

希望语义相近的词语能归到一个词组中 → 语义相似 ⇒ 何必泛化

↓  
机器学习模型中  
高本固底的参数模型  
进行修改并使用

(分割)

# Chinese Word Segmentation

- Although great progress has been made in improving the performance of **Chinese word segmentation** systems, Chinese word segmentation remains an **extremely difficult problem**.
- There is no existing approach that is able to reliably segment **unfamiliar types of texts** before fine-tuning with massive training data, not to mention the **open texts** mixed with the texts from different topics, genres, and regions.
- **Word segmentation** still remains a **main research topic** in the field of Chinese language processing, which indicates that there may be unresolved theoretical and processing issues.

# Word Boundary Ambiguities

歧义

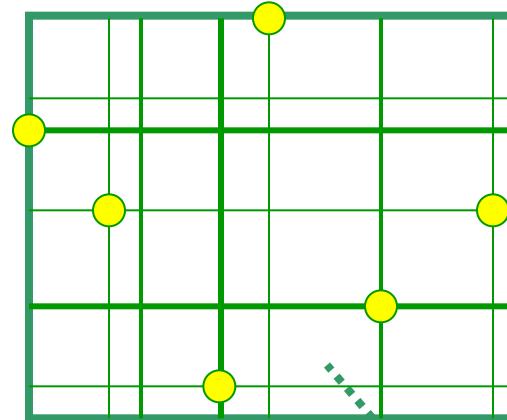
- Combination Ambiguity

以 / 我 / **个人** / 的名义  
他 / 一 / **个** / 人 / 在家

~~并列~~

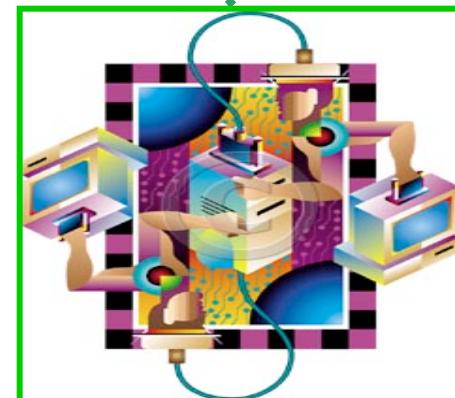
- Overlap Ambiguity

从 / 小学 / 到 / 中学  
从小 / 学 / 计算机



- Global Ambiguity

美国 / 会 / 采取行动  
美 / 国会 / 采取行动



# Unknown words

他从小学画画

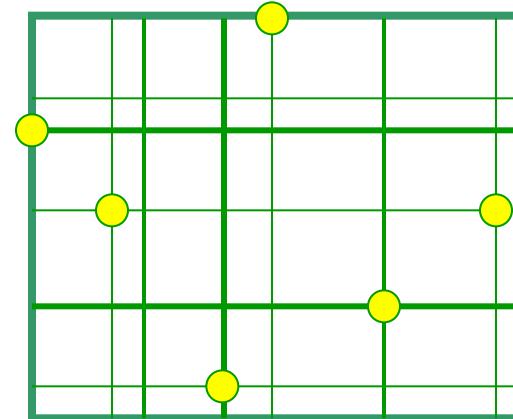
他从小学开始画画

他从小学着画画

他 / 从 / 小 / 学 / 画 / 画

