1

### **Mixtures of Gaussians**

#### Sargur Srihari srihari@cedar.buffalo.edu

### 9. Mixture Models and EM

0. Mixture Models Overview
1.K-Means Clustering
2.Mixtures of Gaussians
3.An Alternative View of EM
4.The EM Algorithm in General

# **Topics in Mixtures of Gaussians**

- Goal of Gaussian Mixture Modeling
- Latent Variables
- Maximum Likelihood
- EM for Gaussian Mixtures

### GMMs and Latent Variables

- A GMM is a linear superposition of Gaussian components
  - Provides a richer class of density models than the single Gaussian
- We formulate a GMM in terms of discrete latent variables
  - This provides deeper insight into this distribution
  - Serves to motivate the EM algorithm
    - Which gives a maximum likelihood solution to no. of components and their means/covariances

# Latent Variable Representation

- Gaussian mixture distribution is written as
  - a linear superposition of K Gaussian components:

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

- Represent K by a K-dimensional binary variable z
  - Using 1-of-*K* representation (one-hot vector)

• Let 
$$z = z_1, ..., z_K$$
 whose elements are  $z_k \in \{0,1\}$  and  $\sum_k z_k = 1$ 



1	k = 1	2
z	10	01
$oldsymbol{\pi}_k$	0.4	0.6
$\mu_{_k}$	-28	1.86
$oldsymbol{\sigma}_{_k}$	0.48	0.88



k	1	2	3
$\boldsymbol{z}$	100	010	001

# Goal of Gaussian Mixture Modeling

• A linear superposition of Gaussians in the form

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

- Goal of Modeling:
  - Find maximum likelihood parameters  $\pi_k$ ,  $\mu_k$ ,  $\Sigma_k$
  - Examples of data sets and models



Each data point is associated with a subclass k with probability  $\pi_k$ 



### Joint Distribution

 Define joint distribution of latent variable and observed variable

$$-p(\boldsymbol{x},\boldsymbol{z}) = p(\boldsymbol{x}|\boldsymbol{z}) \bullet p(\boldsymbol{z})$$

- -x is observed variable
- -z is the hidden or missing variable
- Marginal distribution p(z)
- Conditional distribution  $p(\boldsymbol{x}|\boldsymbol{z})$

### **Graphical Representation of Mixture Model**

- The joint distribution  $p(\pmb{x},\pmb{z})$  is represented in the form  $p(\pmb{z})p(\pmb{x}|\pmb{z})$ 



– We now specify marginal p(z) and conditional p(x|z)

• Using them we specify  $p(\textbf{\textit{x}})$  in terms of observed and latent variables

Specifying the marginal p(z)• Associate a probability with each component  $z_k$ - Denote  $p(z_k = 1) = \pi_k$  where parameters  $\{\pi_k\}$  satisfy  $0 \le \pi_k \le 1$  and  $\sum \pi_k = 1$  $p(\boldsymbol{z})$ 

• Because z uses 1-of-K it follows that

- since  $z_k \in \{0,1\}$  and components of z are mutually exclusive and hence are independent

With one component  $p(z_1) = \pi_1^{z_1}$ 

With two components  $p(z_1, z_2) = \pi_1^{z_1} \pi_2^{z_2}$ 

# Specifying the Conditional $p(\boldsymbol{x}|\boldsymbol{z})$

- For a particular component (value of z)  $p(\boldsymbol{x} \mid z_k = 1) = N(\boldsymbol{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- Thus  $p(\boldsymbol{x}|\boldsymbol{z})$  can be written in the form

$$p(\boldsymbol{x} \mid \boldsymbol{z}) = \prod_{k=1}^{K} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})^{z_{k}}$$



– Due to the exponent  $z_k$  all product terms except for one equal one

# Marginal distribution $p(\mathbf{x})$

- The joint distribution  $p(\mathbf{x}, \mathbf{z})$  is given by  $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$
- Thus marginal distribution of x is obtained by summing over all possible states of z to give

$$p(\boldsymbol{x}) = \sum_{\boldsymbol{z}} p(\boldsymbol{z}) p(\boldsymbol{x} \mid \boldsymbol{z}) = \sum_{\boldsymbol{z}} \prod_{k=1}^{K} \pi^{z_{k}} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})^{z_{k}} = \sum_{k=1}^{K} \pi_{k} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$
  
