# 數位語音處理概論
## Introduction to Digital Speech Processing

## 4.0 More about Hidden Markov Models

### References for 4.0
1. 6.1-6.6, Rabiner and Juang, 2. 4.4.1 of Huang

授課教師：國立臺灣大學 電機工程學系 李琳山 教授

- **Markov Model (Markov Chain)**
  - First-order Markov chain of N states is a triplet (S,**A**,π)
    - S is a set of $N$ states
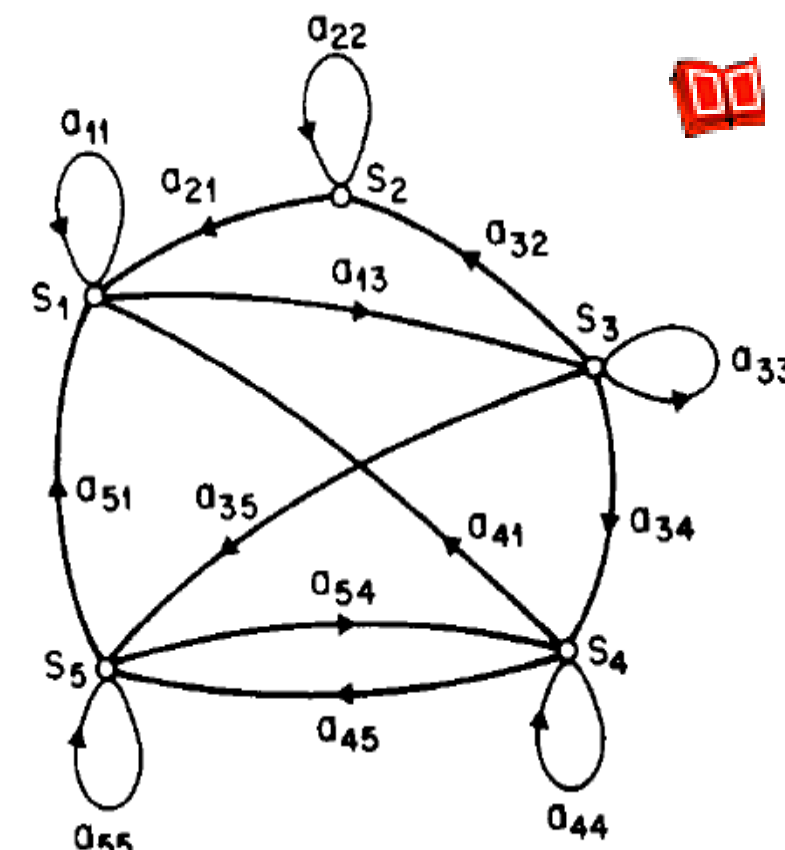    - **A** is the $N \times N$ matrix of state transition probabilities
      $$P(q_t=j|q_{t-1}=i, q_{t-2}=k, \ldots\ldots)=P(q_t=j|q_{t-1}=$$
    - π is the vector of initial state probab
      $$\pi_j=P(q_0=j)$$
  - The output for any given state is an observable event (deterministic)
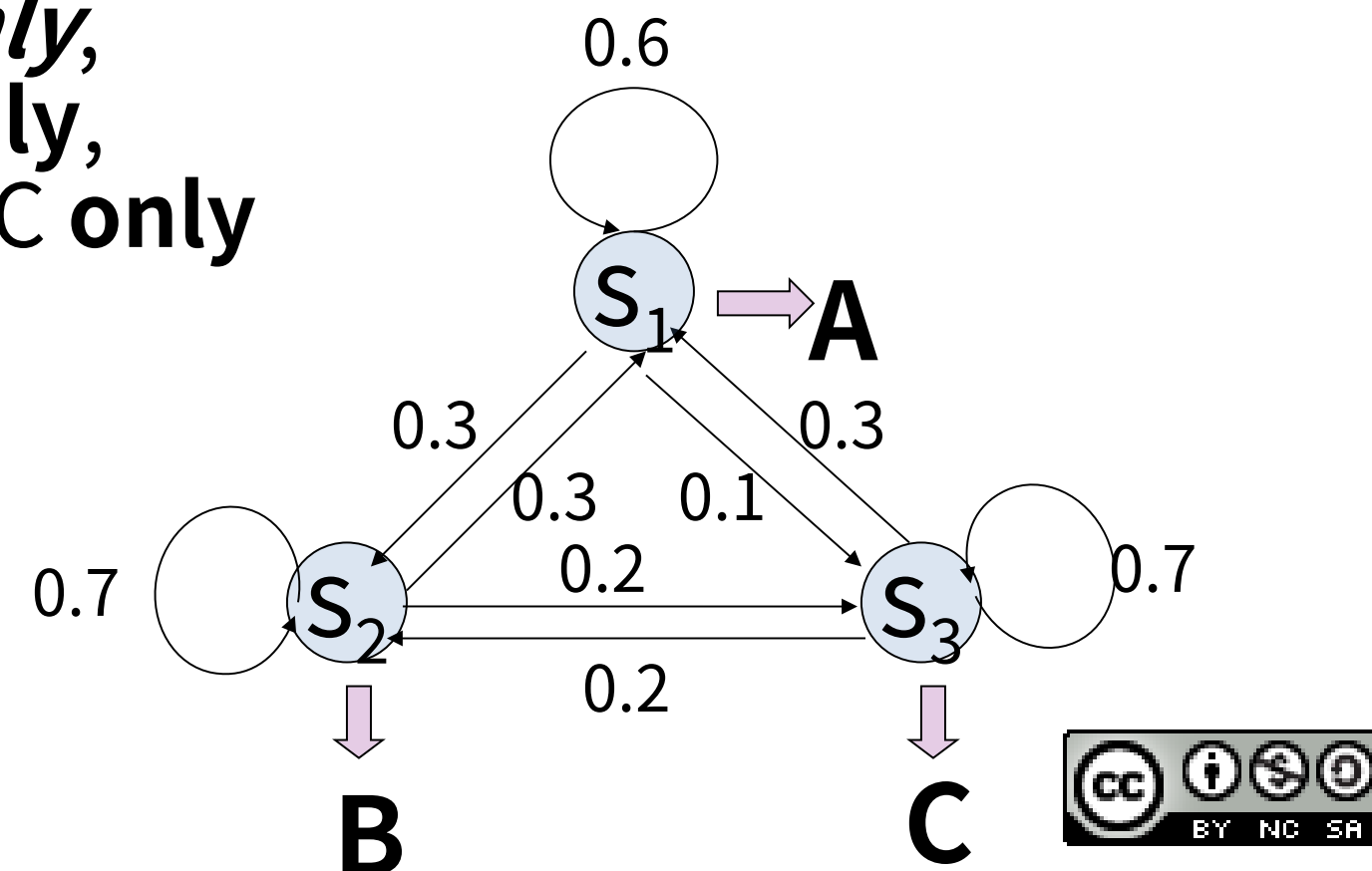  - The output of the process is a seque of observable events



A Markov chain with 5 states (labeled $S_1$ to $S_5$) with state transitions.