他 / 从 / 小 / 学 / 开 / 始 / 画 / 画

他 / 从 / 小 / 学 / 着 / 画 / 画



**数量、时间**：“2012年5月22日”、“一市斤”、“356克”

**人名、机构名、地名**：“李维汉”、“阿凡提”、“上海浦东张江镇”

**外文翻译及缩写**：“阿诺德·施瓦辛格”、“萨默维尔”、“FDA”

**专用名称的缩写**：“复旦”、“中航”、“央行”

**新词**：“安家费”、“三通”、“菜鸟”、“八卦”

# Word Formation

- **重叠词**

**AA式**：爸爸、宝宝、星星

**AABB式**：大大咧咧、形形色色、漂漂亮亮

- **派生词**

**前缀 + 词根**：阿姨、老虎、老婆

**词根 + 后缀**：椅子、木头、鸟儿

**词根 + 中缀 + 词根**：对得起、来得及

- **词转化为短语**

**动宾结构式**：鞠躬 → 鞠个躬 → 鞠个九十度的躬

**偏正或联合结构**：同学 → 同过三年学

**补充结构式**：达到 → 达得/不到

- **短语转化为词**

**ABCD → AC式**：土地改革 → 土改、地下铁道 → 地铁

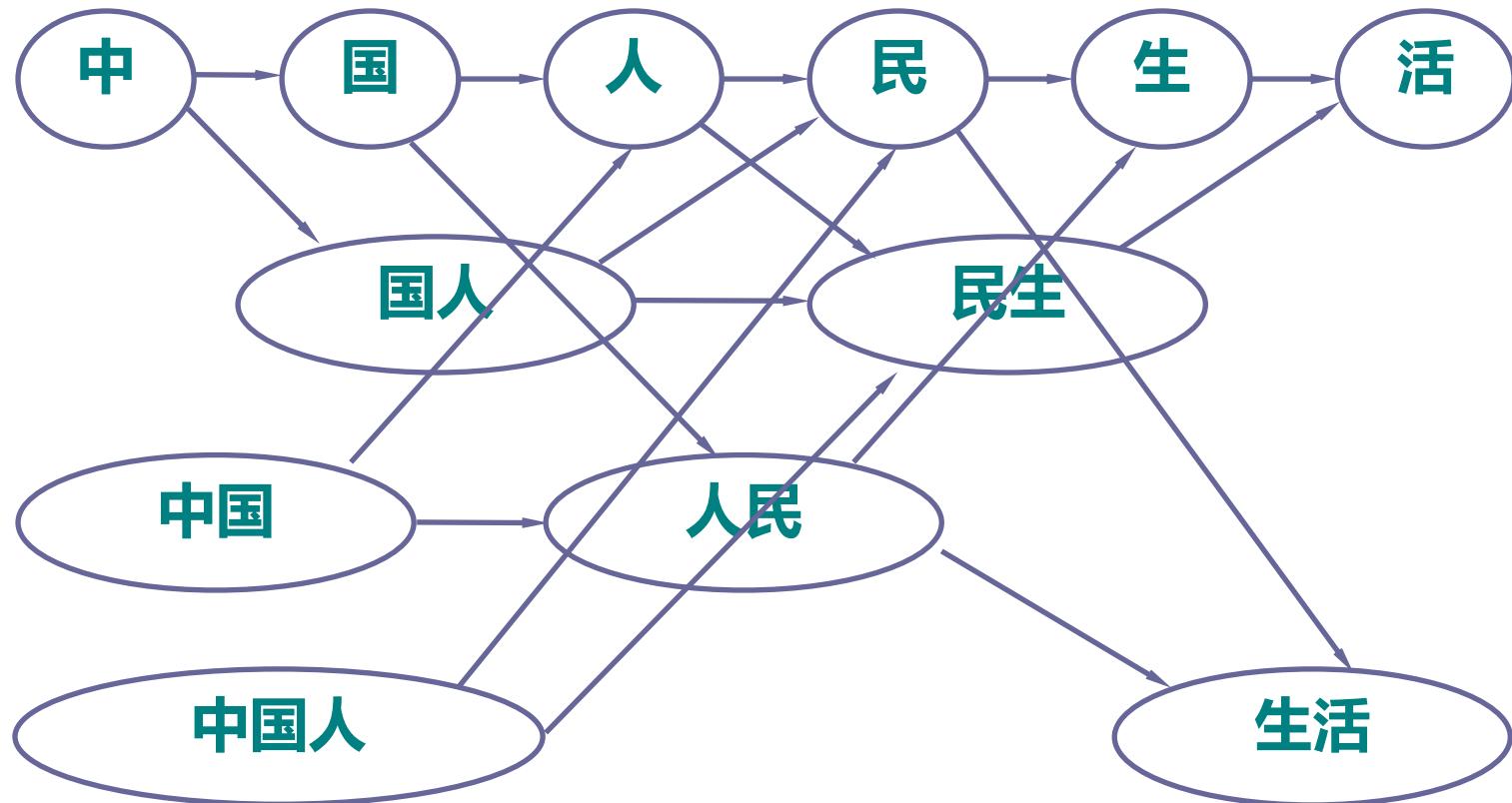
**ABCD → AD式**：空中小姐 → 空姐、高等院校 → 高校

**截段简缩**：中国南极长城站 → 长城站、复旦大学 → 复旦

**综合简缩**：联合国安全理事会 → 安全理事会 → 安理会

# Chinese word segmentation

---



# Character-based Tagging Solution

Sentence: 从 小 学 计 算 机 。

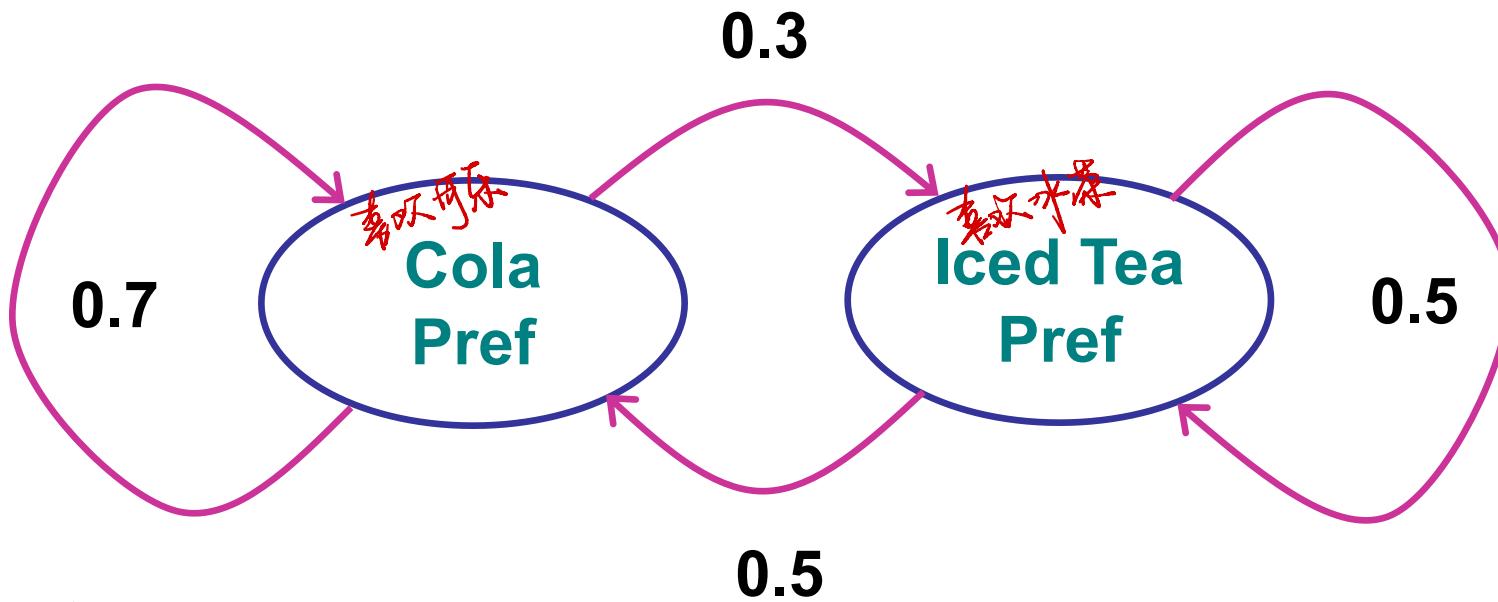
Tags:      B E S B I E S

词开始 词中 词中 词中  
词中 词中 词中 词中

*Character-based models treat known and unknown words equally and no other exception processing is necessary.*

随机游走模型

# The crazy soft drink machine



	Cola	Iced tea	Lemonade
Cola Pref (CP)	0.6	0.1	0.3
Iced Tea Pref (IP)	0.1	0.7	0.2

隐藏马尔可夫元组

