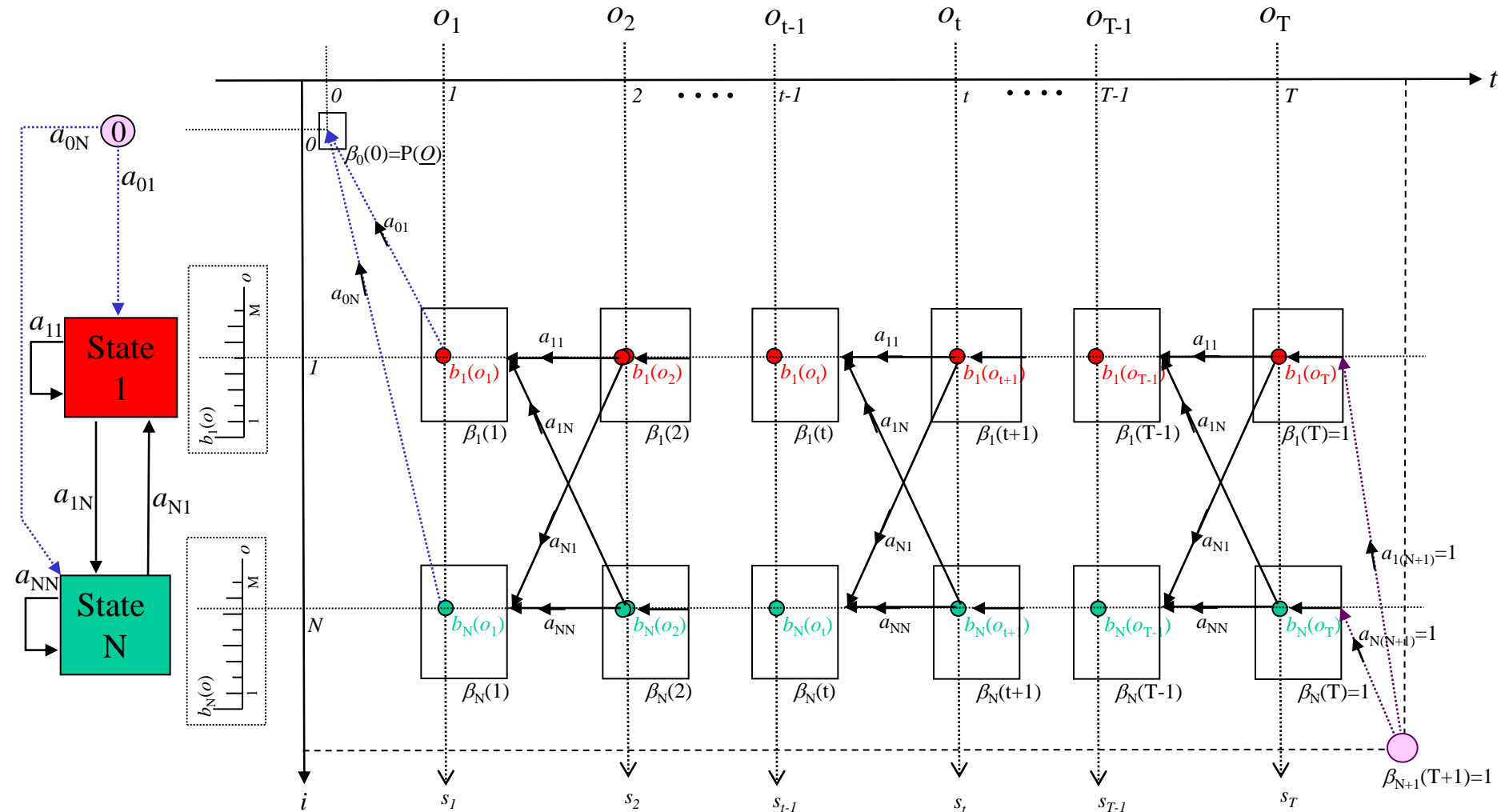
Trellis view of the *forward* algorithm



$$\alpha_{s_t}(t) = \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \quad , s_t \in [1..N]$$

Trellis view of the *backward* algorithm

$$\beta_{s_t}(t) = \sum_{s_{t+1}=1}^N \beta_{s_{t+1}}(t+1) \cdot b_{s_{t+1}}(o_{t+1}) \cdot a_{s_t s_{t+1}}$$



The Estimation (Learning) Problem of HMM

- Given a HMM, with parameters $\{\pi, \mathbf{A}, \underline{\mathbf{B}}\}$, and observation sequence(s)

$$O_1, O_2, O_3, \dots, O_t, \dots, O_T$$

How can we find “better” or even the “best” or “most likely” (new) parameters $\{\pi^*, \mathbf{A}^*, \underline{\mathbf{B}}^*\}$, to help the HM**Machine** generate such an observation sequence with larger or even the maximal probability?
This problem can be reformulate as follows:

$$(\bar{\pi}^*, \underline{\underline{A}}^*, \underline{\underline{B}}^*) = \underset{\forall(\bar{\pi}, \underline{\underline{A}}, \underline{\underline{B}})}{\mathit{Argmax}} P(o_1, o_2, o_3, \dots, o_t, \dots, o_T \mid \bar{\pi}, \underline{\underline{A}}, \underline{\underline{B}}) = ?$$

The solution to the HMM Learning problem

$$\underline{S} \equiv s_1, s_2, \dots, s_{t-1}, s_t, \dots, s_T$$

$$\underline{O} \equiv o_1, o_2, \dots, o_{t-1}, o_t, \dots, o_T$$

$$\text{where } \begin{cases} t \in [1..T] \\ s_t \in [1..N] \\ o_t \in [1..M] \\ \text{when } t = 0, s_0 \equiv 0 \equiv \text{'Begin state'} \end{cases}$$

$$a_{s_{t-1}s_t} \equiv P(s_t | s_{t-1})$$

$$b_{s_t}(o_t) \equiv P(o_t | s_t)$$

$$\alpha_{s_t}(t) \equiv P(o_1, o_2, \dots, o_t, s_t) = \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t)$$

$$\beta_{s_t}(t) \equiv P(o_{t+1}, o_{t+2}, \dots, o_T | s_t) = \sum_{s_{t+1}=1}^N a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{s_{t+1}}(t+1)$$

$$\alpha_{s_0}(0) \equiv 1$$

$$\beta_{s_T}(T) \equiv 1, \begin{cases} \text{in Huang's Textbook, } \beta_{s_T}(T) \equiv 1/N; \\ \text{in HTK, } \beta_{s_T}(T) \equiv a_{s_T(N+1)} \\ \quad = P(S(T+1) = N+1 \equiv \text{'Exit state'} | S(T) = s_T) \end{cases}$$

The following are new : (Prove them)

$$\eta_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t, \underline{O}) = \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \cdot \beta_{s_t}(t)$$

$$\eta_{s_t}(t) \equiv P(s_t, \underline{O}) = \sum_{s_{t-1}=1}^N \eta_{s_{t-1}s_t}(t) = \alpha_{s_t}(t) \cdot \beta_{s_t}(t)$$

$$\eta_{s_{t-1}}(t-1) \equiv P(s_{t-1}, \underline{O}) = \begin{cases} \sum_{s_t=1}^N \eta_{s_{t-1}s_t}(t) \\ \sum_{s_{t-2}=1}^N \eta_{s_{t-2}s_{t-1}}(t-1) \end{cases} = \alpha_{s_{t-1}}(t-1) \cdot \beta_{s_{t-1}}(t-1)$$

$$P(\underline{O}) = \sum_{s_T=1}^N \alpha_{s_T}(T) = \beta_{s_0}(0) = \sum_{s_1=1}^N \eta_{s_1}(1)$$

$$\gamma_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t | \underline{O}) = \frac{P(s_{t-1}, s_t, \underline{O})}{P(\underline{O})} = \frac{\eta_{s_{t-1}s_t}(t)}{P(\underline{O})}$$

$$\gamma_{s_{t-1}}(t-1) \equiv P(s_{t-1} | \underline{O}) = \frac{P(s_{t-1}, \underline{O})}{P(\underline{O})} = \frac{\eta_{s_{t-1}}(t-1)}{P(\underline{O})}$$

$$\gamma_{s_t}(t) \equiv P(s_t | \underline{O}) = \frac{P(s_t, \underline{O})}{P(\underline{O})} = \frac{\eta_{s_t}(t)}{P(\underline{O})}$$

$$\hat{a}_{s_{t-1}s_t} | \underline{O} \equiv (\text{Estimation of } a_{s_{t-1}s_t}, \text{ given } \underline{O})$$

