# PATTERN RECOGNITION
### AND MACHINE LEARNING

## CHAPTER 9: MIXTURE MODELS AND EM

# Learning Objectives

1、What are the differences between supervised and unsupervised learning schemes?

2、What is K-means clustering?

3、What are Gaussian Mixture Models?

4、What are Bernoulli Mixture Models?

5、What is the EM learning scheme?

6、How to understand EM from the perspective of likelihood?

7、How to generalize the EM scheme via decomposition?

# Outlines

➤ Supervised vs Unsupervised Learning

➤ K-means Clustering

➤ Gaussian Mixture Model

➤ Expectation and Maximization

➤ GMM Revisited

➤ Bernoulli Mixture Model

➤ EM Generalization

# Supervised vs Unsupervised Learning

☐ Supervised learning

- ✓ Training data have labels (complete data)
- ✓ To learn the mapping between data and labels
- ✓ Regression, classification
- ✓ Detection, semantic/instance segmentation
- ✓ KNNs, SVMs, decision trees, neural networks
- ✓ Deep neural networks are good at supervised learning

# Supervised vs Unsupervised Learning

☐ Unsupervised learning

- ✓ Training data have no labels (incomplete data)

- ✓ To learn the intrinsic structures of data

- ✓ Clustering, data dimension reduction

- ✓ Segmentation, compression

- ✓ K-means, GMMs, PCA, ICA, NMF

- ✓ GAN is a kind of unsupervised learning

# Unsupervised Learning

# Outlines

# K-means Clustering (I)

◻ 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 (II)

☐ Two-stage optimization

✓ In the 1$^{st}$ stage: minimizing $\mathcal{J}$ with respect to the $r_{nk}$, keeping the $\mu_k$ fixed

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

✓ In the 2$^{nd}$ stage: minimizing $\mathcal{J}$ with respect to the $\mu_k$, keeping $r_{nk}$ fixed
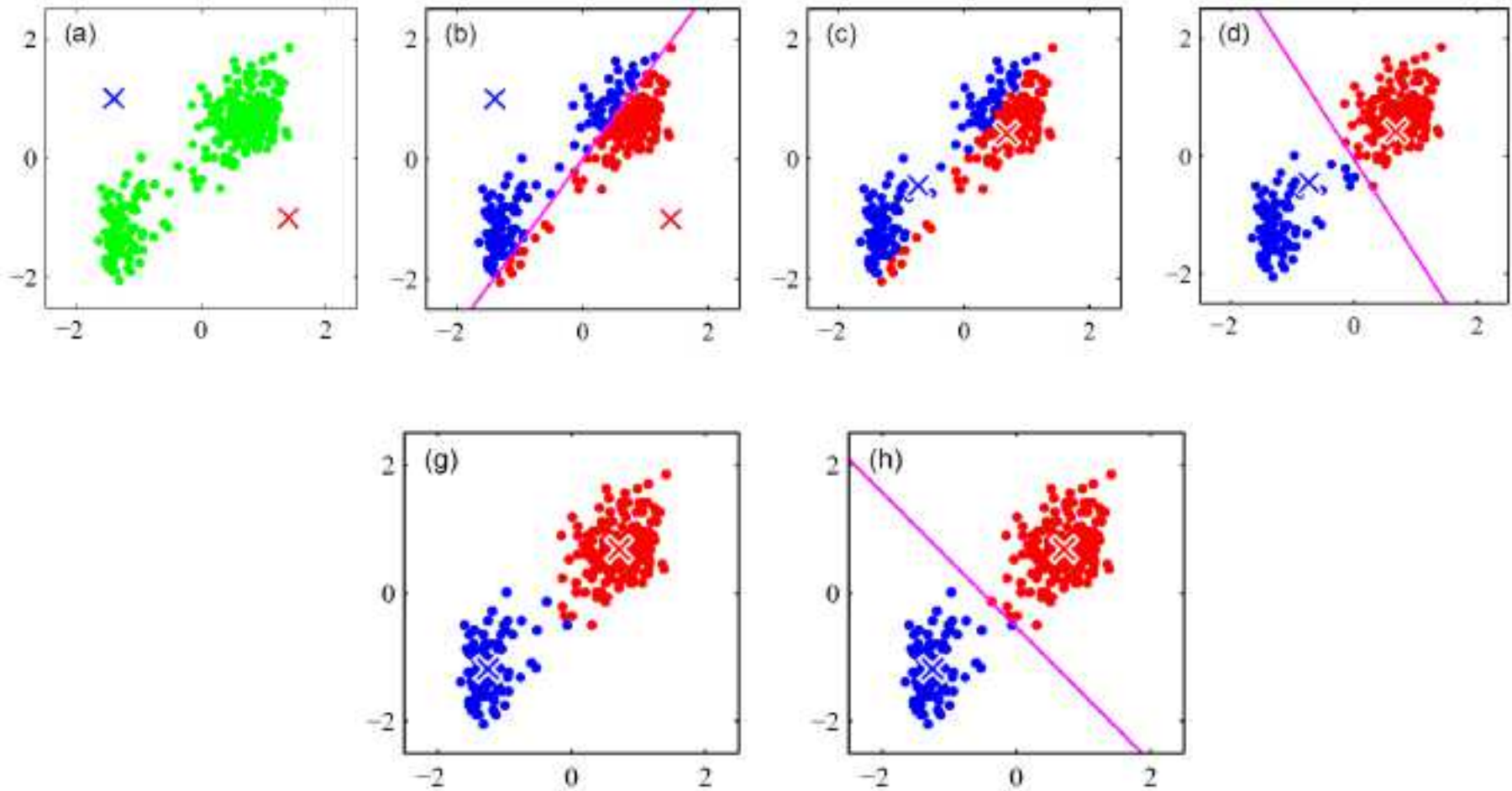
$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} \quad \Longleftarrow \quad \boxed{2 \sum_{n=1}^{N} r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0}$$