Since  $z_{k} \in \{0, 1\}$ 

• This is the standard form of a Gaussian mixture

### Value of Introducing Latent Variable

- If we have observations  $\pmb{x}_1,..,\pmb{x}_N$
- Because marginal distribution is in the form

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

- It follows that for every observed data point  $x_n$  there is a corresponding latent vector  $z_n$ , i.e., its sub-class
- Thus we have found a formulation of Gaussian mixture involving an explicit latent variable
  - We are now able to work with joint distribution p(x,z) instead of marginal p(x)
- Leads to significant simplification through introduction of expectation maximization

#### Another conditional probability (Responsibility)

- In EM  $p(\boldsymbol{z} \mid \boldsymbol{x})$  plays a role
- The probability  $p(z_k=1|\mathbf{x})$  is denoted  $\gamma(z_k)$ 
  - From Bayes theorem

- View  $p(z_k = 1) = \pi_k$  as prior probability of component kand  $\gamma(z_k) = p(z_k = 1 | \mathbf{x})$  as the posterior probability it is also the responsibility that component k takes for explaining the observation  $\mathbf{x}$ 

### Plan of Discussion

- Next we look at
  - 1. How to get data from a mixture model synthetically and then
  - 2. Given a data set  $\{x_1, ..., x_N\}$  how to model the data using a mixture of Gaussians

X

# Synthesizing data from mixture

- Use ancestral sampling
  - Start with lowest numbered node and draw a sample,
    - Generate sample of  $\pmb{z}$ , called  $\hat{\pmb{z}}$
    - move to successor node and draw a sample given the parent value, etc.
    - Then generate a value for  $\pmb{x}$  from conditional  $p(\pmb{x} | \ \hat{\pmb{z}} \ )$
- Samples from p(x,z) are plotted according to value of x and colored with value of z



Complete Data set

(a)



# Illustration of responsibilities

- Evaluate for every data point
  - Posterior probability of each component
  - Responsibilities  $\gamma(z_{nk})$  associated with data point  $x_n$  and component k



- If for a data point  $\gamma(z_{n1}) = 1$  it is colored red
- If for another point  $\gamma(z_{n2}) = \gamma(z_{n3}) = 0.5$  it has equal blue and green and will appear as cyan



# Maximum Likelihood for GMM

- We wish to model data set  $\{x_1, .., x_N\}$  using a mixture of Gaussians (*N* items each of dimension *D*)
- Represent by  $N \times D$  matrix X
  - $n^{\mathsf{th}}$  row is given by  $oldsymbol{x}_n{}^T$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_N \end{bmatrix}$$

- Represent N latent variables with  $N \times K$  matrix Z
  - $n^{\mathsf{th}}$  row is given by  $oldsymbol{z}_n{}^T$

$$Z = \begin{bmatrix} z_1 \\ z_2 \\ z_N \end{bmatrix}$$

- Goal is to state the likelihood function
  - so as to estimate the three sets of parameters
  - by maximizing the likelihood

### Graphical representation of GMM

- For a set of i.i.d. data points {*x<sub>n</sub>*} with corresponding latent points {*z<sub>n</sub>*} where *n*=1,...,*N*
- Bayesian Network for  $p(\pmb{X,Z})$  using plate notation
  - $-N \mathbf{x} D \mathbf{matrix} \mathbf{X}$
  - $-N \mathbf{x} K \mathbf{matrix} \mathbf{Z}$



### Likelihood Function for GMM

Mixture density function is

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\boldsymbol{z}) p(\boldsymbol{x} \mid \boldsymbol{z}) = \sum_{k=1}^{K} \pi_{k} N(\boldsymbol{x} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

Since z has values  $\{z_k\}$  with probabilities  $\{\pi_k\}$ 

Therefore Likelihood function is

$$p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^{N} \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\}$$