# General form of an HMM

An **Hidden Markov Model** (HMM) is specified by a five-tuple  $(S, K, \Pi, A, B)$

Set of states

状态空间分布  
 $S = \{S_1, \dots, S_N\}$

Output alphabet

$$K = \{K_1, \dots, K_M\} = \{1, \dots, M\}$$

Initial state probabilities

$$\Pi = \{\pi_i\}, i \in S$$

State transition probabilities

$$A = \{a_{ij}\}, i, j \in S$$

Symbol emission probabilities

$$B = \{b_{ijk}\}, i, j \in S, k \in K$$

State sequence

$$X = \{X_1, \dots, X_{T+1}\} \quad X_t: S \rightarrow \{1, \dots, N\}$$

Output sequence

$$O = \{O_1, \dots, O_T\} \quad o_t \in K$$

$$P(O_t = k | X_t = s_i, X_{t+1} = s_j) = b_{ijk}$$

# Three fundamental questions for HMMs

---

- Given a model  $\mu = \{A, B, \Pi\}$ , how do we efficiently compute how likely a certain observations is, that is  $P(O|\mu)$ ?
- Given the observation sequence  $O$  and a model  $\mu$ , how do we choose a state sequence  $\{X_1, \dots, X_{T+1}\}$  that best explains the observations?
- Given an observation sequence  $O$ , and a space of possible models found by varying the model parameters  $\mu = \{A, B, \Pi\}$ , how do we find the model that best explains the observed data?

