

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



【本著作除另有註明外，採取 [創用 CC](#)
[「姓名標示—非商業性—相同方式分享」臺灣 3.0](#)
[版授權釋出](#)】

- **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 $M \times N$ matrix of state transition probabilities

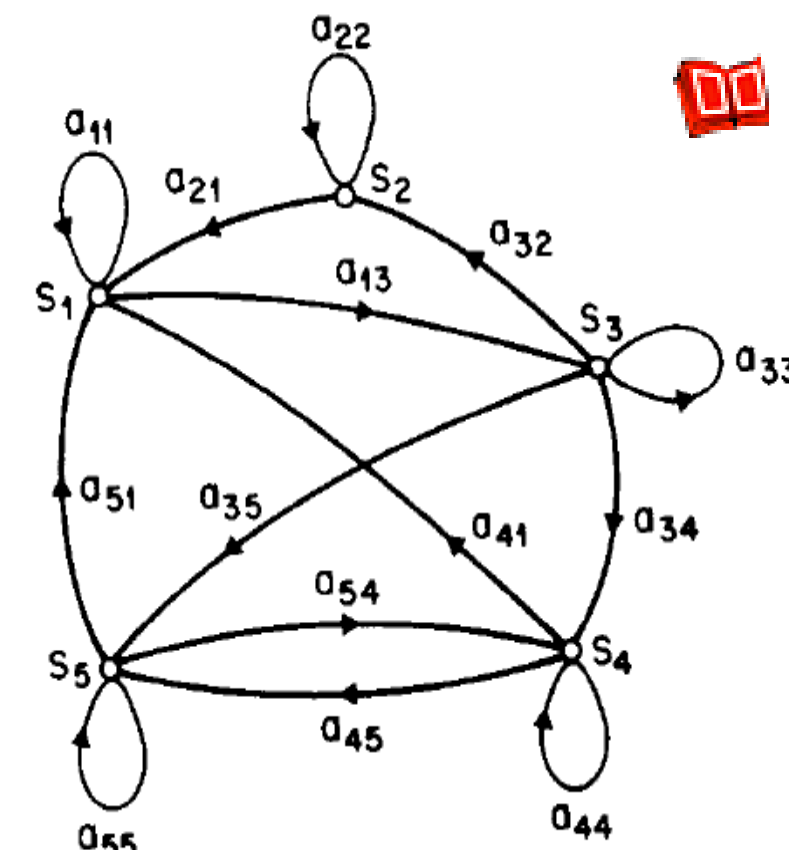
$$P(q_t=j | q_{t-1}=i, q_{t-2}=k, \dots) = P(q_t=j | q_{t-1}=i)$$

- π is the vector of initial state probal

$$\pi_j = P(q_0=j)$$

- The output for any given state is an observable event (deterministic)

- The output of the process is a sequence of observable events



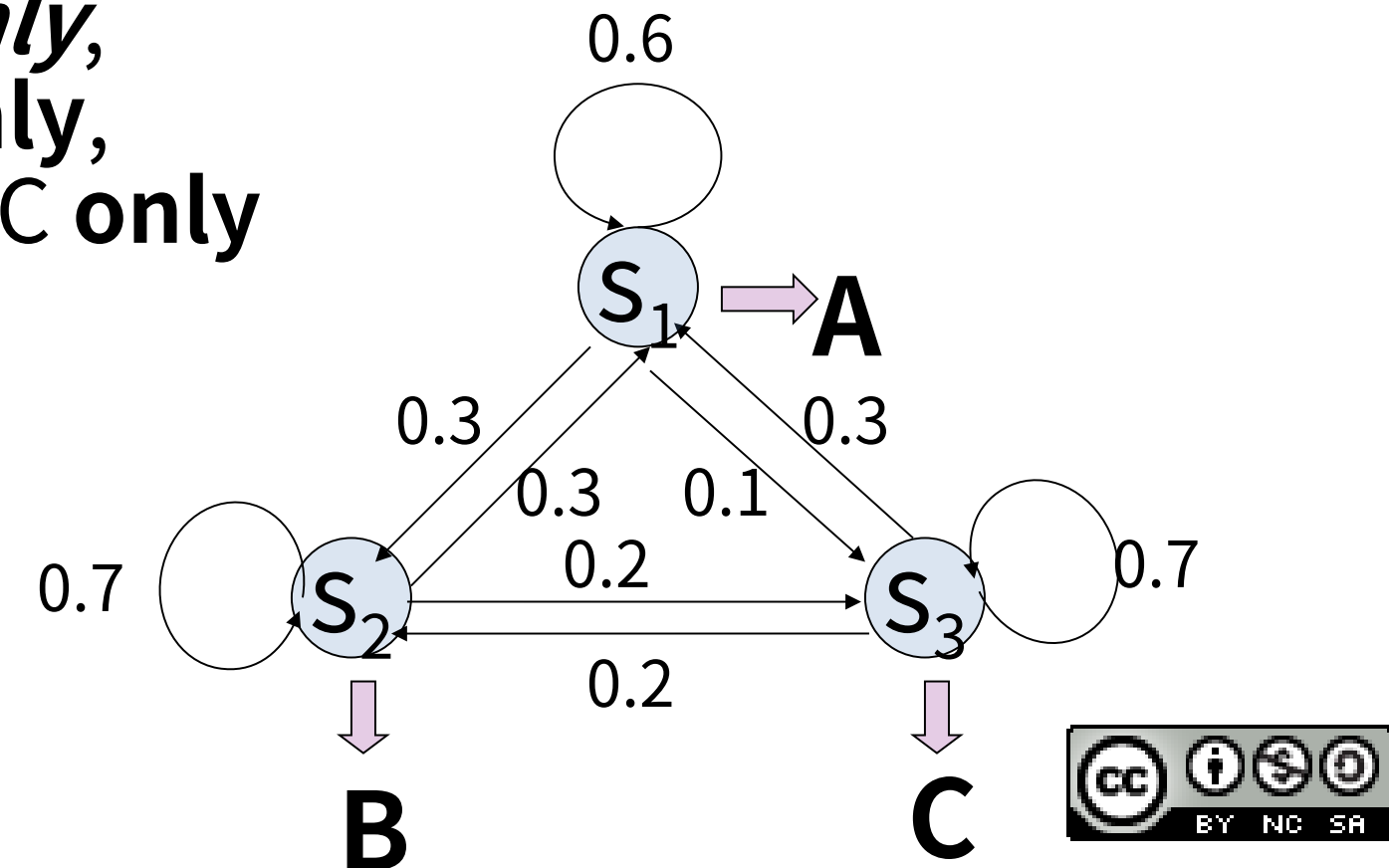
A Markov chain with 5 states (labeled S_1 to S_5) with state transitions.

- **An example : a 3-state Markov Chain λ**

- State 1 generates symbol A **only**,
State 2 generates symbol B **only**,
and State 3 generates symbol C **only**

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

$$\pi = [0.4 \quad 0.5 \quad 0.1]$$



- Given a sequence of observed symbols $\mathbf{O}=\{\text{CABBCABC}\}$, the **only one** corresponding state sequence is

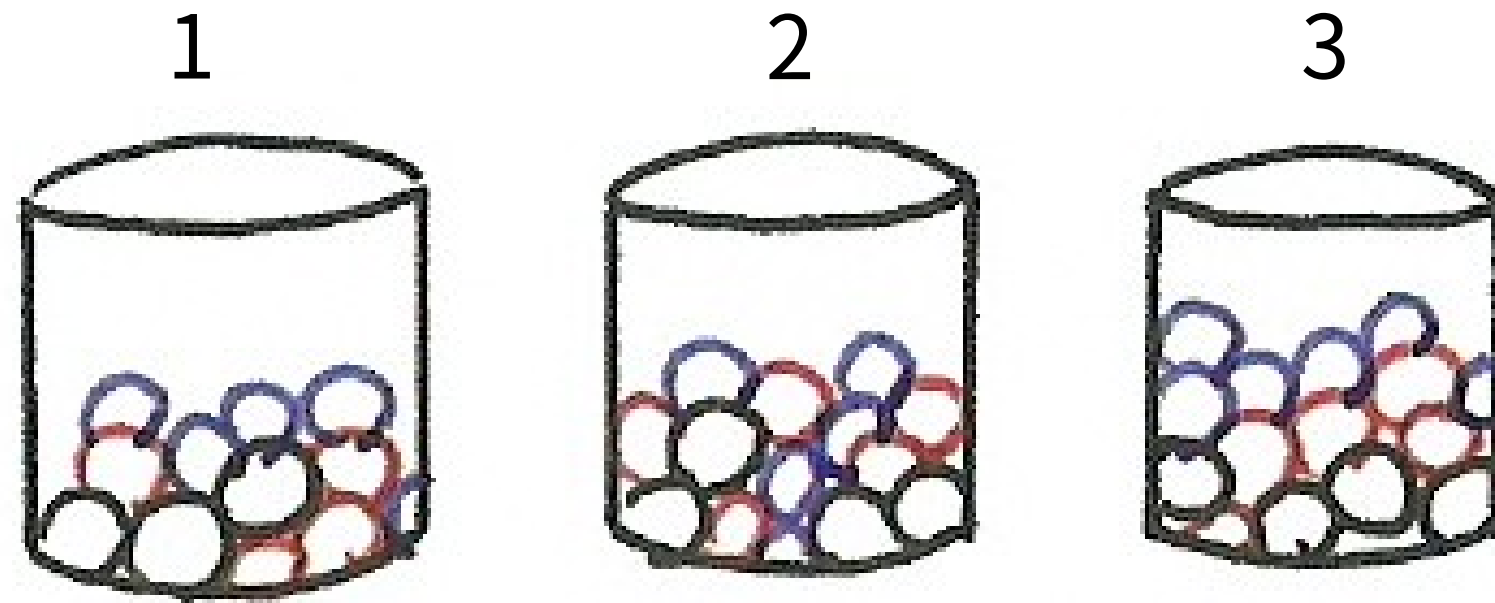
$\{S_3 S_1 S_2 S_2 S_3 S_1 S_2 S_3\}$, and the corresponding probability is

$$\begin{aligned} P(\mathbf{O}|\lambda) &= P(q_0=S_3) P(S_1/S_3) P(S_2/S_1) P(S_2/S_2) P(S_3/S_2) P(S_1/S_3) P(S_2/S_1) P(S_3/S_2) \\ &= 0.1 \times 0.3 \times 0.3 \times 0.7 \times 0.2 \times 0.3 \times 0.3 \times 0.2 = 0.00002268 \end{aligned}$$

- 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 $N \times N$ matrix of state transition probabilities
 - B is a set of N probability functions, each describing the observation probability with respect to a state
 - π is the vector of initial state probabilities

Simplified HMM

RGBGGBBGRRR.....



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

Discrete and finite observations :

- The observations that **all** distinct states generate are finite in number

$$\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_M\}, \mathbf{v}_k \in \mathbf{R}^D$$

- the set of observation probability distributions $B = \{b_j(\mathbf{v}_k)\}$ is defined as $b_j(\mathbf{v}_k) = P(\mathbf{o}_t = \mathbf{v}_k | \mathbf{q}_t = j)$, $1 \leq k \leq M$, $1 \leq j \leq N$

\mathbf{o}_t : observation at time t , \mathbf{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 **all** 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 $b_j(\mathbf{v}) = P(\mathbf{o}_t = \mathbf{v} | \mathbf{q}_t = j)$, $1 \leq j \leq N$

⊠ $b_j(\mathbf{v})$ is a **continuous probability density function** and is often assumed to be a mixture of Gaussian distributions

$$b_j(\mathbf{v}) = \sum_{k=1}^K c_{jk} \frac{1}{(\sqrt{2\pi})^{|\Sigma_{jk}|}} \exp\left(-\frac{1}{2}(\mathbf{v} - \boldsymbol{\mu}_{jk})^T \Sigma_{jk}^{-1} (\mathbf{v} - \boldsymbol{\mu}_{jk})\right) = \sum_{k=1}^K c_{jk} b_{jk}(\mathbf{v})$$

- An example : a 3-state discrete HMM λ

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

$$b_1(A)=0.3, b_1(B)=0.2, b_1(C)=0.5$$

$$b_2(A)=0.7, b_2(B)=0.1, b_2(C)=0.2$$

$$b_3(A)=0.3, b_3(B)=0.6, b_3(C)=0.1$$

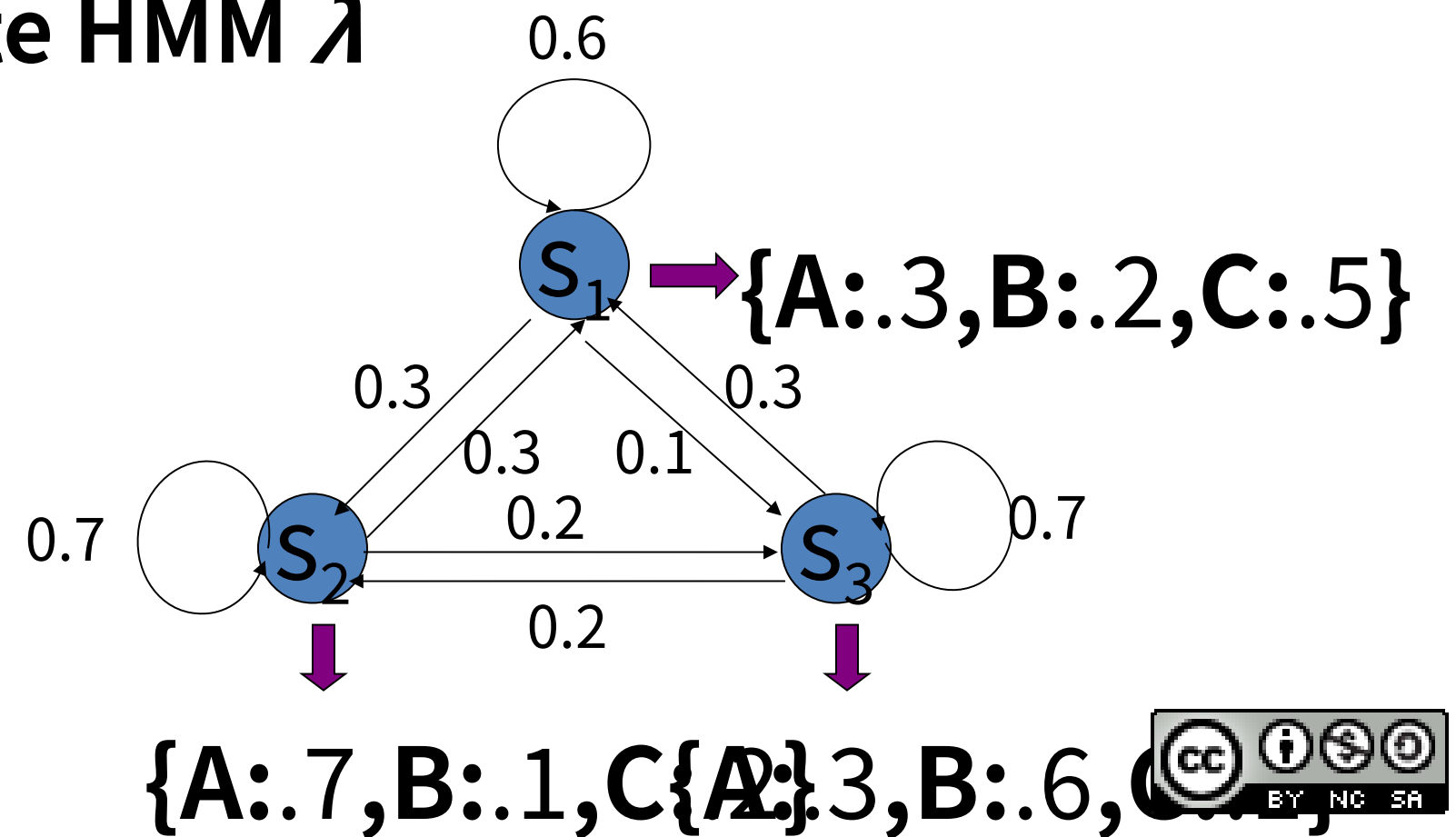
$$\pi = [0.4 \quad 0.5 \quad 0.1]$$

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

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

$$\text{e.g. when } \mathbf{q}_i = \{s_2 s_2 s_3\}, P(\bar{O}|\mathbf{q}_i, \lambda) = P(A|s_2)P(B|s_2)P(C|s_3) = 0.7 * 0.1 * 0.1 = 0.007$$

$$P(\mathbf{q}_i|\lambda) = P(q_0 = s_2)P(s_2|s_2)P(s_3|s_2) = 0.5 * 0.7 * 0.2 = 0.07$$



- Three Basic Problems for HMMs

Given an observation sequence $\bar{O}=(o_1,o_2,\dots,o_T)$, and an HMM

$\lambda=(A,B,\pi)$

- Problem 1:

How to *efficiently* compute $P(\bar{O}|\lambda)$?

☒ *Evaluation problem*

- Problem 2:

How to choose an optimal state sequence $\mathbf{q}=(q_1,q_2,\dots,q_T)$?

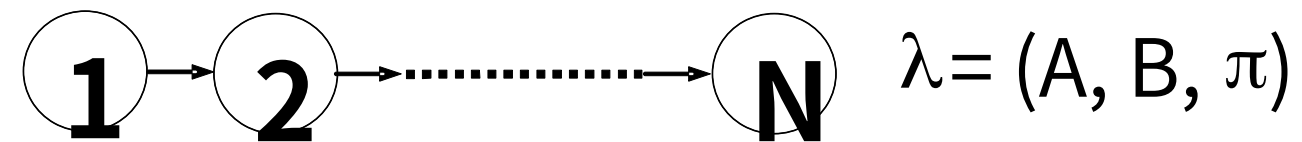
☒ *Decoding Problem*

- Problem 3:

Given some observations \bar{O} for the HMM λ , how to adjust the model parameter $\lambda=(A,B,\pi)$ to maximize $P(\bar{O}|\lambda)$?

