



# Expectation-Maximization Algorithm

DCS310

Sun Yat-sen University

# Outline

---

- The Concerned Problem
- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models

# General Form of the Concerned Problem

---

- Given the joint distribution

$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}),$$

where  $\mathbf{x}$  is the observed variable and  $\mathbf{z}$  is the latent variable, we need to maximize the log likelihood w.r.t.  $\boldsymbol{\theta}$ , that is,

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}),$$

where

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$$

What we have is the joint pdf  $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$ , but what we need to optimize is the marginal pdf  $p(\mathbf{x}; \boldsymbol{\theta})$

# Outline

---

- The Concerned Problem
- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models

# EM Algorithm

---

- Algorithm

*E-step:* Evaluating the expectation

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}[\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]$$

*M-step:* Updating the parameter

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$$

- Key ingredient in EM

- 1) The posteriori distribution of latent variables  $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$
- 2) The expectation of joint distribution  $\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$  w.r.t. the posteriori
- 3) Maximization

# Outline

---

- The Concerned Problem
- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models

# Re-representing the Log-likelihood

- The log-likelihood can be reformulated as

$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}) &= \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} \\&= \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}) q(\mathbf{z})} \\&= \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})}}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})}}_{KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}))} \\&= \mathcal{L}(q, \boldsymbol{\theta}) + KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})), \quad \text{for } \forall \boldsymbol{\theta}, q(\mathbf{z})\end{aligned}$$

*Remark:* The KL-divergence is used to *measure the distance* between two distributions  $q$  and  $p$ , which is defined as

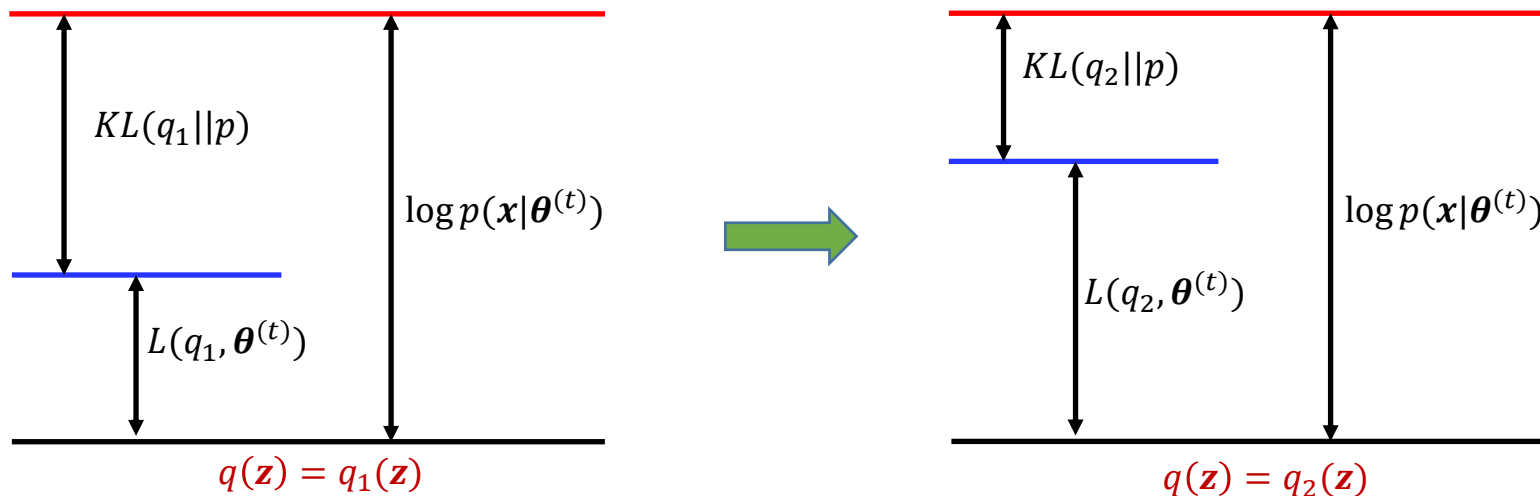
$$KL(q||p) \triangleq \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \geq 0$$

