# The EM Algorithm

$o_1,o_2,\ldots\ldots,o_T$

$\lambda$

0.6

$s_1$ → {A:.3,B:.2,C:.5}

0.3   0.3

0.3   0.1

0.7   $s_2$   0.2   $s_3$   0.7

0.2

{A:.7,B:.1,C:.2}   {A:.3,B:.6,C:.1}

$p(\boldsymbol{O}|\lambda)$

$p(\boldsymbol{O}|\overline{\lambda}) > p(\boldsymbol{O}|\lambda)$

A   B

Observed data : $\boldsymbol{O}$ : "ball sequence"
Latent data : $\boldsymbol{S}$ : "bottle sequence"

Parameters to be estimated to maximize $\log P(\boldsymbol{O}|\lambda)$
$\lambda = \{P(A),P(B),P(B|A),P(A|B),P(R|A),P(G|A),P(R|B),P(G|B)\}$

57

# The EM Algorithm

- **Introduction of EM (Expectation Maximization):**
  - Why EM?
    - Simple optimization algorithms for likelihood function relies on the intermediate variables, called latent (隱藏的)data
      In our case here*, **the state sequence** is the latent data*
    - Direct access to the data necessary to estimate the parameters is impossible or difficult
      In our case here, it is almost impossible to estimate $\{A, B, \pi\}$ without consideration of the **state sequence**
  - Two Major Steps :
    - **E** : expectation with respect to the latent data using the current estimate of the parameters and conditioned on the observations $E\left[\bullet\right]_{s\,|\,\lambda,\,o}$
    - **M**: provides a new estimation of the parameters according to Maximum likelihood (ML) or Maximum A Posterior (MAP) Criteria

# The EM Algorithm

## ML and MAP

- Estimation principle based on observations:

$$x = (x_1, x_2, ..., x_n) \iff X = \{X_1, X_2, ..., X_n\}$$

  - **The Maximum Likelihood (ML) Principle**
  find the model parameter $\Phi$ so that the likelihood $p(x|\Phi)$ is maximum
  *for example, if $\Phi = \{\mu, \Sigma\}$ is the parameters of a multivariate normal distribution, and **X** is i.i.d. (independent, identically distributed), then the ML estimate of $\Phi = \{\mu, \Sigma\}$ is*

  $$\mu_{ML} = \frac{1}{n} \sum_{i=1}^{n} x_i \quad , \quad \Sigma_{ML} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu_{ML})(x_i - \mu_{ML})^t$$

  - **The Maximum A Posteriori (MAP) Principle**
  find the model parameter $\Phi$ so that the likelihood $p(\Phi|x)$ is maximum

# The EM Algorithm

- The EM Algorithm is important to HMMs and other learning techniques
  - Discover new model parameters to maximize the log-likelihood of incomplete data $\log P(O|\lambda)$ by iteratively maximizing the expectation of log-likelihood from complete data $\log P(O,S|\lambda)$

- Using scalar random variables to introduce the EM algorithm
  - The observable training data $O$
    - We want to maximize $P(O|\lambda)$, $\lambda$ is a parameter vector
  - The hidden (unobservable) data $S$
    - E.g. the component densities of observable data $O$, or the underlying state sequence in HMMs

# The EM Algorithm

- Assume we have $\lambda$ and estimate the probability that each $s$ occurred in the generation of $o$
- Pretend we had in fact observed a complete data pair $(o, s)$ with frequency proportional to the probability $P(o, s | \lambda)$, to computed a new $\bar{\lambda}$, the maximum likelihood estimate of $\lambda$
- Does the process converge?
- **Algorithm** unknown model setting

$$P(o, s | \bar{\lambda}) = P(s | o, \bar{\lambda}) P(o | \bar{\lambda}) \quad \text{Bayes' rule}$$

complete data likelihood          incomplete data likelihood

  - **Log-likelihood expression** and expectation taken over $S$

$$\log P(o | \bar{\lambda}) = \log P(o, s | \bar{\lambda}) - \log P(s | o, \bar{\lambda})$$

take expectation over $S$

$$\log P(o | \bar{\lambda}) = \sum_s \left[ P(s | o, \lambda) \log P(o | \bar{\lambda}) \right]$$

$$= \sum_s \left[ P(s | o, \lambda) \log P(o, s | \bar{\lambda}) \right] - \sum_s \left[ P(s | o, \lambda) \log P(s | o, \bar{\lambda}) \right]$$