☒ *Learning/Training Problem*

Basic Problem 1 for HMM



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

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

⊠ **Problem 1:** Given λ and \bar{O} ,
find $P(\bar{O}|\lambda) = \text{Prob}[\text{observing } \bar{O} \text{ given } \lambda]$

⊠ **Direct Evaluation: considering all possible state sequences**

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

$$P(\bar{O}|\bar{q}, \lambda) = \sum_{\text{all } \bar{q}} ([b_{q_1}(o_1) \cdot b_{q_2}(o_2) \cdot \cdots \cdot b_{q_T}(o_T)] \cdot [\pi_{q_1} \cdot a_{q_1 q_2} \cdot a_{q_2 q_3} \cdot \cdots \cdot a_{q_{T-1} q_T}])$$

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

total number of different $\bar{q} : T^N$
huge computation requirements

Basic Problem 1 for HMM

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

$$\begin{aligned}\alpha_t(i) &= P(o_1 o_2 \cdots o_t, q_t = i | \lambda) \\ &= \text{Prob}[\text{observing } o_1 o_2 \cdots o_t, \text{ state } i \text{ at time } t | \lambda]\end{aligned}$$

- Initialization

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

- Induction

$$\begin{aligned}\alpha_{t+1}(j) &= \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}) \\ &\quad 1 \leq j \leq N \\ &\quad 1 \leq t \leq T-1\end{aligned}$$

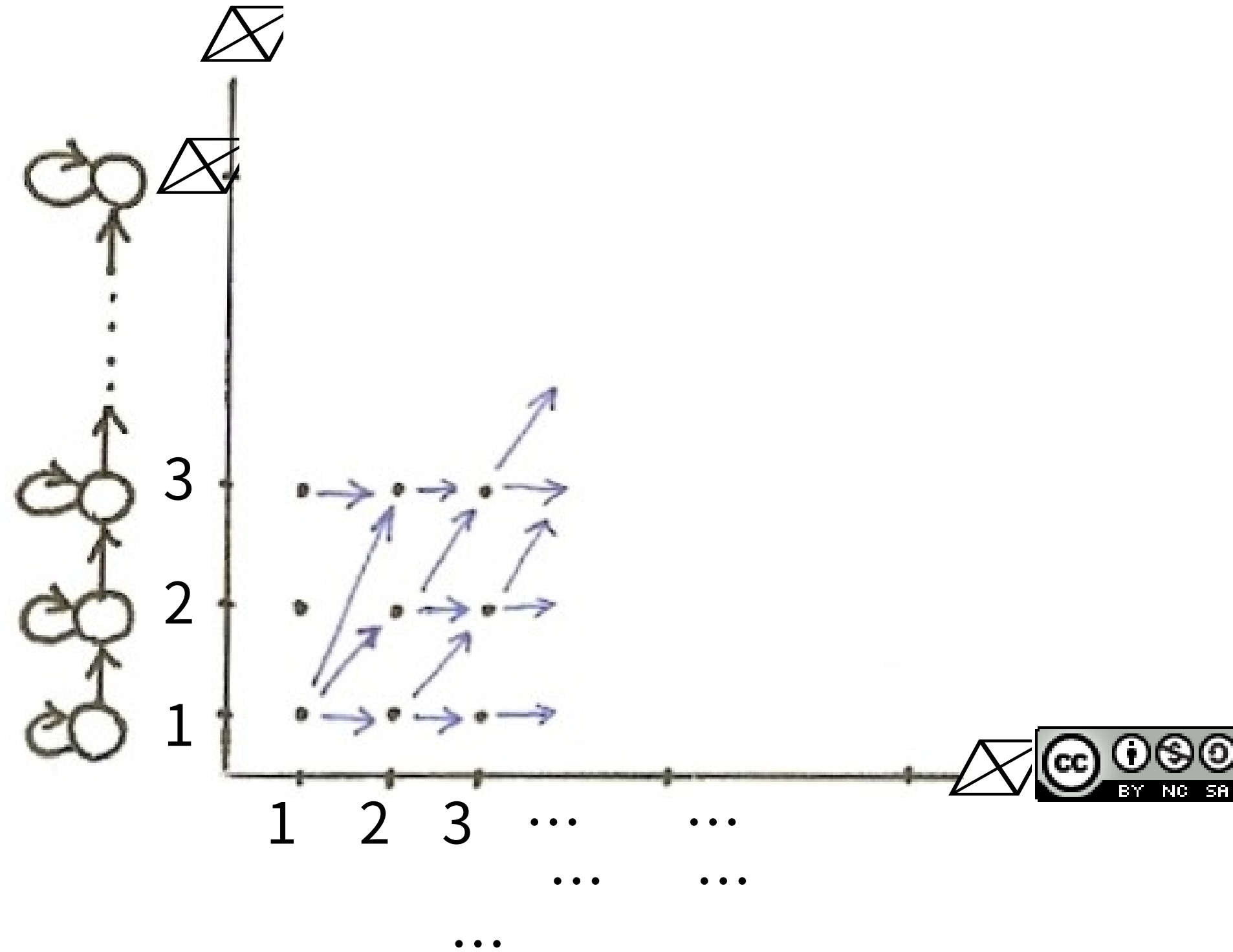
- Termination

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

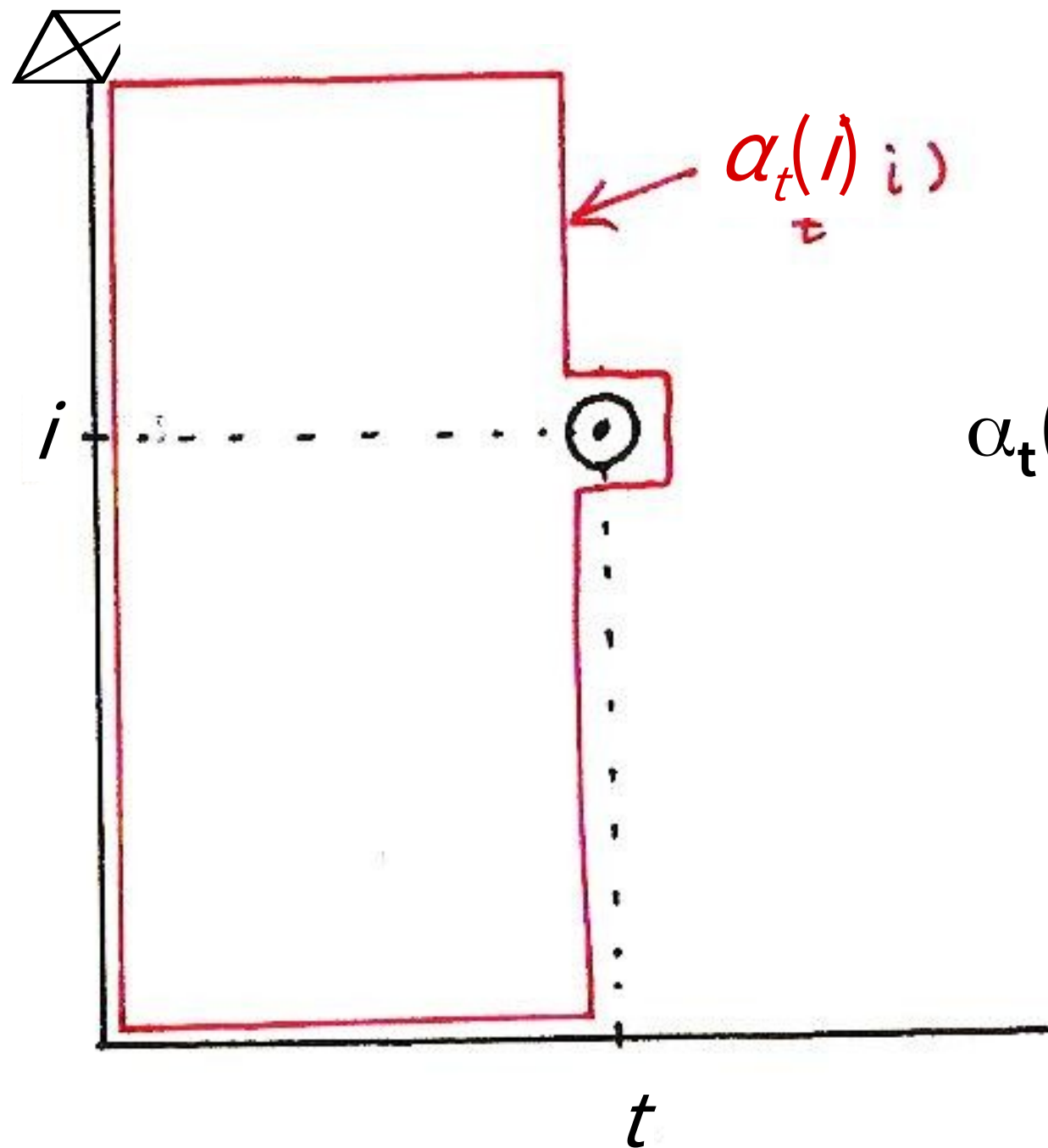
See Fig. 6.5 of Rabiner and Juang

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

Basic Problem 1



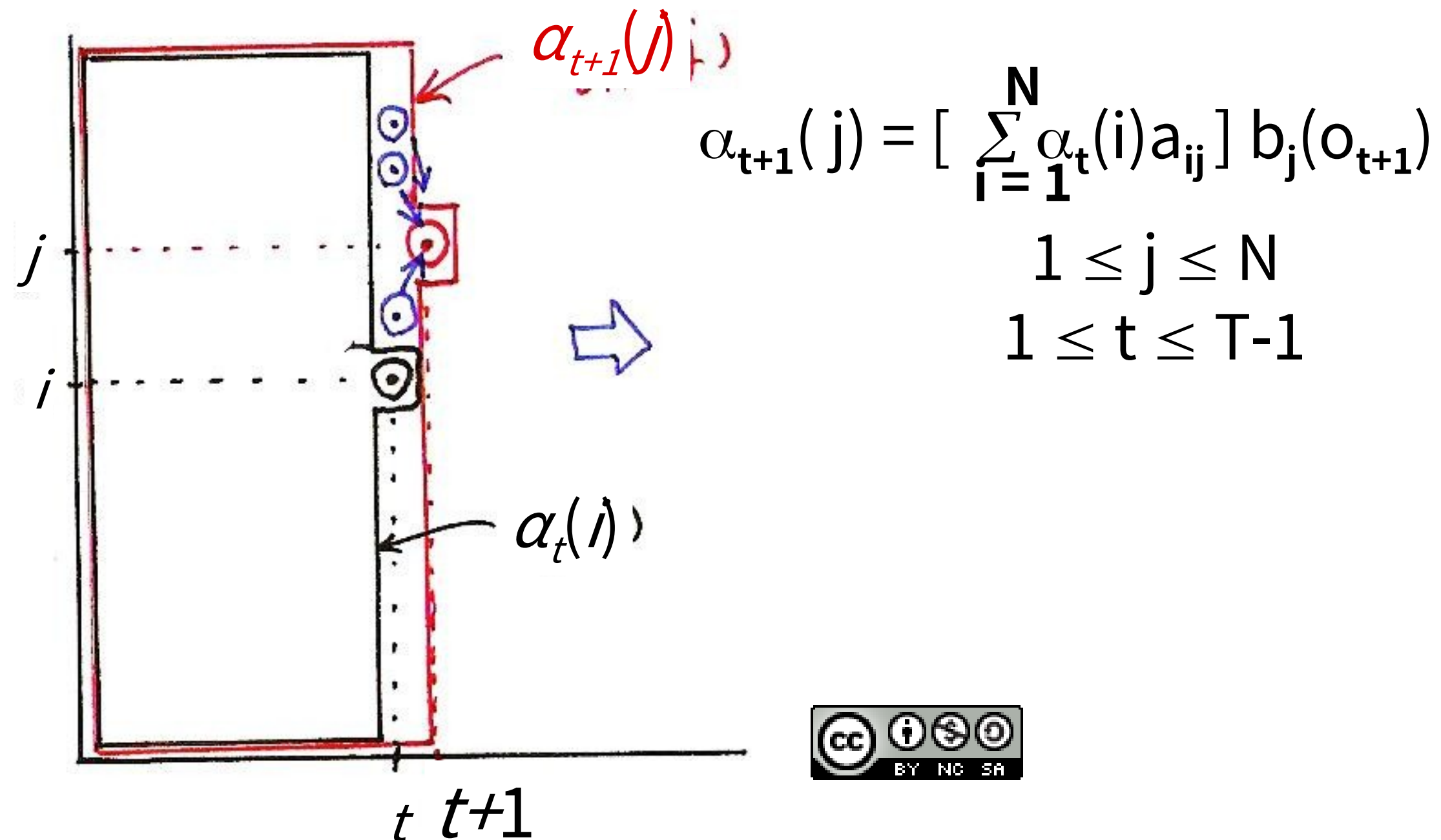
Basic Problem 1



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



Basic Problem 1



Forward Algorithm

Basic Problem 2 for HMM

- **Problem 2:** Given λ and $\bar{O} = o_1 o_2 \cdots o_T$, find a best state sequence $\bar{q} = q_1 q_2 \cdots q_T$
- **Backward Algorithm:** defining a backward variable $\beta_t(i)$

$$\begin{aligned}\beta_t(i) &= P(o_{t+1}, o_{t+2}, \cdots, o_T | q_t = i, \lambda) \\ &= \text{Prob}[\text{observing } o_{t+1}, o_{t+2}, \cdots, o_T | \text{state } i \text{ at time } t, \lambda]\end{aligned}$$

- Initialization

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

- Induction

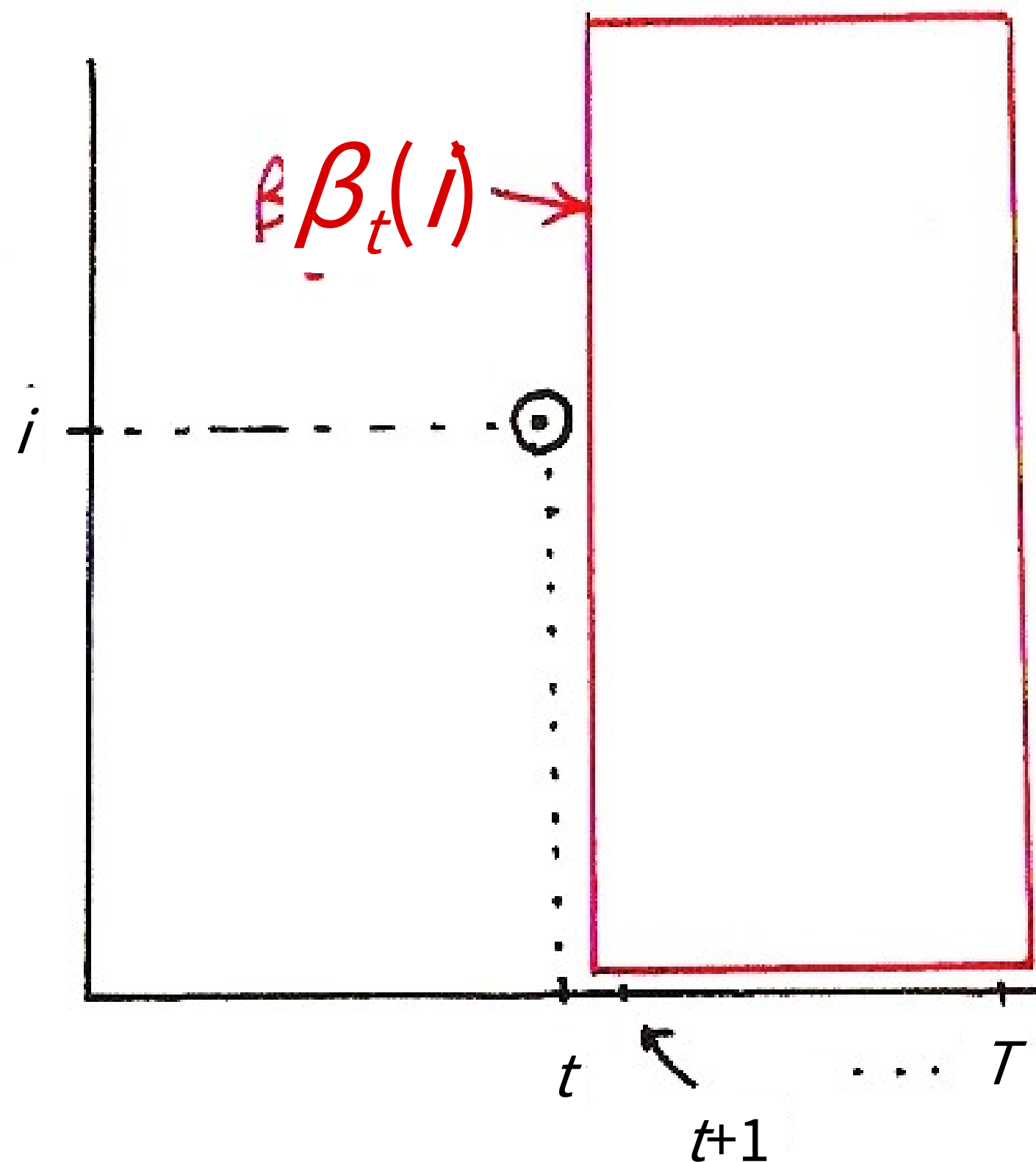
$$\begin{aligned}\beta_t(i) &= \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j) \\ t &= T-1, T-2, \cdots, 2, 1, \quad 1 \leq i \leq N\end{aligned}$$

See Fig. 6.6 of Rabiner and Juang

- **Combining Forward/Backward Variables**

$$\begin{aligned}P(\bar{O}, q_t = i | \lambda) &= \text{Prob}[\text{observing } o_1, o_2, \cdots, o_t, \cdots, o_T, q_t = i | \lambda] \\ &= \alpha_t(i) \beta_t(i) \\ P(\bar{O} | \lambda) &= \sum_{i=1}^N P(\bar{O}, q_t = i | \lambda) = \sum_{i=1}^N [\alpha_t(i) \beta_t(i)]\end{aligned}$$