$$\begin{aligned} &= \frac{\langle \gamma_{s_{t-1}s_t}(t) \rangle}{\langle \gamma_{s_{t-1}}(t-1) \rangle} = \frac{\frac{1}{T} \sum_{t=1}^T \gamma_{s_{t-1}s_t}(t)}{\frac{1}{T} \sum_{t=1}^T \gamma_{s_{t-1}}(t-1)} = \frac{\sum_{t=1}^T \gamma_{s_{t-1}s_t}(t)}{\sum_{t=1}^T \gamma_{s_{t-1}}(t-1)} \\ &= \frac{\sum_{t=1}^T \eta_{s_{t-1}s_t}(t)}{\sum_{t=1}^T \eta_{s_{t-1}}(t-1)} \end{aligned}$$

$$\hat{b}_{s_t}(o) | \underline{O} \equiv (\text{Estimation of } b_{s_t}(o), \text{ given } \underline{O})$$

$$\begin{aligned} &= \frac{\langle \gamma_{s_t}(t) \cdot \delta(o_t - o) \rangle}{\langle \gamma_{s_t}(t) \rangle} = \frac{\frac{1}{T} \sum_{t=1}^T \gamma_{s_t}(t) \cdot \delta(o_t - o)}{\frac{1}{T} \sum_{t=1}^T \gamma_{s_t}(t)} = \frac{\sum_{t=1}^T \gamma_{s_t}(t) \cdot \delta(o_t - o)}{\sum_{t=1}^T \gamma_{s_t}(t)} \\ &= \frac{\sum_{t=1}^T \eta_{s_t}(t) \cdot \delta(o_t - o)}{\sum_{t=1}^T \eta_{s_t}(t)} \end{aligned}$$

$$\text{where } \begin{cases} \langle f(t) \rangle \text{ denotes the time average of } f(t) \\ \delta(o_t - o) = \begin{cases} 1, & \text{if } o_t = o \\ 0, & \text{otherwise} \end{cases} \end{cases}$$

It can be shown that

$$P(\underline{O} | \hat{\pi}, \hat{\underline{A}}, \hat{\underline{B}}) \geq P(\underline{O} | \bar{\pi}, \underline{A}, \underline{B})$$

Ideally, by several iterations,

$$(\hat{\pi}, \hat{\underline{A}}, \hat{\underline{B}}) \rightarrow (\bar{\pi}^*, \underline{A}^*, \underline{B}^*) = \underset{\forall (\bar{\pi}, \underline{A}, \underline{B})}{\text{Argmax}} P(\underline{O} | \bar{\pi}, \underline{A}, \underline{B})$$

Underflow problem of Kernel computing in HMM

$$\begin{aligned}\alpha_{s_t}(t) &\equiv P(o_1, o_2, \dots, o_t, s_t) &= \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \\ \beta_{s_t}(t) &\equiv P(o_{t+1}, o_{t+2}, \dots, o_{T-1}, o_T | s_t) &= \sum_{s_{t+1}=1}^N a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{s_{t+1}}(t+1) \\ \eta_{s_{t-1}s_t}(t) &\equiv P(s_{t-1}, s_t, \underline{Q}) &= \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \cdot \beta_{s_t}(t) \\ \eta_{s_t}(t) &\equiv P(s_t, \underline{Q}) = \sum_{s_{t-1}=1}^N \eta_{s_{t-1}s_t}(t) &= \alpha_{s_t}(t) \cdot \beta_{s_t}(t) \\ P(\underline{Q}) &= \sum_{s_T=1}^N \alpha_{s_T}(T) = \beta_{s_0}(0) = \sum_{s_1=1}^N \eta_{s_1}(1)\end{aligned}$$

when t is large, the values of all the above functions go to 0 (Underflow!)

To avoid the underflow problem, taking \log for all the above functions.

Log Summation

Given x_i be very small positive numbers ($\approx 10^{-10000}$)

$$y_i = \log(x_i)$$

will behave much better when processed (or stored) in a fixed - precision computer

It is easy when $f(x_1, \dots, x_i, \dots, x_N) = \prod_{i=1}^N x_i$ need to be computed,

because

$$\log(f(x_1, \dots, x_i, \dots, x_N)) = \log\left(\prod_{i=1}^N x_i\right) = \sum_{i=1}^N \log(x_i) = \sum_{i=1}^N y_i$$

However, it is not trivial when $h(x_1, \dots, x_i, \dots, x_N) = \sum_{i=1}^N x_i$ need to be computed.

Question :

$$h(x_1, \dots, x_i, \dots, x_N) = \sum_{i=1}^N x_i$$

$$\log(h(x_1, \dots, x_i, \dots, x_N)) = \log\left(\sum_{i=1}^N x_i\right) = ?$$

The answer must be in terms of y_i

[Sol] :

$$\begin{aligned} \log(h(x_1, \dots, x_i, \dots, x_N)) &= \log\left(\sum_{i=1}^N x_i\right) \\ &= \log\left(\sum_{i=1}^N e^{\log(x_i)}\right) = \log\left(\sum_{i=1}^N e^{\log(x_i)} \frac{e^{\log(x_{i^*})}}{e^{\log(x_{i^*})}}\right) \\ &= \log\left(\sum_{i=1}^N e^{\log(x_{i^*})} \frac{e^{\log(x_i)}}{e^{\log(x_{i^*})}}\right) = \log\left(e^{\log(x_{i^*})} \cdot \sum_{i=1}^N e^{(\log(x_i) - \log(x_{i^*}))}\right) \\ &= \log\left(e^{\log(x_{i^*})}\right) + \log\left(\sum_{i=1}^N e^{(\log(x_i) - \log(x_{i^*}))}\right) \\ &= \log(x_{i^*}) + \log\left(\sum_{i=1}^N e^{(\log(x_i) - \log(x_{i^*}))}\right) \\ &= y_{i^*} + \log\left(\sum_{i=1}^N e^{(y_i - y_{i^*})}\right) \end{aligned}$$

where $i^* = \arg \max_{i \in [1..N]} \{x_i\} = \arg \max_{i \in [1..N]} \{\log(x_i)\} = \arg \max_{i \in [1..N]} \{y_i\}$

$$\begin{aligned}
\alpha_{s_t}(t) &= \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \\
\log(\alpha_{s_t}(t)) &= \log\left(\sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t)\right) \\
&= \log\left(\sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t}\right) + \log(b_{s_t}(o_t)) \\
&= \log\left(\alpha_{s_{t-1}^*}^*(t-1) \cdot a_{s_{t-1}^*s_t}^*\right) + \log\left(\sum_{s_{t-1}=1}^N e^{\left(\log(\alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t}) - \log(\alpha_{s_{t-1}^*}^*(t-1) \cdot a_{s_{t-1}^*s_t}^*)\right)}\right) + \log(b_{s_t}(o_t)) \\
&= \log(\alpha_{s_{t-1}^*}^*(t-1)) \\
&\quad + \log(a_{s_{t-1}^*s_t}^*) \\
&\quad + \log\left(\sum_{s_{t-1}=1}^N e^{\left(\log(\alpha_{s_{t-1}}(t-1)) + \log(a_{s_{t-1}s_t}) - \log(\alpha_{s_{t-1}^*}^*(t-1)) - \log(a_{s_{t-1}^*s_t}^*)\right)}\right) \\
&\quad + \log(b_{s_t}(o_t))
\end{aligned}$$

where $s_{t-1}^* = \operatorname{argmax}_{s_{t-1} \in [1..N]} \{\log(\alpha_{s_{t-1}}(t-1)) + \log(a_{s_{t-1}s_t})\}$

Log version of Kernel computing in HMM

$$\begin{aligned}
 \alpha_{s_t}(t) &= \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \\
 \log(\alpha_{s_t}(t)) &= \log(\alpha_{s_{t-1}^*}(t-1)) \\
 &\quad + \log(a_{s_{t-1}^*s_t}) \\
 &\quad + \log\left(\sum_{s_{t-1}=1}^N e^{\left(\log(\alpha_{s_{t-1}}(t-1)) + \log(a_{s_{t-1}s_t}) - \log(\alpha_{s_{t-1}^*}(t-1)) - \log(a_{s_{t-1}^*s_t})\right)}\right) \\
 &\quad + \log(b_{s_t}(o_t)) \\
 \text{where } s_{t-1}^* &= \operatorname{argmax}_{s_{t-1} \in [1..N]} \{\log(\alpha_{s_{t-1}}(t-1)) + \log(a_{s_{t-1}s_t})\}
 \end{aligned}$$

$$\begin{aligned}
 \beta_{s_t}(t) &= \sum_{s_{t+1}=1}^N a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{s_{t+1}}(t+1) \\
 \log(\beta_{s_t}(t)) &= ? \\
 \eta_{s_{t-1}s_t}(t) &= \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \cdot \beta_{s_t}(t) \\
 \log(\eta_{s_{t-1}s_t}(t)) &= ? \\
 \eta_{s_t}(t) &= \alpha_{s_t}(t) \cdot \beta_{s_t}(t) \\
 \log(\eta_{s_t}(t)) &= ? \\
 P(\underline{O}) &= \sum_{s_T=1}^N \alpha_{s_T}(T) = \beta_{s_0}(0) = \sum_{s_1=1}^N \eta_{s_1}(1) \\
 \log(P(\underline{O})) &= ?
 \end{aligned}$$

$$\gamma_{s_{t-1}s_t}(t) = \frac{\eta_{s_{t-1}s_t}(t)}{P(\underline{Q})} = \frac{e^{\log(\eta_{s_{t-1}s_t}(t))}}{e^{\log(P(\underline{Q}))}} = e^{(\log(\eta_{s_{t-1}s_t}(t)) - \log(P(\underline{Q})))}$$

$$\gamma_{s_{t-1}}(t-1) = \frac{\eta_{s_{t-1}}(t-1)}{P(\underline{Q})} = e^{(\log(\eta_{s_{t-1}}(t-1)) - \log(P(\underline{Q})))}$$

$$\gamma_{s_t}(t) = \frac{\eta_{s_t}(t)}{P(\underline{Q})} = e^{(\log(\eta_{s_t}(t)) - \log(P(\underline{Q})))}$$