61

# The EM Algorithm

– Algorithm (Cont.)
  - We can thus express $\log P\left(\boldsymbol{O}|\bar{\lambda}\right)$ as follows

$$\log P\left(\boldsymbol{O}|\bar{\lambda}\right)$$
$$= \sum_{S}\left[P(\boldsymbol{S}|\boldsymbol{O},\lambda)\log P\left(\boldsymbol{O},\boldsymbol{S}|\bar{\lambda}\right)\right] - \sum_{S}\left[P(\boldsymbol{S}|\boldsymbol{O},\lambda)\log P\left(\boldsymbol{S}|\boldsymbol{O},\bar{\lambda}\right)\right]$$
$$= Q\left(\lambda,\bar{\lambda}\right) - H\left(\lambda,\bar{\lambda}\right)$$

where

$$Q\left(\lambda,\bar{\lambda}\right) = \sum_{S}\left[P(\boldsymbol{S}|\boldsymbol{O},\lambda)\log P\left(\boldsymbol{O},\boldsymbol{S}|\bar{\lambda}\right)\right]$$
$$H\left(\lambda,\bar{\lambda}\right) = \sum_{S}\left[P(\boldsymbol{S}|\boldsymbol{O},\lambda)\log P\left(\boldsymbol{S}|\boldsymbol{O},\bar{\lambda}\right)\right]$$

  - We want $\log P\left(\boldsymbol{O}|\bar{\lambda}\right) \geq \log P\left(\boldsymbol{O}|\lambda\right)$

$$\log P\left(\boldsymbol{O}|\bar{\lambda}\right) - \log P\left(\boldsymbol{O}|\lambda\right)$$
$$= \left[Q\left(\lambda,\bar{\lambda}\right) - H\left(\lambda,\bar{\lambda}\right)\right] - \left[Q\left(\lambda,\lambda\right) - H\left(\lambda,\lambda\right)\right]$$
$$= Q\left(\lambda,\bar{\lambda}\right) - Q\left(\lambda,\lambda\right) - H\left(\lambda,\bar{\lambda}\right) + H\left(\lambda,\lambda\right)$$

# The EM Algorithm

- $-H(\lambda, \bar{\lambda}) + H(\lambda, \lambda)$ has the following property

$$-H(\lambda, \bar{\lambda}) + H(\lambda, \lambda)$$

$$= -\sum_S \left[ P(S|O, \lambda) \log \frac{P(S|O, \bar{\lambda})}{P(S|O, \lambda)} \right]$$

Kullbuack-Leibler (KL) distance

$$\geq \sum_S \left[ P(S|O, \lambda) \left( 1 - \frac{P(S|O, \bar{\lambda})}{P(S|O, \lambda)} \right) \right] \quad (\because \log x \leq x - 1)$$

Jensen's inequality

$$= \sum_S \left[ P(S|O, \lambda) - P(S|O, \bar{\lambda}) \right]$$

$$= 0$$

$$\therefore -H(\lambda, \bar{\lambda}) + H(\lambda, \lambda) \geq 0$$

- Therefore, for maximizing $\log P(O|\bar{\lambda})$, we only need to maximize the $Q$-function (auxiliary function)

$$Q(\lambda, \bar{\lambda}) = \sum_S \left[ P(S|O, \lambda) \log P(O, S|\bar{\lambda}) \right]$$

Expectation of the complete data log likelihood with respect to the latent state sequences

63

# EM Applied to Discrete HMM Training

- Apply EM algorithm to iteratively refine the HMM parameter vector $\lambda = (A, B, \pi)$
  - By maximizing the auxiliary function

$$Q(\lambda, \bar{\lambda}) = \sum_{S} \left[ P(S|O, \lambda) \log P(O, S|\bar{\lambda}) \right]$$

$$= \sum_{S} \left[ \frac{P(O, S|\lambda)}{P(O|\lambda)} \log P(O, S|\bar{\lambda}) \right]$$

