

# Unsupervised Methods: EM

## Lecture 7b

Jason Corso

SUNY at Buffalo

26 March 2009

# Gaussian Mixture Models

- Recall the Gaussian distribution:

$$\mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1)$$

# Gaussian Mixture Models

- Recall the Gaussian distribution:

$$\mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1)$$

- It forms the basis for the important Mixture of Gaussians density.

# Gaussian Mixture Models

- Recall the Gaussian distribution:

$$\mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1)$$

- It forms the basis for the important Mixture of Gaussians density.
- The Gaussian mixture is a **linear superposition of Gaussians** in the form:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) . \quad (2)$$

# Gaussian Mixture Models

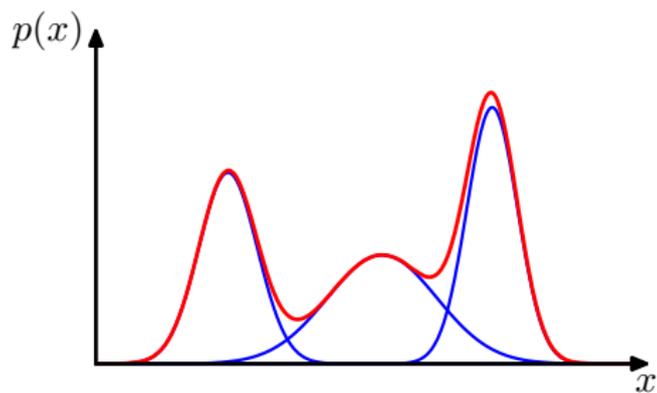
- Recall the Gaussian distribution:

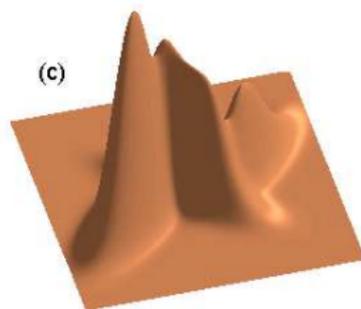
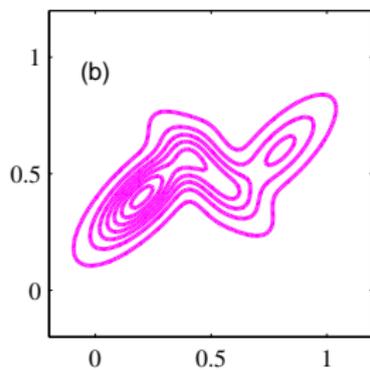
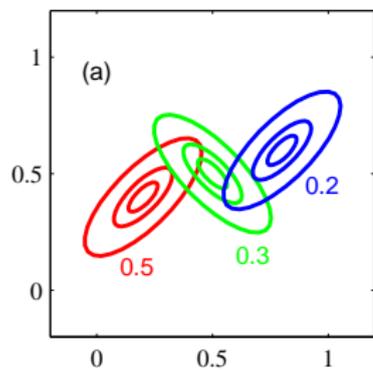
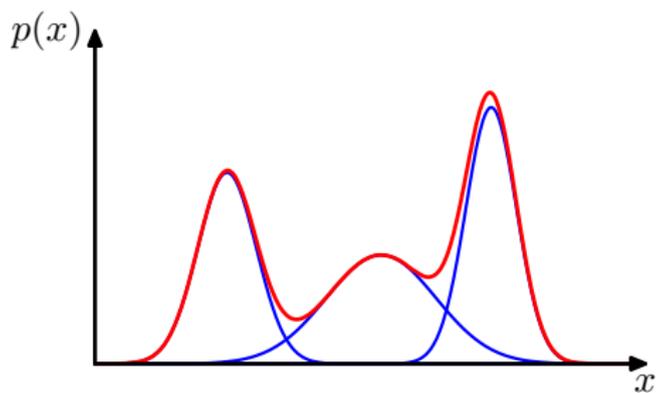
$$\mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1)$$

- It forms the basis for the important Mixture of Gaussians density.
- The Gaussian mixture is a **linear superposition of Gaussians** in the form:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) . \quad (2)$$

- The  $\pi_k$  are non-negative scalars called **mixing coefficients** and they govern the relative importance between the various Gaussians in the mixture density.  $\sum_k \pi_k = 1$ .





# Introducing Latent Variables

- Define a  $K$ -dimensional binary random variable  $\mathbf{z}$ .

