

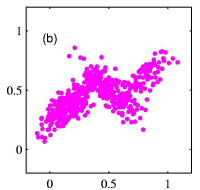
# Latent Variable View of EM

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Main Idea of EM

- Goal of EM is:
  - find maximum likelihood solutions for models having latent variables
- In case of mixture models
  - We have a complex distribution of observed variables  $\mathbf{x}$
  - We wish to estimate its parameters
- Introduce latent variables  $\mathbf{z}$  so that
  - the joint distribution  $p(\mathbf{x}, \mathbf{z})$  of observed and latent variables is more tractable (since we know forms of components)
  - Complicated distribution is formed from simpler components
- The original distribution is obtained by marginalizing the joint distribution



# Alternative View of EM

- This view recognizes key role of latent variables

- Observed data 
$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

- Latent Variables 
$$Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix}$$

- where  $n^{\text{th}}$  row represents  $x_n^T = [x_{n1} \ x_{n2} \ \dots \ x_{nD}]$
- with corresponding row  $z_n^T = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$

- Goal of EM algorithm is to find maximum likelihood solutions for models having latent variables

# Likelihood Function involving Latent Variables

- Joint likelihood function is  $p(X, Z | \theta)$  where  $\theta$  is the set of all model parameters
  - E.g., means, covariances, responsibilities
- Marginal likelihood function of observed data
  - From sum rule

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

- Log likelihood function is

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

# Latent Variables in EM

- Log likelihood function is

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

Summation inside brackets  
due to marginalization  
Not due to log-likelihood

- Key Observation:

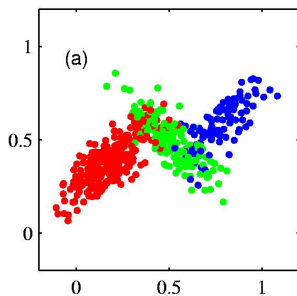
– Summation over latent variables appears inside logarithm

- Even if joint distribution  $p(X, Z | \theta)$  belongs to exponential family the marginal distribution  $p(X | \theta)$  does not
- Results in complicated expressions for maximum likelihood solution

# Complete and Incomplete Data Sets

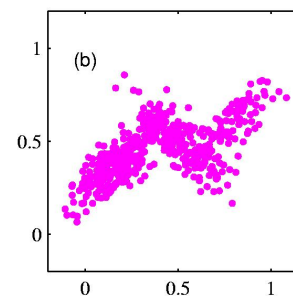
## Complete Data $\{X, Z\}$

- For each observation in  $X$  we know corresponding value of latent variable  $Z$
- Log-likelihood has the form  $\ln p(X, Z | \theta)$ 
  - maximization is straightforward



## Incomplete Data $\{X\}$

- Actual data set
- Log likelihood function is
 
$$\ln p(X | \theta) = \ln \left\{ \sum_Z p(X, Z | \theta) \right\}$$
- Maximization is difficult
  - summations inside logarithm



# Expectation of log-likelihood

- Since we don't have the complete data set  $\{X, Z\}$  we evaluate the expectation of the log-likelihood  $\ln p(X, Z | \theta)$
- Since we are given  $X$ , our state of knowledge of  $Z$  is given only by the posterior distribution of the latent variables  $p(Z | X, \theta)$
- Thus the expected value of the log-likelihood of the complete data set is

$$\sum_Z p(Z | X, \theta) \underbrace{\ln p(X, Z | \theta)}$$

Summation is due to expectation not sum rule!

We maximize this.

Note that the logarithm acts on the joint-- which is tractable

# E and M Steps

- E Step: Use current parameter value  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(Z|X, \theta^{\text{old}})$
- Use this posterior to find *expectation* of complete-data log-likelihood for some general parameter value  $\theta$ . This expectation is

$$Q(\theta, \theta^{\text{old}}) = \sum_z p(Z|X, \theta) \ln p(X, Z | \theta)$$

- M Step: Determine revised parameter estimate  $\theta^{\text{new}}$  by *maximizing*

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$$

# General EM Algorithm

- Given joint distribution  $p(\mathbf{X}, \mathbf{Z} | \theta)$  over observed variables  $\mathbf{X}$  and latent variables  $\mathbf{Z}$  governed by parameters  $\theta$   
goal is to maximize likelihood function  $p(\mathbf{X} | \theta)$

