



Expectation-Maximization Algorithm

Qinliang Su (苏勤亮)

Sun Yat-sen University

suqliang@mail.sysu.edu.cn

General Form of the 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. \mathbf{x} , 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

- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models
- EM Variants

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

Outline

- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models
- EM Variants

Re-representing the Log-likelihood

- The log-likelihood can be reformulated as

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

\forall distribution $q(\mathbf{z})$

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

$$= \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})}}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})}}_{KL(q || p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}))}$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + KL(q || p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})), \quad \text{for } \forall \boldsymbol{\theta}, q(\mathbf{z})$$

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

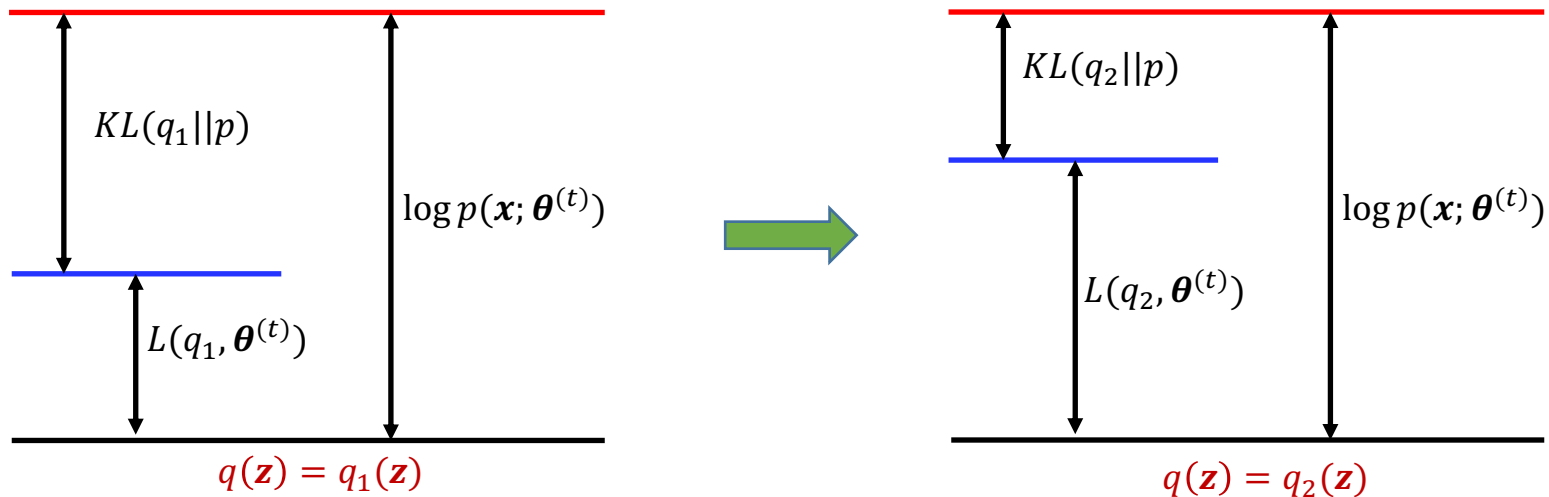
$$KL(q || p) \triangleq \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \geq 0$$