# Introducing Latent Variables

- Define a  $K$ -dimensional binary random variable  $\mathbf{z}$ .
- $\mathbf{z}$  has a 1-of- $K$  representation such that a particular element  $z_k$  is 1 and all of the others are zero. Hence:

$$z_k \in \{0, 1\} \quad (3)$$

$$\sum_k z_k = 1 \quad (4)$$

# Introducing Latent Variables

- Define a  $K$ -dimensional binary random variable  $\mathbf{z}$ .
- $\mathbf{z}$  has a 1-of- $K$  representation such that a particular element  $z_k$  is 1 and all of the others are zero. Hence:

$$z_k \in \{0, 1\} \quad (3)$$

$$\sum_k z_k = 1 \quad (4)$$

- The marginal distribution over  $\mathbf{z}$  is specified in terms of the mixing coefficients:

$$p(z_k = 1) = \pi_k \quad (5)$$

And, recall,  $0 \leq \pi_k \leq 1$  and  $\sum_k \pi_k = 1$ .

- Since  $\mathbf{z}$  has a 1-of- $K$  representation, we can also write this distribution as

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (6)$$

- Since  $\mathbf{z}$  has a 1-of- $K$  representation, we can also write this distribution as

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (6)$$

- The conditional distribution of  $\mathbf{x}$  given  $\mathbf{z}$  is a Gaussian:

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (7)$$

or

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (8)$$

- We are interested in the marginal distribution of  $\mathbf{x}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) \quad (9)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \quad (10)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k^{z_k} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (11)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12)$$

- We are interested in the marginal distribution of  $\mathbf{x}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) \quad (9)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \quad (10)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k^{z_k} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (11)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12)$$

- So, given our latent variable  $\mathbf{z}$ , the marginal distribution of  $\mathbf{x}$  is a Gaussian mixture.

- We are interested in the marginal distribution of  $\mathbf{x}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) \quad (9)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \quad (10)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k^{z_k} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (11)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12)$$

- So, given our latent variable  $\mathbf{z}$ , the marginal distribution of  $\mathbf{x}$  is a Gaussian mixture.
- If we have  $N$  observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , then because of our chosen representation, it follows that we have a latent variable  $\mathbf{z}_n$  for each observed data point  $\mathbf{x}_n$ .

# Component Responsibility Term

- We need to also express the conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$ .

# Component Responsibility Term

- We need to also express the conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$ .
- Denote this conditional  $p(z_k = 1|\mathbf{x})$  as  $\gamma(z_k)$ .

# Component Responsibility Term

- We need to also express the conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$ .
- Denote this conditional  $p(z_k = 1|\mathbf{x})$  as  $\gamma(z_k)$ .
- We can derive this value with Bayes' theorem:

$$\gamma(z_k) \doteq p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \quad (13)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (14)$$

# Component Responsibility Term

- We need to also express the conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$ .
- Denote this conditional  $p(z_k = 1|\mathbf{x})$  as  $\gamma(z_k)$ .
- We can derive this value with Bayes' theorem:

$$\gamma(z_k) \doteq p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \quad (13)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (14)$$

- View  $\pi_k$  as the prior probability of  $z_k = 1$  and the quantity  $\gamma(z_k)$  as the corresponding posterior probability once we have observed  $\mathbf{x}$ .

# Component Responsibility Term

- We need to also express the conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$ .
- Denote this conditional  $p(z_k = 1|\mathbf{x})$  as  $\gamma(z_k)$ .
- We can derive this value with Bayes' theorem:

$$\gamma(z_k) \doteq p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \quad (13)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (14)$$

- View  $\pi_k$  as the prior probability of  $z_k = 1$  and the quantity  $\gamma(z_k)$  as the corresponding posterior probability once we have observed  $\mathbf{x}$ .
- $\gamma(z_k)$  can also be viewed as the responsibility that component  $k$  takes for explaining the observation  $\mathbf{x}$ .

# Sampling from the GMM

- To sample from the GMM, we can first generate a value for  $\mathbf{z}$  from the marginal distribution  $p(\mathbf{z})$ . Denote this sample  $\hat{\mathbf{z}}$ .