- **An example : a 3-state Markov Chain $\lambda$**
  - State 1 generates symbol A *only*,
    State 2 generates symbol B **only**,
    and State 3 generates symbol C **only**

$$A = \begin{bmatrix} 0.6 & 0.3 & 0.1 \\ 0.1 & 0.7 & 0.2 \\ 0.3 & 0.2 & 0.5 \end{bmatrix}$$

$$\pi = \begin{bmatrix} 0.4 & 0.5 & 0.1 \end{bmatrix}$$

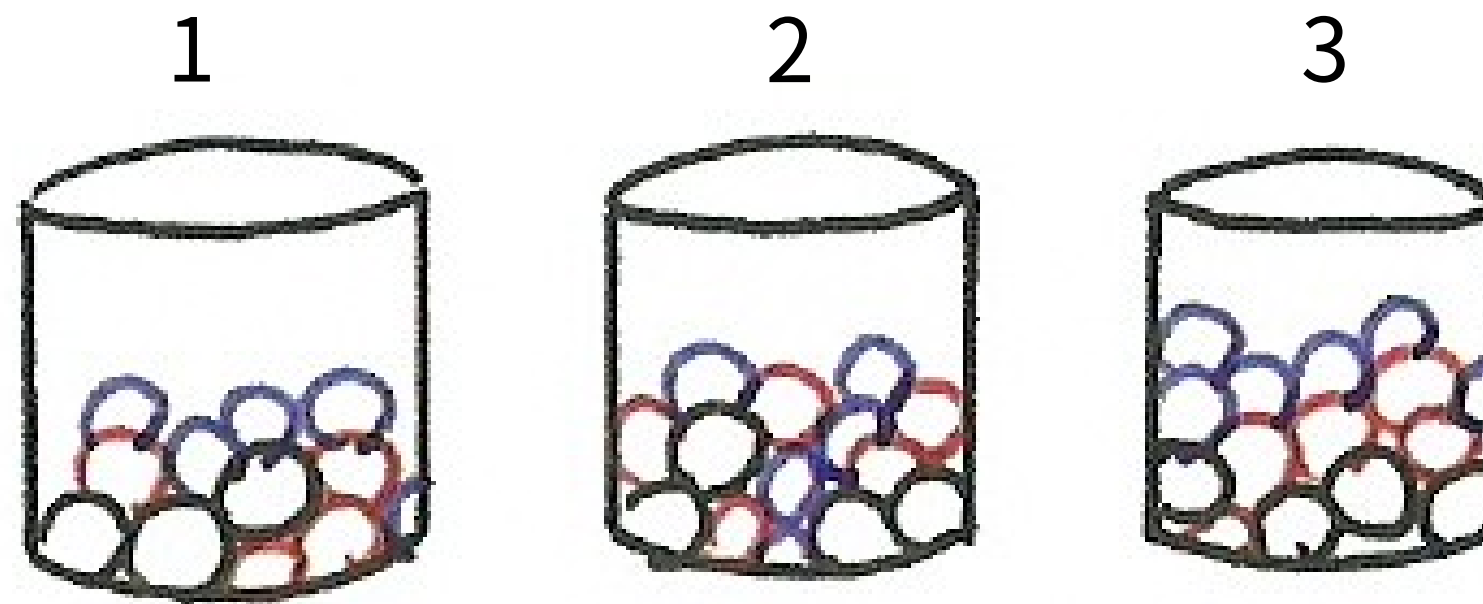  - Given a sequence of observed symbols **O**={CABBCABC}, the **only one** corresponding state sequence is $\{S_3 S_1 S_2 S_2 S_3 S_1 S_2 S_3\}$, and the corresponding probability is
    $P(\mathbf{O}|\lambda) = P(q_0 = S_3) \ P(S_1/S_3) P(S_2/S_1) P(S_2/S_2) P(S_3/S_2) P(S_1/S_3) P(S_2/S_1) P(S_3/S_2)$
    $=0.1 \square 0.3 \square 0.3 \square 0.7 \square 0.2 \square 0.3 \square 0.3 \square 0.2 = 0.00002268$

3

# Hidden Markov Model

- **HMM, an extended version of Markov Model**
  - The observation is **a probabilistic function (discrete or continuous) of a state** instead of an one-to-one correspondence of a state
  - The model is a doubly embedded stochastic process with an underlying stochastic process that is not directly observable (hidden)
    - What is hidden? *The State Sequence*
      *According to the observation sequence, we never know which state sequence generates it*
- **Elements of an HMM {S,A,B,$\pi$}**
  - S is a set of $N$ states
  - **A** is the $M \times N$ matrix of state transition probabilities
  - B is a set of $N$ probability functions, each describing the observation probability with respect to a state
  - $\pi$ is the vector of initial state probabilities

4

# Simplified HMM

1　　　2　　　3



RGBGGBBGRRR......

