

---

# **PATTERN RECOGNITION AND MACHINE LEARNING**

## **CHAPTER 9: MIXTURE MODELS AND EM**

---

# Outlines

---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-

# K-means Clustering (1/3)

---

- Problem of identifying groups, or clusters, of data points in a multidimensional space
  - Partitioning the data set into some number  $K$  of clusters
  - Cluster: a group of data points whose inter-point distances are small compared with the distances to points outside of the cluster
  - Goal: an assignment of data points to clusters such that the sum of the squares of the distances to each data point to its closest vector (the center of the cluster) is a minimum

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

# K-means Clustering (2/3)

---

- Two-stage optimization
  - In the 1<sup>st</sup> stage: minimizing  $\mathcal{J}$  with respect to the  $r_{nk}$ , keeping the  $\mu_k$  fixed

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j ||\mathbf{x}_n - \mu_j||^2 \\ 0 & \text{otherwise} \end{cases}$$

- In the 2<sup>nd</sup> stage: minimizing  $\mathcal{J}$  with respect to the  $\mu_k$ , keeping  $r_{nk}$  fixed

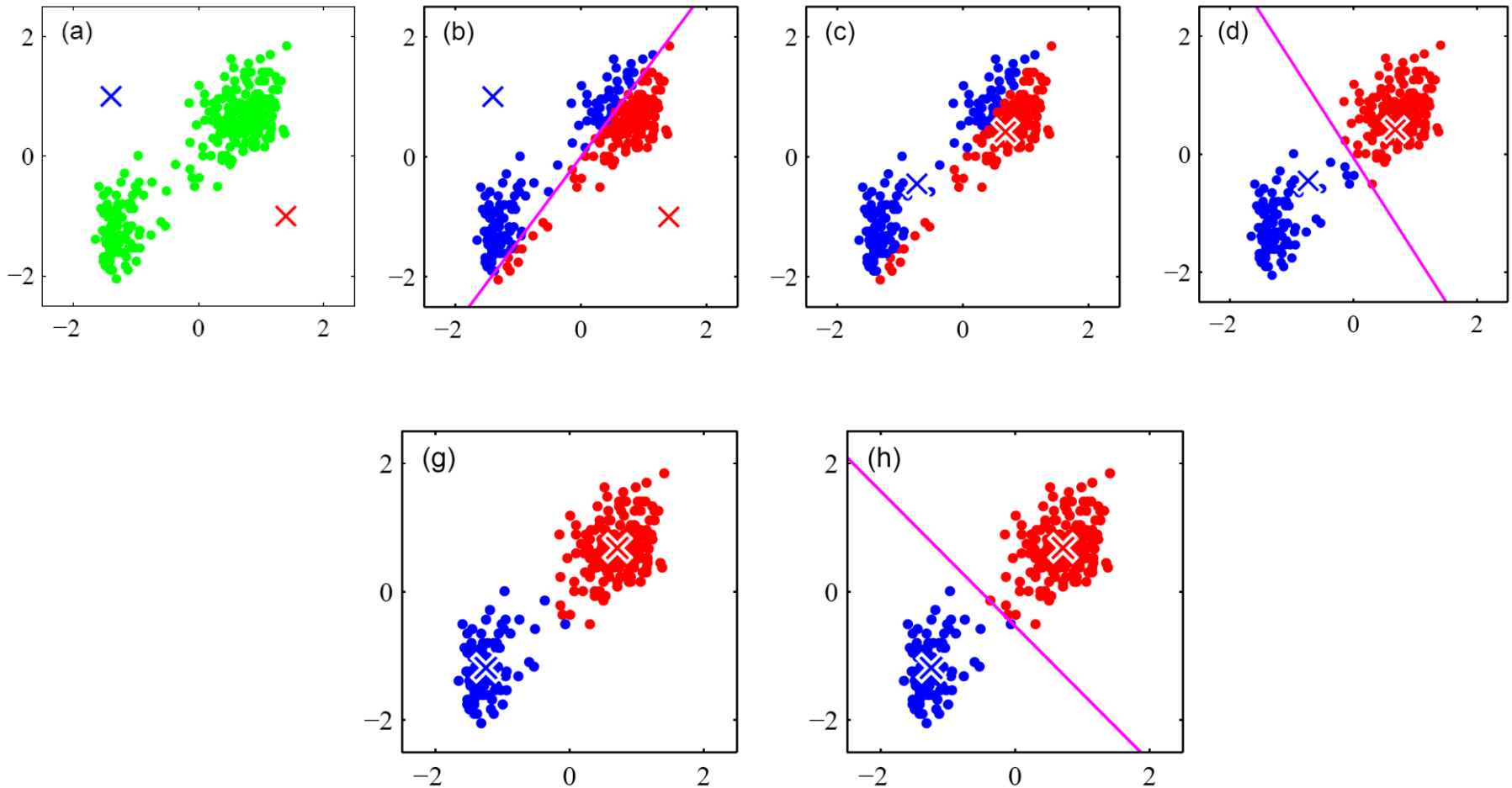
$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$

The mean of all of the data points assigned to cluster k

---

# K-means Clustering (3/3)

---



# Outlines

---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-

# Gaussian Mixture Model (1/4)

---

- Gaussian mixture distribution can be written as a linear superposition of Gaussian

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

- **A** binary random variable  $\mathbf{z}$  having a 1-of-K representation

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

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

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K N(\mathbf{x} | \mu_k, \Sigma_k)^{z_k}$$