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

# K-means Clustering (III)

# Outlines

➢ Supervised vs Unsupervised Learning

➢ K-means Clustering

➢ Gaussian Mixture Model

➢ Expectation and Maximization

➢ GMM Revisited

➢ Bernoulli Mixture Model

➢ EM Generalization

# Gaussian Mixture Model (I)

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

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

☐ random variable **z** having a 1-of-K distribution

$$p(\mathbf{z}) = \prod_{k=1}^{K} \pi_k^{z_k} \qquad \sum_{k=1}^{K} \pi_k = 1 \qquad 0 \leqslant \pi_k \leqslant 1 \qquad p(z_k = 1) = \pi_k$$

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^{K} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \qquad p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

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

# Gaussian Mixture Model (II)

◻ An equivalent formulation of the Gaussian mixture involving an explicit latent variable

✓ Graphical representation of a mixture model

✓ The marginal distribution of **x** is a Gaussian mixture (for every observed data point **x**$_n$, there is a corresponding latent variable **z**$_n$, that is, the cluster label)

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

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$$
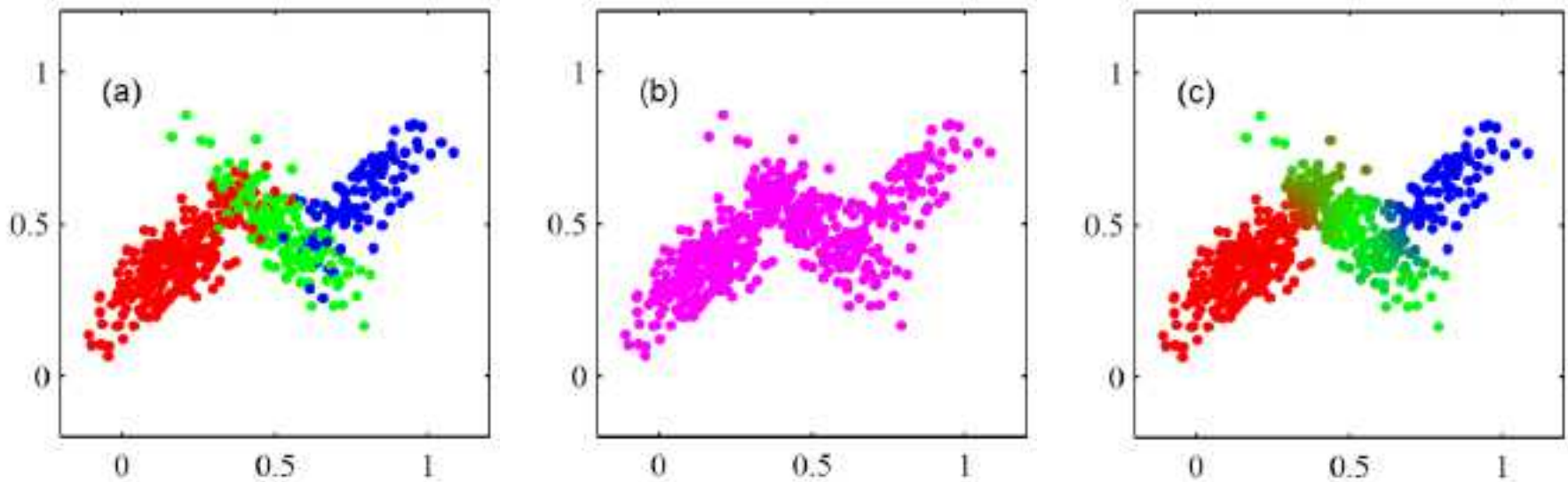
# Gaussian Mixture Model (III)