- **Step 1:** Choose an initial setting for the parameters  $\theta^{\text{old}}$
- **Step 2:** E Step: Evaluate  $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$
- **Step 3:** M Step: Evaluate  $\theta^{\text{new}}$  given by
$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$$

where  $Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$

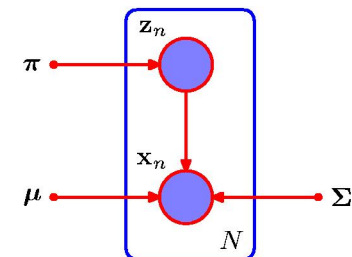
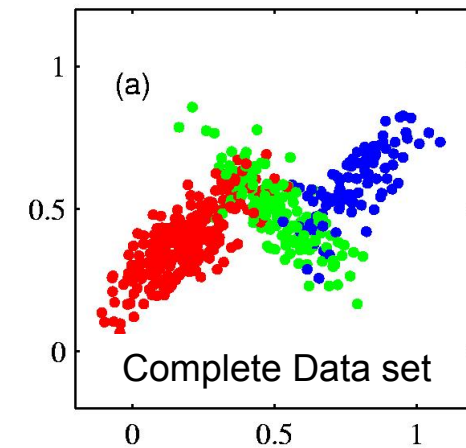
- Check for convergence
  - of either log-likelihood or parameter values
- If not satisfied then let  $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$
- Return to **Step 2**

# Missing Variables

- EM has been described for maximum likelihood function when there are discrete latent variables
- It can also be applied when there are unobserved variables corresponding to missing values in data set
  - Take the joint distribution of all variables and then marginalize over missing ones
  - EM is then used to maximize corresponding likelihood function
- Method is valid when data is missing at random
  - Not if missing value depends on unobserved values
  - E.g., if quantity exceeds some threshold

# Gaussian Mixtures Revisited

- Application of latent variable view to GMM
- Maximizing the log likelihood function for the complete data set  $\{X, Z\}$  leads to the same EM algorithm for Gaussian mixtures as before



# Relation to K-means

- EM for Gaussian mixtures has close similarity to K-means
- K-means performs a hard assignment of data points to clusters
  - Each data point is associated uniquely with one cluster
- EM makes a soft assignment based on posterior probabilities
- K-means does not estimate the covariances of the clusters but only the cluster means

# Mixtures of Bernoulli Distributions

- Previously considered distributions over continuous variables
- Now consider mixtures of discrete binary variables described by Bernoulli distributions
- Sets a foundation for HMMs over discrete variables

# Multivariate Bernoulli

- Set of  $D$  independent binary variables  $x_i$ ,  $i=1,\dots,D$ 
  - E.g., a set of  $D$  coins with heads and tails
- Each governed by parameter  $\mu_i$
- Multivariate distribution

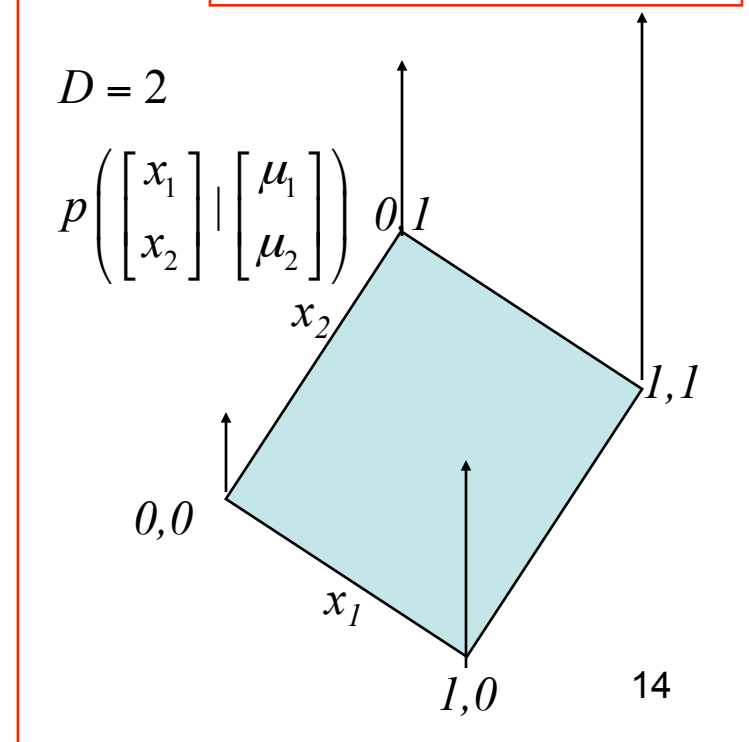
$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{(1-x_i)}$$