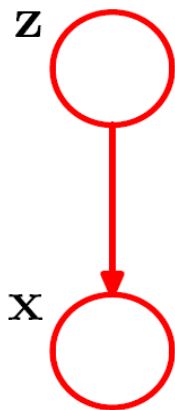
$$p(z_k = 1) = \pi_k \sum_{k=1}^K \pi_k = 1$$

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

# Gaussian Mixture Model (2/4)

---

- An equivalent formulation of the Gaussian mixture involving an explicit latent variable
  - Graphical representation of a mixture model
  - The marginal distribution of  $\mathbf{x}$  is a Gaussian mixture of the form (\*) ( $\rightarrow$  for every observed data point  $\mathbf{x}_n$ , there is a corresponding latent variable  $\mathbf{z}_n$ )



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



# Mixtures of Gaussians (3/4)

---

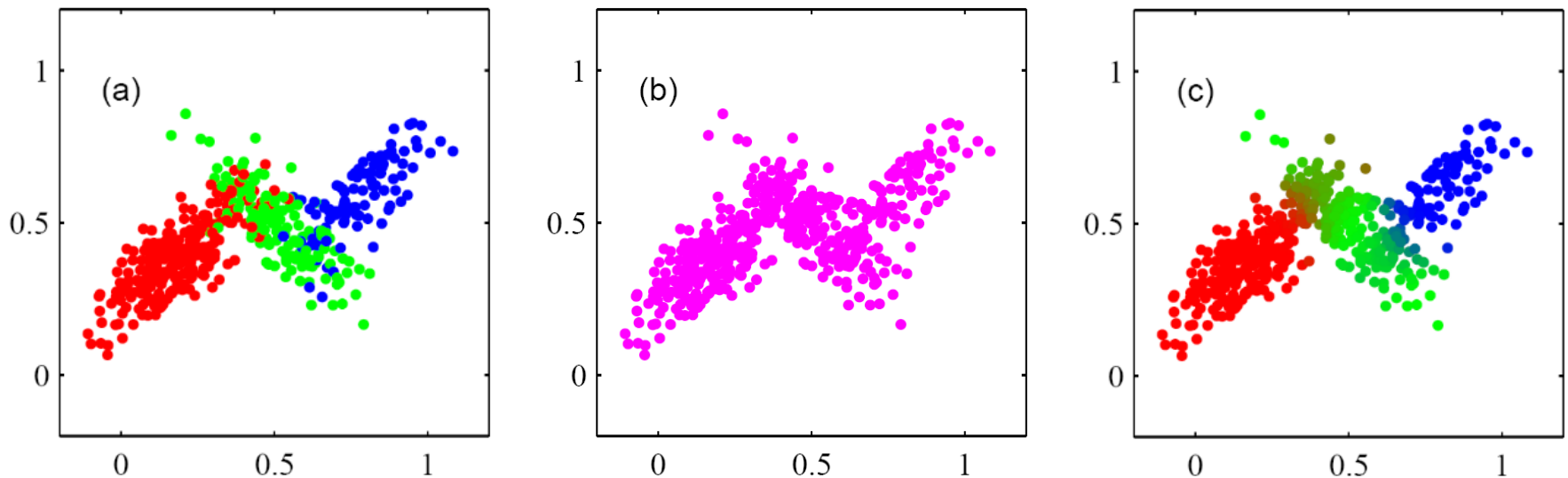
$$\begin{aligned}\gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} \\ &= \frac{\pi_k N(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x} | \mu_j, \Sigma_j)}\end{aligned}$$

- $\gamma(z_k)$  can also be viewed as the **responsibility** that component  $k$  takes for explaining the observation  $\mathbf{x}$

# Mixtures of Gaussians (4/4)

---

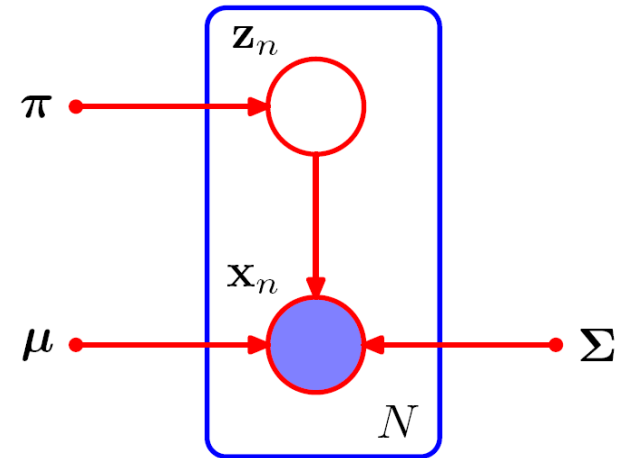
- Generating random samples distributed according to the Gaussian mixture model
  - Generating a value for  $\mathbf{z}$ , which denoted as  $\hat{\mathbf{z}}$  from the marginal distribution  $p(\mathbf{z})$  and then generate a value for  $\mathbf{x}$  from the conditional distribution  $p(\mathbf{x}|\hat{\mathbf{z}})$



# Maximum likelihood (1/3)

---

Graphical representation of a Gaussian mixture model for a set of  $N$  i.i.d. data points  $\{\mathbf{x}_n\}$ , with corresponding latent points  $\{z_n\}$