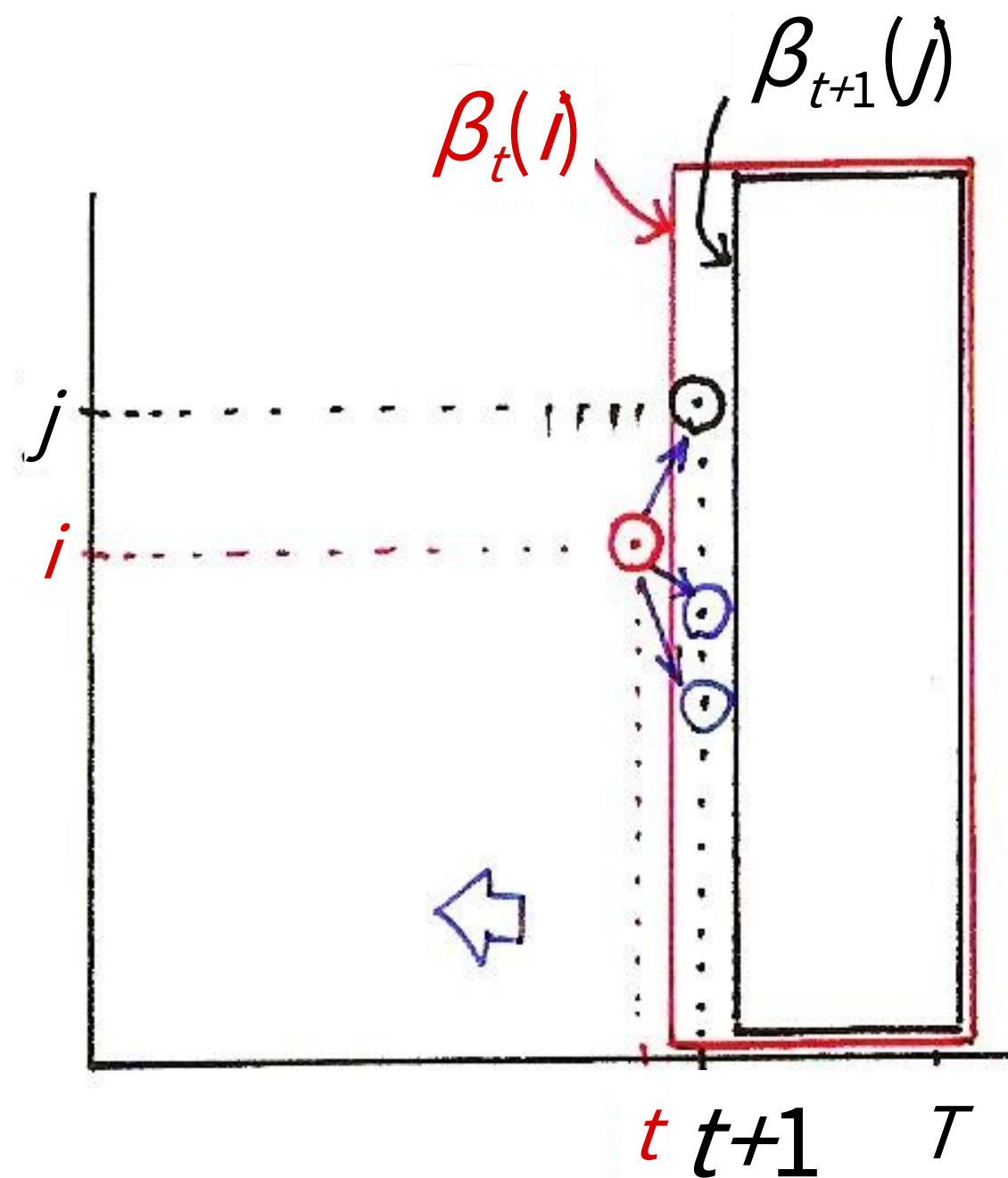
Basic Problem 2



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



Basic Problem 2



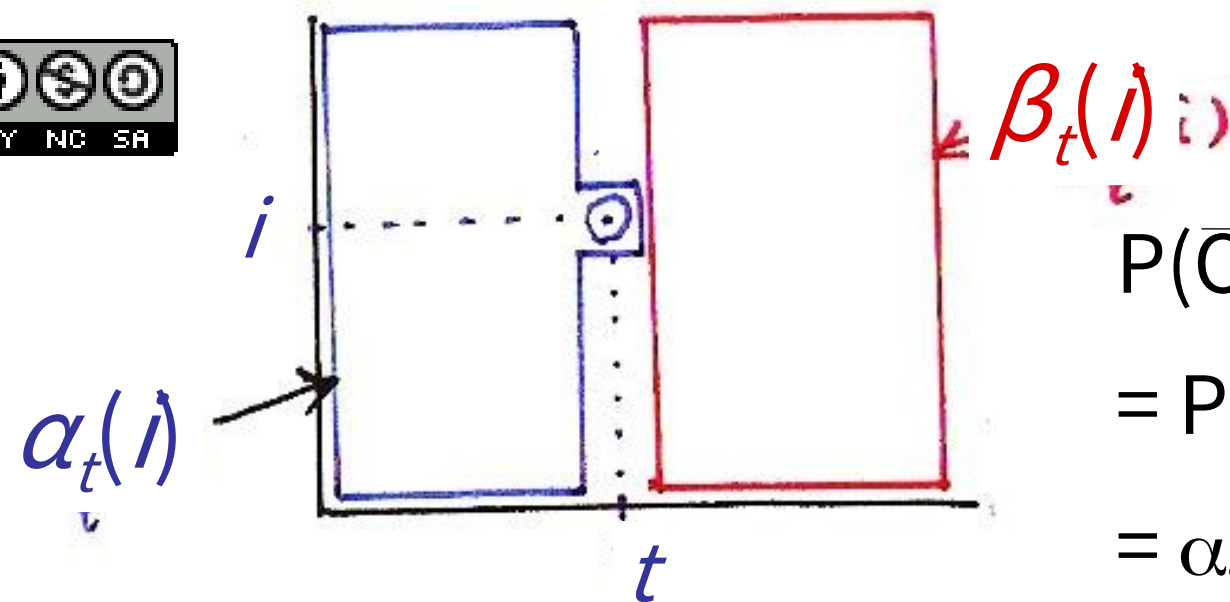
$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

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



Backward Algorithm

Basic Problem 2



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

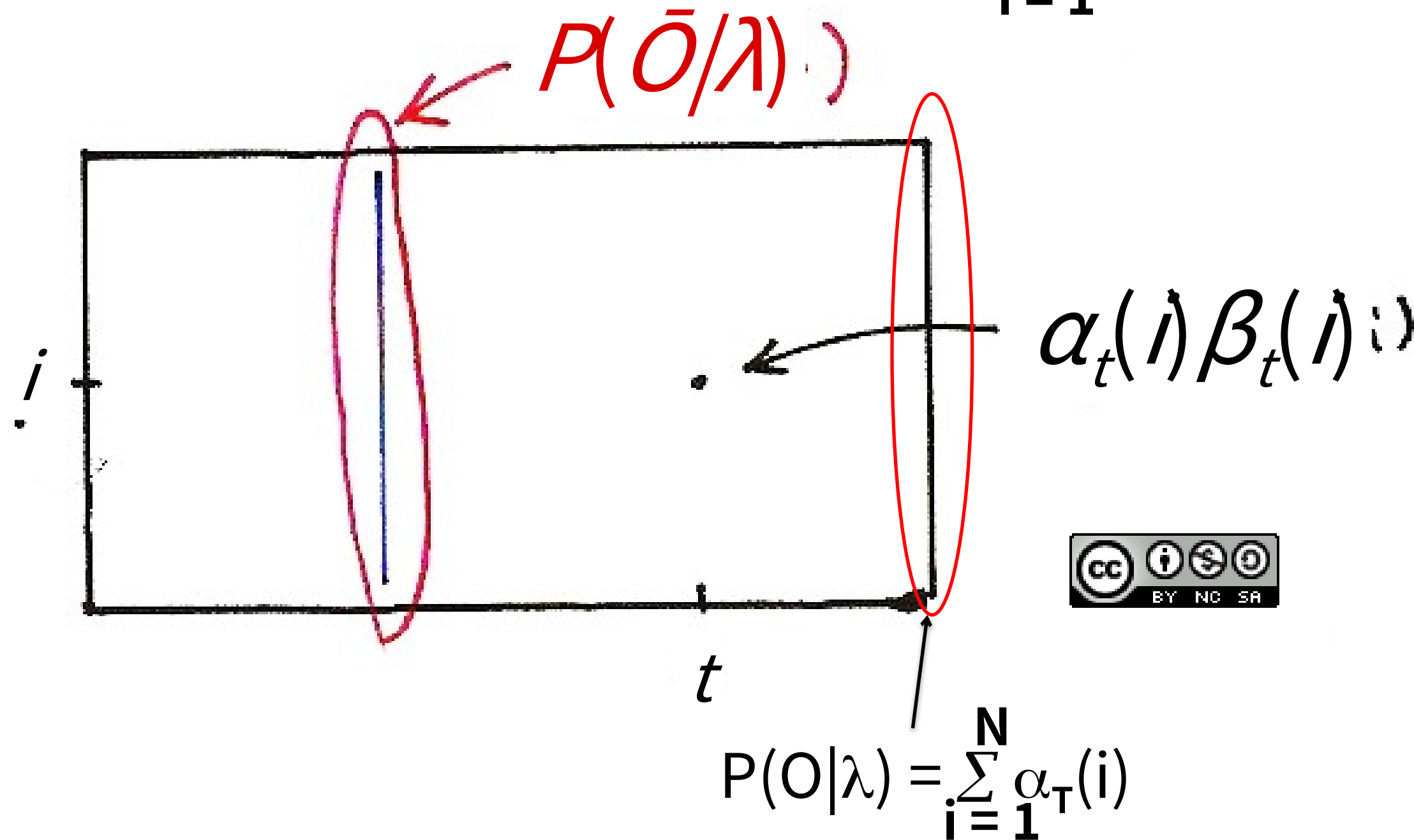
$$= \text{Prob} [\text{observing } o_1, o_2, \dots, o_t, \dots, o_T, q_t = i | \lambda]$$

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

$$\begin{aligned} \alpha_t(i) \beta_t(j) &= \alpha_t(i) \beta_t(j) \quad (\alpha_t(i) \perp \beta_t(j)) \\ \alpha_t(i) \beta_t(j) &= \alpha_t(i) \beta_t(j) \quad (\alpha_t(i) \perp \beta_t(j)) \\ \alpha_t(i) \beta_t(j) &= \alpha_t(i) \beta_t(j) \quad (\alpha_t(i) \perp \beta_t(j)) \end{aligned}$$

Basic Problem 2

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



Basic Problem 2 for HMM

⊗ Approach 1 – Choosing states individually as the most likely state at time t

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

$$e_t(i) = \frac{\prod_{j=1}^t b_j(i) a_{j-1,j}(i)}{\prod_{i=1}^N \prod_{j=1}^t b_j(i) a_{j-1,j}(i)} = \frac{P(\bar{O}, q_t = i \mid \lambda)}{P(\bar{O} \mid \lambda)}$$

- Solution

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

in fact

$$\begin{aligned} q_t^* &= \arg \max_{1 \leq i \leq N} [P(\bar{O}, q_t = i \mid \lambda)] \\ &= \arg \max_{1 \leq i \leq N} [\prod_{j=1}^t b_j(i) a_{j-1,j}(i)] \end{aligned}$$

- Problem

maximizing the probability at each time t

individually may not be a valid sequence

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

Basic Problem 2 for HMM

- Approach 2 – Viterbi Algorithm - finding the single best sequence

$$\bar{q}^* = q_1^* q_2^* \cdots q_T^*$$

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

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

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

- Induction

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

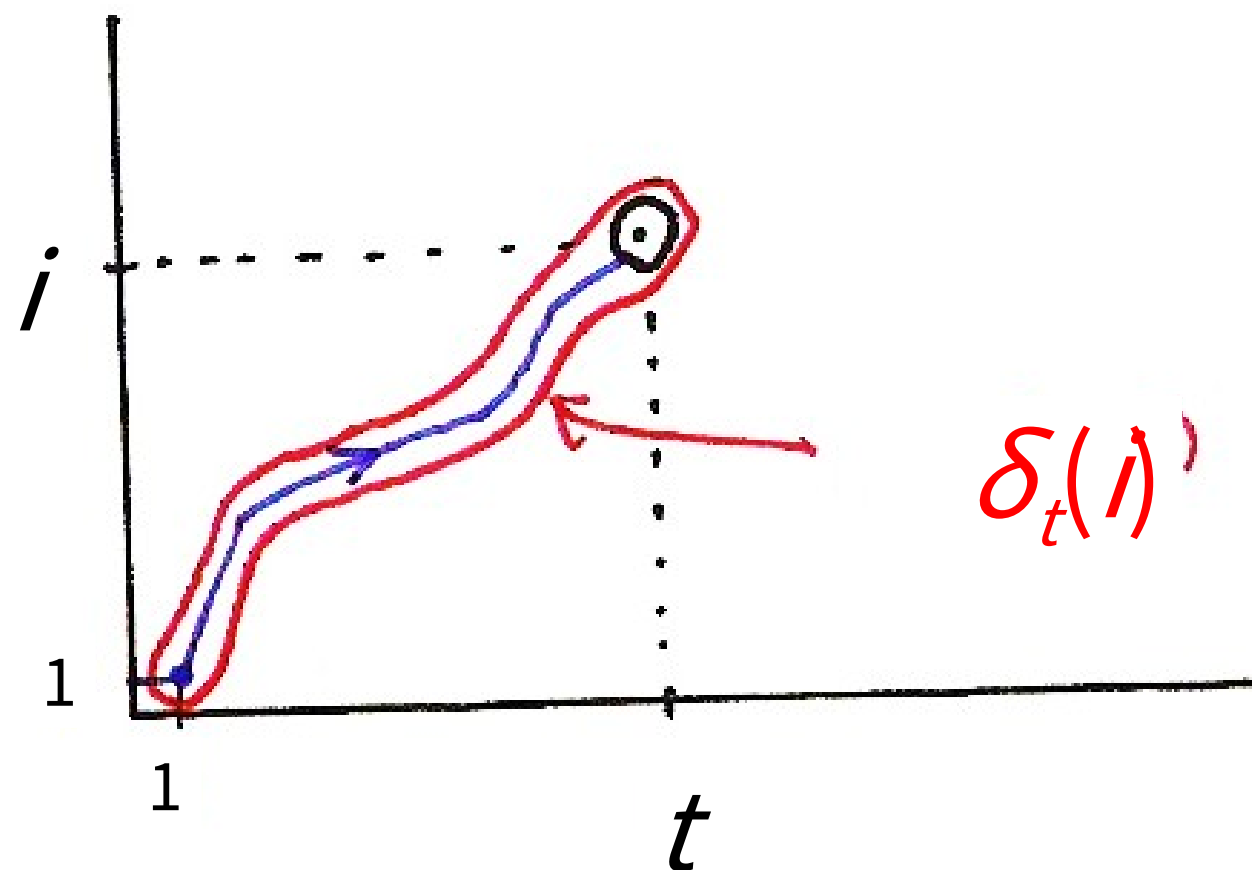
- Backtracking

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

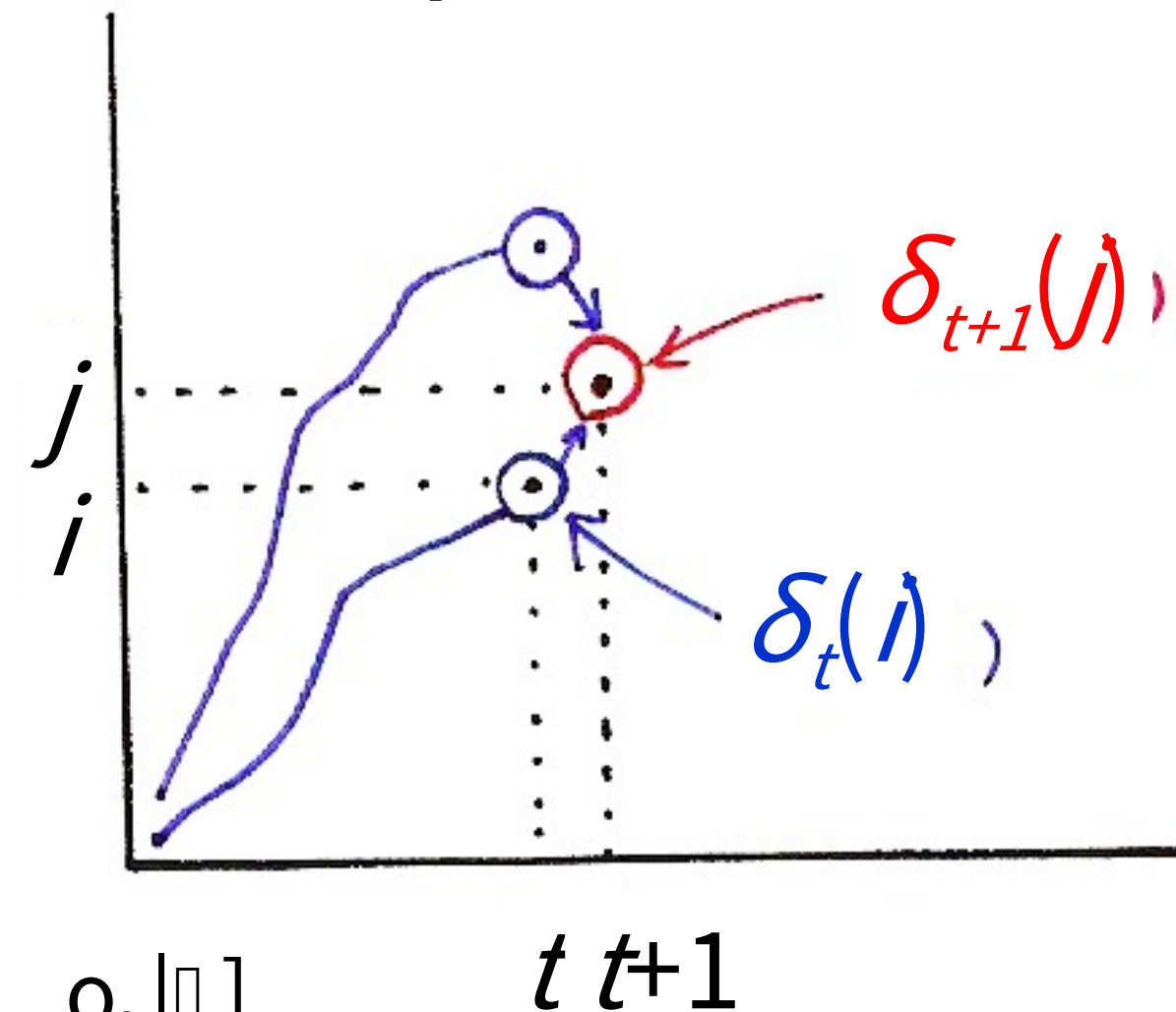
the best previous state at $t-1$ given at state j at time t