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1)p(\mathbf{x}|z_j = 1)}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

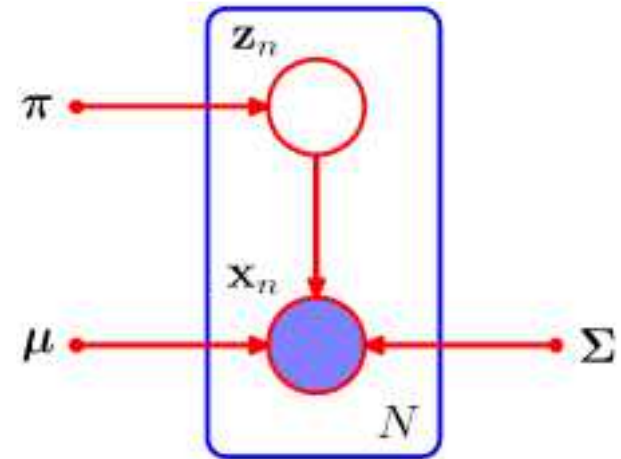☐ $\gamma(z_k)$ can also be viewed as the responsibility that component $k$ takes for explaining the observation **x**

# Gaussian Mixture Model (IV)

□ Generating random samples distributed according to the Gaussian mixture model

✓ Generating a value for **z**, which denoted as $\widehat{\mathbf{z}}$ from the marginal distribution p(**z**) and then generate a value for **x** from the conditional distribution $p(\mathbf{x}|\widehat{\mathbf{z}})$

# Maximum Likelihood (I)

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



□ The log of the likelihood function

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$
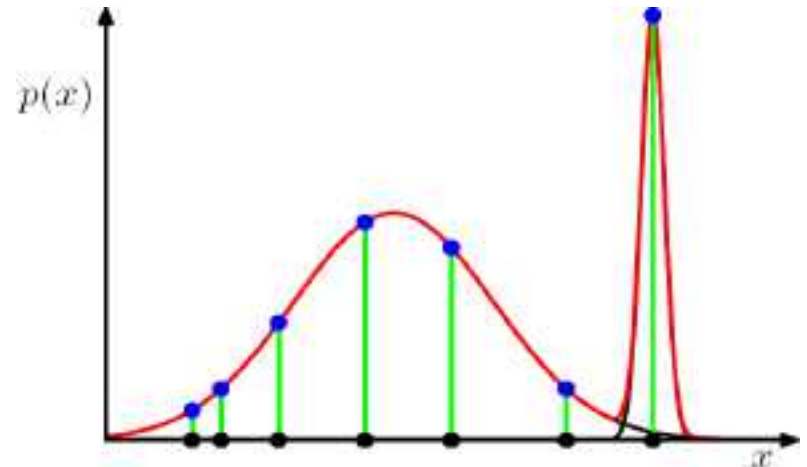
# Maximum Likelihood (II)

◻ For simplicity, consider a Gaussian mixture whose components have covariance matrices given by

$$\Sigma_k = \sigma_k^2 \mathbf{I}$$

✓ 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

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

✓ over-fitting problem

# Maximum Likelihood (III)

- ❑ **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 heuristics to 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 → the presence of the summation over k that appears inside the logarithm gives **no closed form solution** as in the single case

# Outlines

➢ Supervised vs Unsupervised Learning

➢ K-means Clustering

➢ Gaussian Mixture Model

➢ Expectation and Maximization

➢ GMM Revisited

➢ Bernoulli Mixture Model

➢ EM Generalization

# EM for Gaussian Mixtures (I)

① **Initialization**:

Initialize values for means, covariances, and mixing coefficients

② **Expectation or E step**

Using the current values for the parameters to evaluate the posterior probabilities or *responsibilities*