- Thus, with the parameter  $\theta^{(t)}$  at the  $t$ -th iteration, we have

$$\log p(\mathbf{x}; \theta^{(t)}) = \mathcal{L}(q, \theta^{(t)}) + KL(q||p(\mathbf{z}|\mathbf{x}; \theta^{(t)}))$$

This equality holds for any distribution  $q(\mathbf{z})$

- Different  $q(\mathbf{z})$  will lead to different decomposition of  $\log p(\mathbf{x}; \theta^{(t)})$





# Theoretical Justification for EM

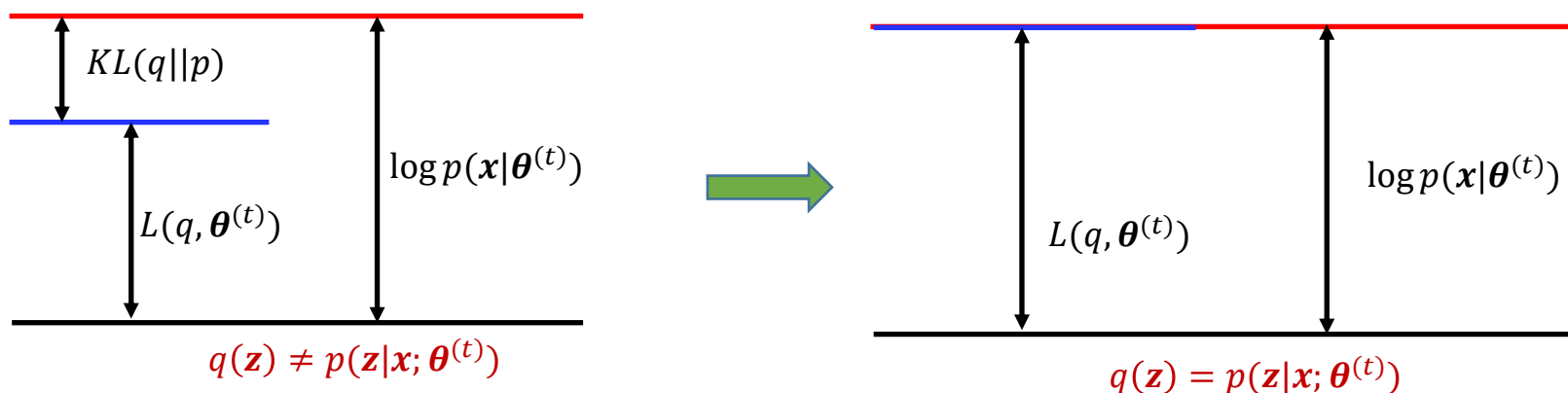
- If we set  $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ , then we have

$$KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})) = 0 \quad \text{i.e., } q \text{ and } p \text{ have no distance}$$

Thus, we have

$$\log p(\mathbf{x}|\boldsymbol{\theta}^{(t)}) = \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)})$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$