  - Where $P(O, S|\lambda)$ and $P(O, S|\bar{\lambda})$ can be expressed as

$$P(O, S|\lambda) = \pi_{s_1} \left[ \prod_{t=1}^{T-1} a_{s_t s_{t+1}} \right] \left[ \prod_{t=1}^{T} b_{s_t}(o_t) \right]$$

$$\log P(O, S|\lambda) = \log \pi_{s_1} + \sum_{t=1}^{T-1} \log a_{s_t s_{t+1}} + \sum_{t=1}^{T} \log b_{s_t}(o_t)$$

$$\log P(O, S|\bar{\lambda}) = \log \bar{\pi}_{s_1} + \sum_{t=1}^{T-1} \log \bar{a}_{s_t s_{t+1}} + \sum_{t=1}^{T} \log \bar{b}_{s_t}(o_t)$$
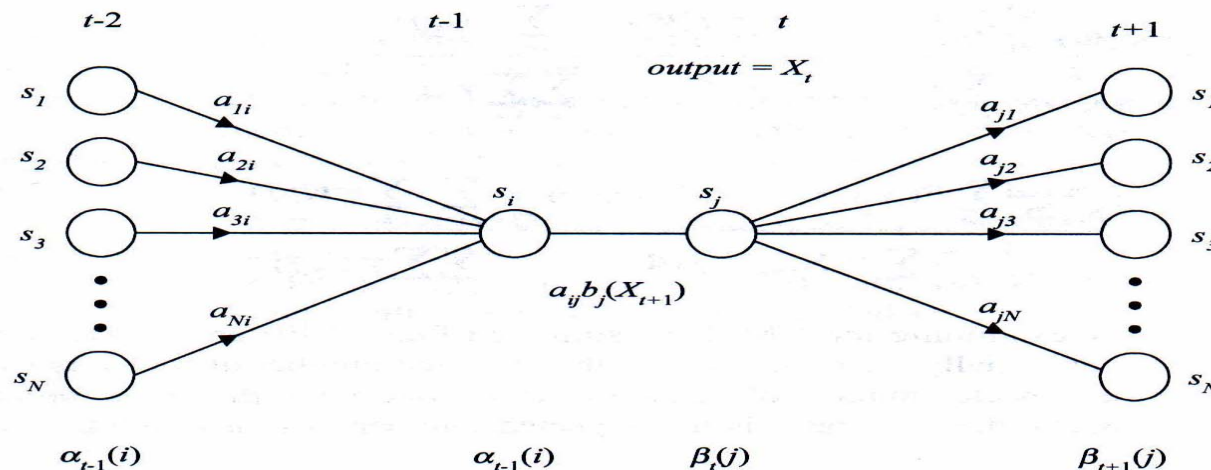
64

# EM Applied to Discrete HMM Training

- Rewrite the auxiliary function as

$$Q\left(\lambda,\overline{\lambda}\right)= Q_{\pi}\left(\lambda,\overline{\pi}\right)+ Q_{a}\left(\lambda,\overline{a}\right)+ Q_{b}\left(\lambda,\overline{b}\right)$$

$w_i$   $y_i$

$$Q_{\pi}\left(\lambda,\overline{\pi}\right)= \sum_{\text{all } S}\left[\frac{P\left(\boldsymbol{O},\boldsymbol{S}\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\log\ \overline{\pi}_{s_1}\right]\ \overset{?}{=}\ \sum_{i=1}^{N}\left[\frac{P\left(\boldsymbol{O},s_1 = i\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\log\ \overline{\pi}_i\right]$$

$$Q_{a}\left(\lambda,\overline{a}\right)= \sum_{\text{all } S}\left[\frac{P\left(\boldsymbol{O},\boldsymbol{S}\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\sum_{t=1}^{T-1}\log\ \overline{a}_{s_t s_{t+1}}\right]\ \overset{?}{=}\ \sum_{i=1}^{N}\sum_{j=1}^{N}\sum_{t=1}^{T-1}\left[\frac{P\left(\boldsymbol{O},s_t = i, s_{t+1} = j\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\log\ \overline{a}_{ij}\right]$$