$\hat{a}_{s_{t-1}s_t} |_{\underline{Q}} \equiv$ (Estimation of $a_{s_{t-1}s_t}$, given \underline{Q})

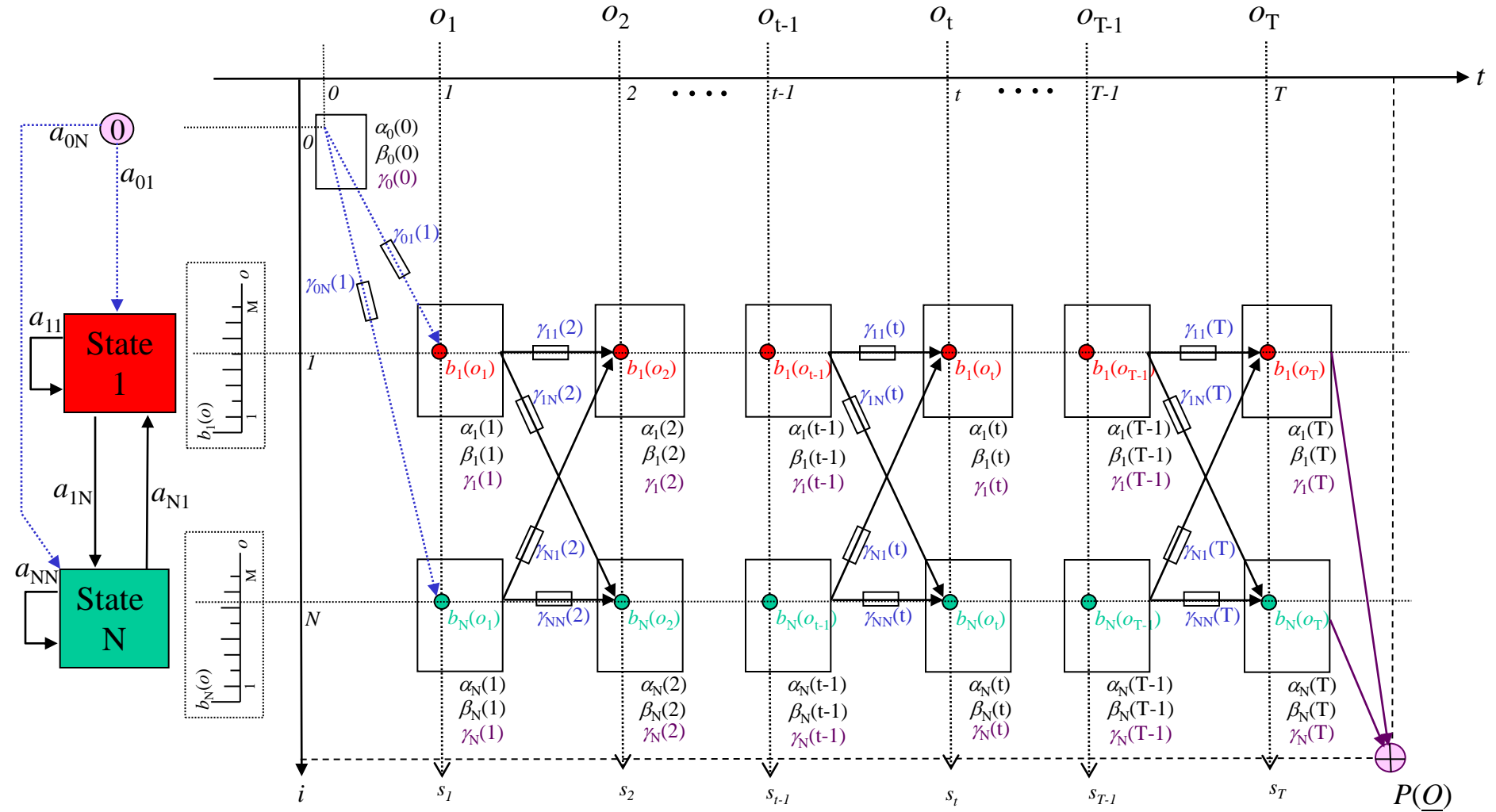
$$= \frac{\langle \gamma_{s_{t-1}s_t}(t) \rangle}{\langle \gamma_{s_{t-1}}(t-1) \rangle} = \frac{\sum_{t=1}^T \gamma_{s_{t-1}s_t}(t)}{\sum_{t=1}^T \gamma_{s_{t-1}}(t-1)} = \frac{\Gamma_{s_{t-1}s_t}}{\Gamma_{s_{t-1}}}$$

$\hat{b}_{s_t}(o) |_{\underline{Q}} \equiv$ (Estimation of $b_{s_t}(o)$, given \underline{Q})

$$= \frac{\langle \gamma_{s_t}(t) \cdot \delta(o_t - o) \rangle}{\langle \gamma_{s_t}(t) \rangle} = \frac{\sum_{t=1}^T \gamma_{s_t}(t) \cdot \delta(o_t - o)}{\sum_{t=1}^T \gamma_{s_t}(t)} = \frac{\Delta_{s_t o}}{\Gamma_{s_t}}$$

$$\text{where } \delta(o_t - o) = \begin{cases} 1, & \text{if } o_t = o \\ 0, & \text{otherwise} \end{cases}$$