where  $\mathbf{x} = (x_1, \dots, x_D)^T$  and

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_D)^T$$

- Mean and covariance are  
 $E[\mathbf{x}] = \boldsymbol{\mu}$ ,  $\text{cov}[\mathbf{x}] = \text{diag}\{\mu_i(1-\mu_i)\}$

A Bernoulli distribution with two variables



# Mixture of multivariate Bernoulli

- Finite mixture of  $K$  Bernoulli distributions
  - E.g.,  $K$  bags of  $D$  coins each where bag  $k$  is chosen with probability  $\pi_k$

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

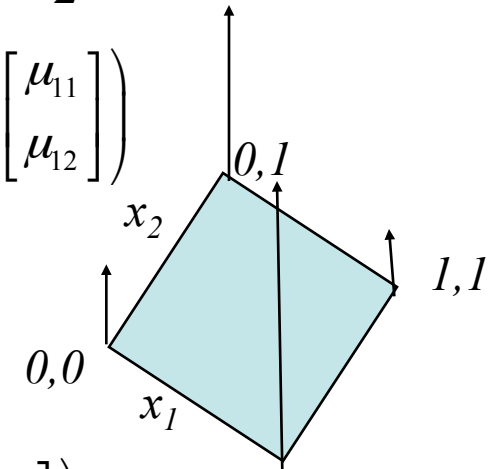
- Where  $\mu = \{\mu_1, \dots, \mu_K\}$ ,  
 $\pi = \{\pi_1, \dots, \pi_K\}$ , and

$$p(\mathbf{x} | \mu_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{(1-x_i)}$$

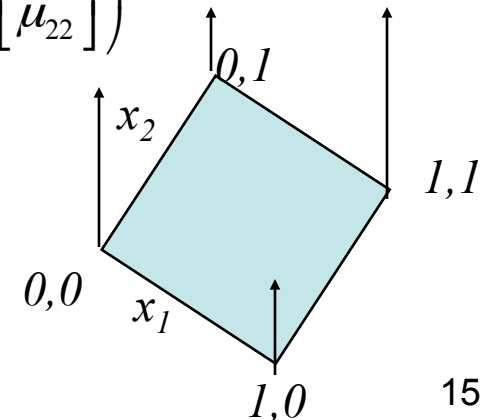
Two Bernoulli distributions with two variables

$$D = 2, K = 2$$

$$p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_{11} \\ \mu_{12} \end{bmatrix}\right)$$



$$p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_{21} \\ \mu_{22} \end{bmatrix}\right)$$



# Log likelihood of Bernoulli mixture

- Given data set  $X = \{x_1, \dots, x_N\}$  log likelihood of model is

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

Summation due to logarithm

Summation due to mixture

- Due to summation inside logarithm there is no closed form m.l.e. solution

# Introduce latent variables

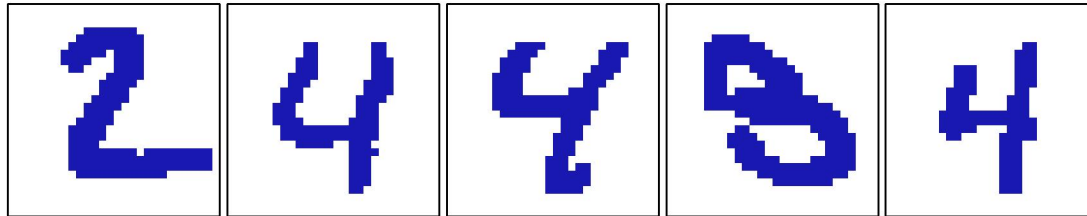
- One of  $K$  representation  $\mathbf{z} = (z_1, \dots, z_K)^T$
- Conditional distribution of  $\mathbf{x}$  given the latent variable is

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

- The EM algorithm is derived by writing the complete data log-likelihood function
- Taking its expectation w.r.t. posterior distribution of the latent variables
- In the E step these responsibilities are evaluated using Bayes theorem
- In the M step we maximize the expected complete data loglikelihood wrt parameters  $\mu_k$  and  $\pi$