$$Q_{b}\left(\lambda,\overline{b}\right)= \sum_{\text{all } S}\left[\frac{P\left(\boldsymbol{O},\boldsymbol{S}\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\sum_{t=1}^{T}\log\ \overline{b}_{s_t}(k)\right]\ \overset{?}{=}\ \sum_{j=1}^{N}\sum_{k}\sum_{t\in o_t = v_k}\left[\frac{P\left(\boldsymbol{O},s_t = j\mid\lambda\right)}{P\left(\boldsymbol{O}\mid\lambda\right)}\log\ \overline{b}_j(k)\right]$$



**Figure 8.7** Illustration of the operations required for the computation of $\gamma_t(i,j)$, which is the probability of taking the transition from state $i$ to state $j$ at time $t$.

65

# EM Applied to Discrete HMM Training

- The auxiliary function contains three independent terms, $\pi_i$, $a_{ij}$ and $b_j(k)$
  - Can be maximized individually
  - All of the same form

$$F(\boldsymbol{y}) = g(y_1, y_2, ...., y_N) = \sum_{j=1}^{N} w_j \, log \, y_j, \quad \text{where } \sum_{j=1}^{N} y_j = 1, \text{ and } y_j \geq 0$$

$$F(\boldsymbol{y}) \text{ has maximum value when}: y_j = \frac{w_j}{\sum_{j=1}^{N} w_j}$$

# EM Applied to Discrete HMM Training

- **Proof**: Apply Lagrange Multiplier

  By applying Lagrange Multiplier $\ell$

  Suppose that $F = \sum\limits_{j=1}^{N} w_j \, log \, y_j = \sum\limits_{j=1}^{N} w_j \, log \, y_j + \ell\left(\sum\limits_{j=1}^{N} y_j - 1\right)$

  **Constraint**

  $$\frac{\partial F}{\partial y_j} = \frac{w_j}{y_j} + \ell = 0 \Rightarrow \ell = -\frac{w_j}{y_j} \; \forall \, j$$

  $$\ell \sum\limits_{j=1}^{N} y_j = -\sum\limits_{j=1}^{N} w_j \Rightarrow \ell = -\sum\limits_{j=1}^{N} w_j$$

  $$\therefore y_j = \frac{w_j}{\sum\limits_{j=1}^{N} w_j}$$

# EM Applied to Discrete HMM Training

- The new model parameter set $\bar{\lambda} = \left( \bar{\pi}, \bar{A}, \bar{B} \right)$ can be expressed as:

$$\bar{\pi}_i = \frac{P\left( \boldsymbol{O}, s_1 = i \mid \lambda \right)}{P\left( \boldsymbol{O} \mid \lambda \right)} = \gamma_1(i)$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} P\left( \boldsymbol{O}, s_t = i, s_{t+1} = j \mid \lambda \right)}{\sum_{t=1}^{T-1} P\left( \boldsymbol{O}, s_t = i \mid \lambda \right)} = \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\bar{b}_i(k) = \frac{\sum_{\substack{t=1 \\ \text{s.t.} \ o_t = v_k}}^{T} P\left( \boldsymbol{O}, s_t = i \mid \lambda \right)}{\sum_{t=1}^{T} P\left( \boldsymbol{O}, s_t = i \mid \lambda \right)} = \frac{\sum_{\substack{t=1 \\ \text{s.t.} \ o_t = v_k}}^{T} \gamma_t(i)}{\sum_{t=1}^{T} \gamma_t(i)}$$

# EM Applied to Continuous HMM Training

- Continuous HMM: the state observation does not come from a finite set, but from a continuous space
  - The difference between the discrete and continuous HMM lies in a different form of state output probability
  - Discrete HMM requires the quantization procedure to map observation vectors from the continuous space to the discrete space