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

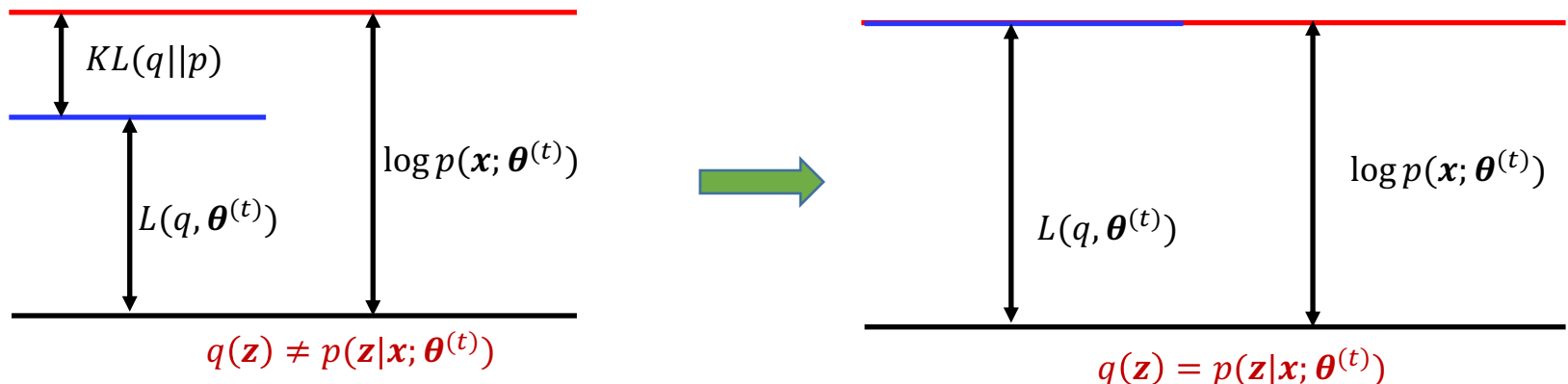
$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) = \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{q(\mathbf{z})} + \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$

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

$$KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})) = 0$$

Thus, we have

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



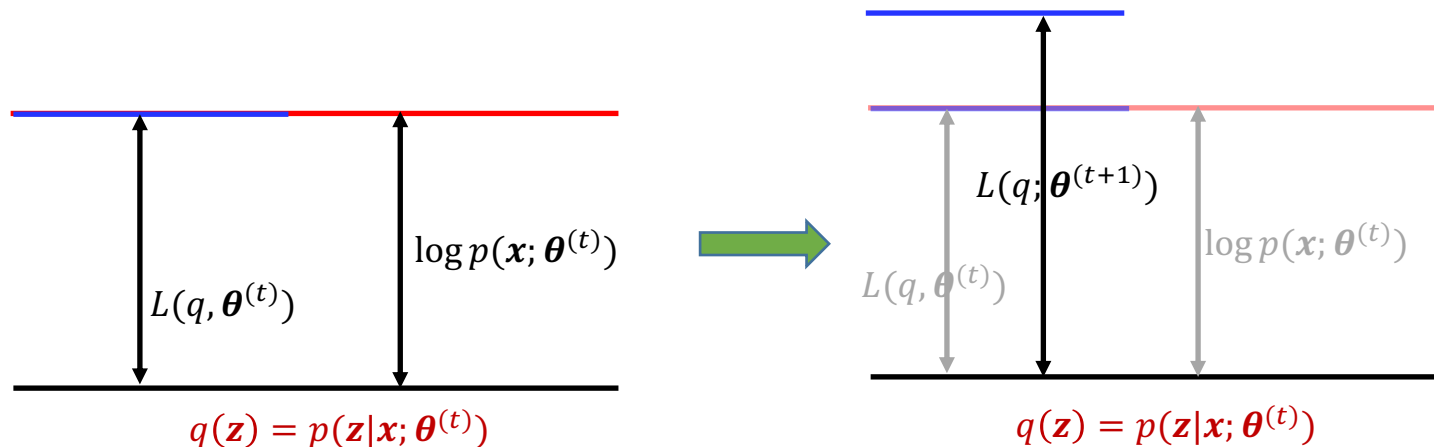
$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) &= \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)}) \\ &= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{p(\mathbf{z}|\mathbf{x};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