# Illustration of Bernoulli Mixture

We are given a set of unlabeled digits 2,3 and 4  
Goal is to use EM to cluster them with  $K=3$



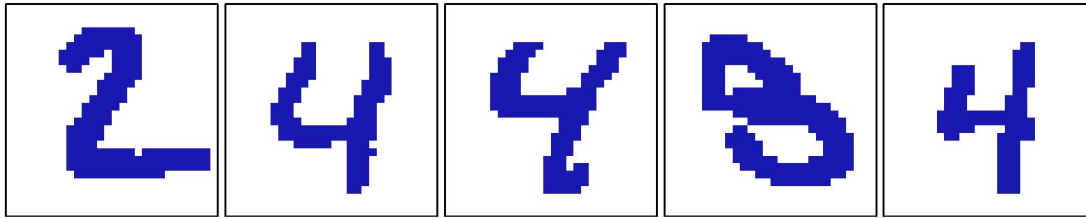
Binary images after grey-scale thresholding at 0.5

$N = 600$

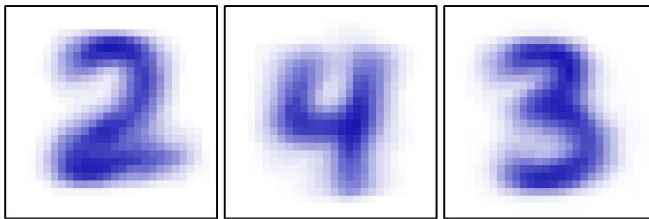
Mixing coefficients initialized with  $\pi_k = 1/K$

Parameters  $\mu_{ki}$  were set to random values chosen uniformly in range  $(0.25, 0.75)$  and normalized so that  $\sum_j \mu_{kj} = 1$

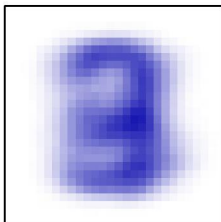
# Result of EM algorithm



EM finds the three clusters



Parameters  $\mu_{ki}$  for each of three components of mixture model



Using single multivariate Bernoulli and maximum likelihood amounts to averaging counts in each pixel

# Bayesian EM for Discrete Case

- Conjugate prior of the parameters of Bernoulli is given by the beta distribution
- Beta prior is equivalent to introducing additional effective observations of  $\mathbf{x}$
- Also introduce priors into the Bernoulli mixture model
- Use EM to maximize posterior probability of distribution
- Can be extended to multinomial discrete variables
  - Introduce Dirichlet priors over model parameters if desired

# EM Algorithm in General

- EM is a general technique for finding maximum likelihood solutions for probabilistic models with latent variables (Dempster 1977)
- EM defined heuristically can be proved to maximize the likelihood function
- Proof involves obtaining lower bound on log-likelihood function

# Decomposition of log-likelihood function

- Given observed variables  $X$  and hidden variables  $Z$ , goal is to maximize likelihood function

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

- Direct optimization of  $p(X | \theta)$  is difficult, whereas complete-data likelihood  $p(X, Z | \theta)$  is easier (since component forms are tractable)
- For any choice of  $q(Z)$  the following holds

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

where we define

$$L(q, \theta) = \sum_Z q(Z) \ln \left\{ \frac{p(X, Z | \theta)}{q(Z)} \right\}$$

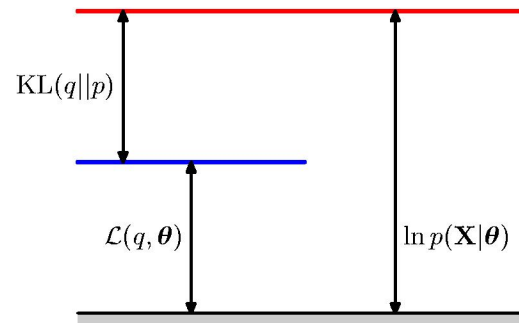
$$KL(q || p) = - \sum_Z q(Z) \ln \left\{ \frac{p(Z | X, \theta)}{q(Z)} \right\}$$

$L$  is a *functional* that takes a function as input and produces a value as output, like entropy