这课(2)讲

# Finding the probability of an observation

Given the observation sequence  $O = \{O_1, \dots, O_T\}$  and a model  $\mu = \{A, B, \Pi\}$ , we wish to know how to efficiently compute  $P(O|\mu)$  – the probability of the observation given the model.

$$\alpha_i(t) = P(o_1 o_2 \dots o_{t-1}, X_t=i | \mu)$$

## Initialization

$$\alpha_i(1) = \pi_i, 1 \leq i \leq N$$

第*t*次  
status  
prob

## Induction

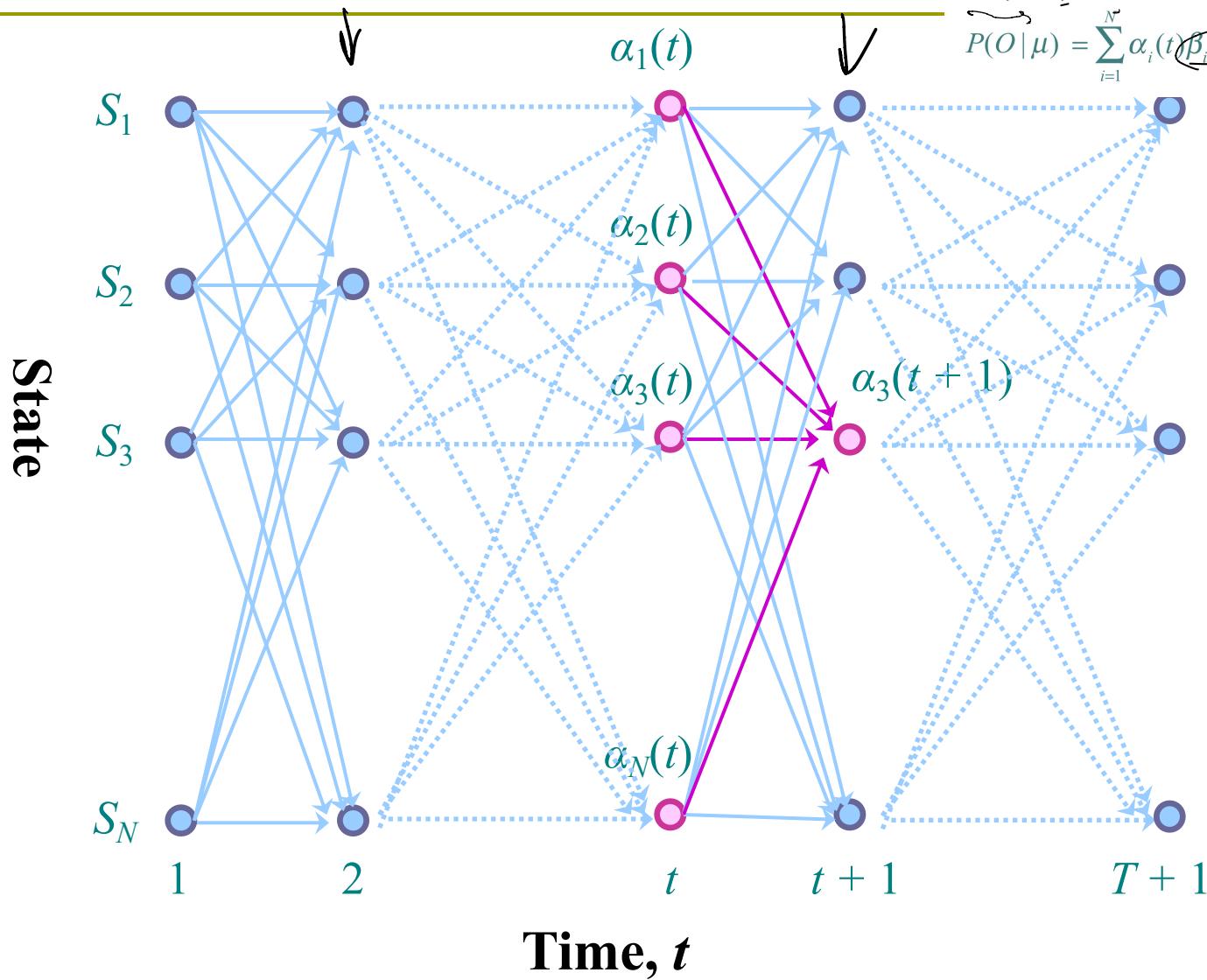
$$\alpha_j(t+1) = \sum_{i=1}^N \alpha_i(t) a_{ij} b_{ij o_t}, 1 \leq t \leq T, 1 \leq j \leq N$$