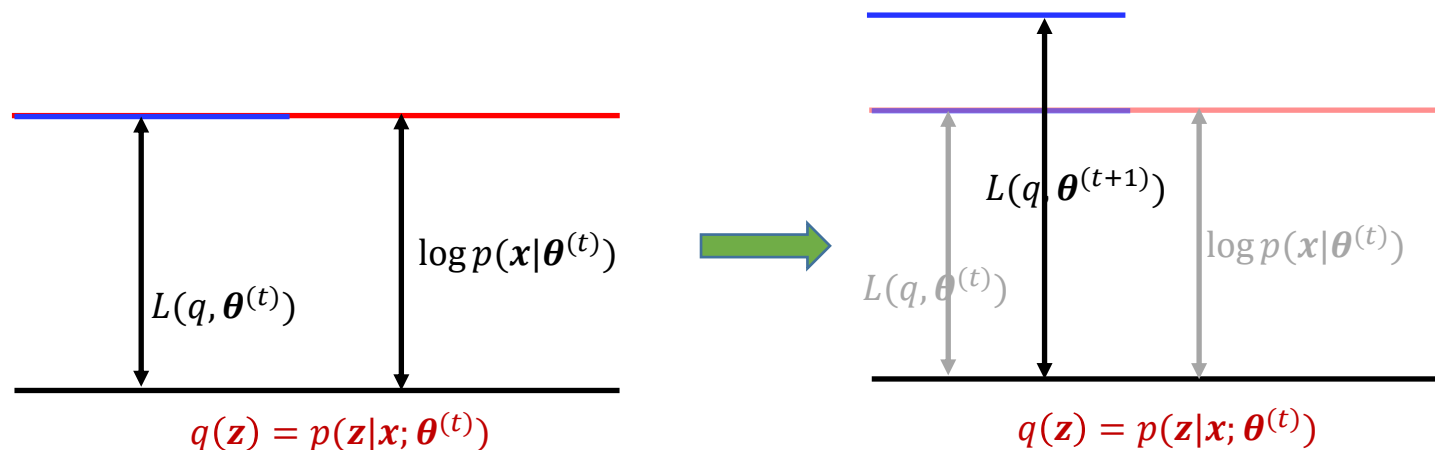
$$\begin{aligned}\log p(\mathbf{x}|\boldsymbol{\theta}^{(t)}) &= \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)}) \\ &= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}\end{aligned}$$

- If we update  $\boldsymbol{\theta}$  as

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}),$$

then we must have the relation

$$\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t+1)}) \geq \underbrace{\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)})}_{=\log p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}$$



- From the nonnegative property of KL-divergence, we know that

$$KL(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) || p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t+1)})) \geq 0$$