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

Viterbi Algorithm

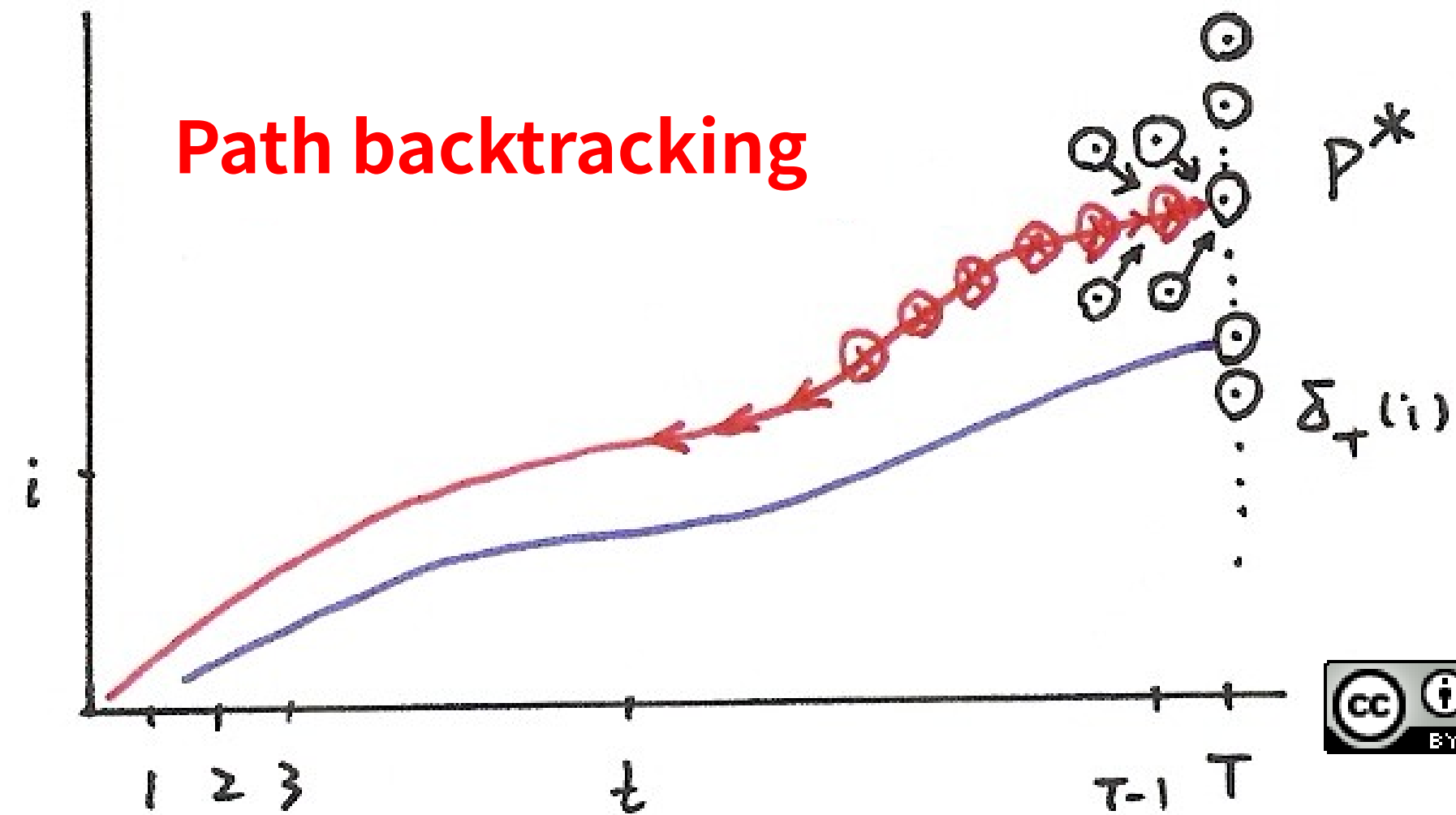


$$\delta_{t+1}(j) = \max_i [\delta_t(i) a_{ij}] \cdot b_j(o_{t+1})$$



$$\begin{aligned} \delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} & P[q_1, q_2, \dots, q_{t-1}, q_t = i, o_1, o_2, \dots, o_t] \end{aligned}$$

Viterbi Algorithm



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

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 2, 1$$

Basic Problem 2 for HMM

☒ Application Example of Viterbi Algorithm

- Isolated word recognition

$$\lambda_0 = (A_0, B_0, \pi_0)$$

$$\lambda_1 = (A_1, B_1, \pi_1)$$

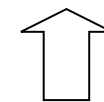
⋮

$$\lambda_n = (A_n, B_n, \pi_n)$$

observation

$$\bar{O} = (o_1, o_2, \dots, o_T)$$

$$k^* = \underset{1 \leq i \leq n}{\operatorname{argmax}} P[\bar{O} | \lambda_i] \approx \underset{1 \leq i \leq n}{\operatorname{argmax}} P^*[\bar{O} | \lambda_i]$$



Basic Problem 1 Basic Problem 2

Forward Algorithm Viterbi Algorithm

- The model with the highest probability for the most probable path (for a single best path) usually also has the highest probability for all possible paths

Basic Problem 3 for HMM

- **Problem 3:** Give \bar{O} and an initial model $\lambda=(A,B,\pi)$, adjust λ to maximize $P(\bar{O}|\lambda)$
 - Baum-Welch Algorithm (Forward-backward Algorithm)
 - Define a new variable

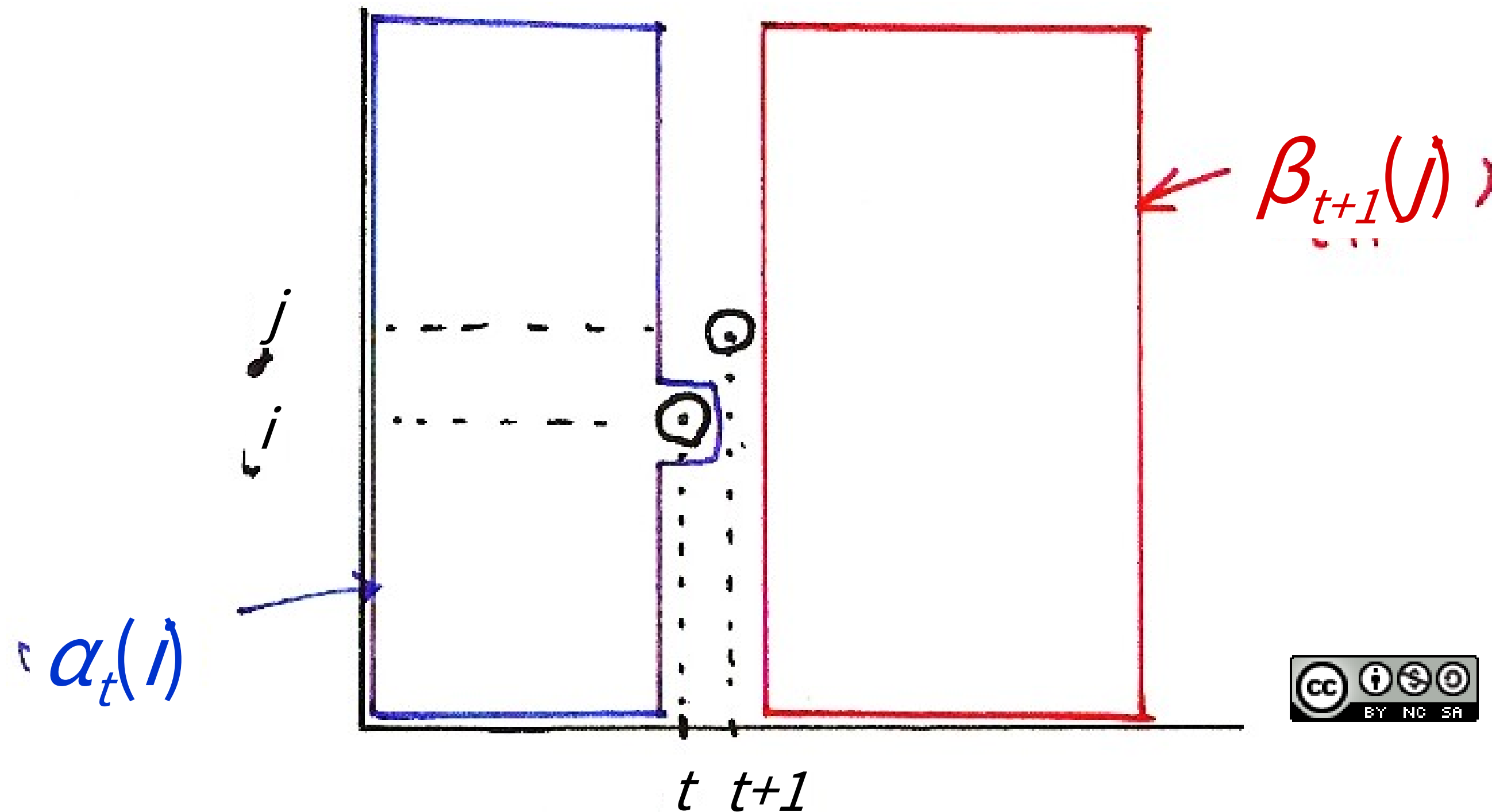
$$\begin{aligned}\varepsilon_t(i, j) &= P(q_t = i, q_{t+1} = j \mid \bar{O}, \lambda) \\ &= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N [\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)]} \\ &= \frac{\text{Prob}[\bar{O}, q_t = i, q_{t+1} = j \mid \lambda]}{P(\bar{O} \mid \lambda)}\end{aligned}$$

See Fig. 6.7 of Rabiner and Juang

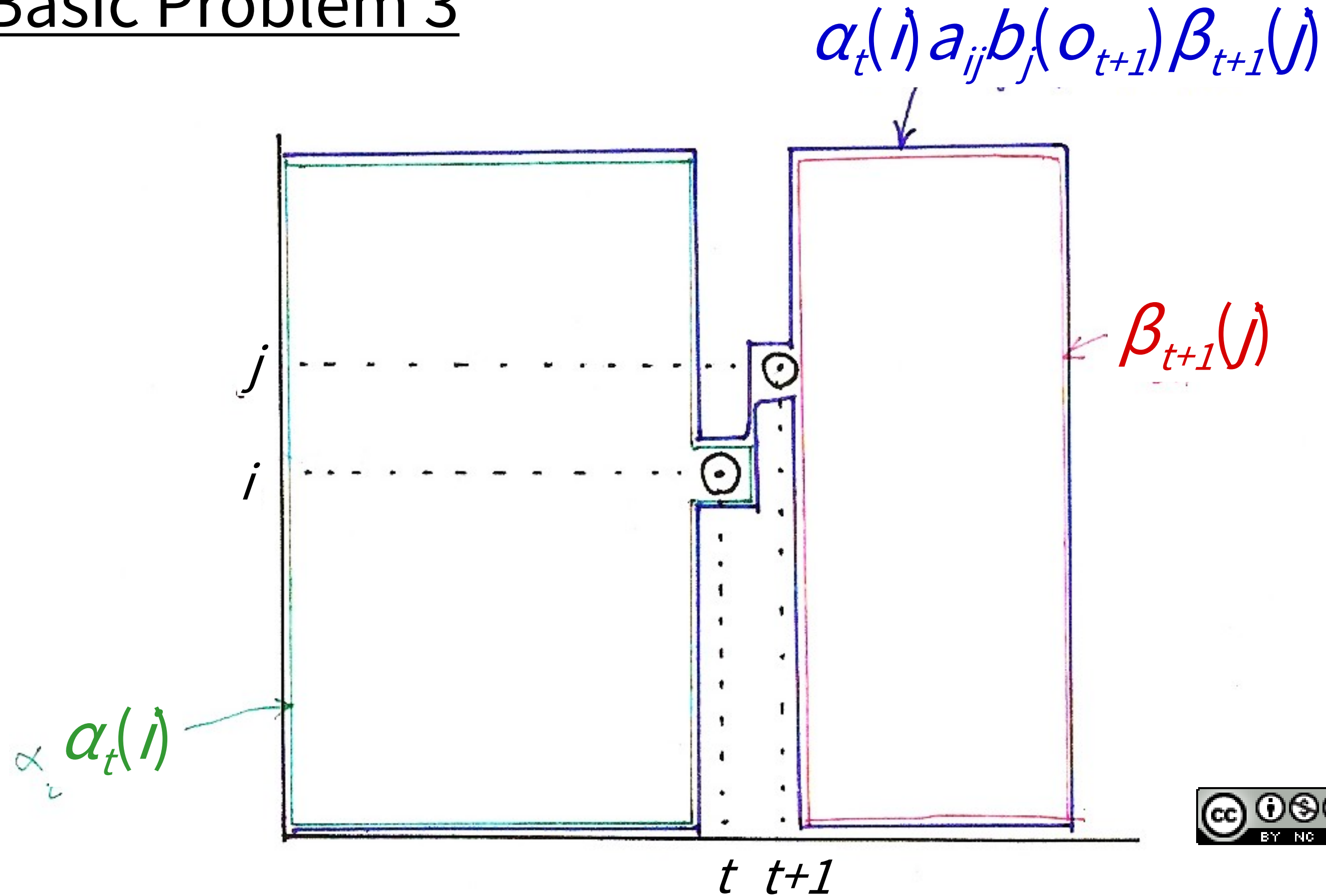
- Recall $\gamma_t(i) = P(q_t = i \mid \bar{O}, \lambda)$
 - $\sum_{t=1}^{T-1} \gamma_t(i) =$ expected number of times that state i is visited in \bar{O} from $t = 1$ to $t = T - 1$
 - $=$ expected number of transitions from state i in \bar{O}
 - $\sum_{t=1}^{T-1} \varepsilon_t(i, j) =$ expected number of transitions from state i to state j in \bar{O}

Basic Problem 3

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



Basic Problem 3



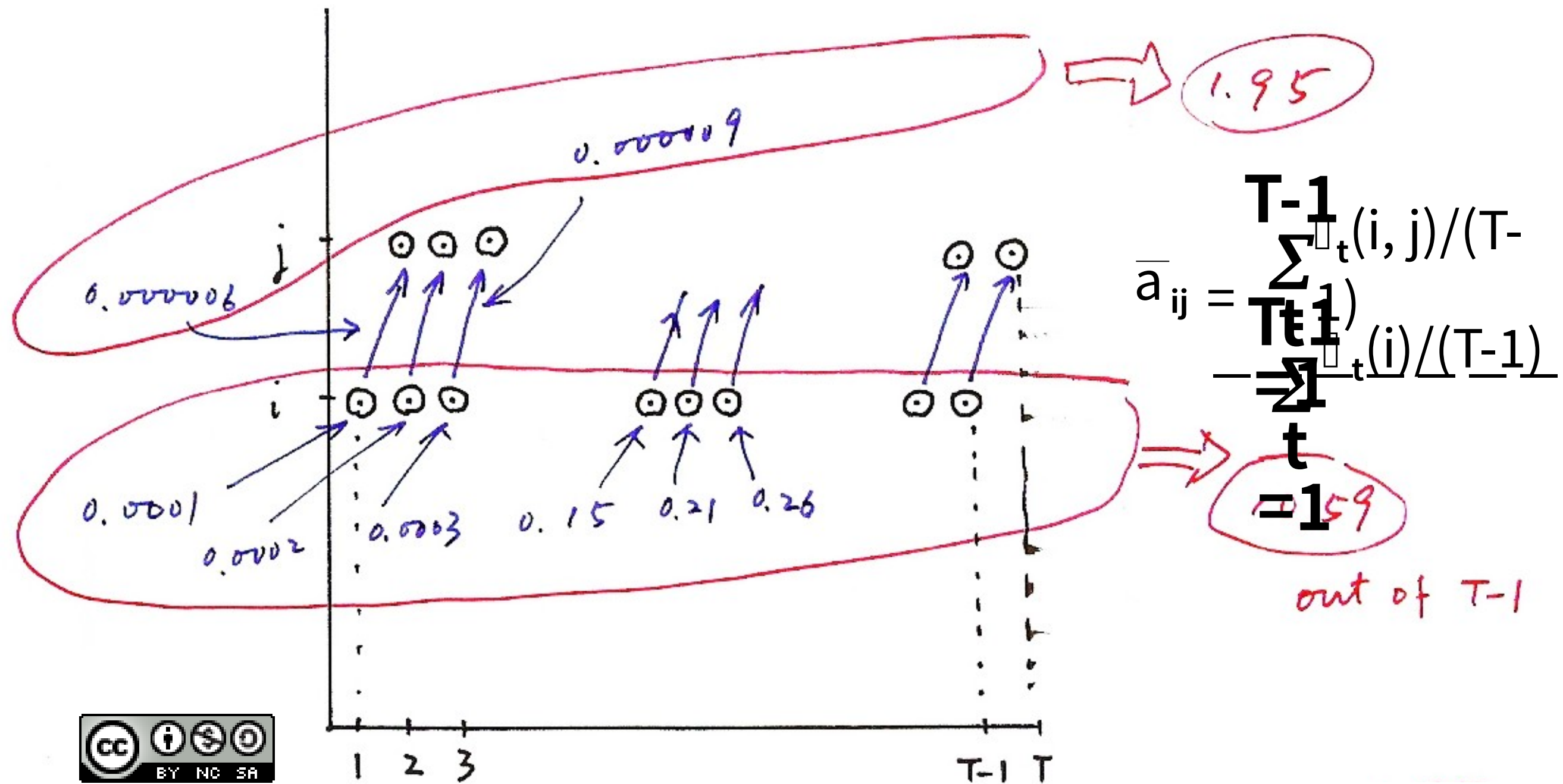
Basic Problem 3

$$\begin{aligned} \varphi_{\varphi}(\varphi) &= \frac{\varphi_{\varphi}(\varphi) \varphi_{\varphi}(\varphi)}{\sum_{\varphi=1} [\varphi \varphi \varphi(\varphi) \varphi_{\varphi}(\varphi)]} = \frac{\varphi(\overline{\varphi}, \varphi_{\varphi} = \varphi | \varphi)}{\varphi(\overline{\varphi} | \varphi)} = \varphi(\varphi_{\varphi} = \varphi | \overline{\varphi}, \varphi) \end{aligned}$$

$$\begin{aligned} \varphi_{\varphi}(\varphi, \varphi) &= \frac{\varphi_{\varphi}(\varphi) \varphi_{\varphi}(\varphi_{\varphi+1}) \varphi_{\varphi+1}(\varphi)}{\sum_{\varphi=1} \sum_{\varphi=1} \varphi_{\varphi}(\varphi) \varphi_{\varphi}(\varphi_{\varphi+1}) \varphi_{\varphi+1}(\varphi)} \end{aligned}$$

$$\varphi \frac{\varphi(\overline{\varphi}, \varphi_{\varphi} = \varphi, \varphi_{\varphi+1} = \varphi | \varphi)}{\varphi(\overline{\varphi} | \varphi)} = \varphi(\varphi_{\varphi} = \varphi, \varphi_{\varphi+1} = \varphi | \overline{\varphi}, \varphi)$$

Basic Problem 3



$$\bar{\Delta}_{\Delta\Delta} = \frac{1.95/69}{10.59/69}$$

Basic Problem 3 for HMM

- Results

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

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

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

$$\bar{b}_j(k) = \text{Prob}[o_t = v_k \mid q_t = j] = \frac{\sum_{t=1}^T \gamma_t(j) \mathbf{1}_{o_t = v_k}}{\sum_{t=1}^T \gamma_t(j)}$$

(for discrete HMM)