- Continuous Mixture HMM
  - The state observation distribution of HMM is modeled by multivariate Gaussian mixture density functions (*M* mixtures)

$$b_j(\boldsymbol{o}) = \sum_{k=1}^{M} c_{jk} b_{jk}(\boldsymbol{o})$$

$$= \sum_{k=1}^{M} c_{jk} N(\boldsymbol{o}; \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}) = \sum_{k=1}^{M} c_{jk} \left( \frac{1}{\left(\sqrt{2\pi}\right)^{L} |\boldsymbol{\Sigma}_{jk}|^{1/2}} \exp\left( -\frac{1}{2} (\boldsymbol{o} - \boldsymbol{\mu}_{jk})^{t} \boldsymbol{\Sigma}_{jk}^{-1} (\boldsymbol{o} - \boldsymbol{\mu}_{jk}) \right) \right)$$

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



Distribution for State *i*

69

# EM Applied to Continuous HMM Training

- Express $b_j(\boldsymbol{o})$ with respect to each single mixture component $b_{jk}(\boldsymbol{o})$

Note:

$$\prod_{t=1}^{T}\left(\sum_{k_t=1}^{M} a_{tk_t}\right)$$

$$= (a_{11} + a_{12} + \ldots + a_{1M})(a_{21} + a_{22} + \ldots + a_{2M})\ldots(a_{T1} + a_{T2} + \ldots + a_{TM})$$

$$= \sum_{k_1=1}^{M}\sum_{k_2=1}^{M}\ldots\sum_{k_T=1}^{M}\prod_{t=1}^{T} a_{tk_t}$$

$$P(\boldsymbol{O}, \boldsymbol{S}|\lambda) = \pi_{s_1}\left\{\prod_{t=1}^{T-1} a_{s_t s_{t+1}}\right\}\left\{\prod_{t=1}^{T} b_{s_t}(\boldsymbol{o}_t)\right\}$$

$$= \pi_{s_1}\left\{\prod_{t=1}^{T-1} a_{s_t s_{t+1}}\right\}\left\{\sum_{k_1=1}^{M}\sum_{k_2=1}^{M}\ldots\sum_{k_T=1}^{M}\prod_{t=1}^{T}\left[c_{s_t k_t} b_{s_t k_t}(\boldsymbol{o}_t)\right]\right\}$$

$$P(\boldsymbol{O}, \boldsymbol{S}, \boldsymbol{K}|\lambda) = \pi_{s_1}\left\{\prod_{t=1}^{T-1} a_{s_t s_{t+1}}\right\}\left\{\prod_{t=1}^{T}\left[c_{s_t k_t} b_{s_t k_t}(\boldsymbol{o}_t)\right]\right\}$$

$\boldsymbol{K}$ : one of the possible mixture component sequence
  along with the state sequence $\boldsymbol{S}$

$$P(\boldsymbol{O}|\lambda) = \sum_{\boldsymbol{S}}\sum_{\boldsymbol{K}} P(\boldsymbol{O}, \boldsymbol{S}, \boldsymbol{K}|\lambda)$$
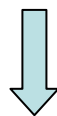
# EM Applied to Continuous HMM Training

- Therefore, an auxiliary function for the EM algorithm can be written as:

$$Q(\lambda, \bar{\lambda}) = \sum_{S} \sum_{K} \left[ P(S, K | O, \lambda) \log P(O, S, K | \bar{\lambda}) \right]$$

$$= \sum_{S} \sum_{K} \left[ \frac{P(O, S, K | \lambda)}{P(O | \lambda)} \log P(O, S, K | \bar{\lambda}) \right]$$

$$\log P(O, S, K | \bar{\lambda}) = \log \bar{\pi}_{s_1} + \sum_{t=1}^{T-1} \log \bar{a}_{s_t s_{t+1}} + \sum_{t=1}^{T} \log \bar{b}_{s_t k_t}(o_t) + \sum_{t=1}^{T} \log \bar{c}_{s_t k_t}$$

$$Q(\lambda, \bar{\lambda}) = Q_{\pi}(\lambda, \bar{\pi}) + Q_{a}(\lambda, \bar{a}) + Q_{b}(\lambda, \bar{b}) + Q_{c}(\lambda, \bar{c})$$

initial probabilities     state transition probabilities     Gaussian density functions     mixture components