# Sampling from the GMM

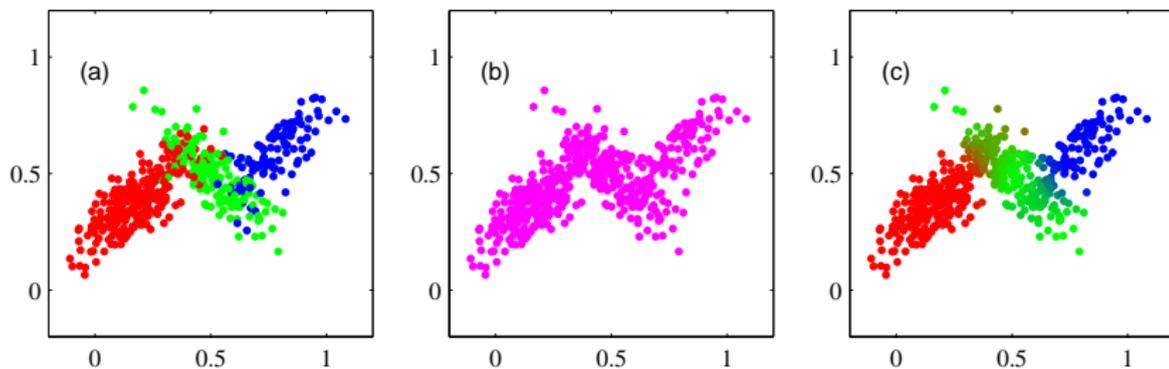
- To sample from the GMM, we can first generate a value for  $\mathbf{z}$  from the marginal distribution  $p(\mathbf{z})$ . Denote this sample  $\hat{\mathbf{z}}$ .
- Then, sample from the conditional distribution  $p(\mathbf{x}|\hat{\mathbf{z}})$ .

# Sampling from the GMM

- To sample from the GMM, we can first generate a value for  $\mathbf{z}$  from the marginal distribution  $p(\mathbf{z})$ . Denote this sample  $\hat{\mathbf{z}}$ .
- Then, sample from the conditional distribution  $p(\mathbf{x}|\hat{\mathbf{z}})$ .
- The figure below-left shows samples from a three-mixture and colors the samples based on their  $\mathbf{z}$ . The figure below-middle shows samples from the marginal  $p(\mathbf{x})$  and ignores  $\mathbf{z}$ . On the right, we show the  $\gamma(z_k)$  for each sampled point, colored accordingly.

# Sampling from the GMM

- To sample from the GMM, we can first generate a value for  $\mathbf{z}$  from the marginal distribution  $p(\mathbf{z})$ . Denote this sample  $\hat{\mathbf{z}}$ .
- Then, sample from the conditional distribution  $p(\mathbf{x}|\hat{\mathbf{z}})$ .
- The figure below-left shows samples from a three-mixture and colors the samples based on their  $\mathbf{z}$ . The figure below-middle shows samples from the marginal  $p(\mathbf{x})$  and ignores  $\mathbf{z}$ . On the right, we show the  $\gamma(z_k)$  for each sampled point, colored accordingly.



# Maximum-Likelihood

- Suppose we have a set of  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that we wish to model with a GMM.

# Maximum-Likelihood

- Suppose we have a set of  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that we wish to model with a GMM.
- Consider this data set as an  $N \times d$  matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row is given by  $\mathbf{x}_n^T$ .

# Maximum-Likelihood

- Suppose we have a set of  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that we wish to model with a GMM.
- Consider this data set as an  $N \times d$  matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row is given by  $\mathbf{x}_n^{\text{T}}$ .
- Similarly, the corresponding latent variables define an  $N \times K$  matrix  $\mathbf{Z}$  with rows  $\mathbf{z}_n^{\text{T}}$ .

# Maximum-Likelihood

- Suppose we have a set of  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that we wish to model with a GMM.
- Consider this data set as an  $N \times d$  matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row is given by  $\mathbf{x}_n^{\text{T}}$ .
- Similarly, the corresponding latent variables define an  $N \times K$  matrix  $\mathbf{Z}$  with rows  $\mathbf{z}_n^{\text{T}}$ .
- The log-likelihood of the corresponding GMM is given by

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]. \quad (15)$$

# Maximum-Likelihood

- Suppose we have a set of  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that we wish to model with a GMM.
- Consider this data set as an  $N \times d$  matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row is given by  $\mathbf{x}_n^{\text{T}}$ .
- Similarly, the corresponding latent variables define an  $N \times K$  matrix  $\mathbf{Z}$  with rows  $\mathbf{z}_n^{\text{T}}$ .
- The log-likelihood of the corresponding GMM is given by

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]. \quad (15)$$

- Ultimately, we want to find the values of the parameters  $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$  that maximize this function.

- However, maximizing the log-likelihood terms for GMMs is much more complicated than for the case of a single Gaussian. Why?

- However, maximizing the log-likelihood terms for GMMs is much more complicated than for the case of a single Gaussian. Why?
- The difficulty arises from the sum over  $k$  inside of the log-term. The log function no longer acts directly on the Gaussian, and no closed-form solution is available.

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.
- Consider simply covariances defined by  $\Sigma_k = \sigma_k^2 \mathbf{I}$ .