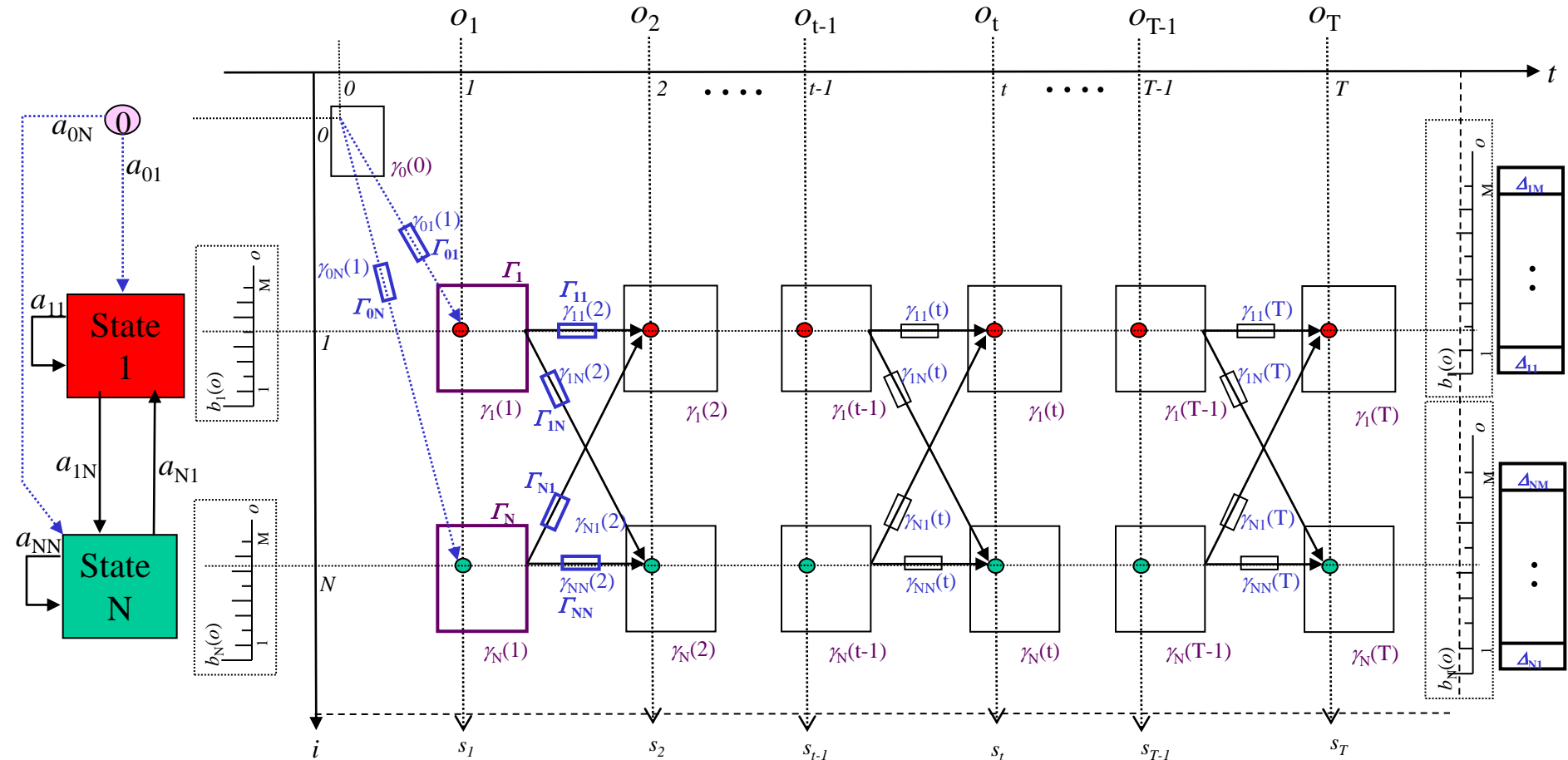
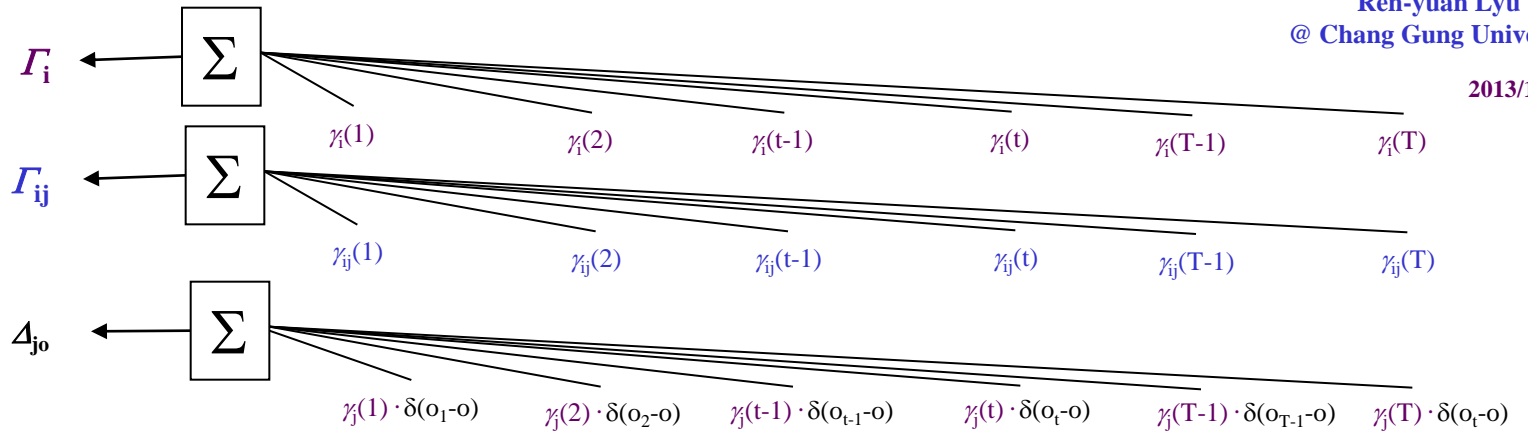
注意: $\gamma_{s_{t-1}s_t}(t)$, $\gamma_{s_{t-1}}(t-1)$, $\gamma_{s_t}(t)$ 本身不需取 log, 它們不會 Underflow



$$\hat{\pi}_j = \hat{a}_{0j} |_{\underline{o}} = \frac{\Gamma_{0j}}{\Gamma_0}$$

$$\hat{a}_{ij} |_{\underline{o}} = \frac{\Gamma_{ij}}{\Gamma_i}$$

$$\hat{b}_j(o) |_{\underline{o}} = \frac{\Delta_{jo}}{\Gamma_j}$$



$$\Gamma_{s_t} = \sum_{t=1}^T \gamma_{s_t}(t)$$

$$\Gamma_{s_{t-1}s_t} = \sum_{t=1}^T \gamma_{s_{t-1}s_t}(t)$$

$$\Delta_{s_t o} = \sum_{t=1}^T \gamma_{s_t}(t) \cdot \delta(o_t - o)$$

$$s_0 \equiv 0$$

$$s_{t-1} = i \in [1..N]$$

$$s_t = j \in [1..N]$$

$$o \in [1..M]$$

$$\hat{a}_{s_{t-1}s_t} |_{\underline{o}} = \frac{\Gamma_{s_{t-1}s_t}}{\Gamma_{s_{t-1}}}$$

$$\hat{b}_{s_t}(o) |_{\underline{o}} = \frac{\Delta_{s_t o}}{\Gamma_{s_t}}$$

$$\Gamma_i = \sum_{t=0}^{T-1} \gamma_i(t)$$

$$\Gamma_{ij} = \sum_{t=1}^T \gamma_{ij}(t)$$

$$\Gamma_0 = \gamma_0(0)$$

$$\Gamma_{0j} = \gamma_{0j}(1)$$

$$\Gamma_j = \sum_{t=1}^T \gamma_j(t)$$

$$\Delta_{jo} = \sum_{t=1}^T \gamma_j(t) \cdot \delta(o_t - o)$$

$$\hat{\pi}_j = \hat{a}_{0j} |_{\underline{o}} = \frac{\Gamma_{0j}}{\Gamma_0}$$

$$\hat{a}_{ij} |_{\underline{o}} = \frac{\Gamma_{ij}}{\Gamma_i}$$

$$\hat{b}_j(o) |_{\underline{o}} = \frac{\Delta_{jo}}{\Gamma_j}$$

Some interpretations about $\alpha, \beta, \eta, \gamma$

$$\alpha_{s_t}(t) \equiv P(o_1, o_2, \dots, o_t, s_t) = \sum_{s_{t-1}=1}^N \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t)$$

$$\beta_{s_t}(t) \equiv P(o_{t+1}, o_{t+2}, \dots, o_{T-1}, o_T | s_t) = \sum_{s_{t+1}=1}^N a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{s_{t+1}}(t+1)$$

$$\eta_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t, \underline{O}) = \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \cdot \beta_{s_t}(t)$$

$$\eta_{s_t}(t) \equiv P(s_t, \underline{O}) = \sum_{s_{t-1}=1}^N \eta_{s_{t-1}s_t}(t) = \alpha_{s_t}(t) \cdot \beta_{s_t}(t)$$

$$P(\underline{O}) = \sum_{s_T=1}^N \alpha_{s_T}(T) = \beta_{s_0}(0) = \sum_{s_1=1}^N \eta_{s_1}(1)$$