- **Two types of HMM's according to the observation functions**

  Discrete and finite observations :
  - The observations that <span style="color:red">all</span> distinct states generate are finite in number
    $\mathbf{V}=\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \ldots\ldots, \mathbf{v}_M\}, \mathbf{v}_k \in \mathbf{R}^D$
  - the set of observation probability distributions B=$\{b_j(\mathbf{v}_k)\}$ is defined as $b_j(\mathbf{v}_k)=P(\mathbf{o}_t=\mathbf{v}_k|\boldsymbol{q}_t=j)$, $1\leq k\leq M, 1\leq j\leq N$
    $\mathbf{o}_t$: *observation at time t*, $\boldsymbol{q}_t$: *state at time t*
    ⬜ *for state j, $b_j(\mathbf{v}_k)$* consists of ***only M probability values***

  Continuous and infinite observations :
  - The observations that <span style="color:red">all</span> distinct states generate are infinite and continuous, $\mathbf{V}=\{\mathbf{v}| \mathbf{v}\in \mathbf{R}^D\}$
  - the set of observation probability distributions B=$\{b_j(\mathbf{v})\}$ is defined as $b_j(\mathbf{v})=P(\mathbf{o}_t=\mathbf{v}|\boldsymbol{q}_t=j)$, $1\leq j\leq N$
    ⬜ *$b_j(\mathbf{v})$ is a **continuous probability density function** and is often assumed to be a mixture of Gaussian distributions*

$$b_j(v)=\sum_{k=1}^{M}c_{jk}\frac{1}{(\sqrt{2\pi})^D|\Sigma_{jk}|^{\frac{1}{2}}}\exp\left(-\frac{1}{2}(v-\mu_{jk})^T\Sigma_{jk}^{-1}(v-\mu_{jk})\right)=\sum_{k=1}^{M}c_{jk}b_{jk}(v)$$

- **An example : a 3-state discrete HMM $\lambda$**
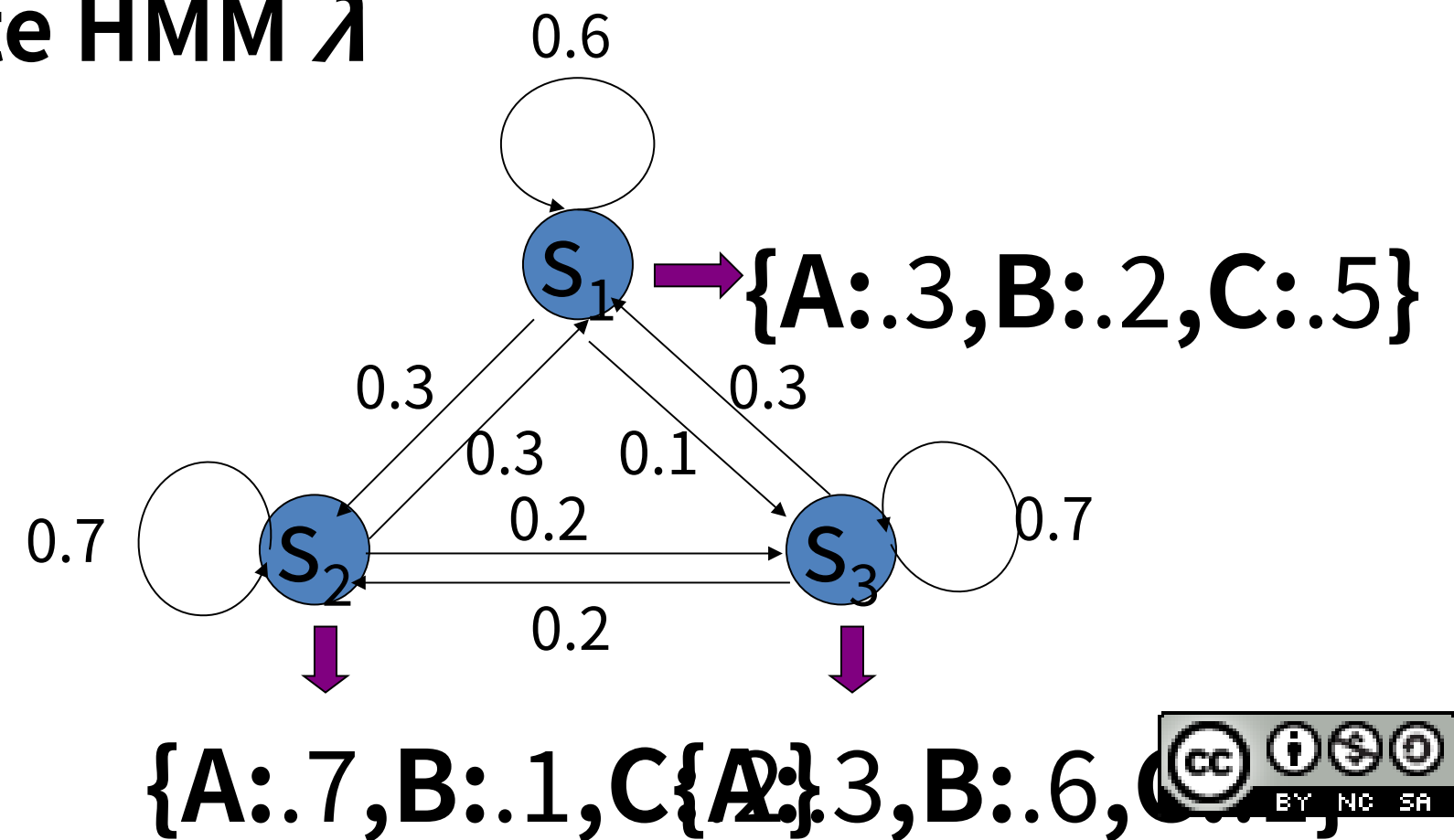
$$\mathbf{A} = \begin{bmatrix} 0.6 & 0.3 & 0.1 \\ 0.1 & 0.7 & 0.2 \\ 0.3 & 0.2 & 0.5 \end{bmatrix}$$

$b_1(\mathbf{A}) = 0.3, b_1(\mathbf{B}) = 0.2, b_1(\mathbf{C}) = 0.5$

$b_2(\mathbf{A}) = 0.7, b_2(\mathbf{B}) = 0.1, b_2(\mathbf{C}) = 0.2$

$b_3(\mathbf{A}) = 0.3, b_3(\mathbf{B}) = 0.6, b_3(\mathbf{C}) = 0.1$

$\pi = \begin{bmatrix} 0.4 & 0.5 & 0.1 \end{bmatrix}$



**{A:.3,B:.2,C:.5}**

**{A:.7,B:.1,C{A}.3,B:.6,C:.1}**

- Given a sequence of observations O={ABC}, there are **27 possible** corresponding state sequences, and therefore the corresponding probability is

$$P(\overline{\mathbf{O}}|\lambda) = \sum_{i=1}^{27} P(\overline{\mathbf{O}}, \mathbf{q}_i|\lambda) = \sum_{i=1}^{27} P(\overline{\mathbf{O}}|\mathbf{q}_i, \lambda) P(\mathbf{q}_i|\lambda), \quad \mathbf{q}_i : \text{state sequence}$$