③ **Maximization or M step**

Using the results of ② 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 → the sample covariances of the clusters found by the K-means algorithm

✓ Mixing coefficients → the fractions of data points assigned to the respective clusters

# EM for Gaussian Mixtures (II)

☐   Goal: 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

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

3.  M step

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n$$

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

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

4.  Evaluate the log likelihood

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

# EM for Gaussian Mixtures (III)

☐ Setting the derivatives of likelihood with respect to the means of the Gaussian components to zero →

$$0 = -\sum_{n=1}^{N} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{\gamma(z_{nk})} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

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

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

☐ Setting the derivatives of likelihood with respect to the covariance of the Gaussian components to zero →

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

✓ 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 (IV)

❑ Setting the derivatives of likelihood with respect to mixing coefficients to zero, subject to their sum equal to 1 →
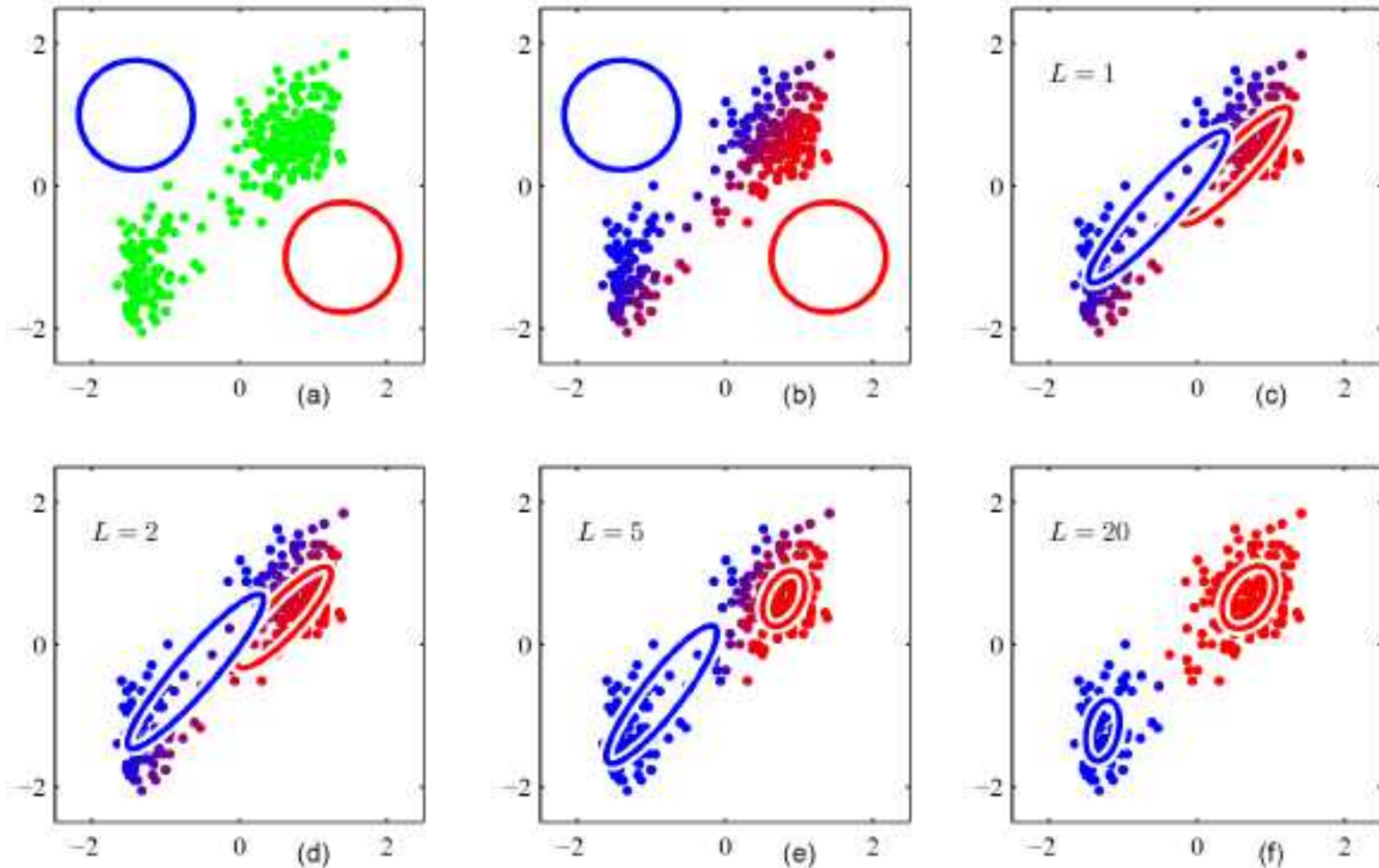
$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right) \qquad N_k = \sum_{n=1}^{N} \gamma(z_{nk})$$