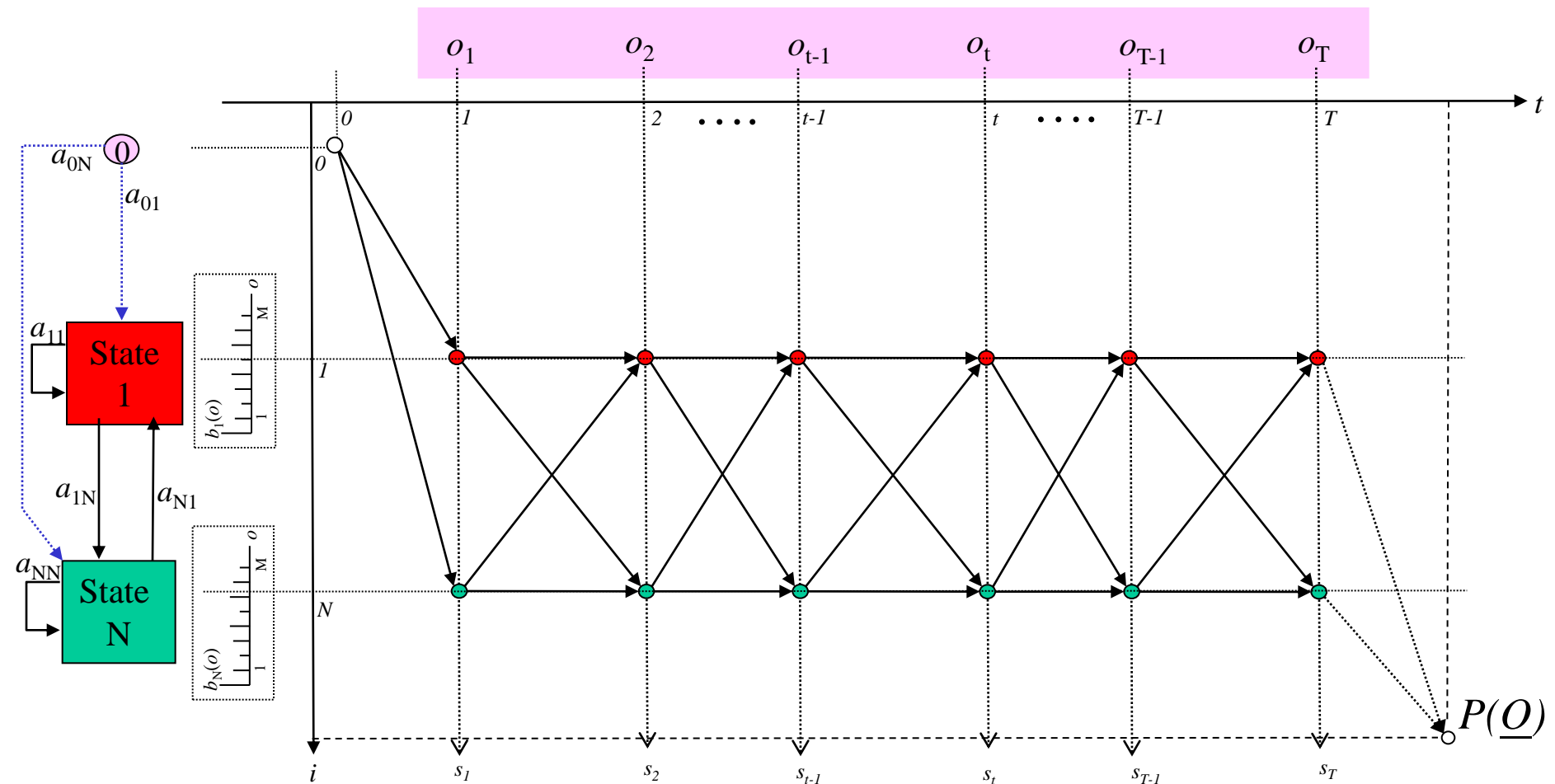
$$\gamma_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t | \underline{O}) = \eta_{s_{t-1}s_t}(t) / P(\underline{O})$$

$$\gamma_{s_t}(t) \equiv P(s_t | \underline{O}) = \eta_{s_t}(t) / P(\underline{O})$$

$$P(\underline{O})$$

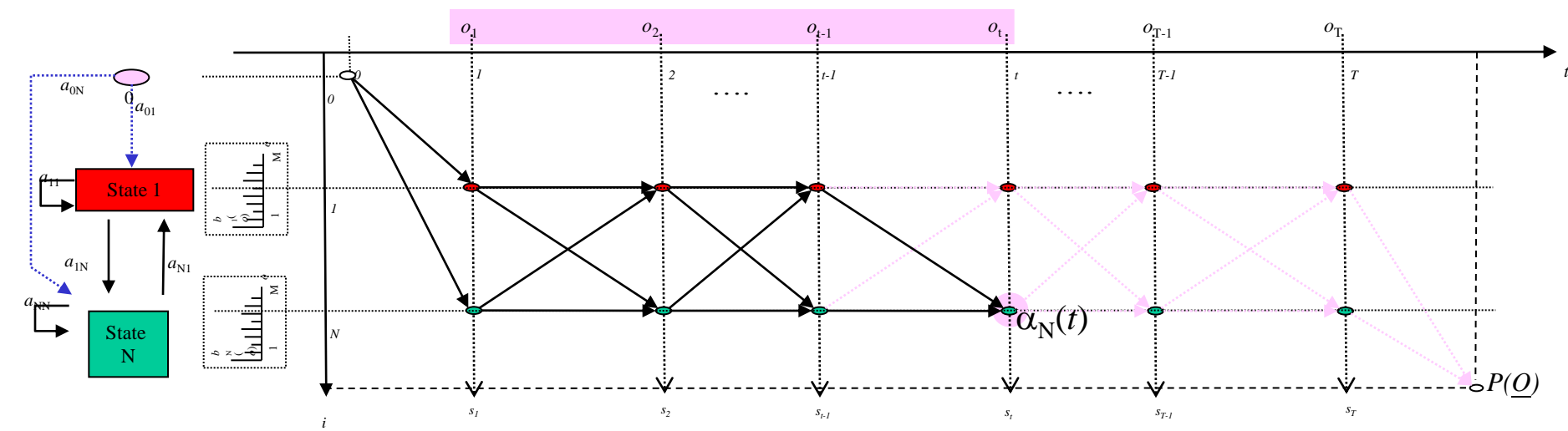
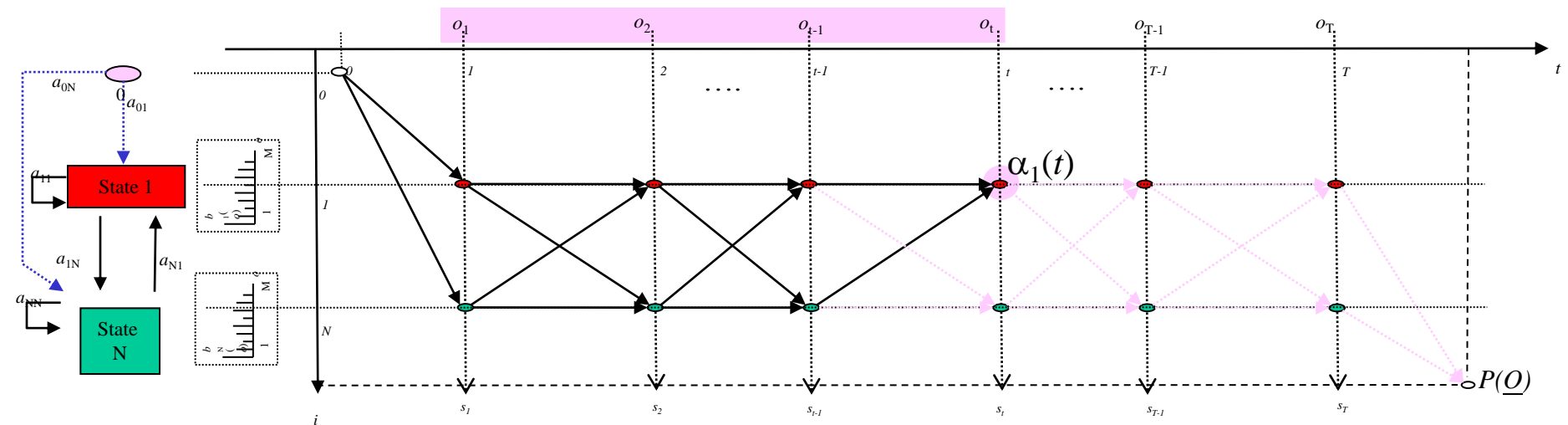
代表所有