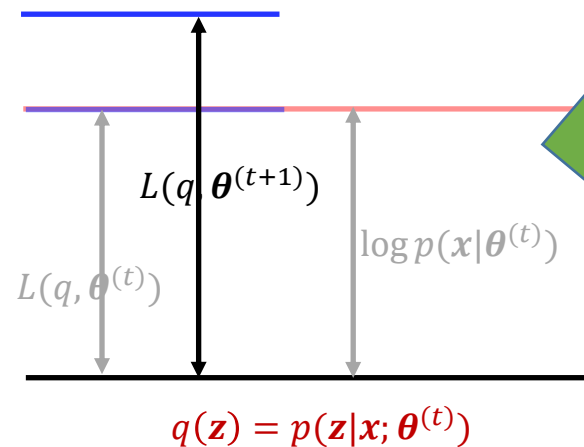
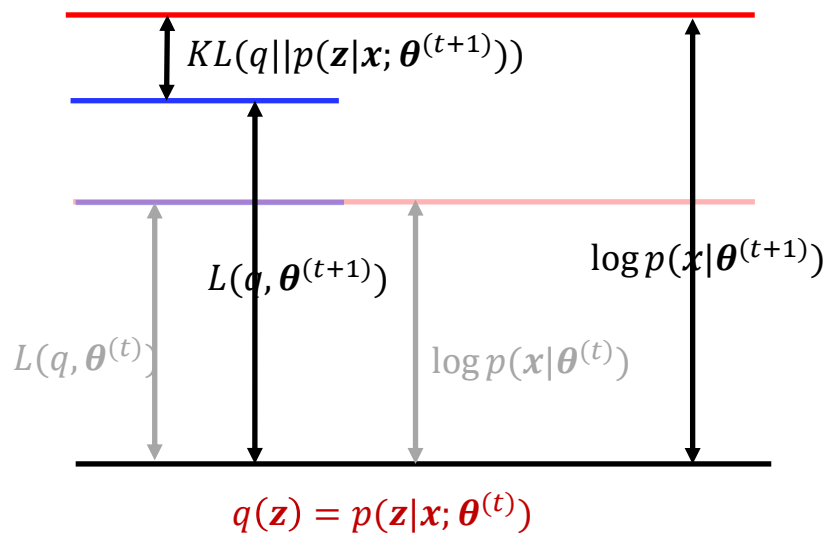
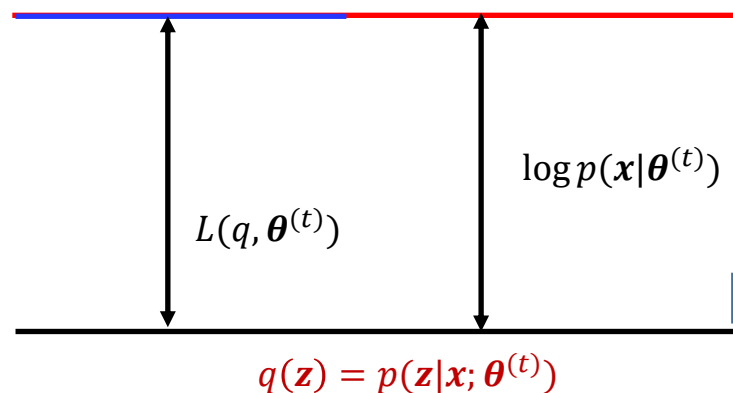
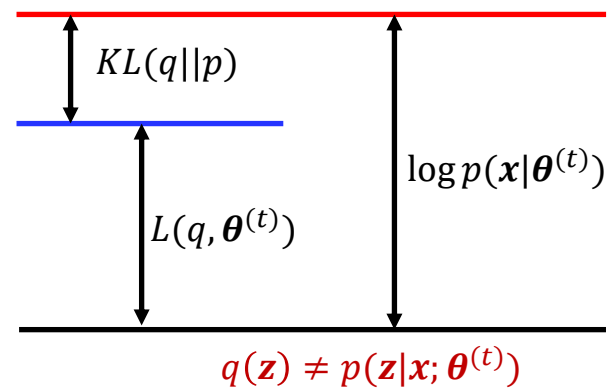
- Because  $\log p(\mathbf{x}; \boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q || p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}))$  holds for any  $q$ , thus we have

$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) = \underbrace{\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t+1)})}_{\geq \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})} + \underbrace{KL(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) || p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t+1)}))}_{\geq 0}$$

- Thus, we can see that

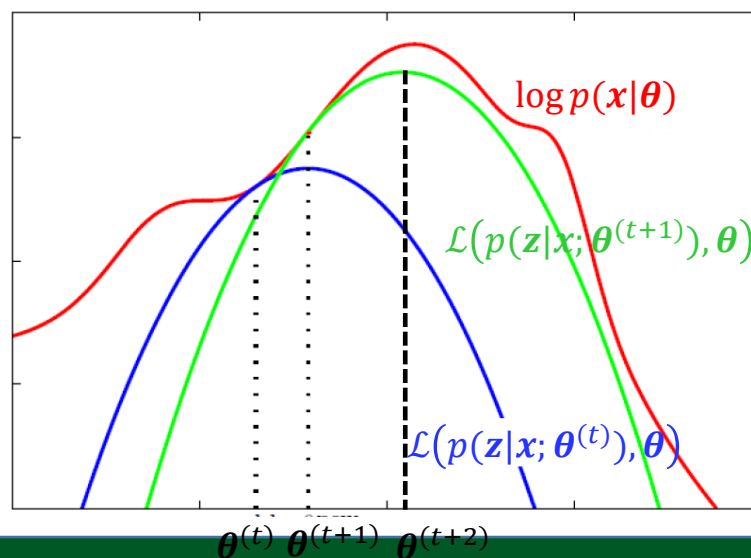
$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) \geq \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})$$