$eg.\ \text{when } \mathbf{q}_i = \{S_2 S_2 S_3\},\ P(\overline{\mathbf{O}}|\mathbf{q}_i, \lambda) = P(\mathbf{A}|S_2) P(\mathbf{B}|S_2) P(\mathbf{C}|S_3) = 0.7 * 0.1 * 0.1 = 0.007$

$P(\mathbf{q}_i|\lambda) = P(q_0 = S_2) P(S_2|S_2) P(S_3|S_2) = 0.5 * 0.7 * 0.2 = 0.07$

7

# Hidden Markov Model

- **Three Basic Problems for HMMs**
  **Given an observation sequence** $O=(o_1,o_2,.....,o_T)$, **and an HMM**
  $\lambda =(A,B,\pi)$

  – Problem *1*:
  How to *efficiently* compute $P(O|\lambda)$ ?
  ⊠ *Evaluation problem*

  – Problem *2*:
  How to choose an optimal state sequence $q=(q_1,q_2,......, q_T)$ ?
  ⊠ *Decoding Problem*

  – Problem *3*:
  Given some observations O for the HMM $\lambda$, how to adjust the model parameter $\lambda =(A,B,\pi)$ to maximize $P(O|\lambda)$?
  ⊠ *Learning /Training Problem*

## Basic Problem 1 for HMM

$$1 \rightarrow 2 \rightarrow \cdots \cdots \rightarrow N \qquad \lambda = (A, B, \pi)$$

$\overline{O} = o_1 o_2 o_3 \cdots \cdots o_t \cdots \cdots o_T$ observation sequence

$\overline{q} = q_1 q_2 q_3 \cdots \cdots q_t \cdots \cdots q_T$ state sequence

⊠**Problem 1:** Given $\lambda$ and $\overline{O}$,

find $P(\overline{O}|\lambda) = $ Prob[observing $\overline{O}$ given $\lambda$]

⊠**Direct Evaluation: considering all possible state sequence** $\overline{q}$

$$P(\overline{O}|\lambda) = \sum_{\text{all } \overline{q}} P(\overline{O}, \overline{q}|\lambda) = \sum_{\text{all } \overline{q}} P(\overline{O}|\overline{q}, \lambda) P(\overline{q}|\lambda)$$

$$P(\overline{O}|\overline{q}, \lambda)$$

⇧

$$P(\overline{O}|\lambda) = \sum_{\text{all } \overline{q}} \left( [b_{q_1}(o_1) \cdot b_{q_2}(o_2) \cdot \cdots \cdots b_{q_T}(o_T)] \cdot \right.$$

$$\left. [\pi_{q_1} \cdot a_{q_1 q_2} \cdot a_{q_2 q_3} \cdot \cdots \cdots a_{q_{T-1} q_T}] \right)$$

⇩

$$P(\overline{q}|\lambda)$$

total number of different $\overline{q}$ : $N^T$

huge computation requirements

# Basic Problem 1 for HMM

- **Forward Algorithm: defining a forward variable $\alpha_t(i)$**

$$\alpha_t(i) = P(o_1 o_2 \cdots\cdots o_t, q_t = i|\lambda)$$

$$= \text{Prob}[\text{observing } o_1 o_2 \cdots o_t, \text{ state } i \text{ at time } t|\lambda]$$

- Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

- Induction

$$\alpha_{t+1}(j) = [\sum_{i=1}^{N} \alpha_t(i)a_{ij}] b_j(o_{t+1})$$
$$1 \leq j \leq N$$
$$1 \leq t \leq T-1$$

- Termination

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

*See Fig. 6.5 of Rabiner and Juang*

- All state sequences, regardless of how long previously, merge to the N state at each time instant t

10

# Basic Problem 1

# Basic Problem 1



$$\alpha_t(i) = P(o_1 o_2 \ldots \ldots o_t, q_t = i | \lambda)$$

# Basic Problem 1



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

$$1 \le j \le N$$
$$1 \le t \le T\text{-}1$$

Forward Algorithm

**Basic Problem 2 for HMM**

- **Problem 2:** Given $\lambda$ and $\overline{O} = o_1 o_2 \cdots o_T$, find a best state sequence $\overline{q} = q_1 q_2 \cdots q_T$

- **Backward Algorithm :** defining a backward variable $\beta_t(i)$

  $\beta_t(i) = P(o_{t+1}, o_{t+2}, \cdots, o_T | q_t = i, \lambda)$

  $\quad\quad = \text{Prob}[\text{observing } o_{t+1}, o_{t+2}, \cdots, o_T | \text{state } i \text{ at time } t, \lambda]$

  - Initialization

    $\beta_T(i) = 1, \quad 1 \leq i \leq N \quad\quad ( \beta_{T-1}(i) = \sum_{j=1}^{N} a_{ij} b_j(o_T) )$

  - Induction

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

    $\quad\quad t = T-1, T-2, \cdots, 2, 1, \quad\quad 1 \leq i \leq N$

    *See Fig. 6.6 of Rabiner and Juang*

- **Combining Forward/Backward Variables**

  $P(\overline{O}, q_t = i | \lambda)$

  $= \text{Prob} [\text{observing } o_1, o_2, \cdots, o_t, \cdots, o_T, q_t = i | \lambda ]$

  $= \alpha_t(i) \beta_t(i)$

  $P(\overline{O}|\lambda) = \sum_{i=1}^{N} P(\overline{O}, q_t = i | \lambda) = \sum_{i=1}^{N} [\alpha_t(i) \beta_t(i)]$

14

# Basic Problem 2

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T \mid q_t = i, \lambda)$$

$i$

$t$    $t+1$    $\cdots$   $T$

15

# Basic Problem 2



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

$$t = T\text{-}1, T\text{-}2, \cdots, 2, 1, \quad 1 \le i \le N$$

Backward Algorithm

# Basic Problem 2

$\beta_t(i)$

$i$

$\alpha_t(i)$

$t$

$P(\overline{O}, q_t = i \mid \lambda)$

$= \text{Prob} [\text{observing } o_1, o_2, \ldots, o_t, \ldots, o_T, q_t = i \mid \lambda]$