$$0 = \sum_{n=1}^{N} \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda \qquad \Longrightarrow \qquad \boxed{\lambda = -N}$$

multiply $\pi_k$ and sum over k

$$\Longrightarrow \qquad \boxed{\pi_k = \frac{N_k}{N}}$$

# EM for Gaussian Mixtures (V)

# Outlines

- ➢ Supervised vs Unsupervised Learning

- ➢ 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 $$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln\left\{\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})\right\}$$ the summation prevents the logarithm from acting directly on the joint distribution

- Instead, the log likelihood function for the complete data set {X, 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( **Z**|**X**, Θ) of the latent variable

- **General EM**

  1. Choose an initial setting for the parameters $\Theta^{old}$

  2. **E step** Evaluate $p(\mathbf{Z}|\mathbf{X},\Theta^{old})$

  3. **M step** Evaluate $\Theta^{new}$ given by
  $\Theta^{new} = \text{argmax}_{\Theta} Q(\Theta, \Theta^{old})$
  $Q(\Theta, \Theta^{old}) = \Sigma_{Z}\, p(\mathbf{Z}|\mathbf{X}, \Theta^{old})\ln p(\mathbf{X}, \mathbf{Z}| \Theta)$

  4. It the covariance criterion is not satisfied, then let $\Theta^{old} \leftarrow \Theta^{new}$

# Expected Complete-Data Log Likelihood

**Expectation of complete-data log likelihood:**

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\mathrm{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\mathrm{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

$$\leq$$

**Log Likelihood:**

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \right\}$$

**Jason Inequality**: as ln (x) is a concave function,
$$E\{\ln f(x)\} \leq \ln( E\{f(x)\} )$$

# An Alternative View of EM for MAP

- ☐ In maximizing the log posterior, ln p(Θ |**X** ) ∝ ln p(**X**|Θ) + ln p(Θ), given the prior p(Θ) the summation prevents the logarithm from acting directly on the joint distribution

- ☐ Instead, the log likelihood function for the complete data set {X, 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( **Z**|**X**, Θ) of the latent variable

- ☐ **General EM for MAP**

  1. Choose an initial setting for the parameters $\Theta^{old}$

  2. **E step** Evaluate $p(\mathbf{Z}|\mathbf{X},\Theta^{old})$

  3. **M step** Evaluate $\Theta^{new}$ given by
     $$\Theta^{new} = \text{argmax}_\Theta \, Q(\Theta, \Theta^{old}) + \ln p(\Theta)$$
     $$Q(\Theta, \Theta^{old}) = \Sigma_Z \, p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}| \Theta)$$

  4. It the covariance criterion is not satisfied, then let $\Theta^{old} \leftarrow \Theta^{new}$

# Gaussian Mixtures Revisited (I)

☐ Maximizing the likelihood for the complete data {X, Z}

$$p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}$$

☐ The logarithm acts directly on the Gaussian distribution → 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 (II)

☐ Unknown latent variables → considering expectation of the complete-data log likelihood with respect to the posterior distribution of the latent variables

✓ Posterior distribution

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \propto \prod_{n=1}^{N} \prod_{k=1}^{K} [\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}$$

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

$$\mathbb{E}[z_{nk}] = \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)]^{z_{nj}}}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \gamma(z_{nk})$$

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

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}$$

# 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:

$$p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp\left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}$$

$$\gamma(z_{nk}) = \frac{\pi_k \exp\left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 / 2\epsilon \right\}}{\sum_j \pi_j \exp\left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon \right\}}$$
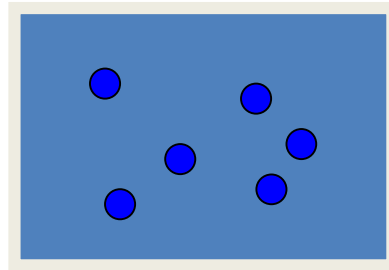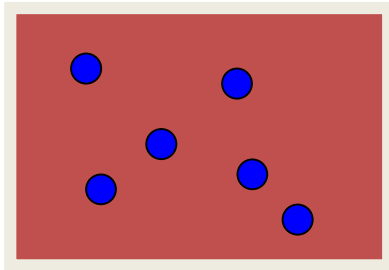
$$\epsilon \to 0 \qquad \gamma(z_{nk}) \to r_{nk} \qquad r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] \to -\frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + \text{const.}$$