- **Continuous Density HMM**

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

$N(\cdot)$: Multi-variate Gaussian

μ_{jk} : mean vector for the k-th mixture component

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

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

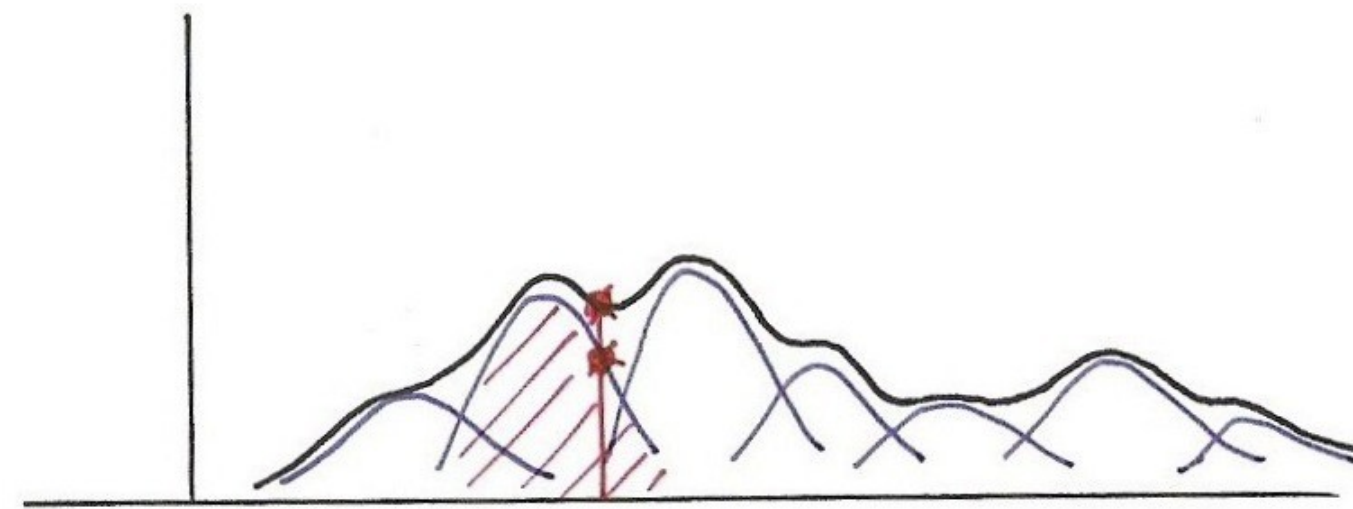
Basic Problem 3 for HMM

- 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

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



- Results

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

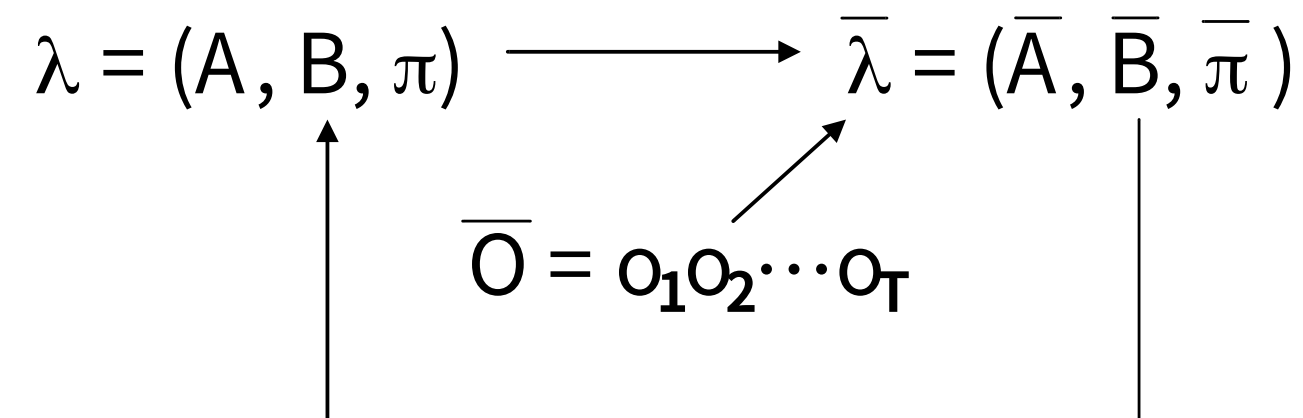
See Fig. 6.9 of Rabiner and Juang

Basic Problem 3 for HMM

- Continuous Density HMM

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

- Iterative Procedure



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

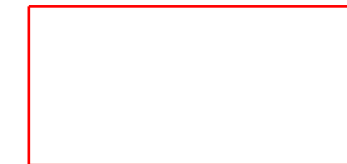
$$P(\bar{O}|\bar{\lambda}) \geq P(\bar{O}|\lambda) \text{ after each iteration}$$

Basic Problem 3

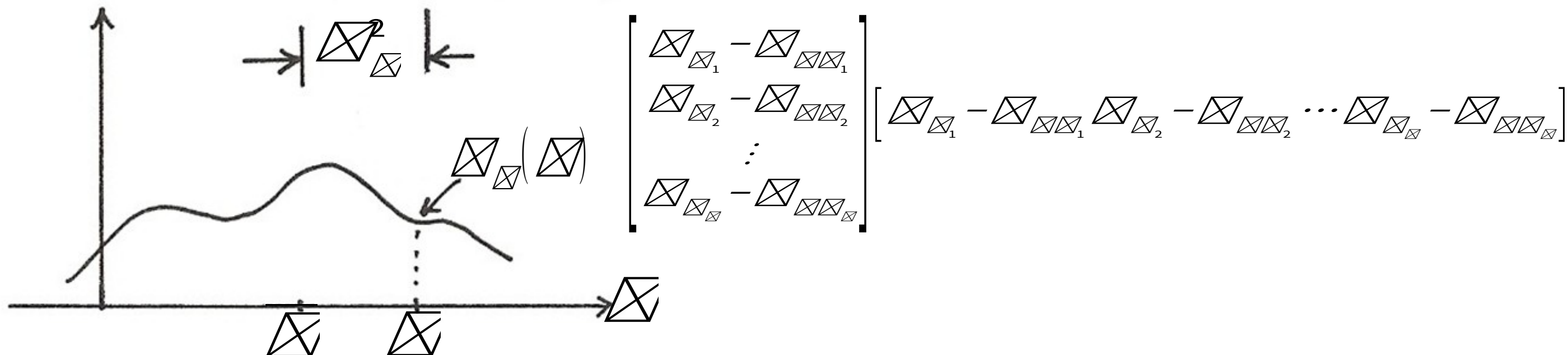
$$\bar{\mu}_{jk} = \frac{\sum_{t=1}^T [\gamma_t(j,k) \cdot o_t]}{\sum_{t=1}^T \gamma_t(j,k)}$$




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



prob. density function



Basic Problem 3

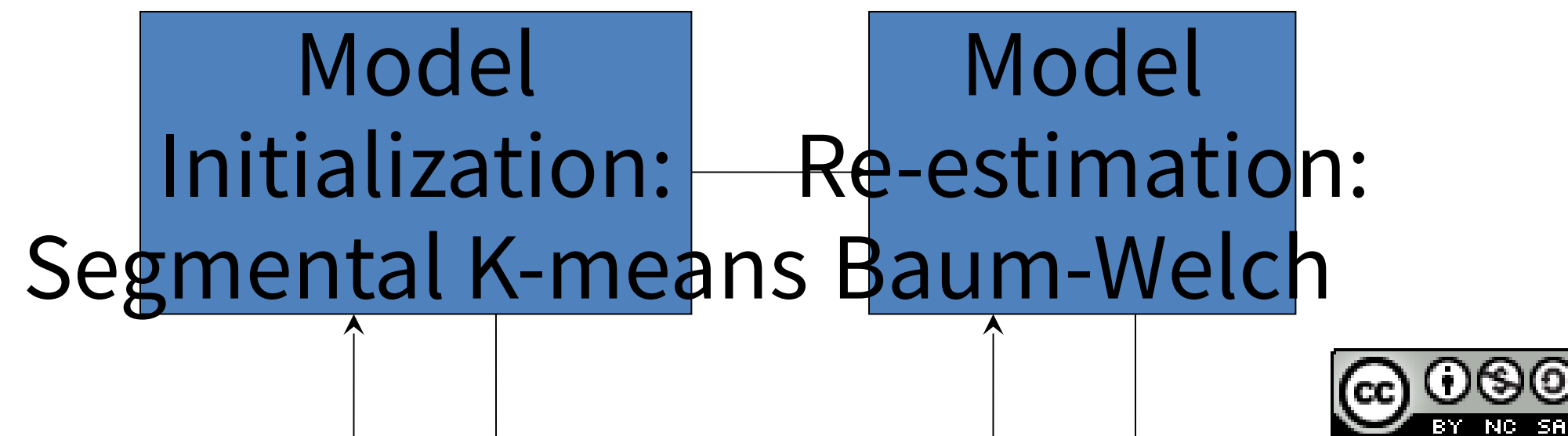


$$\bar{U} = \begin{bmatrix} \diagup & & & & \\ & \ddots & & & \\ & & \bar{\diagup} & & \\ & & & \diagup & \\ & & & & \ddots \end{bmatrix} = E \left(\begin{bmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \\ \vdots \\ \vdots \end{bmatrix} [x_1 - \bar{x}_1, x_2 - \bar{x}_2, \dots] \right)$$

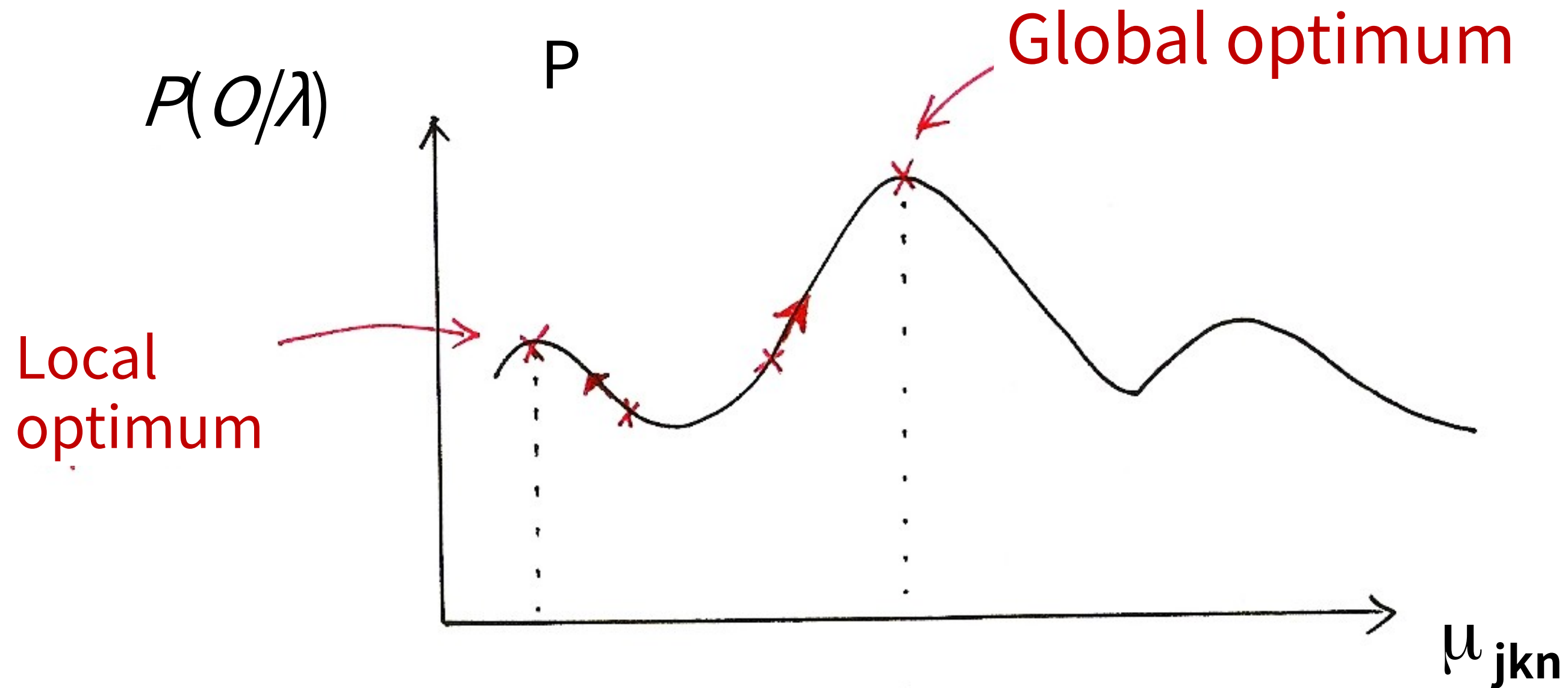
$$\bar{u}_{lm} = E[(x_l - \bar{x}_l)(x_m - \bar{x}_m)]$$

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

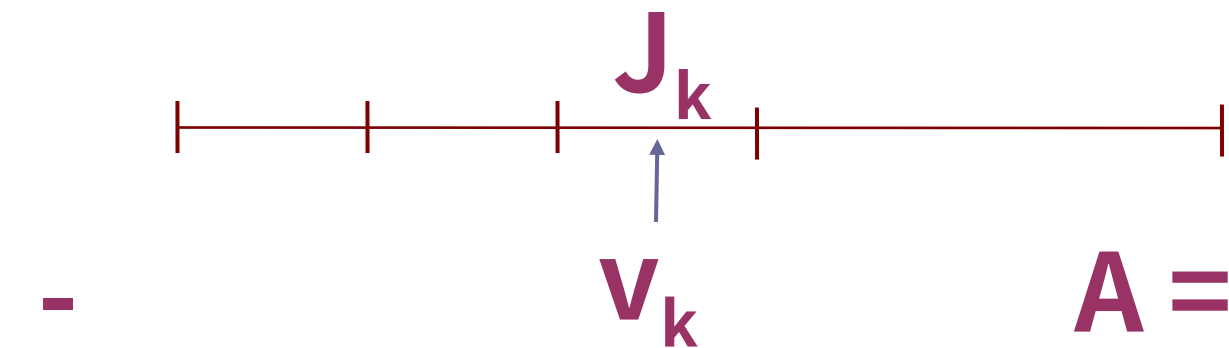


Basic Problem 3



Vector Quantization (VQ)

- **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



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

$$Q: S \rightarrow V$$

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

$$L = 2^R$$

Each v_k represented by an R-bit pattern

–Quantization characteristics (codebook)

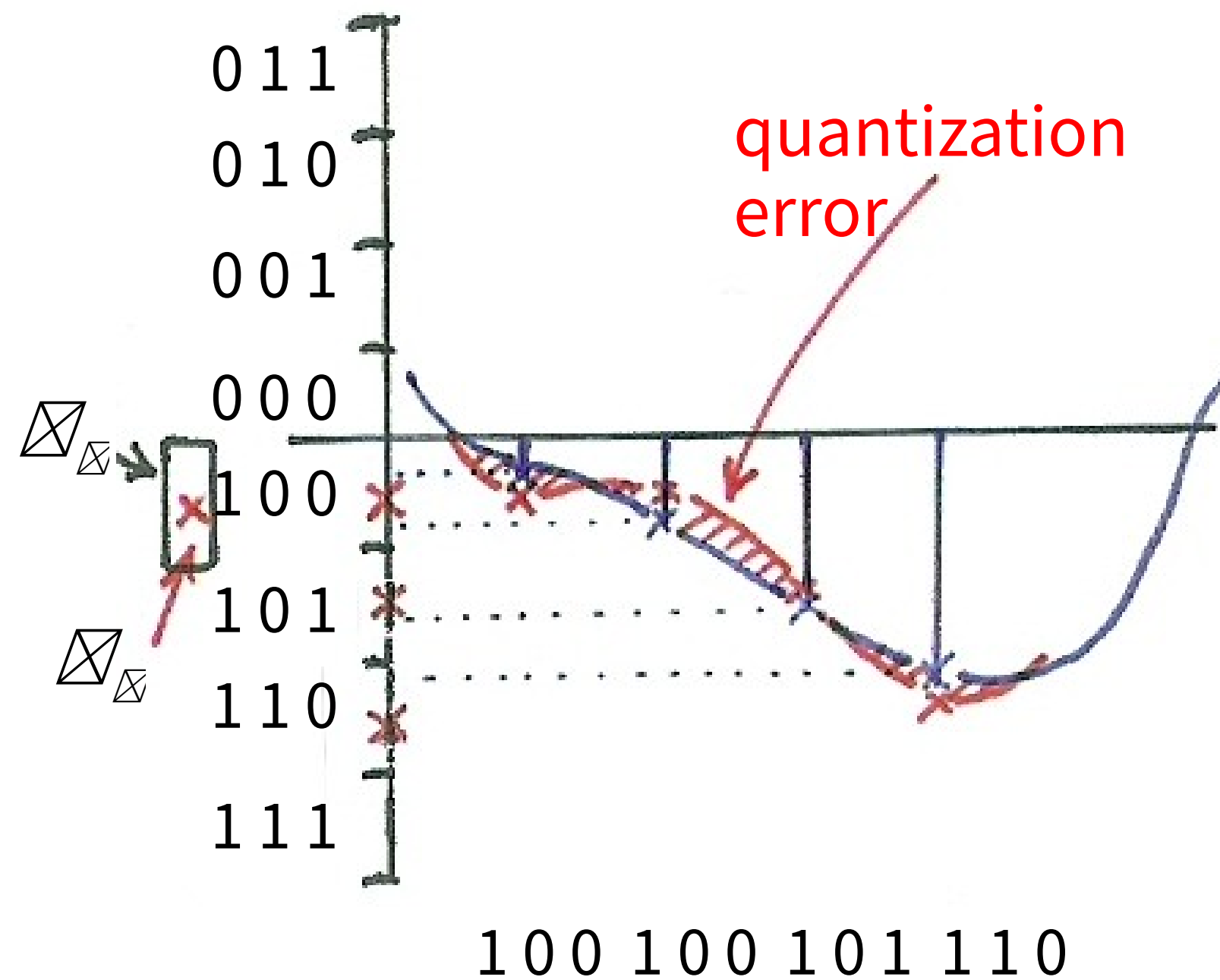
$\{J_1, J_2, \dots, J_L\}$ and $\{v_1, v_2, \dots, v_L\}$ designed considering at least