$= \alpha_t(i)\beta_t(i)$

$$\boxtimes\left(\boxtimes, \boxtimes, \boxtimes\right) = \boxtimes\left(\boxtimes, \boxtimes\right)\boxtimes(\boxtimes \mid \boxtimes, \boxtimes)$$

$$/\!/ \qquad\qquad /\!/ \qquad /\!/ \quad (\boxtimes \perp \boxtimes)$$

$$\boxtimes(\overline{\boxtimes}, \boxtimes_{\boxtimes} = \boxtimes \mid \boxtimes) \boxtimes_{\boxtimes}(\boxtimes) \quad \boxtimes(\boxtimes \mid \boxtimes)$$

$$/\!/$$

$$\boxtimes_{\boxtimes}(\boxtimes)$$

# Basic Problem 2

$$P(\overline{O}|\lambda) = \sum_{i=1}^{N} P(\overline{O}, q_t = i \mid \lambda) = \sum_{i=1}^{N} \left[ \alpha_t(i)\beta_t(i) \right]$$

$P(\overline{O}/\lambda)$ )

$\alpha_t(i)\beta_t(i)$ i)

$i$

$t$

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

☒ **Approach 1 – Choosing state $q_t$ individually as the most likely state at time t**

- Define a new variable $\gamma_t(i) = P(q_t = i \mid \overline{O}, \lambda)$

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{\displaystyle\sum_{i=1}^{N} \alpha_t(i)\beta_t(i)} = \frac{P(\overline{O}, q_t = i \mid \lambda)}{P(\overline{O} \mid \lambda)}$$

- Solution

$$q_t^* = \arg\max_{1 \le i \le N} [\gamma_t(i)], \quad 1 \le t \le T$$

in fact

$$q_t^* = \arg\max_{1 \le i \le N} [P(\overline{O}, q_t = i \mid \lambda)]$$

$$= \arg\max_{1 \le i \le N} [\alpha_t(i)\beta_t(i)]$$

- Problem

maximizing the probability at each time t individually

$q^* = q_1^* q_2^* \cdots q_T^*$ may not be a valid sequence

(e.g. $a_{q_t^* q_{t+1}^*} = 0$)

- **Approach 2 – Viterbi Algorithm - finding the single best sequence $\overline{q}^* = q_1^* q_2^* \cdots q_T^*$**

  - Define a new variable $\delta_t(i)$

    $\delta_t(i) = \max\limits_{q_1,q_2,\cdots q_{t-1}} P[q_1,q_2,\cdots q_{t-1}, q_t = i, o_1,o_2,\cdots,o_t|\lambda]$

    = the highest probability along a certain single path ending at state i at time t for the first t observations, given $\lambda$

  - Induction

    $\delta_{t+1}(j) = \max\limits_{i} [\delta_t(i)a_{ij}] \cdot b_j(o_{t+1})$

  - Backtracking

    $\psi_{t+1}(j) = \arg\max\limits_{1 \leq i \leq N} [\delta_t(i)a_{ij}]$

    the best previous state at t- 1 given at state j at time t

    keeping track of the best previous state for each j and t

20

# Viterbi Algorithm

$$\delta_{t+1}(j) = \max_{i} [\delta_t(i)a_{ij}] \bullet b_j(o_{t+1})$$



$\delta_t(i)$

$\delta_{t+1}(j)$

$\delta_t(i)$

$t \quad t+1$

$$\delta_t(i) = \max_{q_1,q_2,\ldots q_{t-1}} P[q_1,q_2,\ldots q_{t-1}, q_t = i, o_1,o_2,\ldots,o_t | \lambda]$$

# Viterbi Algorithm



**Path backtracking**

# Basic Problem 2 for HMM

- **Complete Procedure for Viterbi Algorithm**

  - Initialization

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

  - Recursion

    $$\delta_{t+1}(j) = \max_{1 \leq i \leq N} [\delta_t(i) a_{ij}] \cdot b_j(o_t)$$

    $$2 \leq t \leq T, \quad 1 \leq j \leq N$$

    $$\psi_{t+1}(j) = \arg\max_{1 \leq i \leq N} [\delta_t(i) a_{ij}]$$

    $$2 \leq t \leq T, \quad 1 \leq j \leq N$$

  - Termination

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

    $$q_T^* = \arg\max_{1 \leq i \leq N} [\delta_T(i)]$$

  - Path backtracking

    $$q_t^* = \psi_{t+1}(q^*_{t+1}) \ , \quad t = T-1, t-2, \cdots..2, 1$$

☒**Application Example of Viterbi Algorithm**
  - Isolated word recognition

$$\lambda_0 = (A_0, B_0, \pi_0)$$
$$\lambda_1 = (A_1, B_1, \pi_1)$$
$$\vdots$$
$$\lambda_n = (A_n, B_n, \pi_n)$$

observation
$$\overline{O} = (o_1, o_2, \ldots o_T)$$
$$k^* = \underset{1 \le i \le n}{\arg\max} P[\overline{O} \mid \lambda_i] \approx \underset{1 \le i \le n}{\arg\max} P[P^* \mid \lambda_i]$$

⇧                    ⇧

Basic Problem 1          Basic Problem 2
Forward Algorithm      Viterbi Algorithm
(for all paths)       (for a single best path)

-The model with the highest probability for the most probable
 usually also has the highest probability for all possible path

24    24

- **Problem 3:** Give $\overline{O}$ and an initial model $\lambda=(A,B,\pi)$, adjust $\lambda$ to maximize $P(\overline{O}|\lambda)$

  - Baum-Welch Algorithm (Forward-backward Algorithm)
  - Define a new variable

$$
\varepsilon_t(\,i,\,j\,) = P(q_t = i,\, q_{t+1} = j \mid \overline{O},\, \lambda)
$$

$$
= \frac{\alpha_t(i)\, a_{ij}\, b_j(o_{t+1})\beta_{t+1}(j)}{\sum\limits_{i=1}^{N}\sum\limits_{j=1}^{N}[\alpha_t(i)a_{ij}\, b_j(o_{t+1})\beta_{t+1}(j)]}
$$

$$
= \frac{\text{Prob}[\overline{O},\, q_t = i,\, q_{t+1} = j|\lambda]}{P(\overline{O}|\lambda)}
$$