71

# EM Applied to Continuous HMM Training

- The only difference we have when compared with Discrete HMM training

$$\gamma_t(j,k)$$

$$Q_b\left(\lambda, \overline{b}\right) = \sum_{t=1}^{T} \left\{ \left[ \sum_{j=1}^{N} \sum_{k=1}^{M} P\left(s_t = j, k_t = k \mid \boldsymbol{O}, \lambda\right) \right] \log \overline{b}_{jk}\left(\boldsymbol{o}_t\right) \right\}$$

$$Q_c\left(\lambda, \overline{c}\right) = \sum_{t=1}^{T} \left\{ \left[ \sum_{j=1}^{N} \sum_{k=1}^{M} P\left(s_t = j, k_t = k \mid \boldsymbol{O}, \lambda\right) \right] \log \overline{c}_{jk}\left(\boldsymbol{o}_t\right) \right\}$$

# EM Applied to Continuous HMM Training

Let $\gamma_t(j,k) = \sum\limits_{k=1}^{M} P(s_t = j, k_t = k | \boldsymbol{O}, \lambda)$

$\bar{b}_{jk}(\boldsymbol{o}_t) = N(\boldsymbol{o}_t; \bar{\boldsymbol{\mu}}_{jk}, \bar{\boldsymbol{\Sigma}}_{jk}) = \dfrac{1}{(2\pi)^{L/2} |\bar{\boldsymbol{\Sigma}}_{jk}|^{1/2}} \exp\left(-\dfrac{1}{2}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})' \bar{\boldsymbol{\Sigma}}_{jk}^{-1}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})\right)$

$\log \bar{b}_{jk}(\boldsymbol{o}_t) = -L/2 \cdot \log(2\pi) + 1/2 \cdot \log|\bar{\boldsymbol{\Sigma}}_{jk}^{-1}| - \left(\dfrac{1}{2}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})' \bar{\boldsymbol{\Sigma}}_{jk}^{-1}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})\right)$

$\dfrac{\partial \log \bar{b}_{jk}(\boldsymbol{o}_t)}{\partial \bar{\boldsymbol{\mu}}_{jk}} = \bar{\boldsymbol{\Sigma}}_{jk}^{-1}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})$

$\dfrac{d(\boldsymbol{x}^T \boldsymbol{C}\boldsymbol{x})}{d\boldsymbol{x}} = (\boldsymbol{C} + \boldsymbol{C}^T)\boldsymbol{x}$

and $\Sigma_{jk}^{-1}$ is symmetric here

$\dfrac{\partial Q_b(\lambda, \bar{b})}{\partial \bar{\boldsymbol{\mu}}_{jk}} = \dfrac{\partial \sum\limits_{t=1}^{T}\left\{\left[\sum\limits_{j=1}^{N}\sum\limits_{k=1}^{M}\gamma_t(j,k)\log \bar{b}_{jk}(\boldsymbol{o}_t)\right]\right\}}{\partial \bar{\boldsymbol{\mu}}_{jk}}$

$\Rightarrow \sum\limits_{t=1}^{T}\left\{\gamma_t(j,k)\bar{\boldsymbol{\Sigma}}_{jk}^{-1}(\boldsymbol{o}_t - \bar{\boldsymbol{\mu}}_{jk})\right\} = 0$

$\Rightarrow \bar{\boldsymbol{\mu}}_{jk} = \dfrac{\sum\limits_{t=1}^{T}[\gamma_t(j,k)\cdot\boldsymbol{o}_t]}{\sum\limits_{t=1}^{T}\gamma_t(j,k)}$

# EM Applied to Continuous HMM Training

$$\log \overline{b}_{jk}(o_t) = -\frac{L}{2} \cdot \log(2\pi) - \frac{1}{2} \cdot \log\left|\overline{\Sigma}_{jk}\right| - \left(\frac{1}{2}(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}(o_t - \overline{\mu}_{jk})\right)$$

$$\frac{\partial \log \overline{b}_{jk}(o_t)}{\partial(\overline{\Sigma}_{jk})} = -\left[\frac{1}{2}\cdot\left|\overline{\Sigma}_{jk}\right|^{-1}\cdot\left|\overline{\Sigma}_{jk}\right|\cdot\overline{\Sigma}_{jk}^{-1} - \left(\overline{\Sigma}_{jk}^{-1}\frac{1}{2}(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}\right)\right]$$