*EM algorithm can guarantee the increase of likelihood at each step*



# A View in the Parameter Space

- 1) E-step ( $t$ ): deriving the expression  $\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta})$  given the model parameter  $\boldsymbol{\theta}^{(t)}$
- 2) M-step ( $t$ ): computing the optimal value  $\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta})$
- 3) E-step ( $t+1$ ): deriving the expression for  $\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t+1)}), \boldsymbol{\theta})$  given the model parameter  $\boldsymbol{\theta}^{(t+1)}$
- 4) Repeating the above process until convergence



# Outline

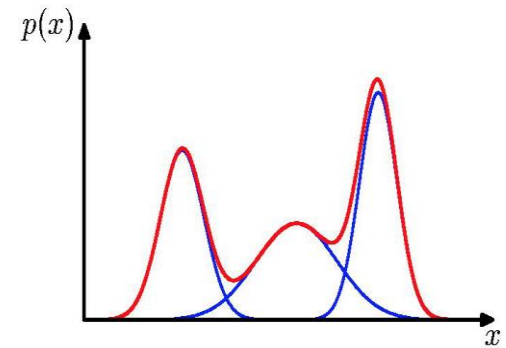
---

- The Concerned Problem
- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models

# Gaussian Mixture Model Review

- For a Gaussian mixture distribution, *i.e.*,

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$



it can be represented as the marginal distribution of the joint distribution

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

$$= \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_k}$$

- $\mathbf{z} = [z_1, z_2, \dots, z_K]$  follows the categorical distribution with parameter  $\boldsymbol{\pi}$

# EM: E-step

---

- The posteriori distribution

$$p(\mathbf{z} = \mathbf{1}_k | \mathbf{x}; \boldsymbol{\theta}) = \frac{p(\mathbf{x}, \mathbf{z} = \mathbf{1}_k; \boldsymbol{\theta})}{\sum_{i=1}^K p(\mathbf{x}, \mathbf{z} = \mathbf{1}_i; \boldsymbol{\theta})}$$

$$= \frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k}{\sum_{i=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \pi_i}$$

- $\mathbf{1}_k$  denotes the one-hot vector with the  $k$ -th element being 1
- The log of the joint distribution  $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_k}$

$$\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{k=1}^K z_k \cdot [\log \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \log \pi_k]$$

*Note that  $\mathbf{z}$  can only be a one-hot vector*



- The expectation

$$\begin{aligned}\mathbb{E}_{p(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}^{(t)})}[\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})] \\ = \sum_{k=1}^K \mathbb{E}_{p(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}^{(t)})}[\mathbf{z}_k][\log \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \log \pi_k]\end{aligned}$$

➤ Due to  $p(\mathbf{z} = \mathbf{1}_k | \mathbf{x}; \boldsymbol{\theta}) = \frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k}{\sum_{i=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \pi_i}$ , we have

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}^{(t)})}[\mathbf{z}_k] = \frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)}) \pi_k}{\sum_{i=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i^{(t)}, \boldsymbol{\Sigma}_i^{(t)}) \pi_i} \triangleq \gamma_k^{(t)}$$

- Therefore, we have

$$\mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \sum_{k=1}^K \gamma_k^{(t)} [\log \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \log \pi_k]$$

- Taking  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$  into  $Q(\cdot)$  gives

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \sum_{k=1}^K \gamma_k^{(t)} \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) - \frac{1}{2}|\boldsymbol{\Sigma}_k| + \log \pi_k \right] + C$$

- $C$  is the constant

- So far, only one data example  $\mathbf{x}$  is considered
- If data  $\mathbf{x}^{(n)}$  for  $n = 1, 2, \dots, N$  are considered,  $Q(\cdot)$  becomes

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}^{(t)} \left[ -\frac{1}{2}(\mathbf{x}^{(n)} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}^{(n)} - \boldsymbol{\mu}_k) - \frac{1}{2}|\boldsymbol{\Sigma}_k| + \log \pi_k \right] + C$$