*See Fig. 6.7 of Rabiner and Juang*

  - Recall $\gamma_t(i) = P(q_t = i \mid \overline{O},\, \lambda)$

$$
\sum_{t=1}^{T-1}\gamma_t(i) = \text{expected number of times that state i is visited in } \overline{O} \text{ from } t = 1 \text{ to } t = T-1
$$

$$
= \text{expected number of transitions from state i in } \overline{O}
$$

$$
\sum_{t=1}^{T-1}\varepsilon_t(i,\,j) = \text{expected number of transitions from state i to state j in } \overline{O}
$$

# Basic Problem 3

$$\alpha_t(i) \, a_{ij} \, b_j(o_{t+1}) \beta_{t+1}(j)$$



$\beta_{t+1}(j)$

$j$

$i$

$\alpha_t(i)$

$t \quad t+1$

# Basic Problem 3

$$\alpha_t(i)a_{ij}b_j(o_{t+1})\beta_{t+1}(j)$$



$\beta_{t+1}(j)$

$j$

$i$

$\alpha_t(i)$

$t \quad t+1$

27

$$\gamma_t(i) = \frac{\alpha_t(i)\,\beta_t(i)}{\sum_{j=1}^{N}[\alpha_t(j)\,\beta_t(j)]} = \frac{P(\overline{O}, q_t = i \mid \lambda)}{P(\overline{O} \mid \lambda)} = P(q_t = i \mid \overline{O}, \lambda)$$

$$\xi_t(i,j) = \frac{\alpha_t(i)\,a_{ij}\,b_j(O_{t+1})\,\beta_{t+1}(j)}{\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_t(i)\,a_{ij}\,b_j(O_{t+1})\,\beta_{t+1}(j)}$$

$$= \frac{P(\overline{O}, q_t = i, q_{t+1} = j \mid \lambda)}{P(\overline{O} \mid \lambda)} = P(q_t = i, q_{t+1} = j \mid \overline{O}, \lambda)$$

# Basic Problem 3

0.000009

1.95

0.000006

0.0001

0.0002

0.0003

0.15    0.21    0.26

1    2    3                    T-1  T

out of T-1

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \epsilon_t(i, j)/(T-1)}{\sum_{t=1}^{T-1} \gamma_t(i)/(T-1)}$$

=159

$$\overline{\otimes}_{\boxtimes\boxtimes} = \frac{1.95/69}{10.59/69}$$

29

- Results

$$\overline{\pi}_i = \gamma_1(i)$$

$$\overline{a}_{ij} = \frac{\displaystyle\sum_{t=1}^{T-1} \varepsilon_t(i, j)}{\displaystyle\sum_{t=1}^{T-1} \gamma_t(i)} \ - \ - \ -$$

$$\overline{b}_j(k) = \text{Prob}[o_t = v_k \,|\, q_t = j\,] = \frac{\displaystyle\sum_{\substack{t=1 \\ o_t = v_k}}^{T} \gamma_t(j)}{\displaystyle\sum_{t=1}^{T} \gamma_t(j)} \ - \ - \ -$$

(for discrete HMM)

- **Continuous Density HMM**

$$b_j(o) = \sum_{k=1}^{M} c_{jk} N(o;\, \mu_{jk},\, U_{jk})$$

N( ): Multi-variate Gaussian

$\mu_{jk}$: mean vector for the k-th mixture component

$U_{jk}$: covariance matrix for the k-th mixture component

$$\sum_{k=1}^{M} c_{jk} = 1 \text{ for normalization}$$

- **Continuous Density HMM**

  - Define a new variable

    $\gamma_t(j, k) = \gamma_t(j)$ but including the probability of $o_t$ evaluated in the k-th mixture component out of all the mixture components

    $$= \left[ \frac{\alpha_t(j)\beta_t(j)}{\sum\limits_{j=1}^{N} \alpha_t(j)\beta_t(j)} \right] \left[ \frac{c_{jk} N(o_t;\ \mu_{jk},\ U_{jk})}{\sum\limits_{m=1}^{M} c_{jm} N(o_t;\ \mu_{jm},\ U_{jm})]} \right]$$



  - Results

    $$\overline{c}_{jk} = \frac{\sum\limits_{t=1}^{T} \gamma_t(j, k)}{\sum\limits_{t=1}^{T} \sum\limits_{k=1}^{M} \gamma_t(j, k)}$$

    *See Fig. 6.9 of Rabiner and Juang*

31

- **Continuous Density HMM**

$$\bar{\mu}_{jk} = \frac{\sum\limits_{t=1}^{T} [\gamma_t(j, k) \cdot o_t]}{\sum\limits_{t=1}^{T} \gamma_t(j, k)}$$

$$\bar{U}_{jk} = \frac{\sum\limits_{t=1}^{T} [\gamma_t(j, k)(o_t - \mu_{jk})(o_t - \mu_{jk})']}{\sum\limits_{t=1}^{T} \gamma_t(j, k)}$$

- **Iterative Procedure**

$$\lambda = (A, B, \pi) \longrightarrow \bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$$

$$\bar{O} = o_1 o_2 \cdots o_T$$

  – It can be shown (by EM Theory (or EM Algorithm))

  $P(\bar{O}|\bar{\lambda}) \geq P(\bar{O}|\lambda)$ after each iteration

32

$$\overline{\mu}_{jk} = \frac{\sum_{t=1}^{T} [\boxed{\gamma_t(j,k)} \cdot o_t]}{\boxed{\sum_{t=1}^{T} \gamma_t(j,k)}}$$

$$\overline{U}_{jk} = \frac{\sum_{t=1}^{T} [\gamma_t(j,k)\boxed{(o_t - \mu_{jk})(o_t - \mu_{jk})'}]}{\sum_{t=1}^{T} \gamma_t(j,k)}$$

prob. density function

# Basic Problem 3

$$\bar{U} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} = E(\begin{bmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \\ \vdots \\ \vdots \end{bmatrix} [x_1 - \bar{x}_1, x_2 - \bar{x}_2, \cdots])$$

$$\bar{u}_{lm} = E[(x_l - \bar{x}_l)(x_m - \bar{x}_m)]$$

# Basic Problem 3 for HMM

- No closed-form solution, but approximated iteratively
- An initial model is needed-model initialization
- May converge to local optimal points rather than global optimal point
  - heavily depending on the initialization
- Model training

| Model Initialization: Segmental K-means | Model Re-estimation: Baum-Welch |
|---|---|