1. error sensitivity

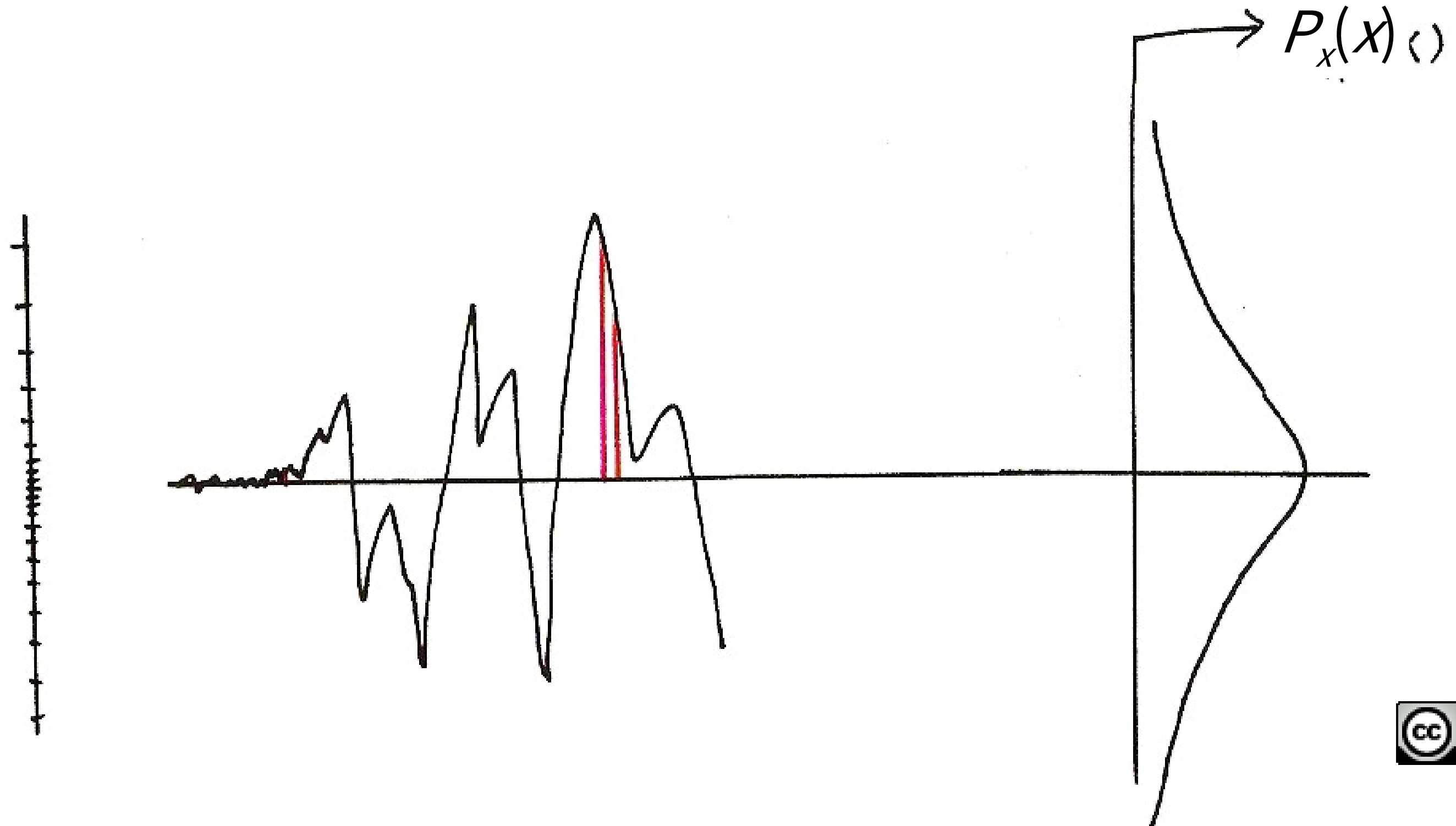
2. probability distribution of $x[n]$

Vector Quantization

Scalar Quantization : Pulse Coded Modulation (PCM)



Vector Quantization



2-dim Vector Quantization (VQ)

Example:

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

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

• VQ

– S divided into L 2-dim regions

$$\{ J_1, J_2, \dots, J_L \}$$

each with a representative

$$\text{vector } \bar{v}_k \in J_k, V = \{ \bar{v}_1, \bar{v}_2, \dots, \bar{v}_L \}$$

– $Q : S \Rightarrow V$

$$Q(\bar{x}_n) = \bar{v}_k \text{ if } \bar{x}_n \in J_k$$

$$L = 2^R$$

each \bar{v}_k represented by an R-bit pattern

– Considerations

1. error sensitivity may depend on $x[n], x[n+1]$ jointly

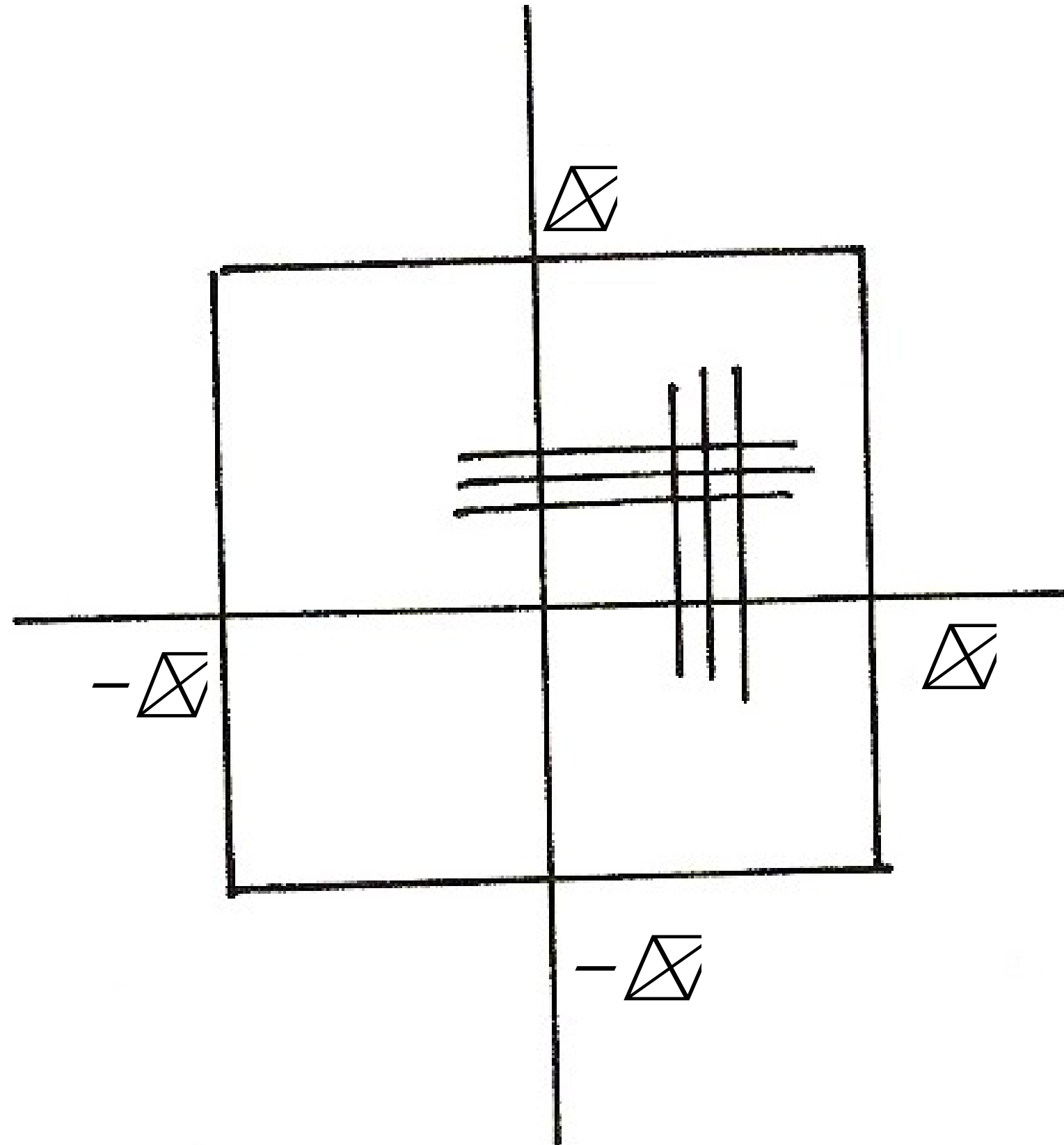
2. distribution of $x[n], x[n+1]$ may be correlated statistically

3. more flexible choice of J_k

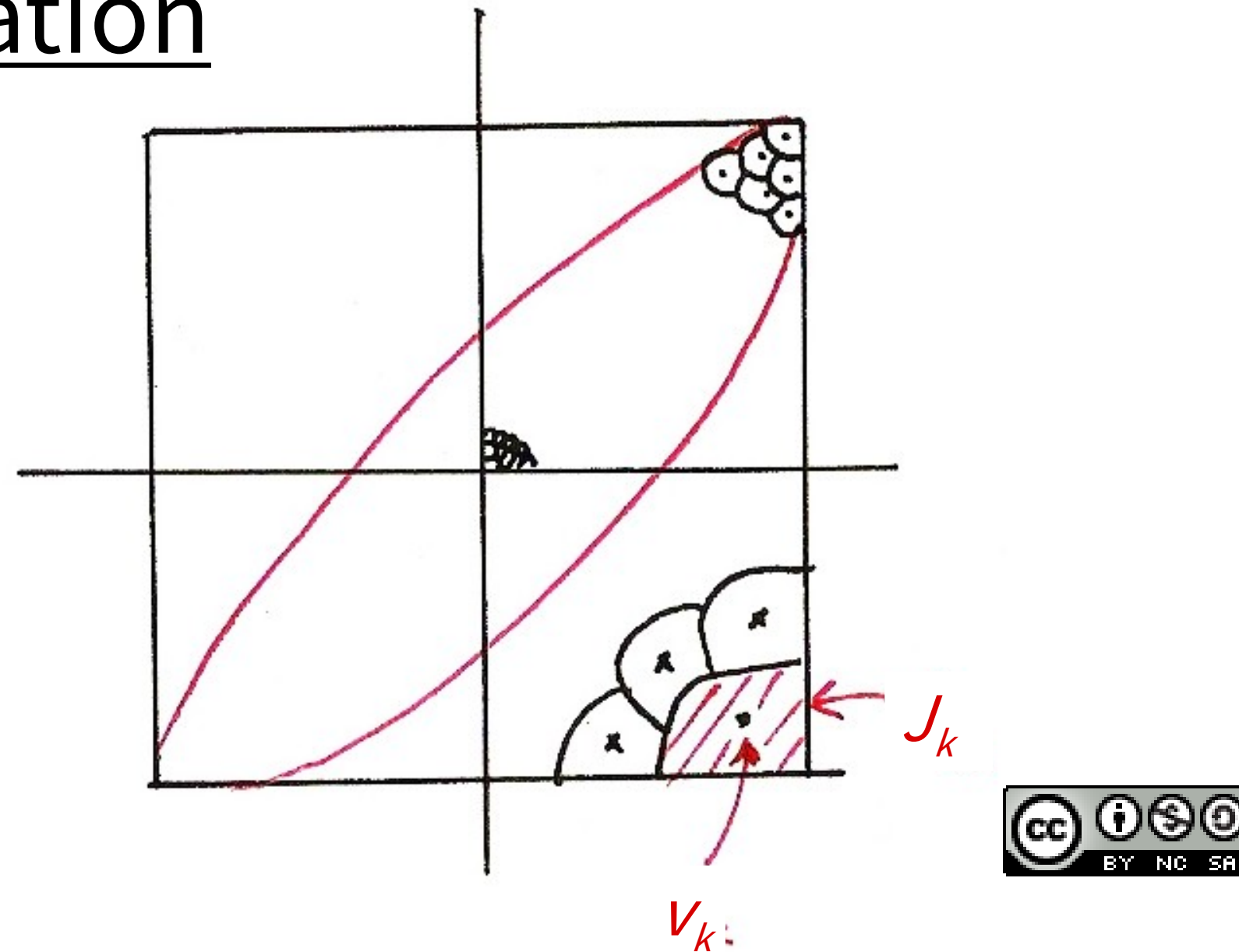
– Quantization Characteristics (codebook)

$$\{ J_1, J_2, \dots, J_L \} \text{ and } \{ \bar{v}_1, \bar{v}_2, \dots, \bar{v}_L \}$$

Vector Quantization



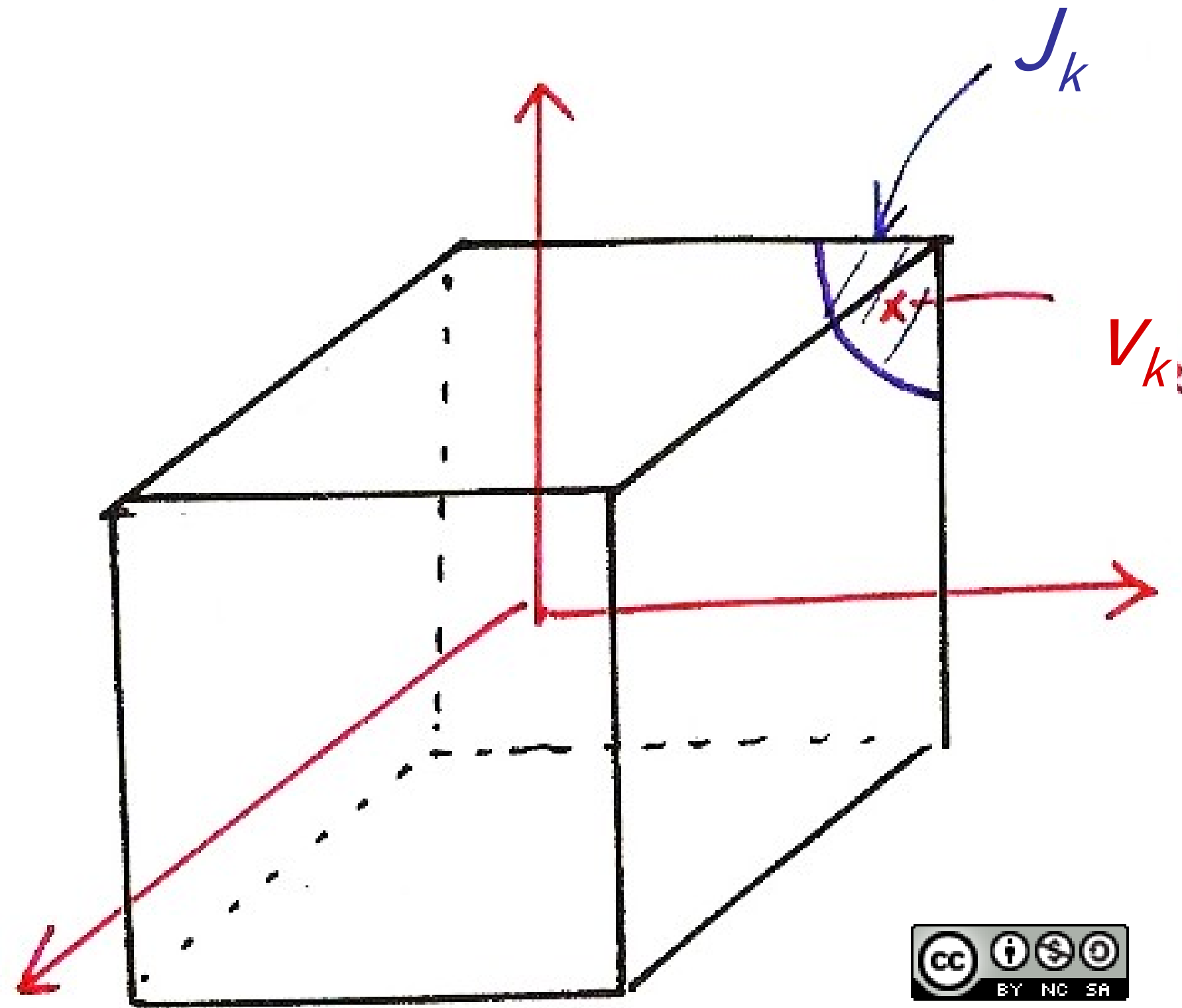
Vector Quantization



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

$$1024 = 2^{10}$$

Vector Quantization



N-dim Vector Quantization

$$x = (x_1, x_2, \dots, x_N)$$
$$S = \{x = (x_1, x_2, \dots, x_N),$$
$$S = \bigcup_{k=1}^L J_k$$
$$|x_k| \leq A, k = 1, 2, \dots, N\}$$

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

$$Q: S \rightarrow V$$

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

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

Codebook Trained by a Large Training Set

• Define distance measure
between

two vectors x, y

$$d(x, y) : S \times S \rightarrow \mathbb{R}^+ \text{ (non-negative}$$

real numbers)

-desired properties

$$d(x, y) \geq 0$$

$$d(x, x) = 0$$

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

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

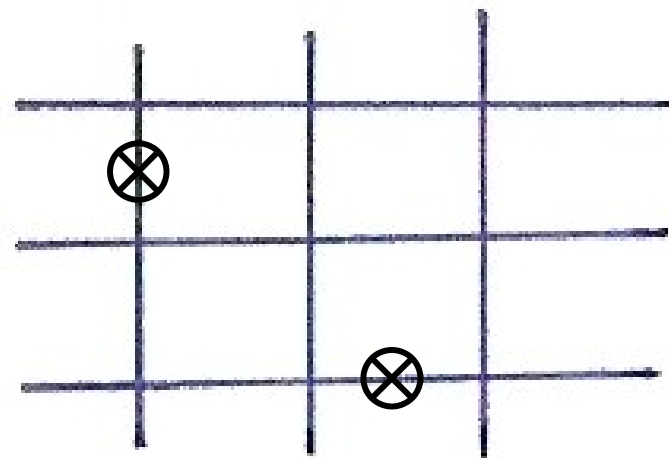
examples :