## Total

$$P(O| \mu) = \left[ \sum_{i=1}^N \alpha_i(T+1) \right]$$

$i \rightarrow j$  prob

# Trellis (lattices)algorithı



$$\begin{aligned}
 & P(O, \boxed{X_t = i} | \mu) \\
 &= P(o_1 \cdots o_T, X_t = i | \mu) \\
 &= P(o_1 \cdots o_{t-1}, X_t = i, o_t \cdots o_T | \mu) \\
 &= P(o_1 \cdots o_{t-1}, X_t = i | \mu) \times P(o_t \cdots o_T | o_1 \cdots o_{t-1}, X_t = i, \mu) \\
 &= P(o_1 \cdots o_{t-1}, X_t = i | \mu) \times P(o_t \cdots o_T | X_t = i, \mu) \\
 &= \alpha_i(t) \beta_i(t) \\
 & P(O | \mu) = \sum_{i=1}^N \alpha_i(t) \beta_i(t), \underline{1 \leq t \leq T+1} \\
 & \text{towards}
 \end{aligned}$$

# The backward procedure

---

Define backward variables

$$\beta_i(t) = P(\underbrace{o_t \dots o_{T+1}}_{}, \underbrace{X_{t+j}}_{}, \mu)$$

**Initialization**

$$\beta_i(T+1) = 1, 1 \leq i \leq N$$

**Induction**

$$\beta_i(t) = \sum_{j=1}^N a_{ij} b_{ij o_t} \beta_j(t+1), 1 \leq t \leq T, 1 \leq j \leq N$$

**Total**

$$P(O | \mu) = \sum_{i=1}^N \pi_i \beta_i(1)$$

# Finding the best state sequence

---

Commonly we want to find the most likely complete path, that is:

$$\arg \max_X P(X | O, \mu)$$

To do this, it is sufficient to maximize for a fixed  $O$ :

$$\arg \max_X P(X, O | \mu)$$

An efficient trellis algorithm for computing this path is the **Viterbi** algorithm

$$\delta_j(t) = \max_{X_1 \cdots X_{t-1}} P(X_1 \cdots X_{t-1}, o_1 \cdots o_{t-1}, X_t = j | \mu)$$

# Viterbi algorithm

---

## Initialization

$$\delta_j(1) = \pi_j, 1 \leq j \leq N$$

## Induction

$$\delta_j(t+1) = \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ijo_t}, 1 \leq j \leq N$$

$$\psi_j(t+1) = \arg \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ijo_t}, 1 \leq j \leq N$$