- 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

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

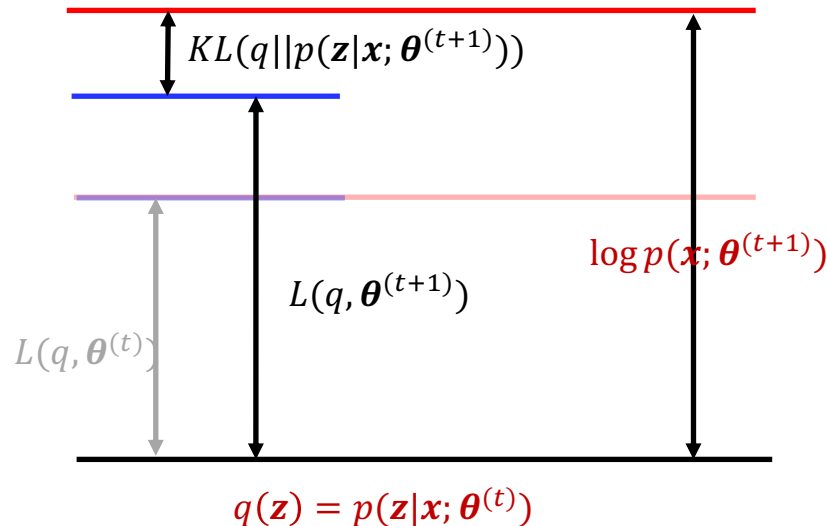


$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) = \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t+1)})}{q(\mathbf{z})} + \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t+1)})}$$

- By setting $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$, we obtain

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

The KL-divergence is always non-negative



- Thus, we can see that

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

$\max_{\boldsymbol{\theta}} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta})$ can guarantee the increase of likelihood at each step

- Equivalence between EM updating

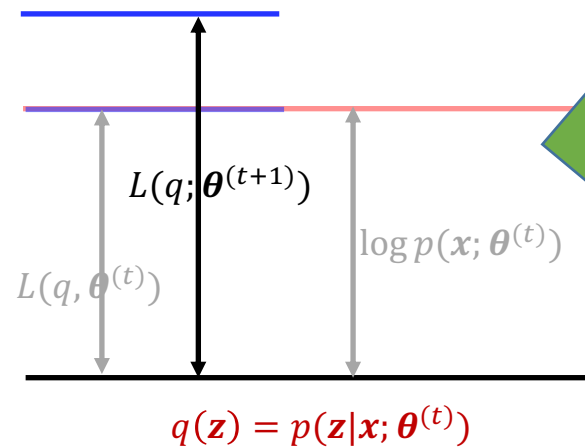
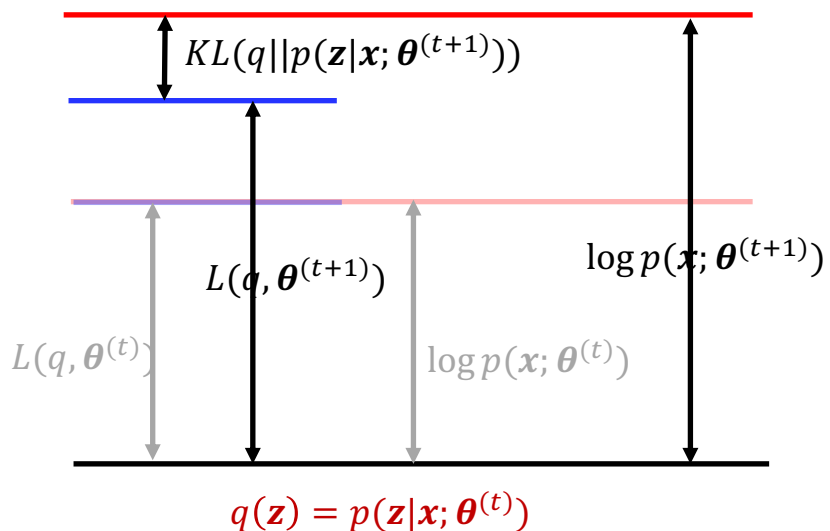
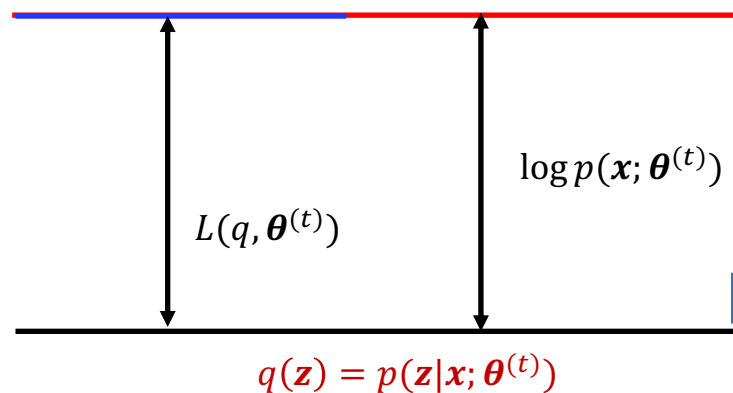
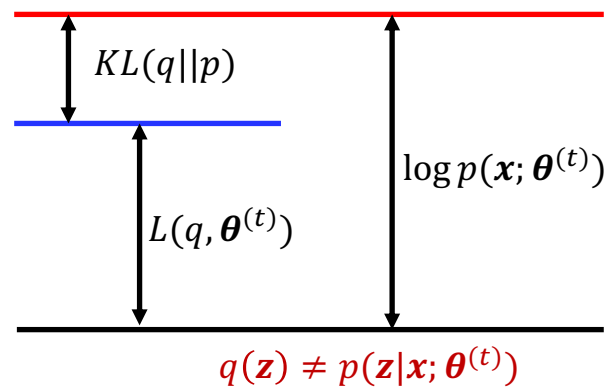
$$\arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) \quad \text{with} \quad Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) \triangleq \mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]$$