Product is over the *N* i.i.d. samples

Therefore log-likelihood function is

$$\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\}$$

Which we wish to maximize A more difficult problem than for a single Gaussian

# Maximization of Log-Likelihood

$$\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\}$$

Goal is to estimate the three sets of parameters

$$\pi_k, \mu_k, \Sigma_k$$

- By taking derivatives in turn w.r.t each while keeping others constant
- But there are no closed-form solutions
  - Task is not straightforward since summation appears in Gaussian and logarithm does not operate on Gaussian
- While a gradient-based optimization is possible, we consider the iterative EM algorithm 20

### Some issues with GMM m.l.e.

- Before proceeding with the m.l.e. briefly mention two technical issues:
  - 1. Problem of singularities with Gaussian mixtures
  - 2. Problem of Identifiability of mixtures

#### Problem of Singularities with Gaussian mixtures

- Consider Gaussian mixture
  - components with covariance matrices  $\Sigma_k = \sigma_k^2 I^{p(x)}$
- Data point that falls on a mean  $\mu_j = x_n$  will contribute to the likelihood function

$$\left| N(\boldsymbol{x}_{n} | \boldsymbol{x}_{n}, \boldsymbol{\sigma}_{j}^{2} I) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\boldsymbol{\sigma}_{j}} \right|$$

since 
$$\exp(\textbf{\textit{x}}_n\textbf{-}\mu_j)^2 = 1$$

• As  $\sigma_j \rightarrow 0$  term goes to infinity

• Therefore maximization of log-likelihood

 $\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\} \text{ is not well-posed}$ 

- Does not happen with a single Gaussian
  - Multiplicative factors go to zero
- Does not happen in the Bayesian approach
- Problem is avoided using heuristics
  - Resetting mean or covariance



One component assigns finite values and other to large value



# Problem of Identifiability

A density  $p(x | \theta)$  is identifiable if  $\theta \neq \theta'$  then there is an x for which  $p(x | \theta) \neq p(x | \theta')$ 

# A *K*-component mixture will have a total of *K*! equivalent solutions

- Corresponding to *K*! ways of assigning *K* sets of parameters to *K* components
  - E.g., for K=3 K!=6: 123, 132, 213, 231, 312, 321
- For any given point in the space of parameter values there will be a further K!-1 additional points all giving exactly same distribution
- However any of the equivalent solutions is as good as the other

Two ways of labeling three Gaussian subclasses



# EM for Gaussian Mixtures

- EM is a method for finding maximum likelihood solutions for models with latent variables
- Begin with log-likelihood function

$$\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\}$$

- We wish to find  $\pi, \mu, \Sigma$  that maximize this quantity
- Task is not straightforward since summation appears in Gaussian and logarithm does not operate on Gaussian
- Take derivatives in turn w.r.t
  - Means  $\mu_k$  and set to zero
  - covariance matrices  $\sum_{k}$  and set to zero
  - mixing coefficients  $\pi_k$  and set to zero

# EM for GMM: Derivative wrt $\mu_k$

Begin with log-likelihood function

$$\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \boldsymbol{\pi}_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right\}$$

- Take derivative w.r.t the means  $\mu_k$  and set to zero
  - Making use of exponential form of Gaussian
  - Use formulas:  $\frac{d}{dx} \ln u = \frac{u'}{u}$  and  $\frac{d}{dx} e^u = e^u u'$

– We get

$$0 = \sum_{n=1}^{N} \frac{\pi_{k} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})}{\sum_{j} \pi_{j} N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{j}, \boldsymbol{\Sigma}_{j})} \sum_{k=1}^{-1} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})$$
  
Inverse of covariance matrix  
 $\gamma(z_{nk})$  the posterior probabilities

### M.L.E. solution for Means

• Multiplying by  $\Sigma_k$  (assuming non-singularity)

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \boldsymbol{\gamma}(\boldsymbol{z}_{nk}) \boldsymbol{x}_{n}$$

Where we have defined

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

Mean of  $k^{\text{th}}$  Gaussian component is the weighted mean of <u>all</u> the points in the data set: where data point  $\boldsymbol{x}_n$  is weighted by the posterior probability that

component k was responsible for generating  $\boldsymbol{x}_n$ 

 Which is the effective number of points assigned to cluster k

### M.L.E. solution for Covariance

- Set derivative wrt  $\Sigma_k$  to zero
  - Making use of mle solution for covariance matrix of single Gaussian

$$\boxed{\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \boldsymbol{\gamma}(\boldsymbol{z}_{nk}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})^{T}}$$