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

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

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

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



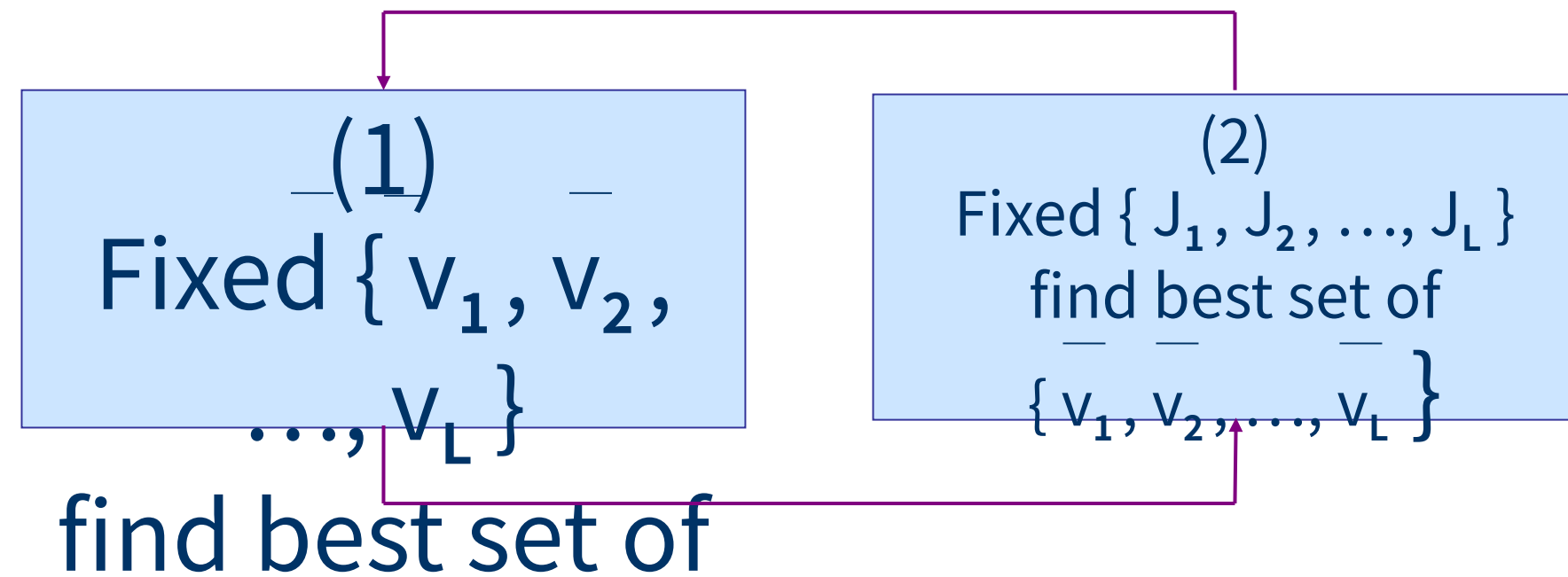
$$d(\bar{x}, \bar{y}) = (\bar{x} - \bar{y})^T \Sigma^{-1} (\bar{x} - \bar{y}) \quad \text{Mahalanobis distance}$$

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

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

Vector Quantization (VQ)

• K-Means Algorithm/Lloyd-Max Algorithm



(1) $J_k = \{ \bar{x} \mid d(\bar{x}, v_k) \leq d(\bar{x}, v_j), j \neq k \}$

(2) Convergence condition

$\rightarrow D = \sum_{\bar{x} \in J_k} d(\bar{x}, Q(\bar{x})) =$

$\min_{\bar{v}_k} \bar{v}_k = \frac{1}{M} \sum_{\bar{x} \in J_k} \bar{x}$

nearest neighbor

$$D^{k+1} = \sum D_k$$

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

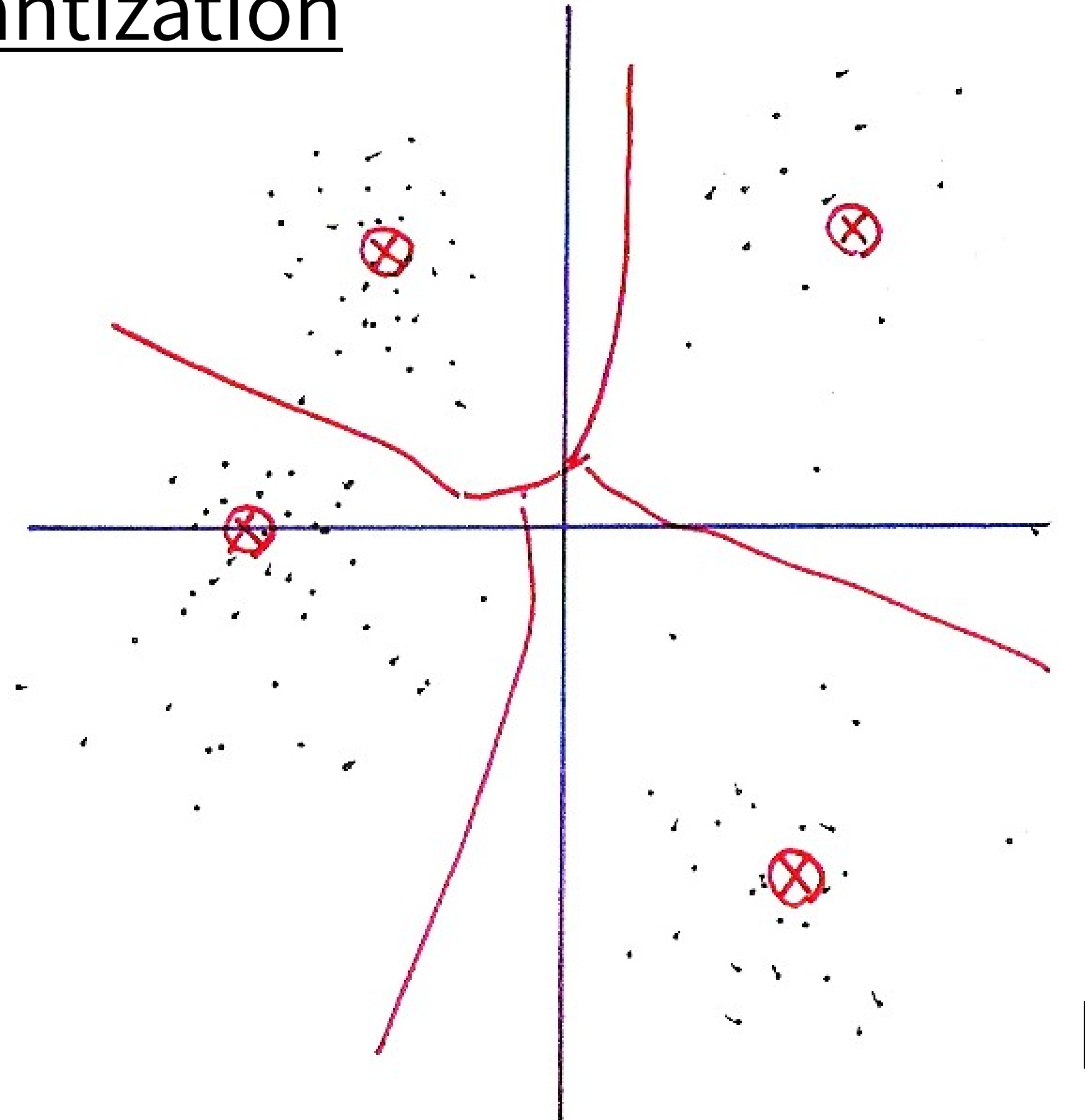
$|D^{(m+1)} - D^{(m)}| < \epsilon, m :$

condition

iteration

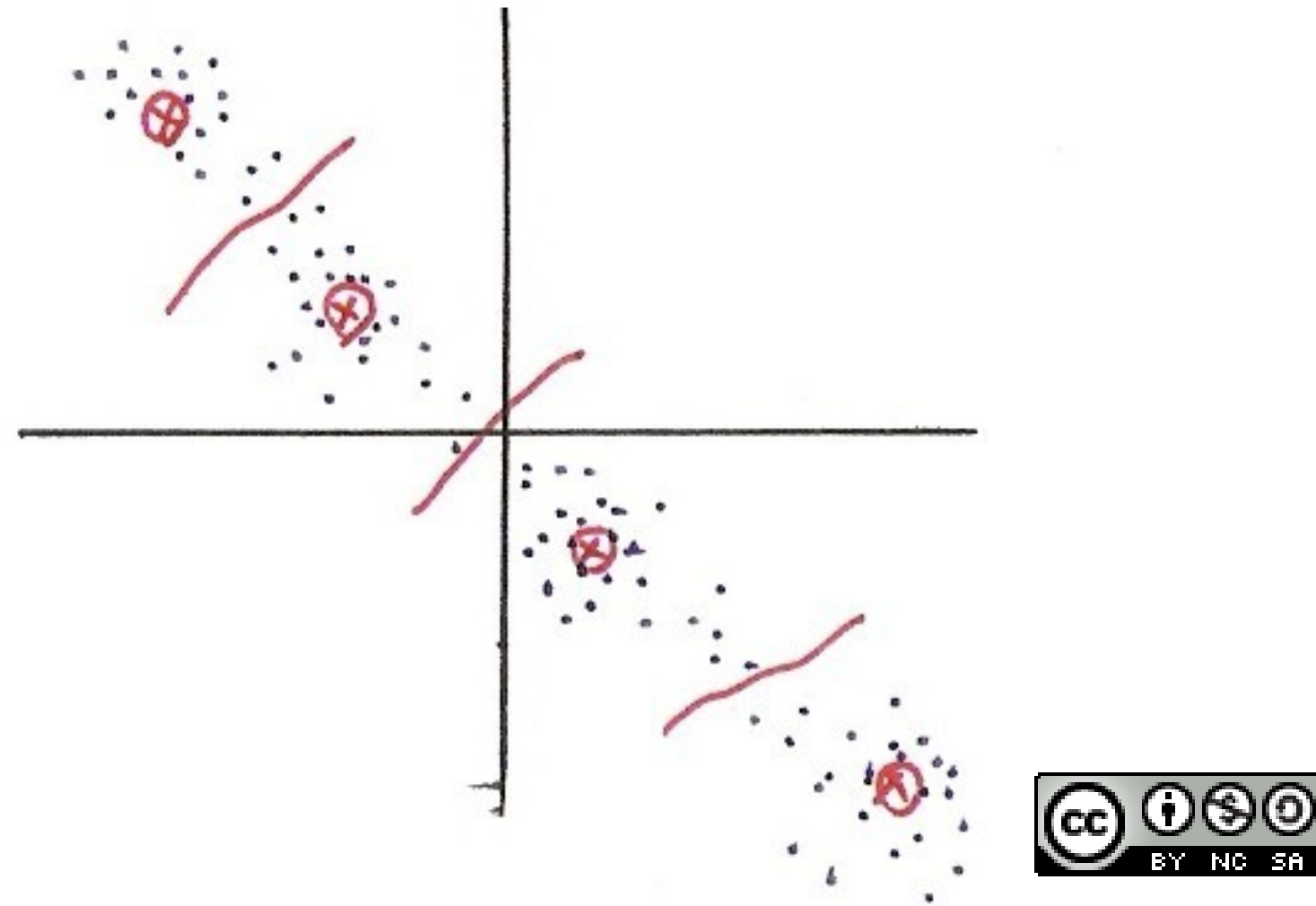
• Iterative Procedure to Obtain Codebook from a Large Training Set

Vector Quantization



- **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. $\bar{V} = \frac{1}{N} \sum_j \mathbf{x}_j$, train a 1-vector VQ codebook
 - 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
 $\bar{V}_k^{(1)} = \bar{V}_k$
 $\bar{V}_k^{(2)}$: 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**

LBG Algorithm



- **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

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

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

⇒ cluster the observation vectors within each state j into a set of M clusters (e.g. with vector quantization)

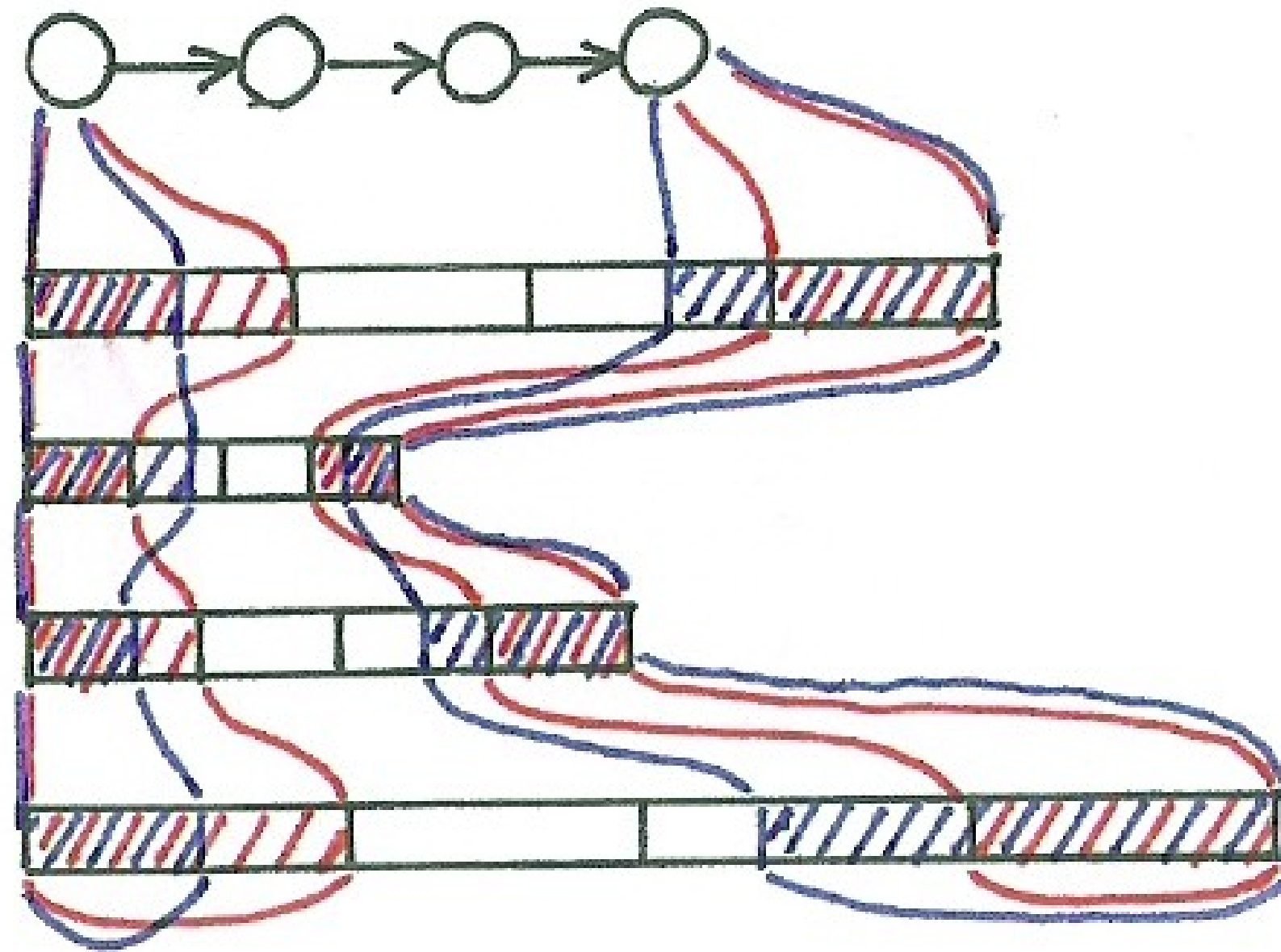
c_{jm} = number of vectors classified in cluster m of state j
divided by number of vectors in state j

μ_{jm} = sample mean of the vectors classified in cluster m of state j

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

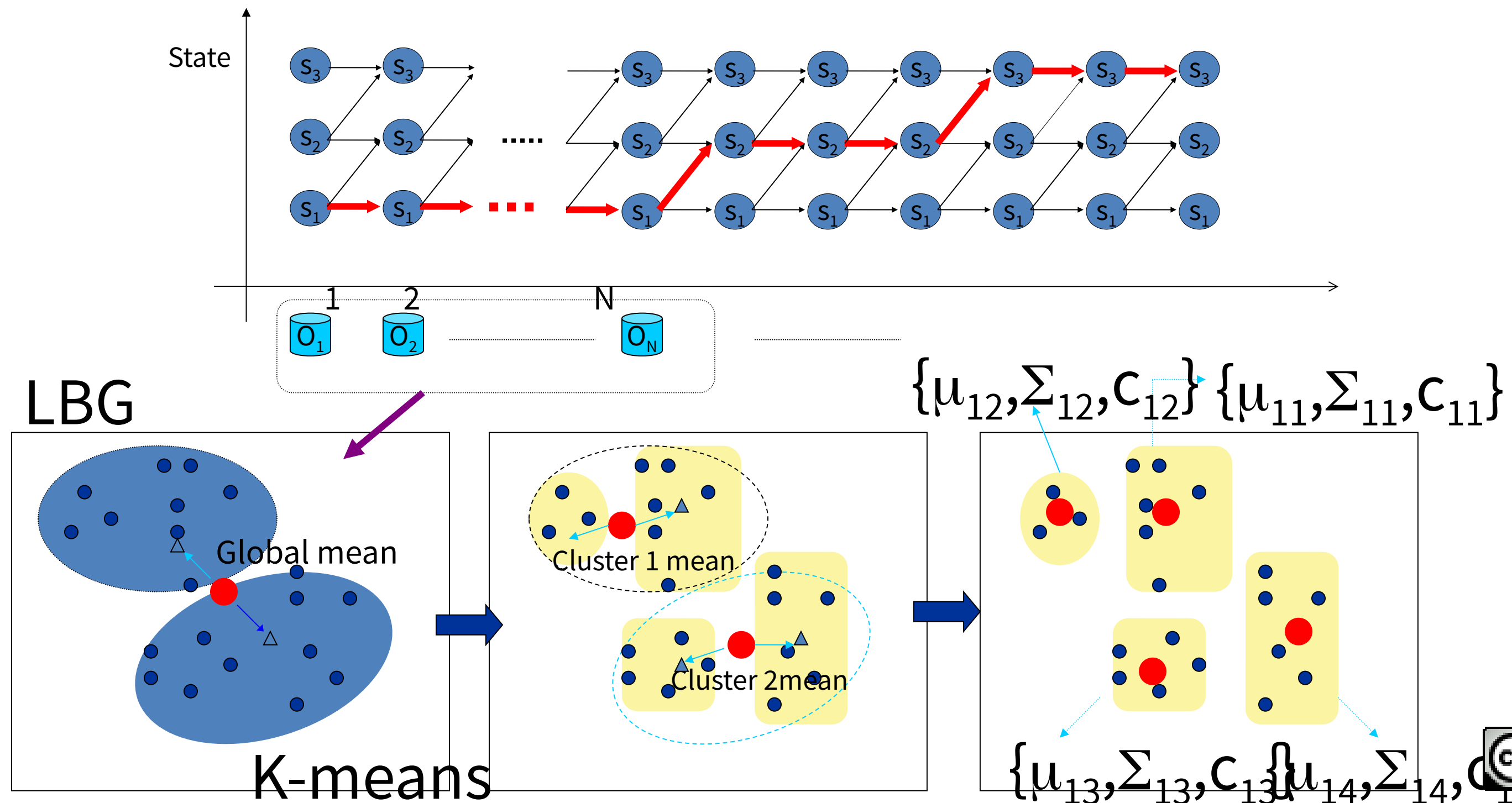
- Step 1 : re-segment the training observation sequences into states based on the initial model by Viterbi Algorithm
- Step 2 : Reestimate the model parameters (same as initial estimation)
- Step 3: Evaluate the model score $P(O|\lambda)$:
If the difference between the previous and current model scores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