The log of the likelihood function

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

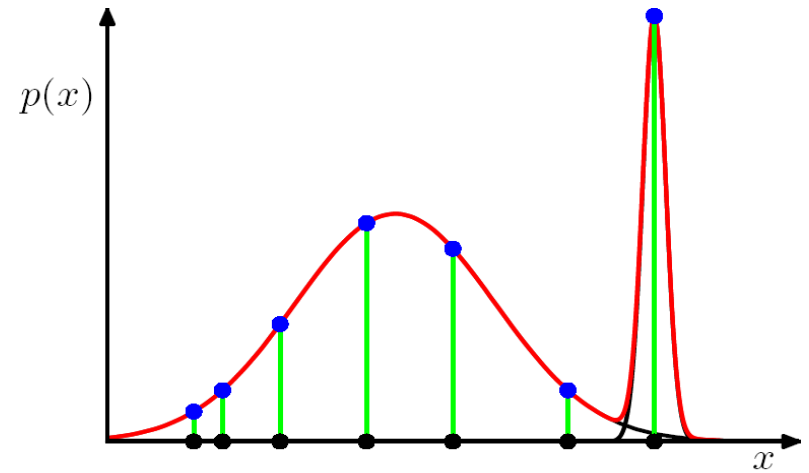
# Maximum likelihood (2/3)

---

- For simplicity, consider a Gaussian mixture whose components have covariance matrices given by  $\Sigma_k = \sigma_k^2 I$ 
  - Suppose that one of the components of the mixture model has its mean  $\mu_j$  exactly equal to one of the data points so that  $\mu_j = \mathbf{x}_n$
  - This data point will contribute a term in the likelihood function of the form

$$N(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 I) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}$$

- Once there are at least two components in the mixture, one of the components can have a finite variance and therefore assign finite probability to all of the data points while the other component can shrink onto one specific data point and thereby contribute an ever increasing additive value to the log likelihood  $\rightarrow$  over-fitting problem



# Maximum likelihood (3/3)

---

- Over-fitting problem
  - Example of the over-fitting in a maximum likelihood approach
  - This problem does not occur in the case of Bayesian approach
  - In applying maximum likelihood to a Gaussian mixture models, there should be steps to avoid finding such pathological solutions and instead seek local minima of the likelihood function that are well behaved
- *Identifiability* problem
  - A K-component mixture will have a total of  $K!$  equivalent solutions corresponding to the  $K!$  ways of assigning K sets of parameters to K components
- Difficulty of maximizing the log likelihood function  $\rightarrow$  the presence of the summation over k that appears inside the logarithm gives **no closed form solution** as in the single case

# Outlines

---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-

# EM for Gaussian Mixtures (1/4)

---

- I. Assign some initial values for the means, covariances, and mixing coefficients
- II. *Expectation* or E step
  - Using the current value for the parameters to evaluate the posterior probabilities or *responsibilities*
- III. *Maximization* or M step
  - Using the result of II to re-estimate the means, covariances, and mixing coefficients
- ❖ It is common to run the K-means algorithm in order to find a suitable initial values
  - The covariance matrices  $\rightarrow$  the sample covariances of the clusters found by the K-means algorithm
  - Mixing coefficients  $\rightarrow$  the fractions of data points assigned to the respective clusters

# EM for Gaussian mixtures (2/4)

---

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters

1. Initialize the means  $\mu_k$ , covariance  $\Sigma_k$  and mixing coefficients  $\pi_k$
2. E step

$$v(z_{nk}) = \frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

3. M step

$$\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N v(z_{nk}) \mathbf{x}_n$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N v(z_{nk}) (\mathbf{x}_n - \mu_k^{new})(\mathbf{x}_n - \mu_k^{new})^T$$

$$\pi_k^{new} = \frac{N_k}{N} \quad N_k = \sum_{n=1}^N v(z_{nk})$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \Pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$



# EM for Gaussian mixtures (3/4)

---

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\} \dots (*2)$$

Setting the derivatives of (\*2) with respect to the means of the Gaussian components to zero  $\rightarrow$

$$0 = - \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)} \Sigma_k (\mathbf{x}_n - \mu_k) \frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)} \leftarrow \text{Responsibility } \gamma(z_{nk})$$

$$\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \begin{array}{l} \text{A weighted mean of all of the points in the data} \\ \text{set} \end{array}$$