- Similar to result for a single Gaussian for the data set but each data point weighted by the corresponding posterior probability
- Denominator is effective no of points in component

### M.L.E. solution for Mixing Coefficients

- Maximize  $\ln p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  w.r.t.  $\pi_k$ 
  - Must take into account that mixing coefficients sum to one
  - Achieved using Lagrange multiplier and maximizing

$$\ln p(X \mid \pi, \mu, \Sigma) + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right)$$

– Setting derivative wrt  $\pi_k$  to zero and solving gives

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

# Summary of m.l.e. expressions

GMM maximum likelihood parameter estimates

Means

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

**Covariance matrices** 

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

**Mixing Coefficients** 

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

- All three are in terms of responsibilities
- and so we have not completely solved the problem

### **EM Formulation**

- The results for  $\mu_k, \Sigma_k, \pi_k$  are not closed form solutions for the parameters
  - Since  $\gamma(z_{nk})$  the responsibilities depend on those parameters in a complex way
- Results suggest an iterative solution
- An instance of EM algorithm for the particular case of GMM

# Informal EM for GMM

- First choose initial values for means, covariances and mixing coefficients
- Alternate between following two updates
  - Called E step and M step
- In E step use current value of parameters to evaluate posterior probabilities, or responsibilities
- In the M step use these posterior probabilities to to reestimate means, covariances and mixing coefficients

# EM using Old Faithful

Data points and Initial mixture model



Initial E step Determine responsibilities



After first M step Re-evaluate Parameters



After 2 cycles

0

(d)

2

-2

-2



After 20 cycles  $\begin{bmatrix} 2 \\ 0 \\ -2 \end{bmatrix}$ 

0

(f)

2

-2

### Comparison with *K*-Means



### Animation of EM for Old Faithful Data

 <u>http://en.wikipedia.org/wiki/</u> <u>File:Em\_old\_faithful.gif</u>

Code in R

#initial parameter estimates (chosen to be deliberately bad) theta <- list( tau=c(0.5,0.5), mu1=c(2.8,75), mu2=c(3.6,58), sigma1=matrix(c(0.8,7,7,70),ncol=2), sigma2=matrix(c(0.8,7,7,70),ncol=2) )

### Practical Issues with EM

- Takes many more iterations than *K*-means
  Each cycle requires significantly more comparison
- Common to run *K*-means first in order to find suitable initialization
- Covariance matrices can be initialized to covariances of clusters found by *K*-means
- EM is not guaranteed to find global maximum of log likelihood function

### Summary of EM for GMM

- Given a Gaussian mixture model
- Goal is to maximize the likelihood function w.r.t. the parameters (means, covariances and mixing coefficients)

**Step1:** Initialize the means  $\mu_k$  covariances  $\Sigma_k$ and mixing coefficients  $\pi_k$  and evaluate initial value of log-likelihood

### **EM** continued

• Step 2: E step: Evaluate responsibilities using current parameter values

$$\gamma(z_k) = \frac{\pi_k N(\boldsymbol{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j N(\boldsymbol{x}_n \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j))}$$

• Step 3: M Step: Re-estimate parameters using current responsibilities

n=1

$$\begin{aligned} \mu_{k}^{\text{new}} &= \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_{n} \\ \overline{\boldsymbol{\Sigma}_{k}^{\text{new}}} &= \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) (\mathbf{x}_{n} - \mu_{k}^{\text{new}}) (\mathbf{x}_{n} - \mu_{k}^{\text{new}})^{T} \\ \overline{\boldsymbol{\pi}_{k}^{\text{new}}} &= \frac{N_{k}}{N} \end{aligned}$$
 where  $N_{k} = \sum_{n=1}^{N} \gamma(z_{nk}) (z_{nk}) (z_{nk} - \mu_{k}^{\text{new}})^{T} \\ \overline{\boldsymbol{\pi}_{k}^{\text{new}}} &= \frac{N_{k}}{N} \end{aligned}$ 

### **EM Continued**

• Step 4: Evaluate the log likelihood

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

- And check for convergence of either parameters or log likelihood
- If convergence not satisfied return to Step 2