# Outlines

➢ Supervised vs Unsupervised Learning

➢ K-means Clustering

➢ Gaussian Mixture Model

➢ Expectation and Maximization

➢ GMM Revisited

➢ Bernoulli Mixture Model

➢ EM Generalization

# Mixtures of Bernoulli Distributions (I)

# Mixtures of Bernoulli Distributions (II)

$$p(\mathbf{x}|\mathbf{z},\boldsymbol{\mu}) = \prod_{k=1}^{K} p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \qquad p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^{K} \pi_k^{z_k}$$

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

# Mixtures of Bernoulli Distributions (III)

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \left\{ \ln \pi_k \right.$$

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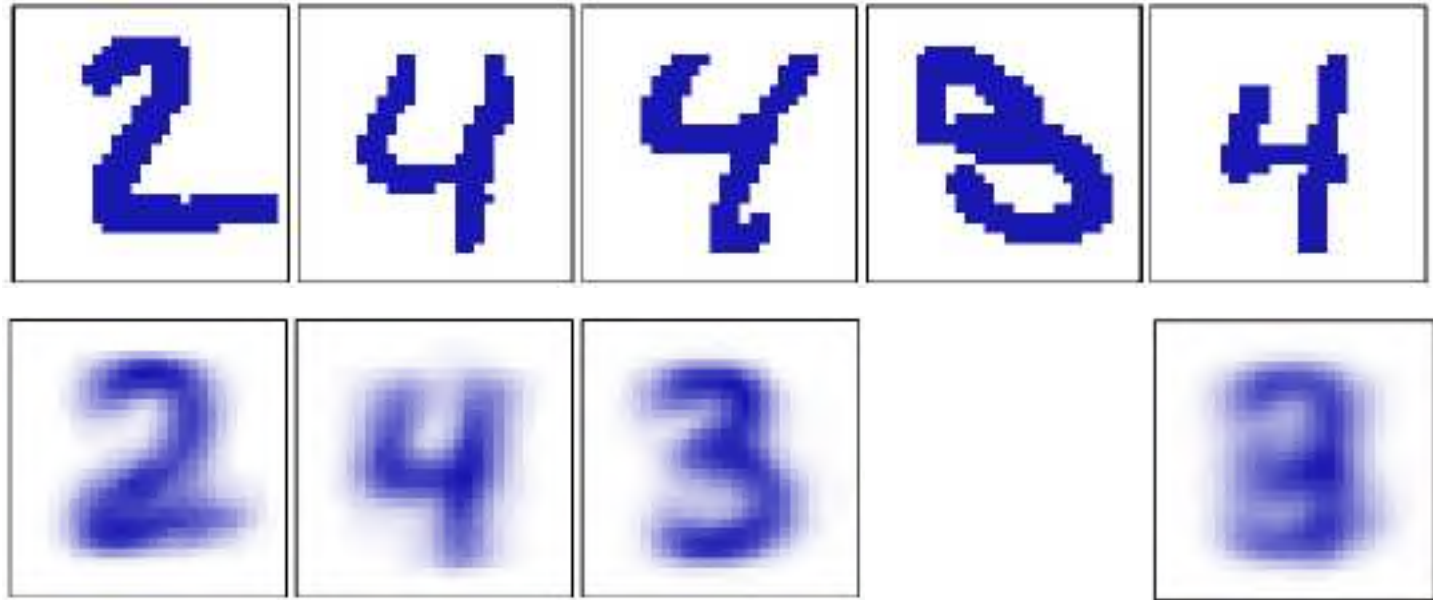
$$\Longrightarrow \quad \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{ \ln \pi_k \right.$$

$$\left. + \sum_{i=1}^{D} [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

# EM for Bernoulli Mixture Models

**E-Step:**

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \frac{\displaystyle\sum_{z_{nk}} z_{nk} \left[\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)\right]^{z_{nk}}}{\displaystyle\sum_{z_{nj}} \left[\pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)\right]^{z_{nj}}}$$

$$= \frac{\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)}{\displaystyle\sum_{j=1}^{K} \pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)}.$$