# Basic Problem 3



$P(O/\lambda)$

P

Global optimum

Local optimum

$\mu_{jkn}$

- **An Efficient Approach for Data Compression**
  - replacing a set of real numbers by a finite number of bits
- **An Efficient Approach for Clustering Large Number of Sample Vectors**
  - grouping sample vectors into clusters, each represented by a single vector (codeword)
- **Scalar Quantization**
  - replacing a single real number by an R-bit pattern
  - a mapping relation

$$J_k$$

$$v_k \qquad A =$$

$$A = \\ m_L$$

$$- \\ A = \\ m$$

$$S = \bigcup_{k=1}^{L} J_k , \ V = \{ v_1 , v_2 , \ldots, v_L \}$$

$$Q : S \to V$$

$$Q(x[n]) = v_k \ \text{if} \ x[n] \in J_k$$

$$L = 2^R$$

Each $v_k$ represented by an R-bit pattern

– Quantization characteristics (codebook)

$\{ J_1 , J_2 , \ldots, J_L \}$ and $\{ v_1 , v_2 , \ldots, v_L \}$

designed considering at least

1. error sensitivity
2. probability distribution of x[n]

37

# Vector Quantization

## Scalar Quantization ：Pulse Coded Modulation (PCM)



quantization error

011
010
001
000
100
101
110
111

100 100 101 110

# Vector Quantization

$P_x(x)$

## *2-dim Vector Quantization (VQ)*

Example:

$\overline{x}_n = ( x[n] , x[n+1] )$

$S = \{\overline{x}_n = (x[n] , x[n+1] ) ; |x[n]| < A, |x[n+1]| < A\}$

- **VQ**
  - S divided into L 2-dim regions
    $\{ J_1 , J_2 , \ldots, J_k , \ldots J_L \}$ $S = \bigcup_{k=1}^{L} J_k$

    each with a representative

    vector $\overline{v}_k \in J_k$, $V = \{ \overline{v}_1 , \overline{v}_2 , \ldots, \overline{v}_L \}$
  - $Q : S \Rightarrow V$
    $Q(\overline{x}_n) = \overline{v}_k$ if $\overline{x}_n \in J_k$
    $L = 2^R$
    each $\overline{v}_k$ represented by an R-bit pattern

- Considerations
  1. error sensitivity may depend on x[n], x[n+1] jointly
  2. distribution of x[n] , x[n+1] may be correlated statistically
  3. more flexible choice of $J_k$
- Quantization Characteristics (codebook)
  $\{ J_1 , J_2 , \ldots, J_L \}$ and $\{ \overline{v}_1 , \overline{v}_2 , \ldots, \overline{v}_L \}$

# Vector Quantization

# Vector Quantization



$J_k$

$V_k$

$(256)^2 = (2^8)^2 = 2^{16}$

$1024 = 2^{10}$

# Vector Quantization

## N-dim Vector Quantization

$x = (x_1, x_2, \ldots, x_N)$

$S = \{x = (x_1, x_2, \ldots, x_N),$

$S = \bigcup_{k=1}^{L} J_k$

$|x_k| \leq A, k = 1, 2, \ldots N\}$

$V = \{v_1, v_2, \ldots, v_L\}$

$Q : S \to V$

$Q(x) = v_k$ if $x \in J_k$

$L = 2^R$, each $v_k$ represented
   by an R-bit pattern

## Codebook Trained by a Large Training Set

· **Define distance measure between**
   **two vectors x, y**

$d(x, y) : S \times S \to R^+$ (non-negative

real numbers)

-desired properties

$d(x, y) \geq 0$

$d(x, x) = 0$

$d(x, y) = d(y, x)$

$d(x, y) + d(y, z) \geq d(x, z)$

examples :

$d(x, y) = \Sigma^i (x_i - y_i)^2$

$d(x, y) = \Sigma | x_i - y_i |$

$d(x, y) = (x-y)^t \Sigma^{-1}(x-y)$

# Distance Measures

$$d(\bar{x}, \bar{y}) = \sum_i \left| x_i - y_i \right| \quad \text{city block distance}$$



$$d(\bar{x}, \bar{y}) = (\bar{x} - \bar{y})^T \Sigma^{-1} (\bar{x} - \bar{y}) \quad \text{Mahalanobis distance}$$

$$\Sigma = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}, \quad d(\bar{x}, \bar{y}) = \sum_i (x_i - y_i)^2$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \cdots & 0 \\ & \sigma_2^2 & \cdots & \sigma_i^2 \\ \vdots & & \ddots \end{bmatrix}, \quad d(\bar{x}, \bar{y}) = \sum_i \frac{(x_i - y_i)^2}{\sigma_i^2}$$

45

- **K-Means Algorithm/Lloyd-Max Algorithm**

$$\boxed{\begin{array}{c}(1)\\ \text{Fixed } \{\, v_1, v_2,\\ \dots, v_L\,\}\\ \text{find best set of}\end{array}} \qquad \boxed{\begin{array}{c}(2)\\ \text{Fixed } \{\, J_1, J_2, \dots, J_L\,\}\\ \text{find best set of}\\ \{\, v_1, v_2, \dots, v_L\,\}\end{array}}$$

(1) $J_k = \{\, x \mid d(x, v_k) \le d(x, v_j),\ j \ne k \,\}$  $\{ v_1, v_2 \dots v_L \}$

$\qquad \to D = \sum_{all\,x} d(x, Q(x)) = $ min  $v_k = \dfrac{1}{M} \sum_{x \in J_k} x$