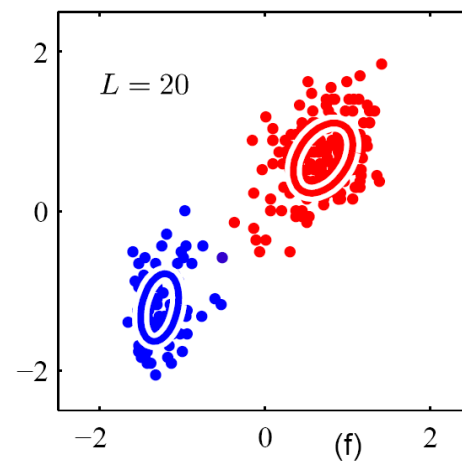
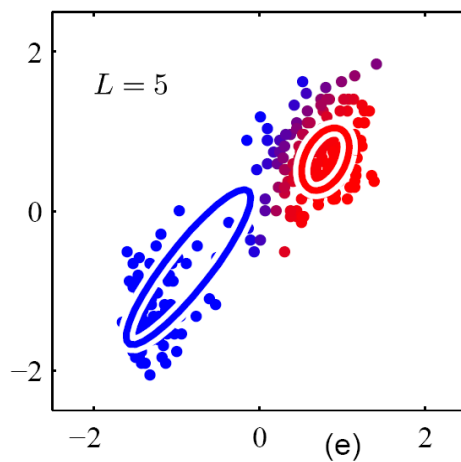
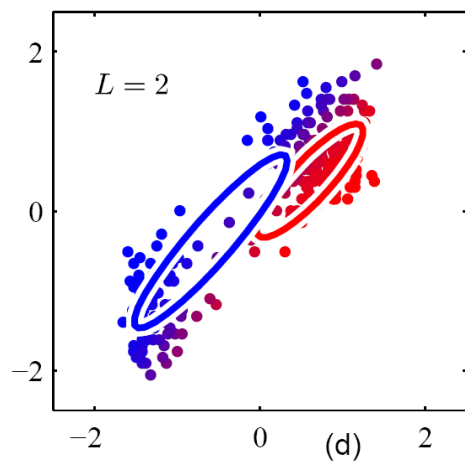
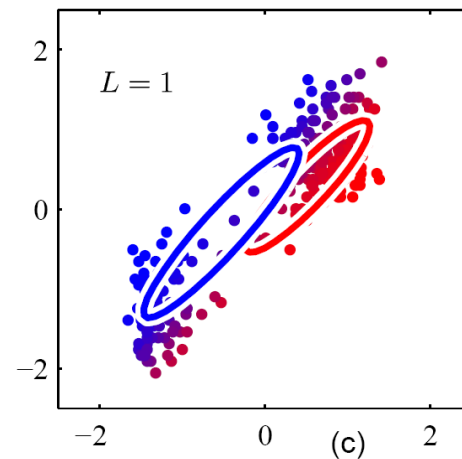
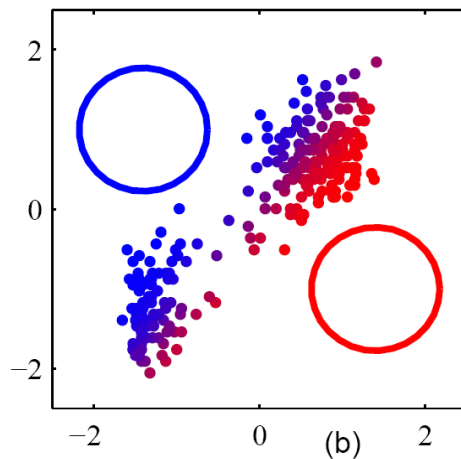
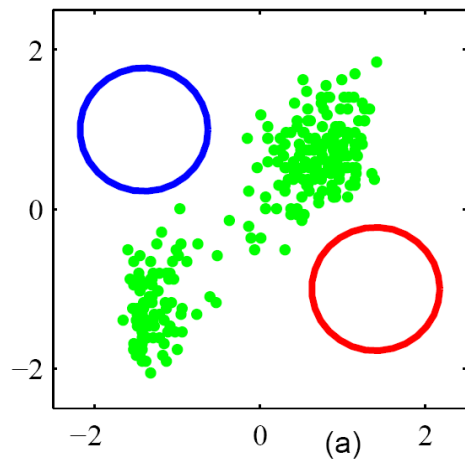
Setting the derivatives of (\*2) with respect to the covariance of the Gaussian components to zero  $\rightarrow$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{new})(\mathbf{x}_n - \mu_k^{new})^T$$

- Each data point weighted by the corresponding posterior probability
- The denominator given by the effective # of points associated with the corresponding component

# EM for Gaussian mixtures (4/4)

---



# Outlines

---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-

# An Alternative View of EM

---

- In maximizing the log likelihood function  $\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} P(\mathbf{X}, \mathbf{Z}|\theta) \right\}$  the **summation** prevents the logarithm from acting directly on the joint distribution
- Instead, the log likelihood function for the **complete** data set  $\{\mathbf{X}, \mathbf{Z}\}$  is straightforward.
- In practice since we are not given the complete data set, we consider instead its **expected** value  $Q$  under the **posterior** distribution  $p(\mathbf{Z}|\mathbf{X}, \theta)$  of the latent variable

- General EM

1. Choose an initial setting for the parameters  $\theta^{\text{old}}$
2. **E step** Evaluate  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$
3. **M step** Evaluate  $\theta^{\text{new}}$  given by
$$\theta^{\text{new}} = \operatorname{argmax}_{\theta} Q(\theta, \theta^{\text{old}})$$
$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$
4. If the covariance criterion is not satisfied, then let  $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$

# Gaussian Mixtures Revisited (1/2)

---

- Maximizing the likelihood for the **complete** data  $\{X, Z\}$

$$p(X, Z \mid \mu, \Sigma, \pi) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k N(X_n \mid \mu_k, \Sigma_k)]^{z_{nk}}$$

$$\ln p(X, Z \mid \mu, \Sigma, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln N(X_n \mid \mu_k, \Sigma_k) \}$$

$$\ln p(X \mid \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(X_n \mid \mu_k, \Sigma_k) \right\}$$

- The logarithm acts directly on the Gaussian distribution  $\rightarrow$  much simpler solution to the maximum likelihood problem
  - ♦ the maximization with respect to a mean or a covariance is exactly as for a single Gaussian (closed form)

# Gaussian Mixtures Revisited (2/2)

---

- Unknown latent variables  $\rightarrow$  considering **expectation** of the complete-data log likelihood with respect to the posterior distribution of the latent variables

- Posterior distribution
$$p(Z | X, \mu, \Sigma, \pi) \propto p(X | Z, \mu, \Sigma, \pi) p(Z) = \prod_{n=1}^N \prod_{k=1}^K [\pi_k N(x_n | \mu_k, \Sigma_k)]^{z_{nk}}$$

(ref. (9.10), (9.11))

- The expected value of the indicator variable under this posterior distribution

$$\begin{aligned} E(z_{nk}) &= p(z_{nk} = 1 | x_n) = \frac{p(z_{nk} = 1) p(x_n | z_{nk} = 1)}{p(x_n)} \\ &= \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)} = \gamma(z_{nk}) \end{aligned}$$