**M-Step:**

$$N_k = \sum_{n=1}^{N} \gamma(z_{nk}) \qquad \pi_k = \frac{N_k}{N}$$

$$\overline{\mathbf{x}}_k = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n \qquad \boldsymbol{\mu}_k = \overline{\mathbf{x}}_k$$

# 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

# Outlines

➤ Supervised vs Unsupervised Learning

➤ K-means Clustering

➤ Gaussian Mixture Model

➤ Expectation and Maximization

➤ GMM Revisited

➤ Bernoulli Mixture Model

➤ EM Generalization

# The EM Algorithm in General (I)

☐ Direct optimization of p (X|θ) is difficult while optimization of complete data likelihood p (X, Z|θ) is significantly easier.

☐ Decomposition of the likelihood p (X|θ)

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta})$$

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

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\mathrm{KL}(q\|p) = -\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\mathrm{KL}(q\|p) \geqslant 0 \implies \mathcal{L}(q, \bar{\boldsymbol{\theta}}) \leqslant \ln p(\mathbf{X}|\boldsymbol{\theta})$$
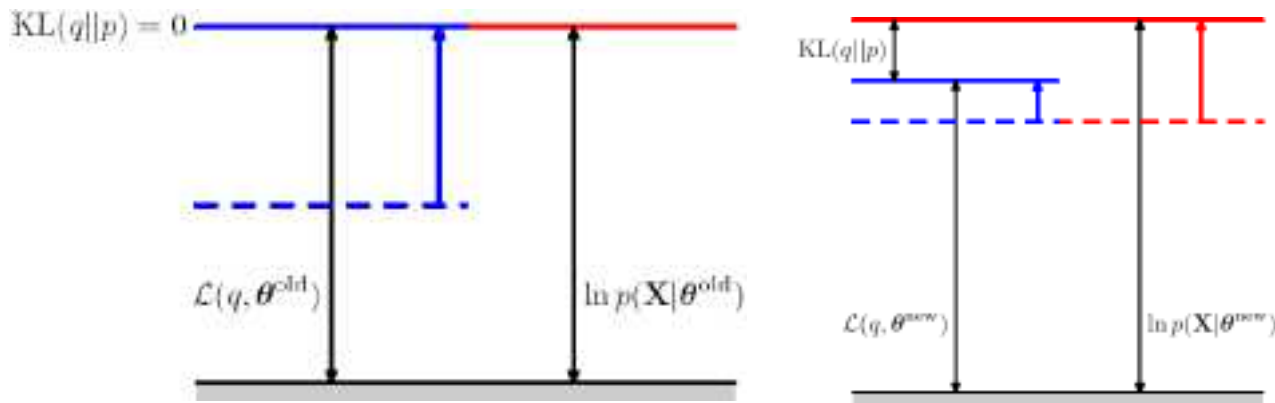
# The EM Algorithm in General (II)

☐ (**E step**) The lower bound $\mathcal{L}$ (q, $\theta_{old}$) is maximized while holding $\theta_{old}$ fixed. Since ln p(X| $\theta$) does not depend on q(Z), $\mathcal{L}$ (q, $\theta_{old}$) will be the largest when KL(q||p) vanishes (i.e. when q(Z) is equal to the posterior distribution p(Z|X, $\theta_{old}$))

☐ (**M step**) q(Z) is fixed and the lower bound $\mathcal{L}$ (q, $\theta_{old}$) is maximized wrt. $\theta$ to $\theta_{new}$ . When the lower bound is increased, $\theta$ is updated making 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
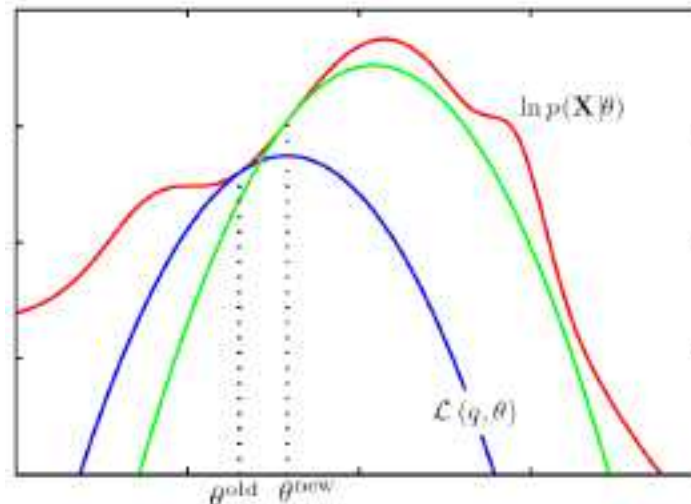
$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$$

$$= \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) + \text{const} \qquad \Longleftarrow \qquad \boxed{q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})}$$