and the updating rule $\arg \max_{\boldsymbol{\theta}} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta})$

$$\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}) = \underbrace{\sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}_{\mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]} - \underbrace{\sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}_{\text{constant}}$$

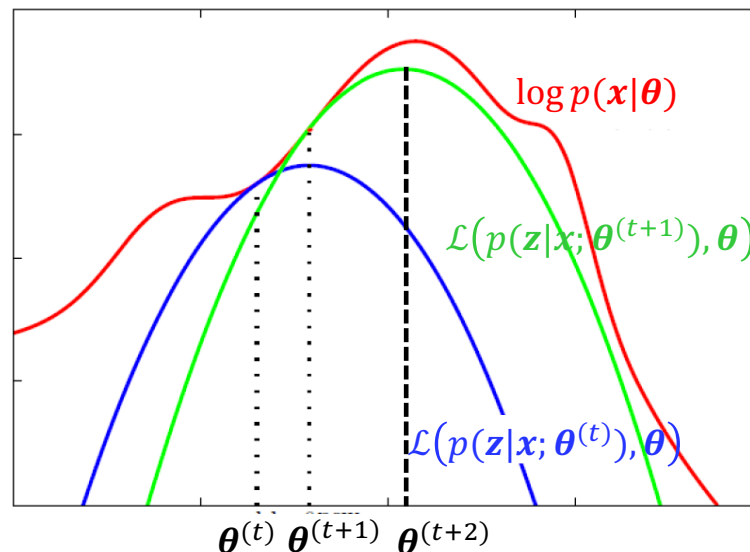
$$\text{Therefore,} \quad \arg \max_{\boldsymbol{\theta}} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}) \Leftrightarrow \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$$

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



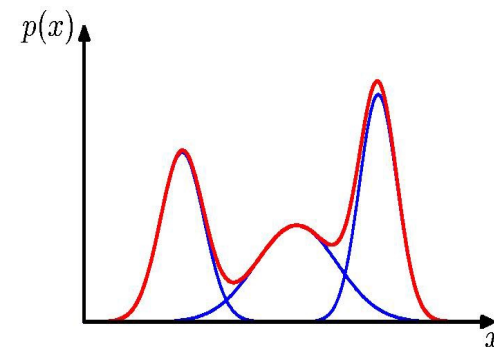
Outline

- EM Algorithm
- Theoretical Guarantees
- **Example: Training Gaussian Mixture Models**
- EM Variants

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 Two Steps

- It is a latent-variable model, thus we can use **EM** to optimize it

Remark: maximizing $\max_{\theta} \mathcal{L}(p(\mathbf{z}|\mathbf{x}; \theta^{(t)}), \theta)$ is equivalent to $\max_{\theta} Q(\theta; \theta^{(t)})$