Segmental K-Means



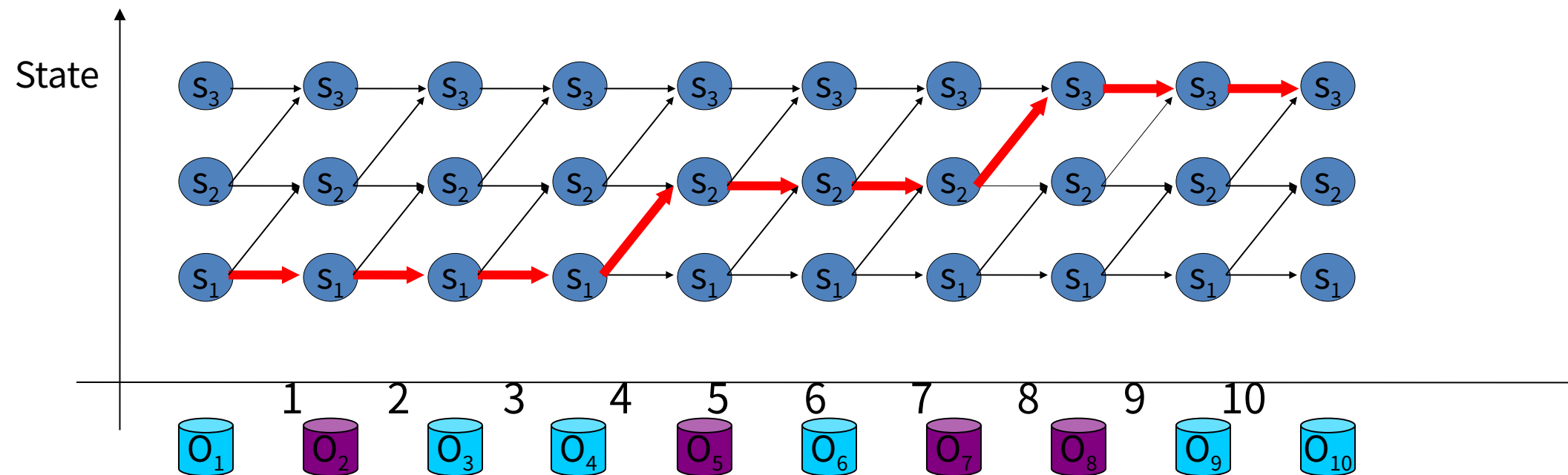
Initialization in HMM Training

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



Initialization in HMM Training



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



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

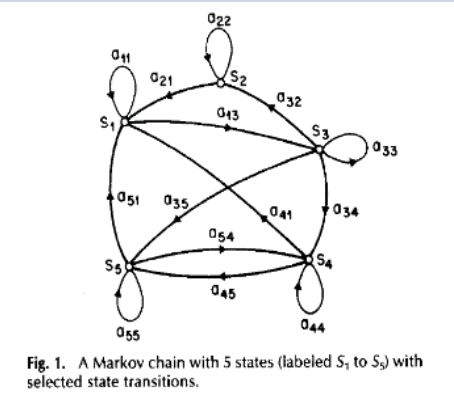

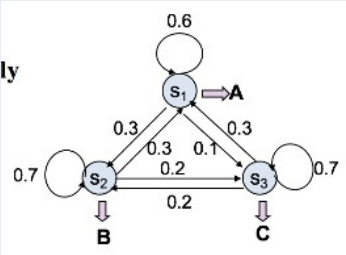

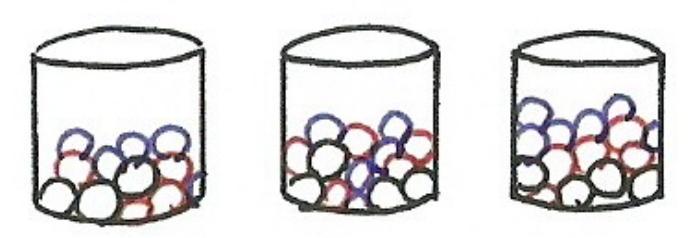
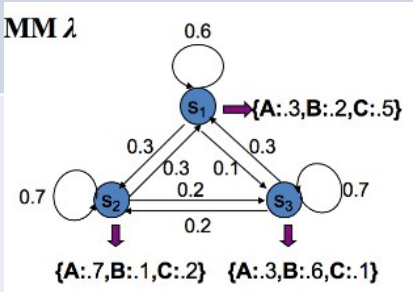

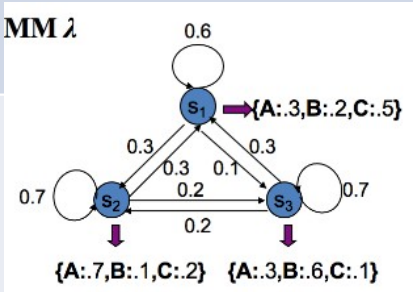

$$b_2(\mathbf{v}_1)=1/3, b_2(\mathbf{v}_2)=2/3$$

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

\mathbf{v}_1 
 \mathbf{v}_2 



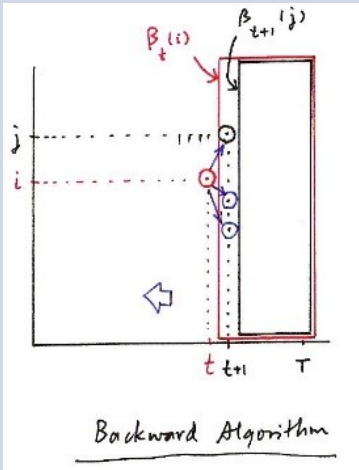

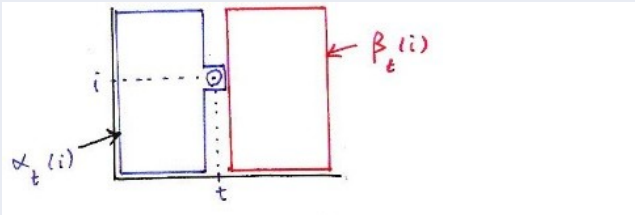

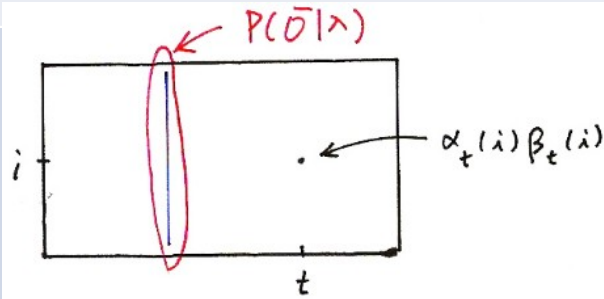
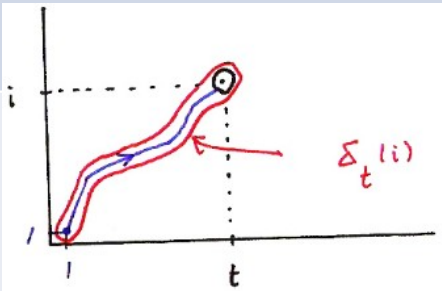


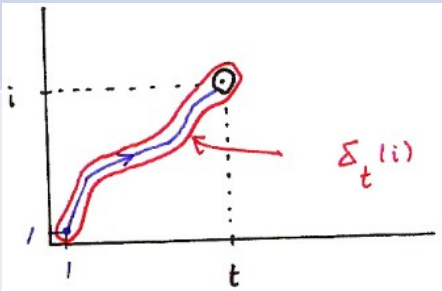

版權聲明

頁碼	作品	版權標示	作者 / 來源
2	 <p>Fig. 1. A Markov chain with 5 states (labeled S_1 to S_5) with selected state transitions.</p>		Lawrence Rabiner, Biing-Hwang Juang / FUNDAMENTALS OF SPEECH RECOGNITION Chap. 6, Sec. 6.2 Discrete-Time Markov Processes, page 323, Prentice- Hall International, Inc.
3			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。
5	 		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。
7			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

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。 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。 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。 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。

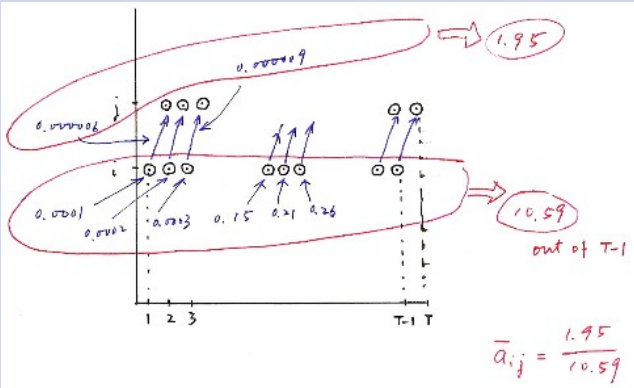

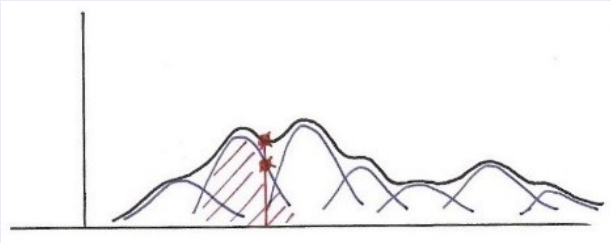

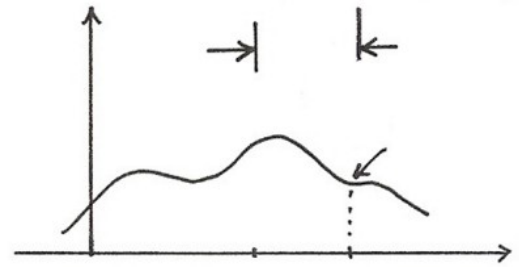
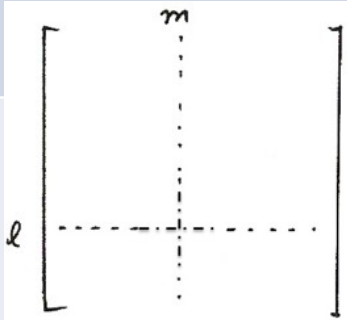

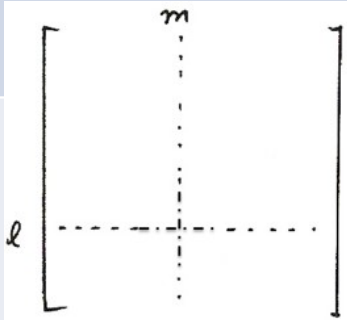

版權聲明

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。

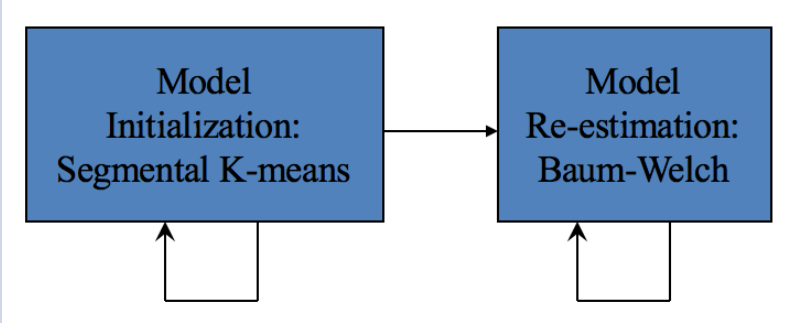

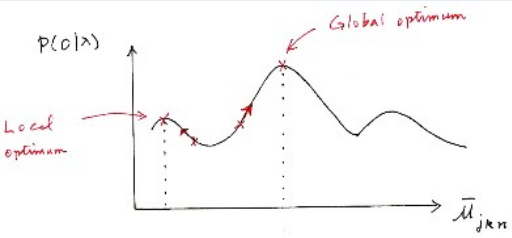

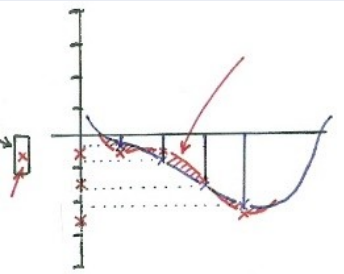
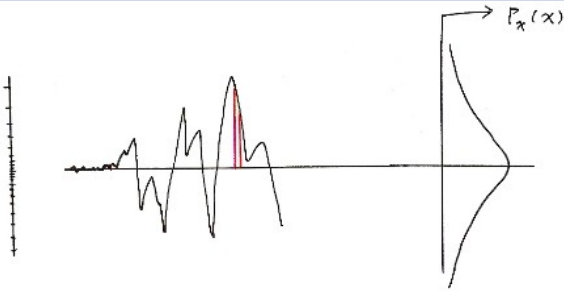


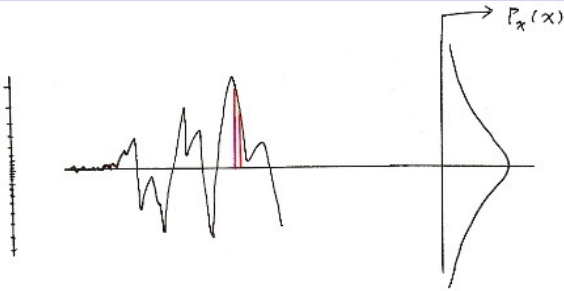

版權聲明

頁碼	作品	版權標示	作者 / 來源
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。
22			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。
26	 	 	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。
27			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。

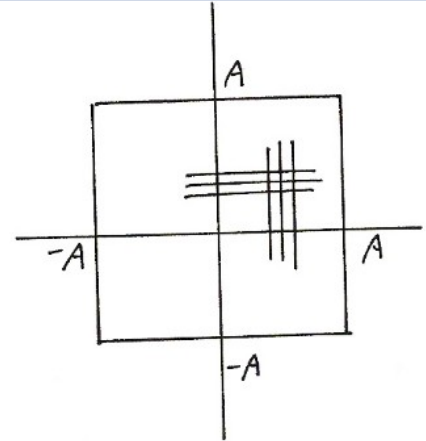

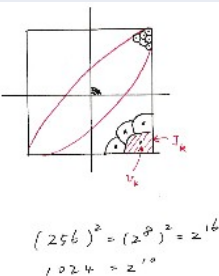

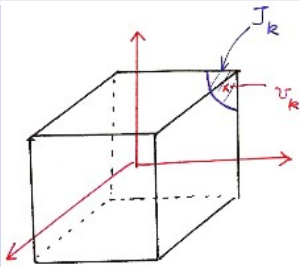
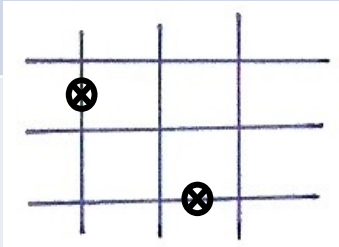



版權聲明

頁碼	作品	版權標示	作者 / 來源
29			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。
31			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。
33	 		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。
34			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。

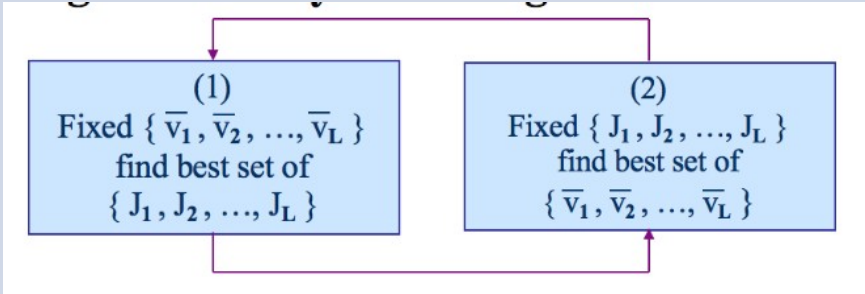

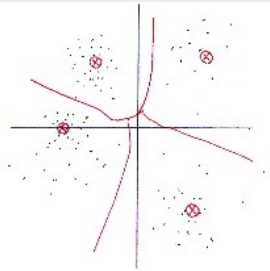

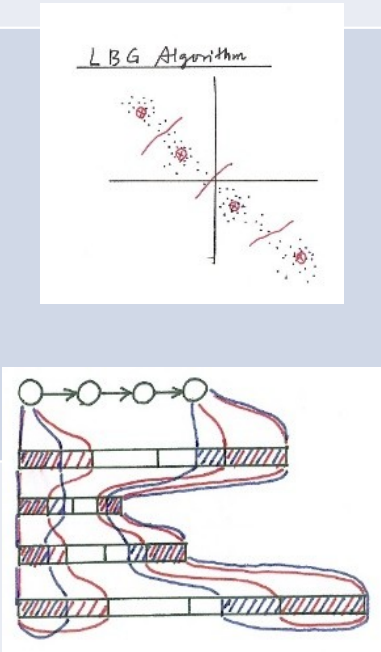


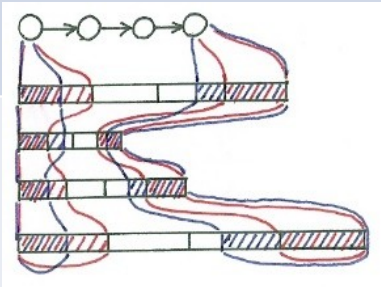

版權聲明

頁碼	作品	版權標示	作者 / 來源
35			<p>國立臺灣大學電機工程學系李琳山 教授。</p> <p>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。</p>
36			<p>國立臺灣大學電機工程學系李琳山 教授。</p> <p>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。</p>
38	 	 	<p>國立臺灣大學電機工程學系李琳山 教授。</p> <p>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。</p>
39			<p>國立臺灣大學電機工程學系李琳山 教授。</p> <p>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。</p>

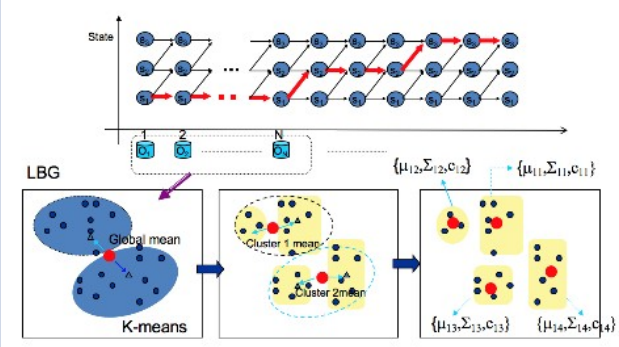

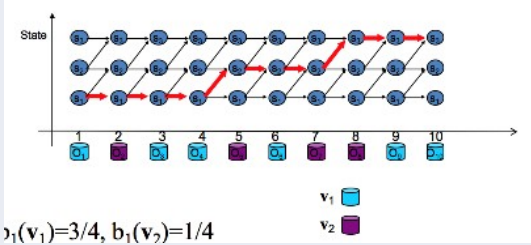

版權聲明

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。 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。 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。 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。 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
46			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。
47			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。
49		 	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。
51			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

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。 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。 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。