從($i=0, t=0$)到($i=N+1, t=T+1$)所有的路徑，
如圖所示總共有 N^T 條



代表所有
從 $(i=0, t=0)$ 到 (i, t) 所有的路徑，
如圖所示總共有 N^{t-1} 條

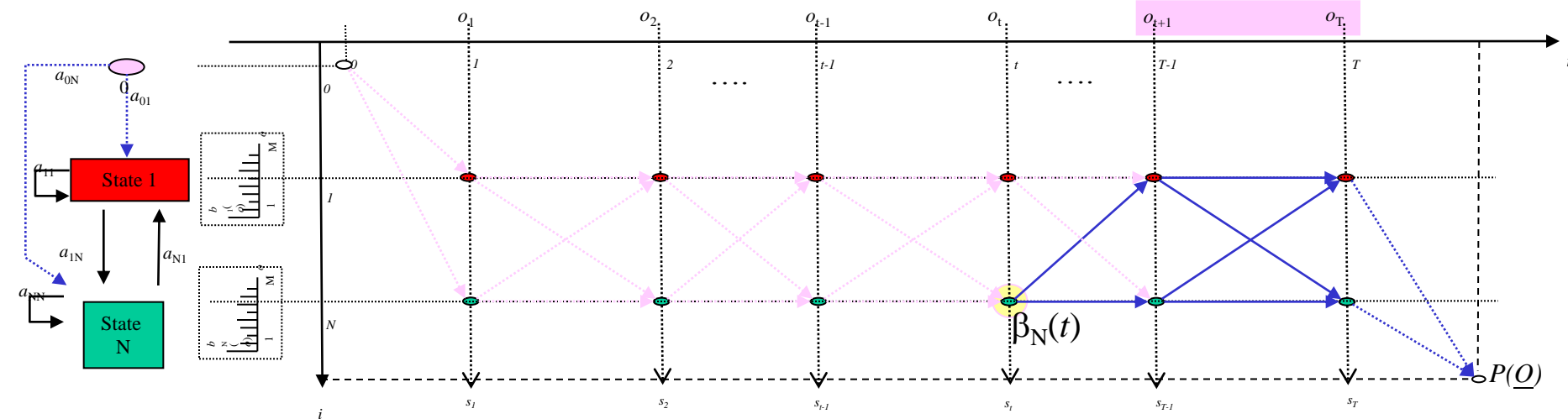
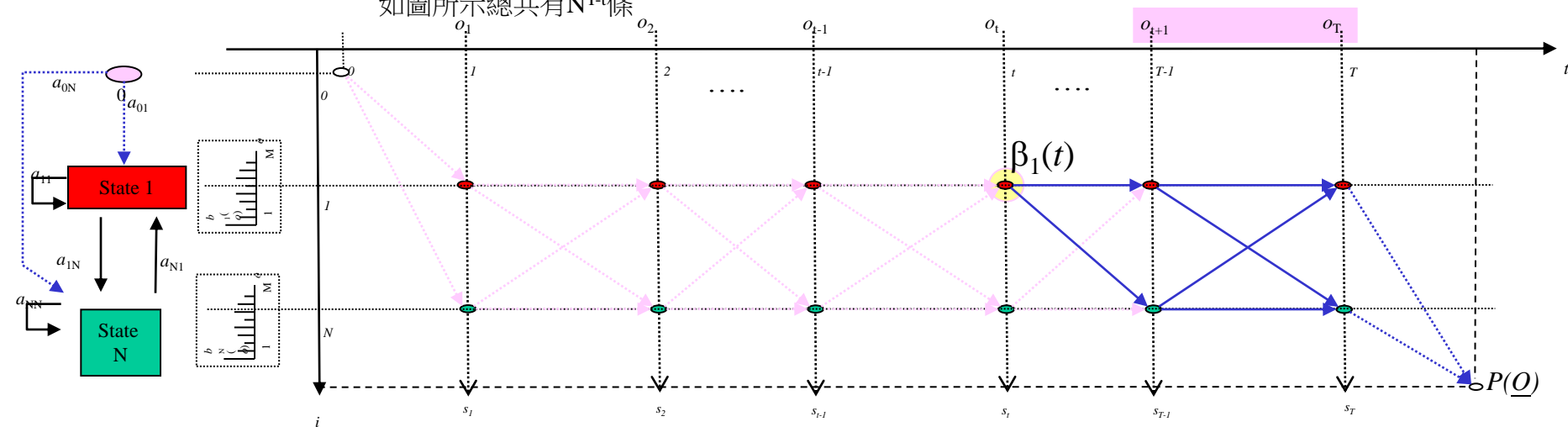
$$\alpha_i(t)$$



代表所有

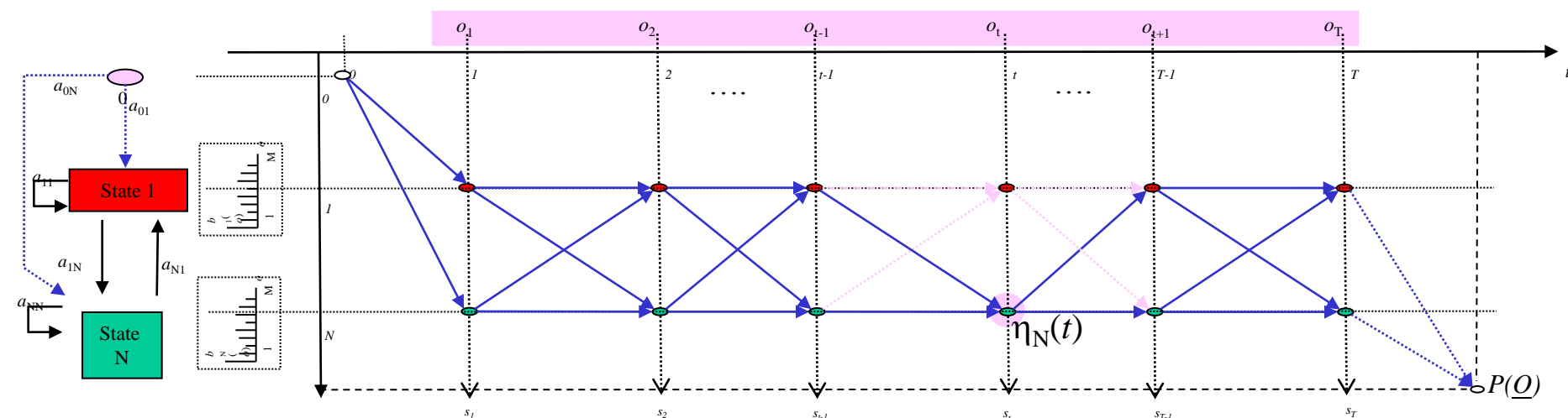
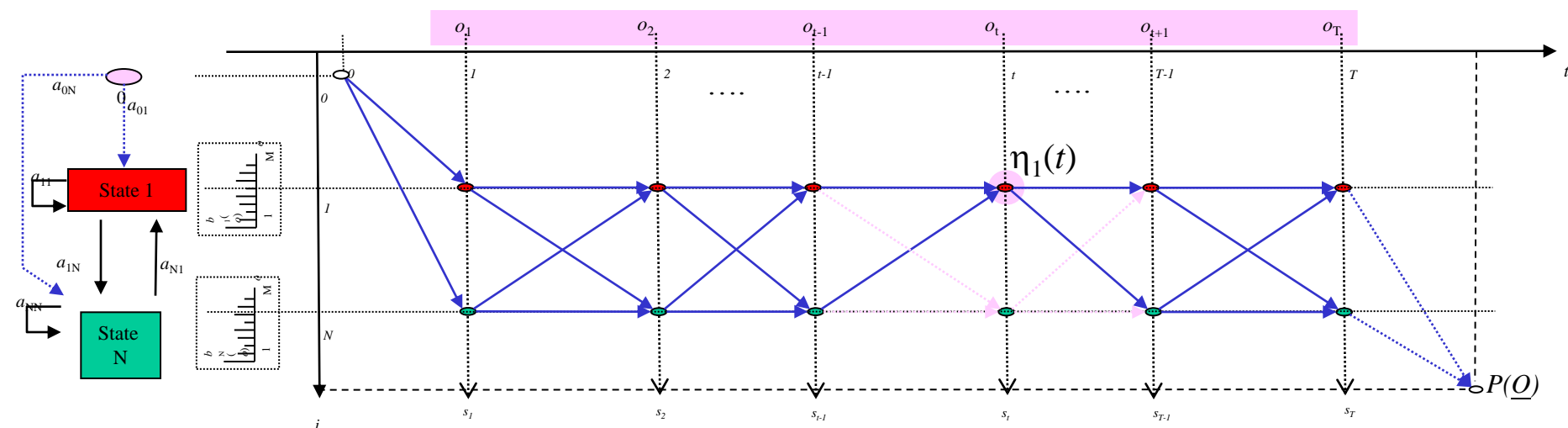
從 (i,t) 到 $(i=N+1, t=T+1)$ 所有的路徑，
如圖所示總共有 N^{T-t} 條