- Reminder:* Key integrant in EM

➤ **E-step:** Expectation *w.r.t.* the posteriori $p(\mathbf{z}|\mathbf{x}; \theta^{(t)})$

$$Q(\theta; \theta^{(t)}) = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}^{(n)}|\mathbf{x}^{(n)}; \theta^{(t)})} [\log p(\mathbf{x}^{(n)}, \mathbf{z}^{(n)}; \theta)]$$

➤ **M-step:** Maximization

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

EM: E-step

- The posteriori distribution

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

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

- $\mathbf{1}_k$ denotes the one-hot vector with the k -th element being 1
- The log of the joint distribution $\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$

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

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

- The expectation

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

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

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

- Therefore, we have

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

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

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

- 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, the $Q(\cdot)$ becomes

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

EM: M-step

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

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

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

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

- For π_k , we need to consider the optimization under constraint $\sum_{k=1}^K \pi_k = 1$, leading to the solution

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

where $N_k = \sum_{n=1}^N \gamma_{nk}^{(t)}$ 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 γ_{nk} as

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

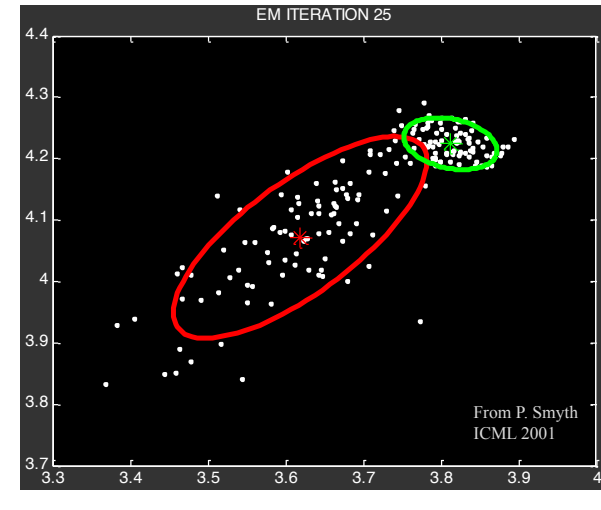
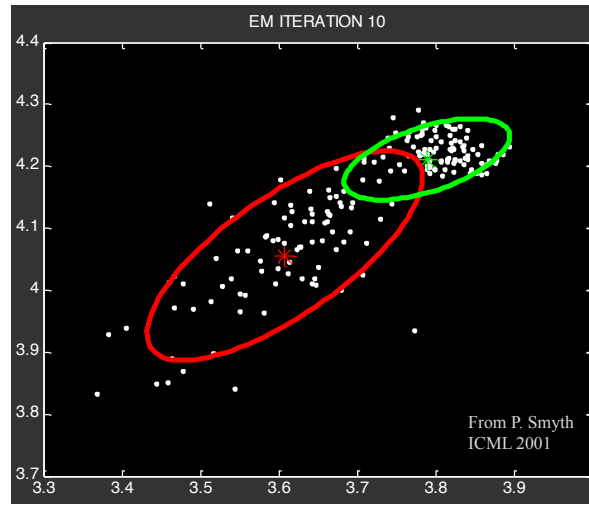
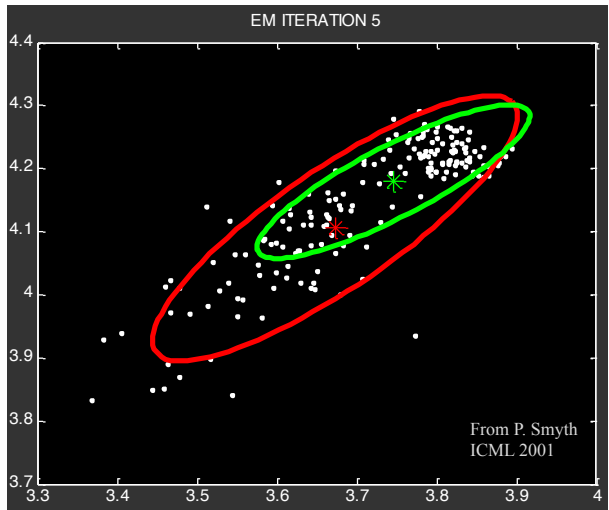
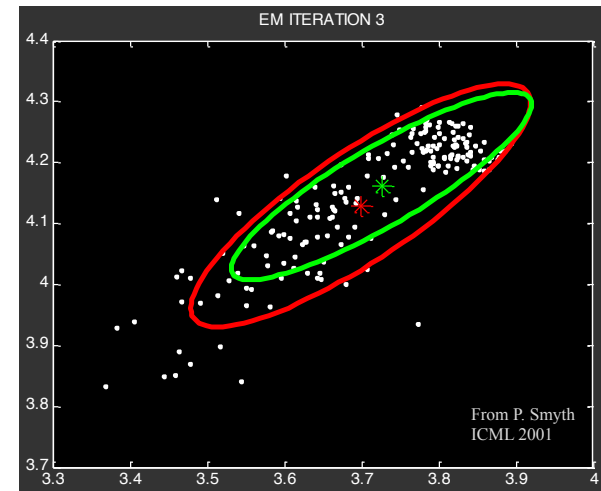
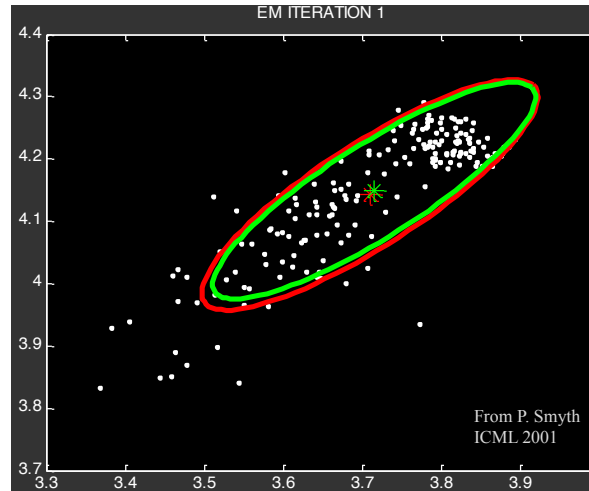
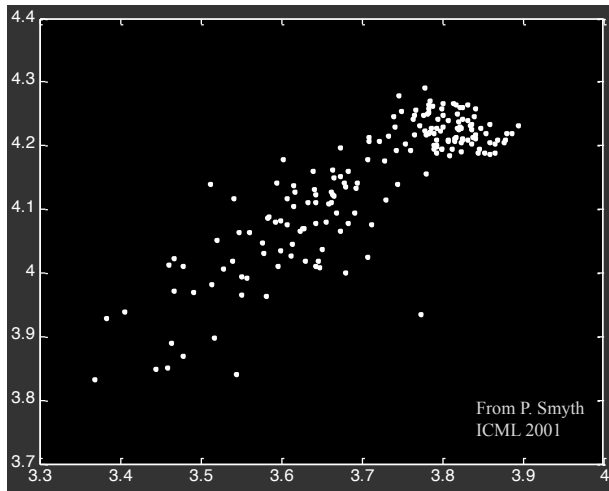
- Given the γ_{nk} , update $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ and π_k as