- The **expected value** of the complete-data log likelihood function

$$E[\ln p(X, Z | \mu, \Sigma, \pi)] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln N(x_n | \mu_k, \Sigma_k) \} \dots (*3)$$

# Relation to K-means

---

- ❖ K-means performs a **hard** assignment of data points to the clusters (each data point is associated **uniquely** with one cluster)
- ❖ EM makes a **soft** assignment based on the posterior probabilities
- ❖ K-means can be derived as a particular limit of EM for Gaussian mixtures:

As epsilon gets smaller, the terms for which  $\|x_n - \mu_j\|^2$  is farthest will go to zero most quickly. Hence the responsibility go all zero except for the term  $k$  for which the responsibility will go to unity

$$p(x_n | \mu_k, \epsilon I_k) = \frac{1}{\sqrt{2\pi\epsilon}} \exp\left\{-\frac{1}{2\epsilon} \|x_n - \mu_k\|^2\right\}, \quad \gamma(z_{nk}) = \frac{\pi_k \exp\{-\|x_n - \mu_k\|^2 / 2\epsilon\}}{\sum_j \pi_j \exp\{-\|x_n - \mu_j\|^2 / 2\epsilon\}}$$

$$E[\ln p(X, Z | \mu, \Sigma, \pi)] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2 + \text{const} \quad J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- ❖ Thus maximizing the expected complete data log-likelihood is equivalent to minimizing the distortion measure  $J$  for the K-means
- ❖ (In Elliptical K-means, the covariance is estimated also.)

# Outlines

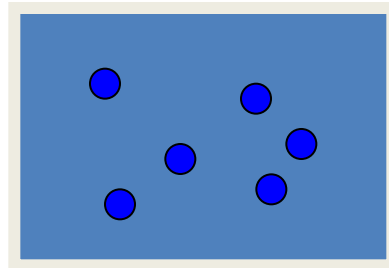
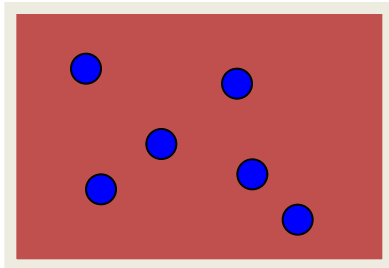
---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-



# Mixtures of Bernoulli Distributions

---



# Mixtures of Bernoulli Distributions (1/3)

---

$$* p(\mathbf{x} \mid \boldsymbol{\mu}) = \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{(1-x_i)} \quad (\text{single component})$$

$$(\mathbf{x} = (x_1, \dots, x_D), E(\mathbf{x}) = \boldsymbol{\mu} = (\mu_1, \dots, \mu_D), \text{Cov}(\mathbf{x}) = \boldsymbol{\Sigma} = \text{diag}\{\mu_i(1 - \mu_i)\})$$

$$* p(\mathbf{x} \mid \boldsymbol{\mu}, \pi) = \sum_{k=1}^K \pi_k p(\mathbf{x} \mid \boldsymbol{\mu}_k), \quad p(\mathbf{x} \mid \boldsymbol{\mu}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{(1-x_i)} \quad (\text{mixture})$$

$$(\boldsymbol{\mu} = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}, \quad \boldsymbol{\mu}_k = (\mu_{k1}, \dots, \mu_{kD}))$$

$$E(\mathbf{x}) = \sum_{k=1}^K \pi_k E(\mathbf{x} \mid \boldsymbol{\mu}_k) = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k,$$

$$\text{Cov}(\mathbf{x}) = E(\mathbf{x}\mathbf{x}^T) - E(\mathbf{x})E(\mathbf{x})^T = \sum_{k=1}^K \pi_k (E(\mathbf{x}\mathbf{x}^T \mid \boldsymbol{\mu}_k)) - E(\mathbf{x})E(\mathbf{x})^T = \sum_{k=1}^K \pi_k \{\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T\} - E(\mathbf{x})E(\mathbf{x})^T$$

$$(\text{Let } E(\mathbf{x}\mathbf{x}^T \mid \boldsymbol{\mu}_k) \equiv (c_{ij(k)}), \text{ then } c_{ii(k)} = E(x_i^2 \mid \boldsymbol{\mu}_k) = \mu_k, \quad c_{ij(k)} = E(x_i x_j \mid \boldsymbol{\mu}_k) = \mu_k^2 \quad (i \neq j))$$

and note that individual variables  $x_i$  are *independent, given  $\boldsymbol{\mu}$* )

# Mixtures of Bernoulli Distributions (2/3)

---

\* [E.g.]  $\mathbf{x} = (1, 0, 0, 1, 1)$

(single component) :

$$p(H) = \mu : p(\mathbf{x}) = \mu^3(1 - \mu)^2, E(\mathbf{x}) = (\mu, \dots, \mu), \text{Cov}(\mathbf{x}) = \Sigma = \mu(1 - \mu)I$$

(mixture):

$$p(H \mid c_1) = \mu_1, p(H \mid c_2) = \mu_2 : p(\mathbf{x}) = \pi_1 \mu_1^3(1 - \mu_1)^2 + \pi_2 \mu_2^3(1 - \mu_2)^2,$$

$$E(\mathbf{x}) = (\pi_1 \mu_1 + \pi_2 \mu_2, \dots, \pi_1 \mu_1 + \pi_2 \mu_2)$$

$$\begin{aligned} \text{Cov}(\mathbf{x}) &= \pi_1 \{ \mu_1(1 - \mu_1)I + \mu_1^2 \mathbf{1} \} + \pi_2 \{ \mu_2(1 - \mu_2)I + \mu_2^2 \mathbf{1} \} - (\pi_1 \mu_1 + \pi_2 \mu_2)^2 \mathbf{1} \\ &= \{ \pi_1 \mu_1(1 - \mu_1) + \pi_2 \mu_2(1 - \mu_2) \} I + \{ \pi_1 \mu_1^2 + \pi_2 \mu_2^2 - (\pi_1 \mu_1 + \pi_2 \mu_2)^2 \} \mathbf{1} \end{aligned}$$

\* Because the covariance matrix  $\text{Cov}(\mathbf{x})$  is no longer diagonal, the mixture distribution can capture correlations between the variables, unlike a single Bernoulli distribution.