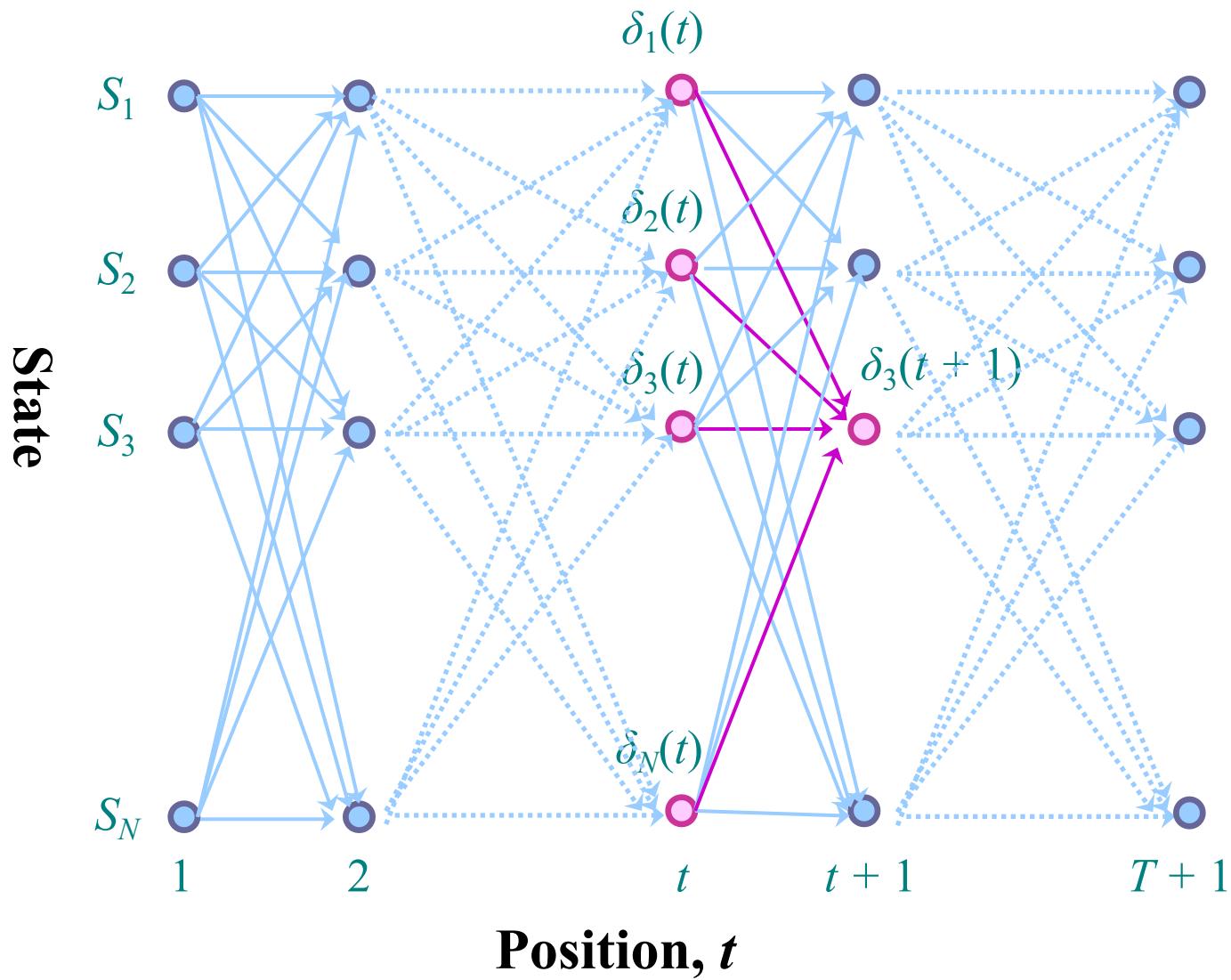
## Termination

$$\hat{X}_{T+1} = \arg \max_{1 \leq i \leq N} \delta_i(T+1)$$

$$\hat{X}_t = \psi_{\hat{X}_{t+1}}(t+1)$$

$$P(\hat{X}) = \max_{1 \leq i \leq N} \delta_i(T+1)$$

# Viterbi algorithm



# The crazy soft drink machine

Transition	Cola Pref (CP)	Iced Tea Pref (IP)
Cola Pref (CP)	0.7	0.3
Iced Tea Pref (IP)	0.5	0.5

$P(O_k | S_i)$  機連事件發生  
 $= P(O_k | S_{i \rightarrow j})$  不一樣

Emission	Cola	Iced tea	Lemonade
Cola Pref (CP)	0.6	0.1	0.3
Iced Tea Pref (IP)	0.1	0.7	0.2

Initial	Cola Pref (CP)	Iced Tea Pref (IP)
$\pi$	0.4	0.6

# Evaluation

---

Times	1	2	3
CP	$0.6 \times 0.7$		
IP	$0.1 \times 0.6$		



Transition	CP	IP
CP	0.7	0.3
IP	0.5	0.5

Emission	Cola	Iced tea	Lemonade
CP	0.6	0.1	0.3
IP	0.1	0.7	0.2

Initial	CP	IP
$\pi$	0.4	0.6

# Evaluation

Times	1	2	3
CP	0.2400	0.0198	0.0149
IP	0.0600	0.0714	0.0083

0.0232



$$\begin{aligned}
 &= 0.2400 \times 0.7 \times 0.1 + \\
 &\quad 0.0600 \times 0.5 \times 0.1 \\
 &= 0.0198
 \end{aligned}$$

Transition	CP	IP
CP	0.7	0.3
IP	0.5	0.5

Emission	Cola	Iced tea	Lemonade
CP	0.6	0.1	0.3
IP	0.1	0.7	0.2

Initial	CP	IP
$\pi$	0.4	0.6

How about



?