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

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

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

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



Relation to Soft K -Means

- When restricting Σ_k to the form $\Sigma_k = \sigma^2 \mathbf{I}$, the EM updating rules for GMM are

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

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

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

- Updates in soft K -means

Setting π_k and β as
 $\pi_k = \frac{1}{K}, \beta = \frac{1}{2\sigma^2}$

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

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

Outline

- EM Algorithm
- Theoretical Guarantees
- Example: Training Gaussian Mixture Models
- **EM Variants**

Review of the EM Algorithms

- To use EM algorithms, the key steps below are required

1) Computing the posteriori distribution

$$p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$$

2) Evaluating the expectation of $\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$ w.r.t. the posteriori $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$, i.e.,

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

3) Maximizing

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

However, not all of them are always achievable

Two Issues in EM Algorithm

- Issue one

The maximization is not achievable

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

- Issue two

- 1) The posteriori $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ cannot be derived analytically
- 2) Even if $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ can be obtained, we still cannot derive the close-form expression for the expectation

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

Generalized EM

- It is quite often in training LVMs that the optimization $\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$ cannot be solved

How to address this issue?

- Maximizing $\mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$ is not necessary. Increasing $\mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$ is sufficient to guarantee the EM algorithm to work
- That is, if we adopt SGD to update the parameter as

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \gamma \cdot \left. \frac{\partial \mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}$$

we can also guarantee the monotonic increase of log-likelihood

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

- Sketch of proof

➤ First, after the SGD update, it can be easily seen that

$$Q(\boldsymbol{\theta}^{(t+1)}; \boldsymbol{\theta}^{(t)}) \geq Q(\boldsymbol{\theta}^{(t)}; \boldsymbol{\theta}^{(t)})$$

➤ From $\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}) = \int p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} d\mathbf{z} = Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) - \int p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) d\mathbf{z}$, we further have

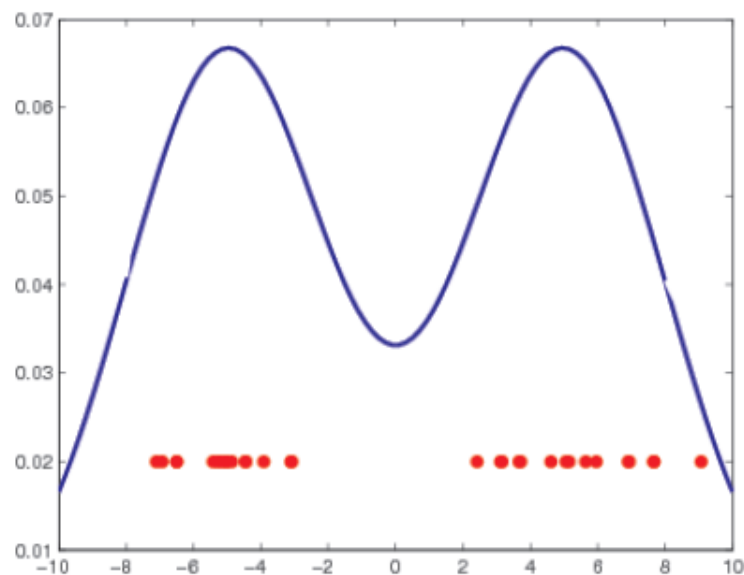
$$\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}); \boldsymbol{\theta}^{(t+1)}) \geq \underbrace{\mathcal{L}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}); \boldsymbol{\theta}^{(t)})}_{=\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$

➤ Due to

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

MCMC EM

- For any probability distributions, we can always draw samples from it, e.g., using **Markov chain Monte Carlo (MCMC)** methods



- Although the exact expression of the posteriori $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ is not known, we can use samples drawn from it to approximate it

- Thus, we can draw lots of samples \mathbf{z}_s for $s = 1, \dots, S$ from the posteriori distribution $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ such that

$$\mathbf{z}_s \sim p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$$

- Then, the expectation $\mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]$ can be approximated as

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) \approx \frac{1}{S} \cdot \sum_{s=1}^S \log p(\mathbf{x}, \mathbf{z}_s; \boldsymbol{\theta})$$

- We can optimize the approximate $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$ with SGD algorithm

The two sub-problems in the Issue Two are both solved. Thus, latent-variable models can always be trained with MCMC EM

- Drawing samples from a distribution is *computationally expensive*
- An alternative approach is to use a simple distribution $q(\mathbf{z}; \boldsymbol{\phi})$ to approximate the exact posterior distribution $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$

How to get the approximate simple distribution $q(\mathbf{z}; \boldsymbol{\phi})$?

- Idea
 - 1) Assuming a simple form for $q(\mathbf{z}; \boldsymbol{\phi})$, e.g.,

$$q(\mathbf{z}; \boldsymbol{\phi}) = \prod \mathcal{N}(z_i; \mu_i, \sigma_i^2)$$

- 2) Finding the best $\boldsymbol{\phi}$ that minimizes the KL-divergence

$$KL(q(\mathbf{z}; \boldsymbol{\phi}) \| p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}))$$

- Steps to update the model parameter θ

1) Finding the best approximate $q(\mathbf{z}; \phi)$ such that

$$\phi^{(t)} = \arg \min_{\phi} KL(q(\mathbf{z}; \phi) \| p(\mathbf{z} | \mathbf{x}; \theta^{(t)}))$$