$$\beta_i(t)$$



$$\gamma_i(t)$$

$$\eta_i(t)$$



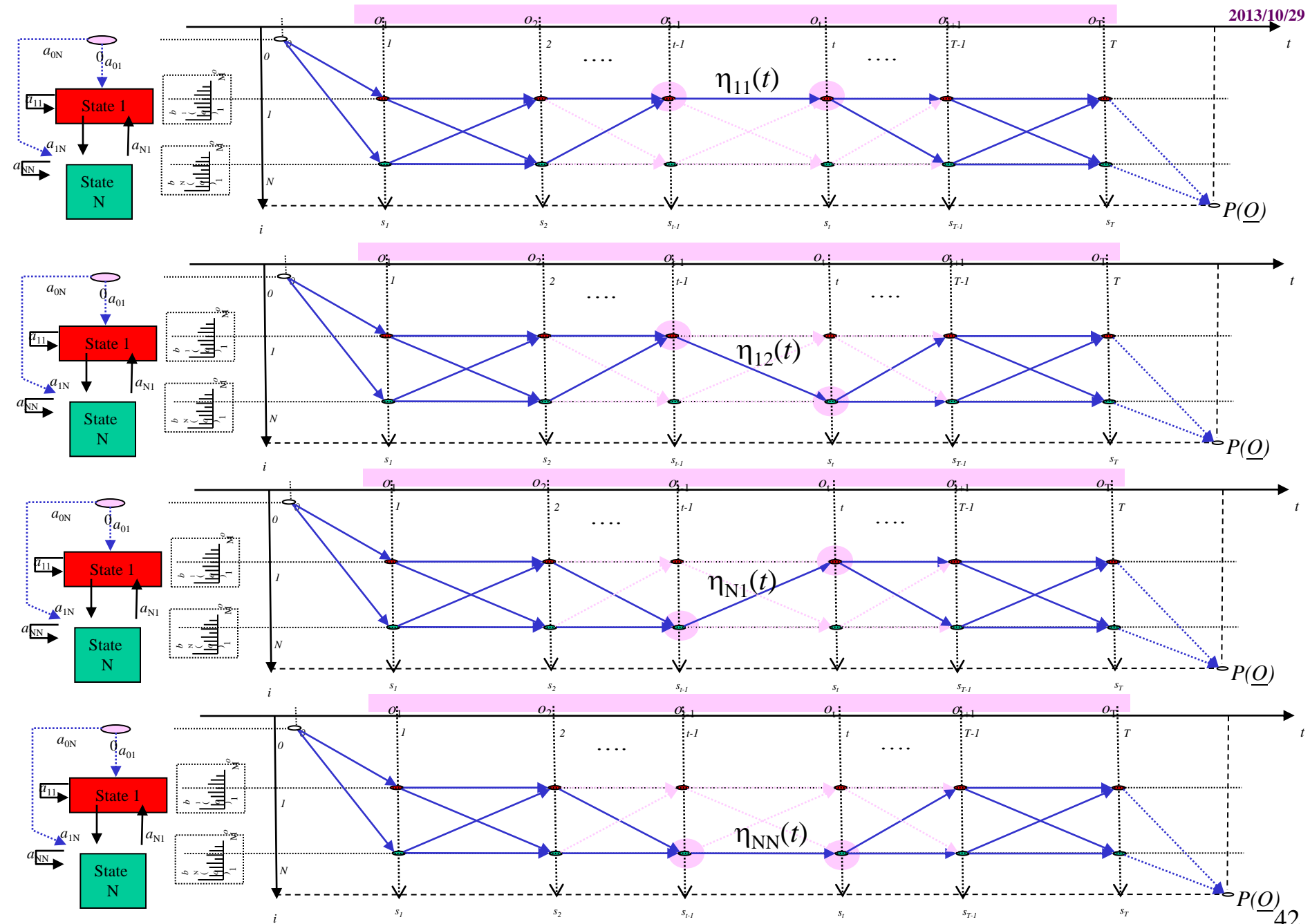
代表所有
從(i=0,t=0)到(I,t)

再從(i,t)到(i=N+1, t=T+1)的
所有的路徑，如圖所示總共有 N^{T-1} 條

這些路徑的共同特色為：它們都經過(t,i)

2013/10/29

$$\eta_{ij}(t), \gamma_{ij}(t)$$



Training Formula & Interpretations

$$\Gamma_i = \sum_{t=0}^{T-1} \gamma_i(t)$$

走完全程後，(state=i)會出現的次數

$$\Gamma_{ij} = \sum_{t=1}^T \gamma_{ij}(t)$$

走完全程後，(state=i)且(next state=j)會出現的次數

$$\Gamma_0 = \gamma_0(0)$$

走完全程後，一開始(t=0時) (state=0)會出現的次數

$$\Gamma_{0j} = \gamma_{0j}(1)$$

走完全程後，一開始(t=0時) (state=0)且(next state=j) 會出現的次數

$$\Gamma_j = \sum_{t=1}^T \gamma_j(t)$$

走完全程後，(state=j)會出現的次數

$$\Delta_{jo} = \sum_{t=1}^T \gamma_j(t) \cdot \delta(o_t - o)$$

走完全程後，(state=j)且(observation=o) 會出現的次數

$$\hat{\pi}_j = \hat{a}_{0j} |_{\underline{o}} = \frac{\Gamma_{0j}}{\Gamma_0}$$

$$\hat{a}_{ij} |_{\underline{o}} = \frac{\Gamma_{ij}}{\Gamma_i}$$

$$\hat{b}_j(o) |_{\underline{o}} = \frac{\Delta_{jo}}{\Gamma_j}$$

The Decoding Problem of HMM

- Given a HMM, with parameters $\{\pi, \mathbf{A}, \mathbf{B}\}$, and an observation sequence

$$O_1, O_2, O_3, \dots, O_t, \dots, O_T$$

What is the “best” or “most likely” state sequence

$$S_1^*, S_2^*, S_3^*, \dots, S_t^*, \dots, S_T^*$$

to help the HM**Machine** generate such an observation sequence?

This problem can be reformulate as follows:

$$(s_1^*, s_2^*, s_3^*, \dots, s_t^*, \dots, s_T^*) = \underset{\forall (s_1, s_2, s_3, \dots, s_t, \dots, s_T)}{\mathbf{Argmax}} \quad P(s_1, s_2, s_3, \dots, s_t, \dots, s_T \mid o_1, o_2, o_3, \dots, o_t, \dots, o_T) = ?$$

Optimal principle: (The Dynamic Programming Principle)

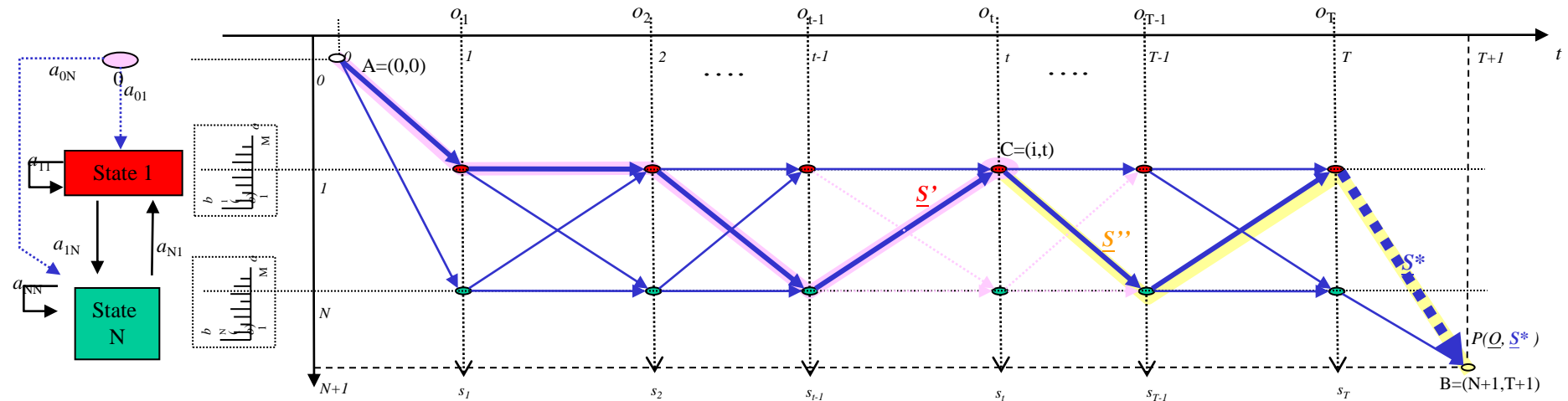
If path \underline{S}^* is the “**optimal**” **path** of all **paths** from $A=(0,0)$ to $B=(N+1,T+1)$, which passes $C=(i, t)$,
Then the **partial** path $\underline{S}' (= \underline{S}^*[0..t])$ of \underline{S}^* from $A=(0,0)$ to $C=(i, t)$
must also be the “**optimal**” **path** of all **paths** from $A=(0,0)$ to $C=(i, t)$,
And the **partial** path $\underline{S}'' (= \underline{S}^*[t..T+1])$ of \underline{S}^* from $C=(i,t)$ to $B=(N+1, T+1)$
must also be the “**optimal**” **path** of all **paths** from $C=(i,t)$ to $B=(N+1, T+1)$.

e.g., in the following figure,

$\underline{S}^* = (0,0) \rightarrow (1,1) \rightarrow (1,2) \rightarrow \dots \rightarrow (N,t-1) \rightarrow (1,t) \rightarrow (N,T-1) \rightarrow (1,T) \rightarrow (N+1) \rightarrow T+1)$

$\underline{S}' = (0,0) \rightarrow (1,1) \rightarrow (1,2) \rightarrow \dots \rightarrow (N,t-1) \rightarrow (1,t)$

$\underline{S}'' = (1,t) \rightarrow (N,T-1) \rightarrow (1,T) \rightarrow (N+1) \rightarrow T+1)$



If you want to find the optimal path from $t=0$ to $t=T+1$,
You should find the optimal partial path from $t=0$ to t for each possible i in $\{1, \dots, N\}$
Because at time t , the optimal path may pass one of points in $\{(1,t), \dots, (N,t)\}$

$$t \in [1..T]$$

$$i \in \Omega_S = \{1..N\}$$

$\underline{S}'_i(t) \equiv$ the optimal partial path from $(i=0, t=0)$ to (i, t)

$$= "(0,0) \sim (s_1^*, 1) \sim (s_2^*, 2) \sim (s_{t-1}^*, t-1) \sim (i, t)"$$

\equiv the optimal partial state sequence from beginning to time $(t-1)$ and $(s_t = i)$

$$= [(s_0 \equiv 0), s_1^*, s_2^*, \dots, s_{t-1}^*, (s_t = i)]$$

$$= \underset{s_1, s_2, \dots, s_{t-1}}{\text{ArgMax}} P(s_1, s_2, \dots, s_{t-1}, (s_t = i), o_1, o_2, \dots, o_t)$$