# Mixtures of Bernoulli Distributions (3/3)

---

$$\ln p(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k p(\mathbf{x}_n \mid \boldsymbol{\mu}_k) \right\}$$

$$p(\mathbf{x} \mid \mathbf{z}, \boldsymbol{\mu}) = \prod_{k=1}^K p(\mathbf{x} \mid \boldsymbol{\mu}_k)^{z_k} \quad (\mathbf{z} = (z_1, \dots, z_K)^T \text{ is a binary indicator variables})$$

$$p(\mathbf{z} \mid \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}$$

(complete - data log likelihood function) :

$$\ln p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

$$E_z[\ln p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

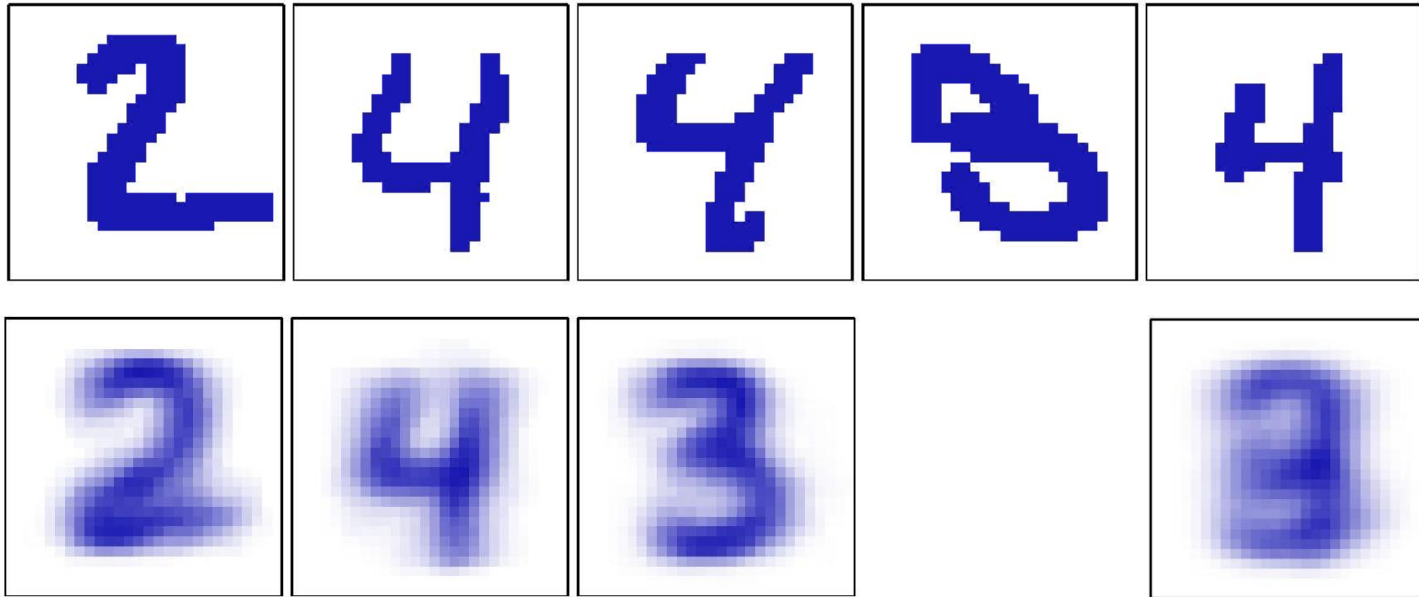
$$(\text{E - step}) \quad \gamma(z_{nk}) = E[z_{nk}] = \frac{\pi_k p(\mathbf{x}_n \mid \boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n \mid \boldsymbol{\mu}_j)}, \quad N_k = \sum_{n=1}^N \gamma(z_{nk}), \quad \bar{\mathbf{x}}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$(\text{M - step}) \quad \boldsymbol{\mu}_k = \bar{\mathbf{x}}_k, \quad \pi_k = \frac{N_k}{N}$$

\* In contrast to the mixture of Gaussians, there are no singularities in which the likelihood goes to infinity

# Mixtures of Bernoulli Distributions

---



- ❖  $N=600$  digit images, 3 mixtures
- ❖ A mixture of  $k=3$  Bernoulli distributions by 10 EM iterations
- ❖ Parameters for each of the three components/single multivariate Bernoulli
- ❖ The analysis of Bernoulli mixtures can be extended to the case of multinomial binary variables having  $M>2$  states (Ex.9.19)

# Outlines

---

- K-means Clustering
  - Gaussian Mixture Model
  - Expectation and Maximization
  - GMM Revisited
  - Bernoulli Mixture Model
  - EM Generalization
-

# The EM Algorithm in General (1/3)

- ❖ Direct optimization of  $p(X|\theta)$  is difficult while optimization of complete data likelihood  $p(X, Z|\theta)$  is significantly easier.
- ❖ Decomposition of the likelihood  $p(X|\theta)$