# EM: M-step

---

- By taking derivatives *w.r.t.*  $\boldsymbol{\mu}_k$ ,  $\boldsymbol{\Sigma}_k$  and  $\pi_k$  and setting them to zero, we obtain the optimal  $\boldsymbol{\theta}$  as

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k^{(t+1)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)})^T$$

$$\pi_k^{(t+1)} = \frac{N_k}{N}$$

where  $N_k = \sum_{n=1}^N \gamma_{nk}$  is the effective number of examples assigned to the  $k$ -th class

# Summary of EM Algorithm

---

- Given the current estimate  $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}_{k=1}^K$ , update  $\gamma_{nk}$  as

$$\gamma_{nk} \leftarrow \frac{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\pi_k}{\sum_{i=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\pi_i}$$

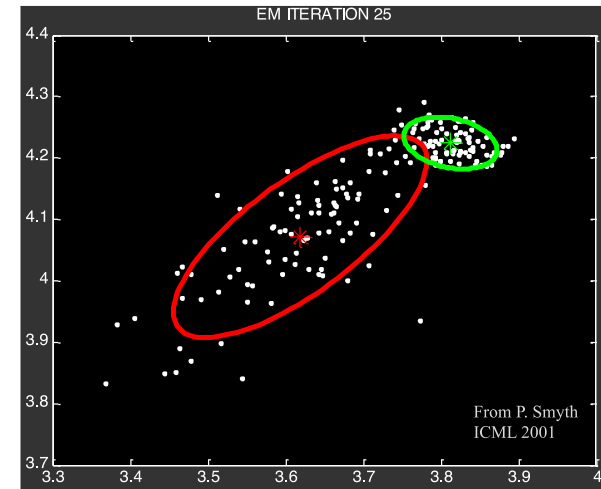
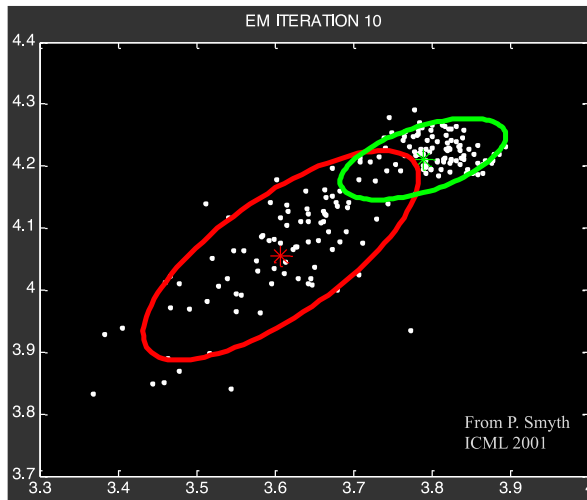
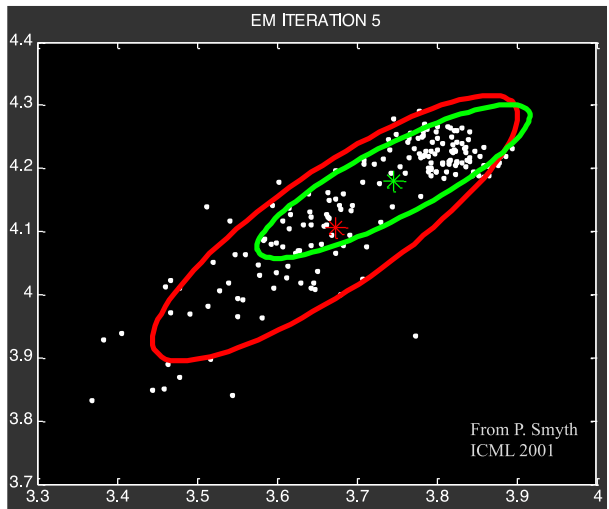
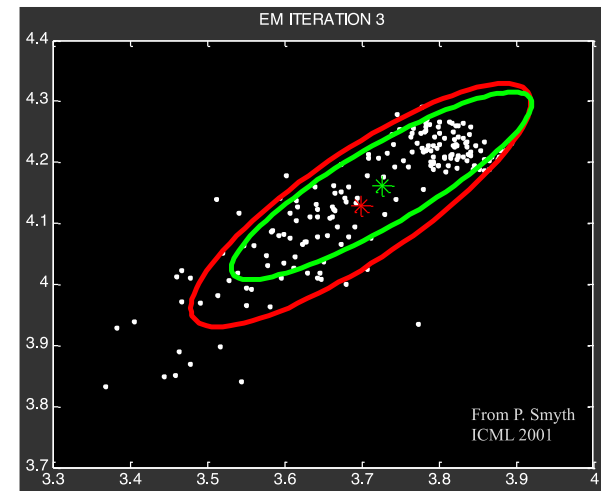