Kullback-Leibler  
Divergence

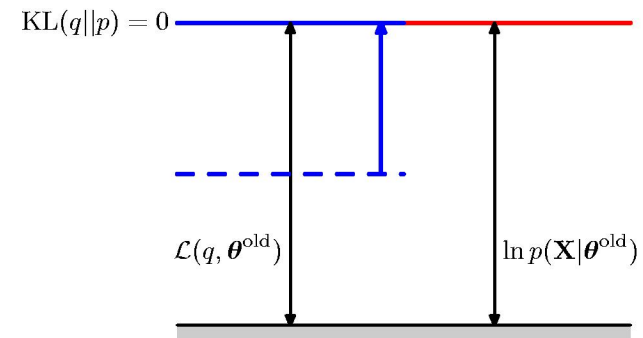
# Bounds on the log-likelihood

- Since  $\text{KL}(q||p) \geq 0$
- It follows that  $\mathcal{L}(q, \theta) \leq \ln p(\mathbf{X}|\theta)$ 
  - Or  $\mathcal{L}(q, \theta)$  is a lower bound on  $\ln p(\mathbf{X}|\theta)$

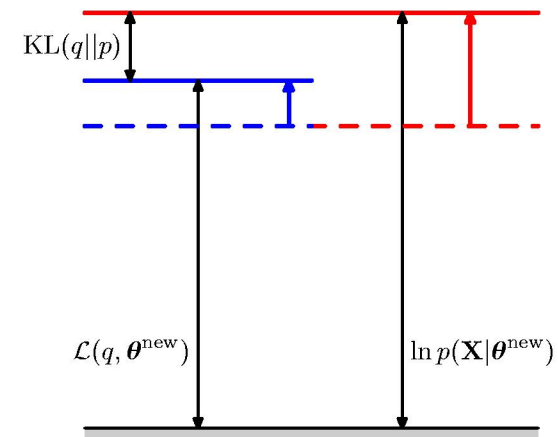


# Illustration of E and M steps

- E step

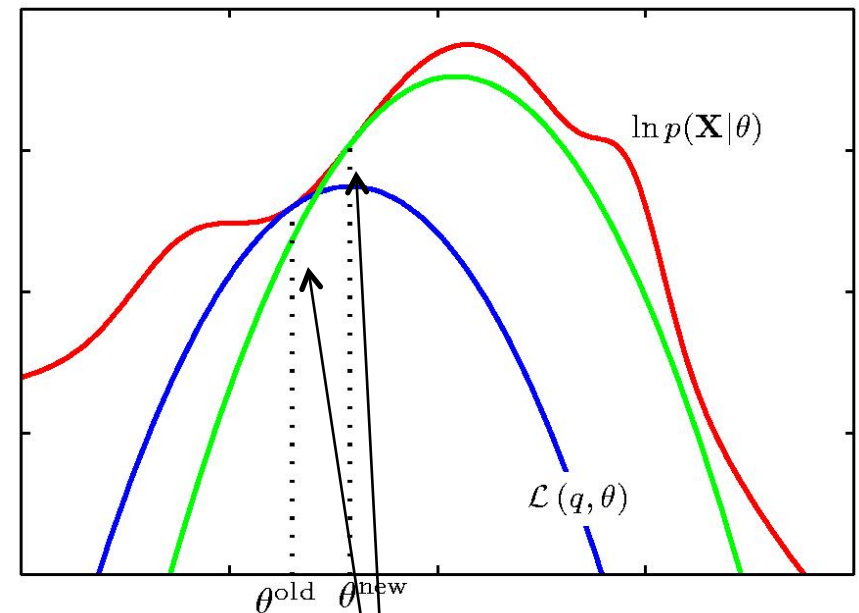


- M step



# View of EM in parameter space

- EM involves alternately computing lower bound on log-likelihood for the current parameter values
- And maximizing this bound to obtain new parameter values



Lower bound on  
log-likelihood function tangential  
at  $\theta^{\text{old}}$  and  $\theta^{\text{new}}$

# Generalized EM (GEM)

- EM breaks own potentially difficult problem of maximizing the likelihood function into two stages, the E step and the M step
- One or both may remain intractable
- GEM addresses the problem of the intractable M step
  - Instead of maximizing  $L(q, \theta)$  wrt  $\theta$  it changes parameters so as to increase its value
- Similar generalization of the E step can be made