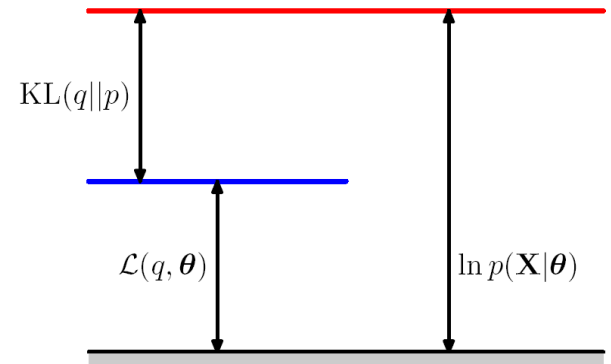
$$p(X|\theta) = \sum_Z p(X, Z|\theta)$$

$$\ln p(X|\theta) = \mathcal{L}(q, \theta) + KL(q||p)$$

$$(\text{where } \mathcal{L}(q, \theta) = \sum_Z q(Z) \ln \left\{ \frac{p(X, Z|\theta)}{q(Z)} \right\}, \quad KL(q||p) = - \sum_Z q(Z) \ln \left\{ \frac{p(Z|X, \theta)}{q(Z)} \right\})$$

(substitute  $\ln p(X, Z|\theta) = \ln p(Z|X, \theta) + \ln p(X|\theta)$  into the expression for  $\mathcal{L}(q, \theta)$ )

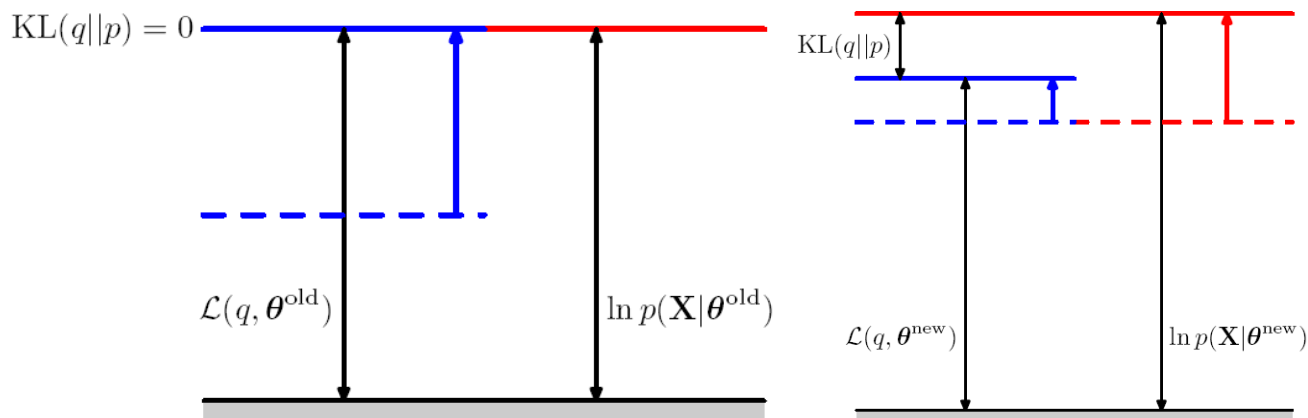
$\mathcal{L}(q, \theta) \leq \ln p(X|\theta) \rightarrow \mathcal{L}(q, \theta)$  is a lower bound on  $\ln p(X|\theta)$



# The EM Algorithm in General (2/3)

- ❖ (E step) The lower bound  $\mathcal{L}(q, \theta_{\text{old}})$  is maximized while holding  $\theta_{\text{old}}$  fixed. Since  $\ln p(\mathbf{X} | \theta)$  does not depend on  $q(\mathbf{Z})$ ,  $\mathcal{L}(q, \theta_{\text{old}})$  will be the largest when  $\text{KL}(q||p)$  vanishes (i.e. when  $q(\mathbf{Z})$  is equal to the posterior distribution  $p(\mathbf{Z} | \mathbf{X}, \theta_{\text{old}})$ )
- ❖ (M step)  $q(\mathbf{Z})$  is fixed and the lower bound  $\mathcal{L}(q, \theta_{\text{old}})$  is maximized wrt.  $\theta$  to  $\theta_{\text{new}}$ . When the lower bound is increased,  $\theta$  is updated making  $\text{KL}(q||p)$  greater than 0. Thus the increase in the log likelihood function is **greater** than the increase in the lower bound.
- ❖ In the M step, the quantity being maximized is the expectation of the **complete-data** log-likelihood

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta_{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) - \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta_{\text{old}}) \ln p(\mathbf{Z} | \mathbf{X}, \theta_{\text{old}}) \\ &= Q(\theta, \theta_{\text{old}}) + \text{const} \quad (\text{where } q(\mathbf{Z}) \text{ is set to be } p(\mathbf{Z} | \mathbf{X}, \theta_{\text{old}}))\end{aligned}$$