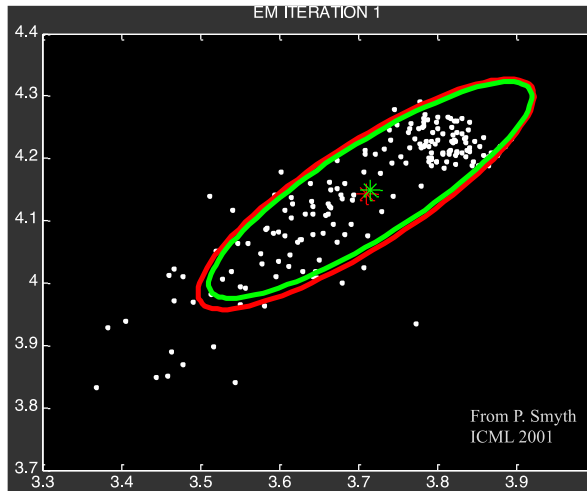
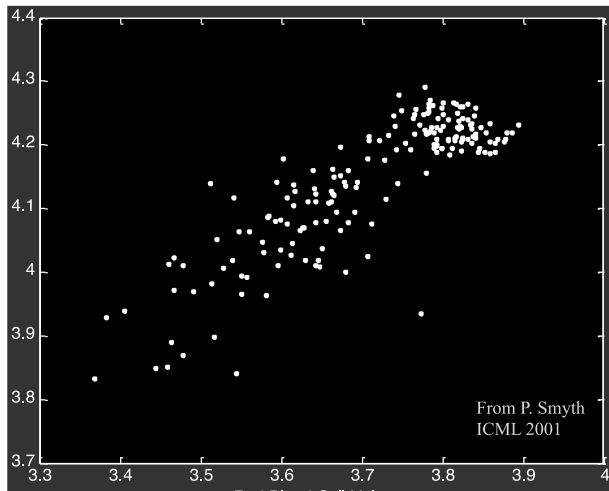
- Given the  $\gamma_{nk}$ , update  $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$  and  $\pi_k$  as

$$N_k \leftarrow \sum_{n=1}^N \gamma_{nk}$$

$$\boldsymbol{\mu}_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\pi_k \leftarrow \frac{N_k}{N}$$



# Relation to Soft $K$ -Means

- When restricting  $\Sigma_k = \sigma^2 \mathbf{I}$ , the updating of GMM becomes

$$\pi_k \leftarrow \frac{\sum_{n=1}^N \gamma_{nk}}{N}$$

$$\gamma_{nk} \leftarrow \frac{e^{-\beta_k \|x^{(n)} - \mu_k\|^2}}{\sum_{i=1}^K e^{-\beta_i \|x^{(n)} - \mu_i\|^2}}$$

$$\mu_k \leftarrow \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}}$$

where  $\beta_i = \frac{\ln \pi_i}{2\sigma^2}$

- Updates in soft  $K$ -means

$$r_{nk} = \frac{e^{-\beta \|x^{(n)} - \mu_k\|^2}}{\sum_{i=1}^K e^{-\beta \|x^{(n)} - \mu_i\|^2}}$$

$$\mu_k \leftarrow \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}}$$