$\qquad \sum_{x \in J_k}$ nearest neighbor condition

(3) Convergence condition

$D = \sum_{k=1}^{L} D_k$

after each iteration D is reduced, but $D \ge 0$

$\mid D^{(m+1)} - D^{(m)} \mid < \in$, m : iteration

- **Iterative Procedure to Obtain Codebook from a Large Training Set**

(2) For each k

46

# Vector Quantization

# Vector Quantization (VQ)

- **K-means Algorithm may Converge to Local Optimal Solutions**
  - depending on initial conditions, not unique in general
- **Training VQ Codebook in Stages— LBG Algorithm**
  - step 1: Initialization. L = 1, train a 1-vector VQ codebook

$$\bar{v} = \frac{1}{N} \sum_j \bar{x}_j$$

  - step 2: Splitting.
        Splitting the L codewords into 2L codewords, L = 2L
    - example 1
      $$\bar{v}_k^{(1)} = \bar{v}_k(1 + \varepsilon)$$
      $$\bar{v}_k^{(2)} = \bar{v}_k(1 - \varepsilon)$$
    - example 2
      $$\bar{v}_k^{(1)} = \bar{v}_k$$
      $$\bar{v}_k^{(2)} : \text{the vector most far apart}$$

  - step 3: k-means Algorithm: to obtain L-vector codebook
  - step 4: Termination. Otherwise go to step 2

- **Usually Converges to Better Codebook**

48

# LBG Algorithm

# Initialization in HMM Training

- **An Often Used Approach— Segmental K-Means**
  - Assume an initial estimate of all model parameters (e.g. estimated by segmentation of training utterances into states with equal length)
    - For discrete density HMM

      $$b_j(k) = \frac{number\ of\ vectors\ in\ state\ j\ associated\ with\ codeword\ k}{total\ number\ of\ vectors\ in\ state\ j}$$

    - For continuous density HMM (M Gaussian mixtures per state)

      $\Rightarrow$ cluster the observation vectors within each state $j$ into a set of M clusters

      (e.g. with vector quantization)

      $c_{jm}$ = number of vectors classified in cluster $m$ of state $j$

      divided by number of vectors in state $j$

      $\mu_{jm}$ = sample mean of the vectors classified in cluster $m$ of state $j$

      $\Sigma_{jm}$ = sample covariance matrix of the vectors classified in cluster $m$ of state $j$

  - Step 1 : re-segment the training observation sequences into states based on the initial model by Viterbi Algorithm
  - Step 2 : Reestimate the model parameters (same as initial estimation)
  - Step 3: Evaluate the model score P($\underline{O}$|$\lambda$):
    If the difference between the previous and current model scores exceeds a threshold, go back to Step 1, otherwise stop and the initial model is obtained

50

# Segmental K-Means

- **An example for Continuous HMM**
  - 3 states and 4 Gaussian mixtures per state

- **An example for discrete HMM**
  - 3 states and 2 codewords



$$b_1(\mathbf{v}_1)=3/4, \ b_1(\mathbf{v}_2)=1/4$$

$$b_2(\mathbf{v}_1)=1/3, \ b_2(\mathbf{v}_2)=2/3$$

$$b_3(\mathbf{v}_1)=2/3, \ b_3(\mathbf{v}_2)=1/3$$

# 版權聲明

# 版權聲明

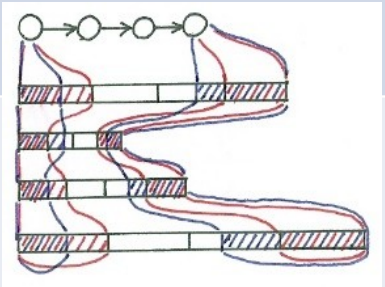| 頁碼 | 作品 | 版權標示 | 作者 / 來源 |
|---|---|---|---|
| 11 |  |  | Lawrence Rabiner, Biing-Hwang Juang / FUNDAMENTALS OF SPEECH RECOGNITION Chap. 6, Sec. 6.2 Discrete-Time Markov Processes, page 323, Prentice-Hall International, Inc. |
| 12 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 13 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 15 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |

# 版權聲明

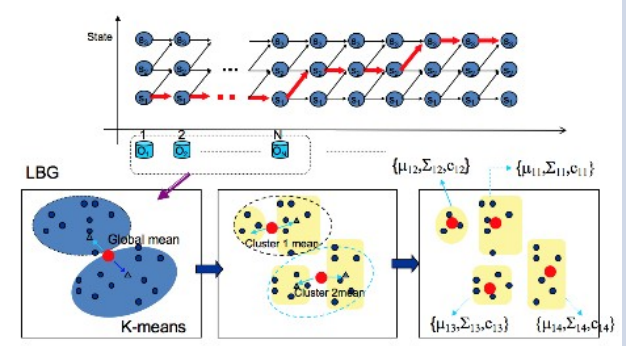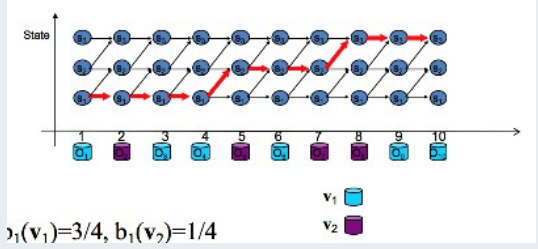| 頁碼 | 作品 | 版權標示 | 作者 / 來源 |
|---|---|---|---|
| 16 |  |  | Lawrence Rabiner, Biing-Hwang Juang / FUNDAMENTALS OF SPEECH RECOGNITION Chap. 6, Sec. 6.2 Discrete-Time Markov Processes, page 323, Prentice-Hall International, Inc. |
| 17 |  |  | 國立臺灣大學電機工程學系李琳山 教授。本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 18 |  |  | 國立臺灣大學電機工程學系李琳山 教授。本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 21 |  |  | 國立臺灣大學電機工程學系李琳山 教授。本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |

# 版權聲明

# 版權聲明

# 版權聲明

# 版權聲明

| 頁碼 | 作品 | 版權標示 | 作者 / 來源 |
|------|------|----------|-------------|
| 41 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 42 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 43 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 45 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |

# 版權聲明

# 版權聲明

| 頁碼 | 作品 | 版權標示 | 作者 / 來源 |
|------|------|----------|-------------|
| 52 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |
| 53 |  |  | 國立臺灣大學電機工程學系李琳山 教授。<br>本作品採用創用 CC 「姓名標示 - 非商業性 - 相同方式分享 3.0 臺灣」許可協議。 |