2) Using $q(\mathbf{z}; \phi^{(t)})$ to compute expectation $\mathbb{E}_{p(\mathbf{z} | \mathbf{x}; \theta^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \theta)]$ approximately as

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

3) Obtaining the new value $\theta^{(t+1)}$ as

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

- The two optimization problems can be equivalently written as

$$\min_{\phi} KL(q(\mathbf{z}; \phi) \| p(\mathbf{z}|\mathbf{x}; \theta^{(t)})) \Leftrightarrow \max_{\phi} \int q(\mathbf{z}; \phi) \log \frac{p(\mathbf{z}|\mathbf{x}; \theta^{(t)})}{q(\mathbf{z}; \phi)} d\mathbf{z}$$

$$\Leftrightarrow \max_{\phi} \int q(\mathbf{z}; \phi) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta^{(t)})}{q(\mathbf{z}; \phi)} d\mathbf{z}$$

$$\max_{\theta} \mathbb{E}_{q(\mathbf{z}; \phi^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \theta)] \Leftrightarrow \max_{\theta} \int q(\mathbf{z}; \phi^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q(\mathbf{z}; \phi^{(t)})} d\mathbf{z}$$

- The algorithm to optimize θ and ϕ can be understood as solving the following optimization problem **in an alternative way**

$$\max_{\phi, \theta} \mathcal{L}(\mathbf{x}; \theta, \phi)$$

with

$$\mathcal{L}(\mathbf{x}; \theta, \phi) \triangleq \int q(\mathbf{z}; \phi) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q(\mathbf{z}; \phi)} d\mathbf{z}$$

- Instead of updating θ, ϕ alternatively, we can also update them *simultaneously* with the SGD algorithm, that is,

$$\theta^{(t+1)} = \theta^{(t)} + \gamma \cdot \left. \frac{\partial \mathcal{L}(x; \theta, \phi)}{\partial \theta} \right|_{\theta=\theta^{(t)}}$$

$$\phi^{(t+1)} = \phi^{(t)} + \gamma \cdot \left. \frac{\partial \mathcal{L}(x; \theta, \phi)}{\partial \phi} \right|_{\phi=\phi^{(t)}}$$

The method is dubbed *variational Bayesian EM* (VB-EM)

- In general, we optimize $\mathcal{L}(x; \theta, \phi)$ w.r.t. the two parameters θ, ϕ simultaneously

- Actually, it can be proved that $\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$ is a *lower bound* of the log-likelihood $\ln p(\mathbf{x}; \boldsymbol{\theta})$ for any $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$, that is,

$$\ln p(\mathbf{x}; \boldsymbol{\theta}) \geq \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$$

(Proof can be found in the next slide)

When the log-likelihood $\ln p(\mathbf{x}; \boldsymbol{\theta})$ cannot be directly maximized, we can seek to optimize its lower bound

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) \triangleq \int q(\mathbf{z}; \boldsymbol{\phi}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z}; \boldsymbol{\phi})} d\mathbf{z}$$

where $q(\mathbf{z}; \boldsymbol{\phi})$ can be set to be any distribution form, *e.g.*,

$$q(\mathbf{z}; \boldsymbol{\phi}) = \prod \mathcal{N}(z_i; \mu_i, \sigma_i^2)$$

Proof of $\ln p(\mathbf{x}; \boldsymbol{\theta}) \geq \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$

$$\begin{aligned}\ln p_{\boldsymbol{\theta}}(\mathbf{x}) &= \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln p_{\boldsymbol{\theta}}(\mathbf{x}) d\mathbf{z} = \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) q_{\boldsymbol{\phi}}(\mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z}) p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\&= \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z})} d\mathbf{z} + \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln \frac{q_{\boldsymbol{\phi}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\&= \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z})} d\mathbf{z} + KL(q_{\boldsymbol{\phi}}(\mathbf{z}) || p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})) \\&\geq \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}) \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z})} d\mathbf{z} \\&\triangleq \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})\end{aligned}$$