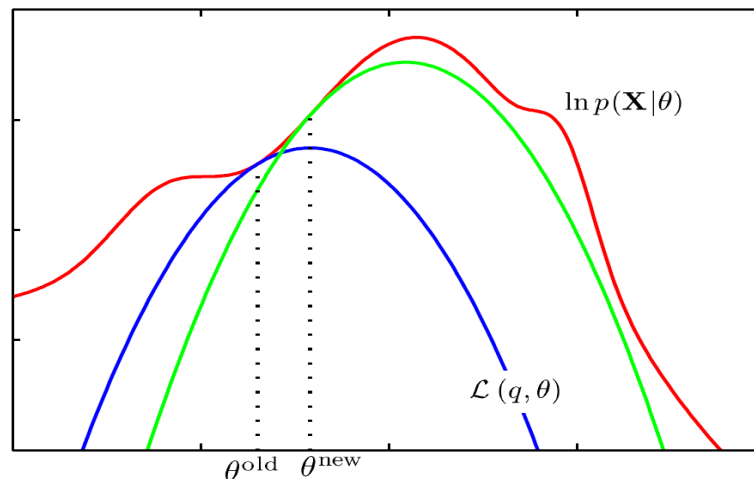




# The EM Algorithm in General (3/3)

---

- ❖ Start with initial parameter value  $\theta_{\text{old}}$
- ❖ In the first E step, evaluation of posterior distribution over latent variables gives rise to a lower bound  $\mathcal{L}(q, \theta_{\text{old}})$  whose value equals the log likelihood at  $\theta_{\text{old}}$  (blue curve)
- ❖ Note that the bound makes a **tangential contact** with the log-likelihood at  $\theta_{\text{old}}$ , so that both curves have the same gradient
- ❖ For mixture components from the exponential family, this bound is a **convex** function
- ❖ In the M step, the bound is maximized giving the value  $\theta_{\text{new}}$  which gives a larger value of log-likelihood than  $\theta_{\text{old}}$ .
- ❖ The subsequent E step constructs a bound tangential at  $\theta_{\text{new}}$  (green curve)



# The EM Algorithm in General (3a/3)

---

- ❖ EM can be also used to maximize the **posterior distribution**  $p(\theta|X)$  over parameters.
- ❖ Optimize the RHS alternatively wrt  $q$  and  $\theta$
- ❖ Optimization wrt  $q$  is the same E step
- ❖ M step required only a small modification through the introduction of the prior term  $\ln p(\theta)$

$$p(\theta|X) = p(\theta, X) / p(X), \quad \ln p(\theta|X) = \ln p(\theta, X) - \ln p(X)$$

$$\text{by decomposition (9.70), } \ln p(\theta|X) = \ln p(X|\theta) + \ln p(\theta) - \ln p(X)$$

$$= \mathcal{L}(q, \theta) + KL(q \| p) + \ln p(\theta) - \ln p(X) \geq \boxed{\mathcal{L}(q, \theta) + \ln p(\theta) + \text{const.}}$$

# The EM Algorithm in General (3b/3)

---

For complex problems, either E step or M step or both are intractable:

- ❖ Intractable M: Generalized EM (GEM), expectation conditional maximization (ECM)
- ❖ Intractable E: Partial E step
- ❖ **GEM**: instead of maximizing  $\mathcal{L}(q, \theta)$  wrt  $\theta$ , it seeks to change the parameters to increase its value.
- ❖ **ECM**: makes several constrained optimization within each M step. For instance, parameters are partitioned into groups and the M step is broken down into multiple steps each of which involves optimizing one of the subset with the remainder held fixed.
- ❖ **Partial (or incremental) EM**: (Note) For any given  $\theta$ , there is a unique maximum  $\mathcal{L}(q^*, \theta)$  wrt  $q$ . Since  $\mathcal{L}(q^*, \theta) = \ln p(X|\theta)$ , there is a  $\theta^*$  for the global maximum of  $\mathcal{L}(q, \theta)$  and  $\ln p(X|\theta^*)$  is a global maximum too. Any algorithm that converges to the global maximum of  $\mathcal{L}(q, \theta)$  will find a value of  $\theta$  that is also a global maximum of the log likelihood  $\ln p(X|\theta)$
- ❖ Each E or M step in partial E step algorithm is increasing the value of  $\mathcal{L}(q, \theta)$  and if the algorithm converges to a local (or global) maximum of  $\mathcal{L}(q, \theta)$ , this will correspond to a local (or global) maximum of the log likelihood function  $\ln p(X|\theta)$ .

# The EM Algorithm in General (3c/3)

---

- ❖ (Incremental EM) For a Gaussian mixture, suppose  $\mathbf{x}_m$  is updated with old and new values of responsibilities  $\gamma^{\text{old}}(z_{mk})$ ,  $\gamma^{\text{new}}(z_{mk})$  in the E-step .

In the M step, the means are updated as,

$$(9.17) \quad \mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$$

$$(9.18) \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\mu_k^{\text{new}} = \mu_k^{\text{old}} + \left( \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} \right) (x_m - \mu_k^{\text{old}})$$

$$N_k^{\text{new}} = N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})$$

$$N_k^{\text{new}} \mu_k^{\text{new}} = N_k^{\text{old}} \mu_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) x_m - \gamma^{\text{old}}(z_{mk}) x_m$$

$$N_k^{\text{new}} \mu_k^{\text{new}} = [N_k^{\text{new}} - \gamma^{\text{new}}(z_{mk}) + \gamma^{\text{old}}(z_{mk})] \mu_k^{\text{old}} + [\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})] x_m$$

$$\mu_k^{\text{new}} = \left( 1 - \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} \right) \mu_k^{\text{old}} + \left( \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} \right) x_m$$

- ❖ Both E step and M step take fixed time **independent** of the total number of data points. Because the parameters are revised after each data point, rather than waiting until after the whole data set is processed, this incremental version can converge faster than the batch version.

# HW5

---

K-means: 9.1

Mixture Model: 9.3 9.7 9.10 9.18 9.19

EM: 9.20 9.21 9.25 9.26

Due Dec. 6

---