$$= -\frac{1}{2}\cdot\left[\overline{\Sigma}_{jk}^{-1} - \overline{\Sigma}_{jk}^{-1}(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}\right]$$

$$\frac{d(a^T X^{-1} b)}{dX} = -X^T a b^T X^T$$

$$\frac{d[\det(X)]}{dX} = \det(X)\cdot X^{-T}$$

and $\Sigma_{jk}$ is symmetric here

$$\frac{\partial Q_b(\lambda, \overline{b})}{\partial(\overline{\Sigma}_{jk})} = \frac{\partial \sum\limits_{t=1}^{T}\left\{\left[\sum\limits_{j=1}^{N}\sum\limits_{k=1}^{M}\gamma_t(j,k)\log \overline{b}_{jk}(o_t)\right]\right\}}{\partial(\overline{\Sigma}_{jk}^{-1})}$$

$$\Rightarrow \sum_{t=1}^{T}\left\{\gamma_t(j,k)\left(-\frac{1}{2}\right)\cdot\left[\overline{\Sigma}_{jk}^{-1} - \overline{\Sigma}_{jk}^{-1}(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}\right]\right\} = 0$$

$$\Rightarrow \sum_{t=1}^{T}\gamma_t(j,k)\overline{\Sigma}_{jk}^{-1} = \sum_{t=1}^{T}\gamma_t(j,k)\overline{\Sigma}_{jk}^{-1}(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}$$

$$\Rightarrow \sum_{t=1}^{T}\gamma_t(j,k)\overline{\Sigma}_{jk}\overline{\Sigma}_{jk}^{-1}\overline{\Sigma}_{jk} = \sum_{t=1}^{T}\gamma_t(j,k)\overline{\Sigma}_{jk}\overline{\Sigma}_{jk}^{-1}(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\overline{\Sigma}_{jk}^{-1}\overline{\Sigma}_{jk}$$

$$\Rightarrow \overline{\Sigma}_{jk} = \frac{\sum\limits_{t=1}^{T}\left[\gamma_t(j,k)\cdot(o_t - \overline{\mu}_{jk})(o_t - \overline{\mu}_{jk})'\right]}{\sum\limits_{t=1}^{T}\gamma_t(j,k)}$$

# EM Applied to Continuous HMM Training

- The new model parameter set for each mixture component and mixture weight can be expressed as:

$$\overline{\boldsymbol{\mu}}_{jk} = \frac{\sum_{t=1}^{T}\left[\dfrac{P(\boldsymbol{O}, s_t = j, k_t = k | \lambda)}{P(\boldsymbol{O}|\lambda)} \boldsymbol{o}_t\right]}{\sum_{t=1}^{T}\dfrac{P(\boldsymbol{O}, s_t = j, k_t = k|\lambda)}{P(\boldsymbol{O}|\lambda)}} = \frac{\sum_{t=1}^{T}[\gamma_t(j,k)\boldsymbol{o}_t]}{\sum_{t=1}^{T}\gamma_t(j,k)}$$

$$\overline{\boldsymbol{\Sigma}}_{jk} = \frac{\sum_{t=1}^{T}\left[\dfrac{P(\boldsymbol{O}, s_t = j, k_t = k|\lambda)}{P(\boldsymbol{O}|\lambda)}(\boldsymbol{o}_t - \overline{\boldsymbol{\mu}}_{jk})(\boldsymbol{o}_t - \overline{\boldsymbol{\mu}}_{jk})^t\right]}{\sum_{t=1}^{T}\dfrac{P(\boldsymbol{O}, s_t = j, k_t = k|\lambda)}{P(\boldsymbol{O}|\lambda)}} = \frac{\sum_{t=1}^{T}\left[\gamma_t(j,k)(\boldsymbol{o}_t - \overline{\boldsymbol{\mu}}_{jk})(\boldsymbol{o}_t - \overline{\boldsymbol{\mu}}_{jk})^t\right]}{\sum_{t=1}^{T}\gamma_t(j,k)}$$

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