# Decoding

Times	1	2	3
CP	0.2400	0.0168	0.0076
IP	0.0600	0.0504	0.0050
Times	1	2	3
CP	CP	CP	IP
IP	CP	CP	IP

Transition	CP	IP
CP	0.7	0.3
IP	0.5	0.5
Emission	Cola	Iced tea
CP	0.6	0.1
IP	0.1	0.7
Initial	CP	IP
$\pi$	0.4	0.6



$$= 0.2400 \times 0.7 \times 0.1 = 0.0168$$

$$= 0.0600 \times 0.5 \times 0.1 = 0.0003$$

States sequence: CP IP CP

# Parameter estimation

Given a certain observation sequence, we want to find the values of the model parameters  $\mu = \{A, B, \Pi\}$  which best explain what we observed. Using **Expectation Maximization** (EM), that means we want to find the values that maximize  $P(O|\mu)$ :

$$\begin{aligned} P(O, \boxed{X_t = i} | \mu) &\xrightarrow{\text{添付資料}} \text{Scribble} \\ &= P(o_1 \cdots o_T, X_t = i | \mu) \\ &= P(o_1 \cdots o_{t-1}, X_t = i, o_t \cdots o_T | \mu) \\ &= P(o_1 \cdots o_{t-1}, X_t = i | \mu) \times P(o_t \cdots o_T | o_1 \cdots o_{t-1}, X_t = i, \mu) \\ &= P(o_1 \cdots o_{t-1}, X_t = i | \mu) \times P(o_t \cdots o_T | X_t = i, \mu) \\ &= \alpha_i(t) \beta_i(t) \end{aligned}$$

$$P(O | \mu) = \sum_{i=1}^N \alpha_i(t) \beta_i(t), 1 \leq t \leq T+1$$

EM 算法

随机游走

→ 随机游走

随机游走  $\rightarrow$  不规则随机，无周期性  $\xrightarrow{P=1}$

随机游走

随机游走

# Parameter estimation

$$\gamma_i(t) = P(X_t = i | O, \mu)$$

$$= \frac{P(X_t = i, O | \mu)}{P(O | \mu)}$$

$$= \frac{\alpha_i(t) \beta_i(t)}{\sum_{j=1}^N \alpha_j(t) \beta_j(t)}$$

$$p_t(i, j) = P(X_t = i, X_{t+1} = j | O, \mu)$$

$$= \frac{P(X_t = i, X_{t+1} = j, O | \mu)}{P(O | \mu)}$$

$$= \frac{\alpha_i(t) \boxed{a_{ij} b_{ijo_t}} \beta_j(t+1)}{\sum_{m=1}^N \sum_{n=1}^N \alpha_m(t) \boxed{a_{mn} b_{mno_t}} \beta_n(t+1)}$$

*Note that:*

$$\underbrace{\gamma_i(t)}_{\text{initially}} = \sum_{j=1}^N p_t(i, j)$$

*Fix status at  $i \rightarrow j$  initially*

# Parameter estimation

$$\hat{\pi}_i = \text{Expected frequency in state } i \text{ at time } t=1$$

$= \gamma_i(1)$

$$\hat{a}_{ij} = \frac{\text{Expected number of transitions from state } i \text{ to } j}{\text{Expected number of transitions from state } i}$$

$$= \frac{\sum_{t=1}^T p_t(i, j)}{\sum_{t=1}^T \gamma_i(t)}$$

← fix status of  $i$   
 ← fix  $j$

$$\hat{b}_{ijk} = \frac{\text{Expected number of transitions from state } i \text{ to } j \text{ with } k \text{ observed}}{\text{Expected number of transitions from state } i \text{ to } j}$$

$$= \frac{\sum_{\{t: o_t = k, 1 \leq t \leq T\}} p_t(i, j)}{\sum_{t=1}^T p_t(i, j)}$$

$P(O | \hat{\mu}) \geq P(O | \mu)$