$s_\tau^* \equiv$ the optimal state at time τ

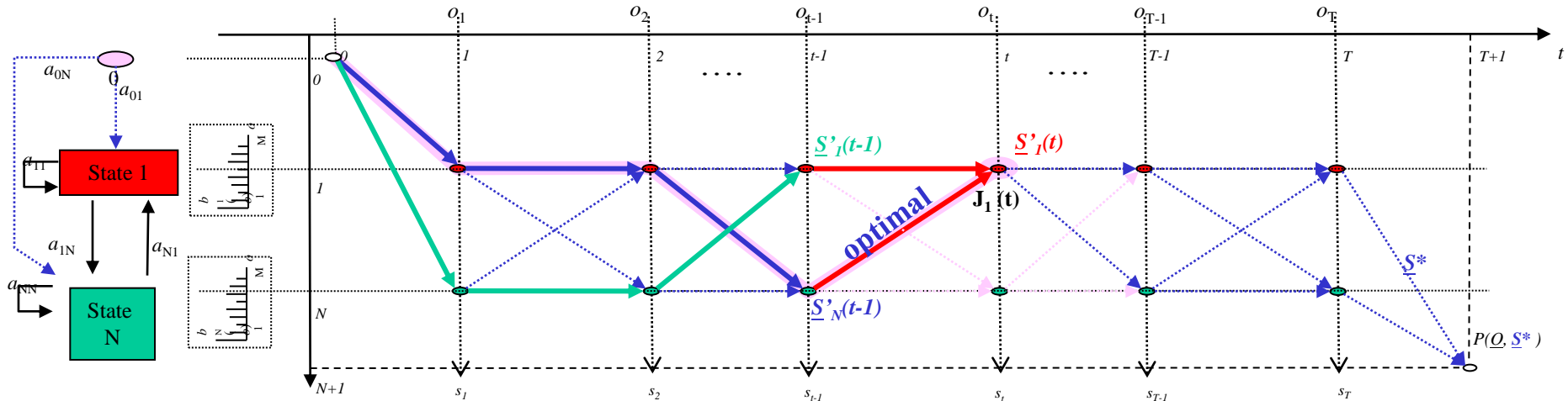
It can be shown by the optimal principle that

$$\underline{S}'_i(t) = \underset{j \in \Omega_S = \{1..N\}}{\text{Opt}} \{ \underline{S}'_j(t-1) \sim (i, t) \}$$

And let's record the optimal argument as follows.

$J_i(t) \equiv$ the optimal pre-state to (i, t)

$$= \underset{j \in \Omega_S = \{1..N\}}{\text{ArgOpt}} \{ \underline{S}'_j(t-1) \sim (i, t) \}$$



$$\underline{S}'_i(1) = (0,0) \sim (i,1) \quad i \in \Omega_S = \{1..N\}$$

$$J_i(1) = 0$$

$$\underline{S}'_i(2) = \underset{j \in \Omega_S = \{1..N\}}{\text{Opt}} \{ \underline{S}'_j(1) \sim (i,2) \}$$

$$J_i(2) = \underset{j \in \Omega_S = \{1..N\}}{\text{ArgOpt}} \{ \underline{S}'_j(1) \sim (i,2) \}$$

$$\underline{S}'_i(t) = \underset{j \in \Omega_S = \{1..N\}}{\text{Opt}} \{ \underline{S}'_j(t-1) \sim (i,t) \}$$

$$J_i(t) = \underset{j \in \Omega_S = \{1..N\}}{\text{ArgOpt}} \{ \underline{S}'_j(t-1) \sim (i,t) \}$$

$$\underline{S}'_i(T) = \underset{j \in \Omega_S = \{1..N\}}{\text{Opt}} \{ \underline{S}'_j(T-1) \sim (i,T) \}$$

$$J_i(T) = \underset{j \in \Omega_S = \{1..N\}}{\text{ArgOpt}} \{ \underline{S}'_j(T-1) \sim (i,T) \}$$

$$\underline{S}'_{N+1}(T+1) = \underset{j \in \Omega_S = \{1..N\}}{\text{optimal}} \{ \underline{S}'_j(T) \sim (N+1, T+1) \}$$

$$J_{N+1}(T+1) = \underset{j \in \Omega_S = \{1..N\}}{\text{ArgOpt}} \{ \underline{S}'_j(T) \sim (N+1, T+1) \}$$

$s^*(t) \equiv$ the optimal state at time t

$$s^*(T+1) \equiv N+1$$

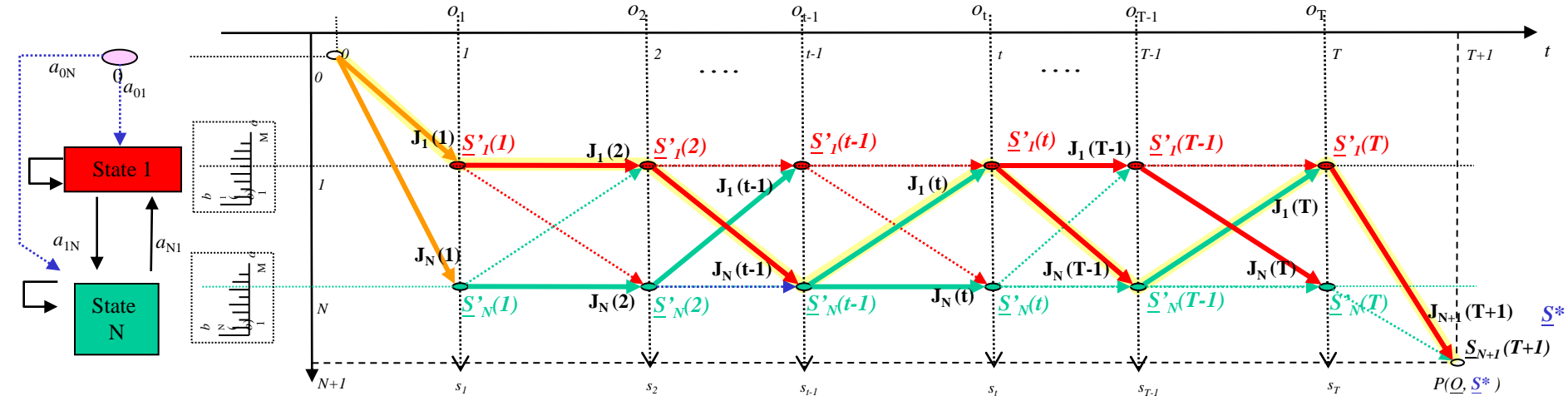
$$s^*(T) = J_{N+1}(T+1) = J_{s^*(T+1)}(T+1)$$

$$s^*(t) = J_{s^*(t+1)}(t+1)$$

$$s^*(1) = J_{s^*(2)}(2)$$

$$s^*(0) = J_{s^*(1)}(1) \equiv 0$$

$$\underline{S}^* = (s_0^* \equiv 0), s^*(1), s^*(2), \dots, s^*(t), \dots, s^*(T), (s^*(T+1) \equiv N+1)$$



An example application



GGRFFGRFFRRGGRRGGRGRGRGRGFRRGFGRRRGGGF

G: Green, down

R: Red, up

F: Flat, level



The observation sequence

$O(t)$ G G G R F F G R F F R R G G G R R G G R G R G R G F R R G F G F R R G G G F

A *possible* state sequence

$S(t)$

空頭市場
 $S=1$

多頭市場
 $S=2$

空頭市場
 $S=1$

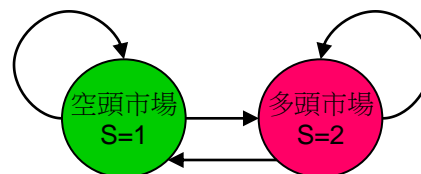
多頭市場
 $S=2$

空頭市場
 $S=1$

多頭市場
 $S=2$

空頭市場
 $S=1$

	空	多	
空	22	3	25
多	3	9	12
			37



	空	多	
G	15	2	
R	6	7	
F	5	3	
	26	12	38

Problems

- (1) Provide an initial estimate of $\{\pi, \mathbf{A}, \underline{\mathbf{B}}\}$
- (2) $P(\underline{O}) = ?$
- (3) The optimal sequence $\underline{S}^* = ?$
- (4) Re-estimate $\{\pi, \mathbf{A}, \underline{\mathbf{B}}\}$ as $\{\pi', \mathbf{A}', \underline{\mathbf{B}}'\}$
such that $P(\underline{O} | \pi', \mathbf{A}', \underline{\mathbf{B}}') > P(\underline{O} | \pi, \mathbf{A}, \underline{\mathbf{B}})$