# The EM Algorithm in General (III)

☐ Start with initial parameter value $\theta_{old}$

☐ In the first **E step**, evaluation of posterior distribution over latent variables gives rise to a lower bound $\mathcal{L}$ (q, $\theta_{old}$) whose value equals the log likelihood at $\theta_{old}$ (blue curve)

☐ Note that the bound makes a tangential contact with the log-likelihood at $\theta_{old}$, 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 $\theta_{new}$ which gives a larger value of log-likelihood than $\theta_{old}$.

☐ The subsequent E step constructs a bound tangential at $\theta_{new}$ (green curve)

# The EM Algorithm for MAP

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

$$\boxed{\ln p(\boldsymbol{\theta}|\mathbf{X}) = \ln p(\boldsymbol{\theta}, \mathbf{X}) - \ln p(\mathbf{X})}$$

$$
\begin{aligned}
\ln p(\boldsymbol{\theta}|\mathbf{X}) &= \mathcal{L}(q, \boldsymbol{\theta}) + \mathrm{KL}(q\|p) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}) \\
&\geqslant \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}).
\end{aligned}
$$

# EM Algorithm Variations

**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|$\theta$).

# Incremental EM Algorithm

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

☐ In the **M step**, the means are updated as,

$$\boldsymbol{\mu}_k^{new} = \boldsymbol{\mu}_k^{old} + \left( \frac{\gamma^{new}(z_{mk}) - \gamma^{old}(z_{mk})}{N_k^{new}} \right) (\mathbf{x}_m - \boldsymbol{\mu}_k^{old})$$

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

☐ 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.

# EM for Linear Regression

☐ The complete-data log likelihood

$$\ln p(\mathbf{t}, \mathbf{w}|\alpha, \beta) = \ln p(\mathbf{t}|\mathbf{w}, \beta) + \ln p(\mathbf{w}|\alpha)$$

☐ E-Step: expectation over **w**

$$\mathbb{E}\left[\ln p(\mathbf{t}, \mathbf{w}|\alpha, \beta)\right] = \frac{M}{2}\ln\left(\frac{\alpha}{2\pi}\right) - \frac{\alpha}{2}\mathbb{E}\left[\mathbf{w}^{\mathrm{T}}\mathbf{w}\right] + \frac{N}{2}\ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2}\sum_{n=1}^{N}\mathbb{E}\left[(t_n - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}_n)^2\right]$$

☐ M-Step

$$\alpha = \frac{M}{\mathbb{E}\left[\mathbf{w}^{\mathrm{T}}\mathbf{w}\right]} = \frac{M}{\mathbf{m}_N^{\mathrm{T}}\mathbf{m}_N + \mathrm{Tr}(\mathbf{S}_N)}$$

$$\frac{1}{\beta} = \frac{1}{N}\sum_{i=1}^{N}(t_n - \mathbf{m}_N^T \boldsymbol{\phi}_n)^2$$

# EM for Sparse Kernel Machines

☐ The complete-data log likelihood

$$\ln p(\mathbf{t}, \mathbf{w}|\alpha, \beta) = \ln p(\mathbf{t}|\mathbf{w}, \beta) + \ln p(\mathbf{w}|\alpha)$$

☐ E-Step: expectation over **w**

$$\mathbb{E}_{\mathbf{w}}\left[\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)\right]$$

☐ M-Step

$$
\begin{aligned}
\alpha_i^{\text{new}} &= \frac{1}{m_i^2 + \Sigma_{ii}} \\
(\beta^{\text{new}})^{-1} &= \frac{\|\mathbf{t} - \mathbf{\Phi}\mathbf{m}_N\|^2 + \beta^{-1}\sum_i \gamma_i}{N}
\end{aligned}
$$

# Summary

- Supervised vs Unsupervised Learning

- K-means Clustering

- Gaussian Mixture Model

- Expectation and Maximization

- GMM Revisited

- Bernoulli Mixture Model

- EM Generalization