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.
- Consider simply covariances defined by  $\Sigma_k = \sigma_k^2 \mathbf{I}$ .
- Suppose that one of the components of the mixture model,  $j$ , has its mean  $\mu_j$  exactly equal to one of the data points so that  $\mu_j = \mathbf{x}_n$  for some  $n$ .

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.
- Consider simply covariances defined by  $\Sigma_k = \sigma_k^2 \mathbf{I}$ .
- Suppose that one of the components of the mixture model,  $j$ , has its mean  $\mu_j$  exactly equal to one of the data points so that  $\mu_j = \mathbf{x}_n$  for some  $n$ .
- This term contributes

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{(1/2)} \sigma_j} \quad (16)$$

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.
- Consider simply covariances defined by  $\Sigma_k = \sigma_k^2 \mathbf{I}$ .
- Suppose that one of the components of the mixture model,  $j$ , has its mean  $\mu_j$  exactly equal to one of the data points so that  $\mu_j = \mathbf{x}_n$  for some  $n$ .
- This term contributes

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{(1/2)} \sigma_j} \quad (16)$$

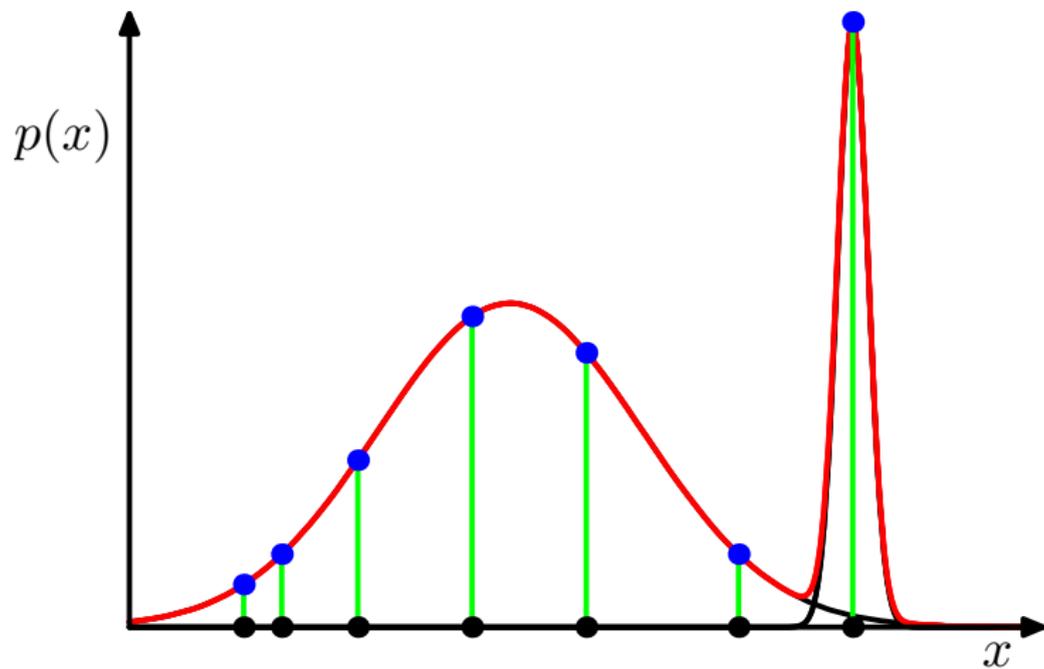
- Consider the limit  $\sigma_j \rightarrow 0$  to see that this term goes to infinity and hence the log-likelihood will also go to infinity.

# Singularities

- There is a significant problem when we apply MLE to estimate GMM parameters.
- Consider simply covariances defined by  $\Sigma_k = \sigma_k^2 \mathbf{I}$ .
- Suppose that one of the components of the mixture model,  $j$ , has its mean  $\mu_j$  exactly equal to one of the data points so that  $\mu_j = \mathbf{x}_n$  for some  $n$ .
- This term contributes

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{(1/2)} \sigma_j} \quad (16)$$

- Consider the limit  $\sigma_j \rightarrow 0$  to see that this term goes to infinity and hence the log-likelihood will also go to infinity.
- **Thus, the maximization of the log-likelihood function is not a well posed problem because such a singularity will occur whenever one of the components collapses to a single, specific data point.**



# Expectation-Maximization for GMMs

- **Expectation-Maximization** or EM is an elegant and powerful method for finding MLE solutions in the case of missing data such as the latent variables  $z$  indicating the mixture component.

# Expectation-Maximization for GMMs

- **Expectation-Maximization** or EM is an elegant and powerful method for finding MLE solutions in the case of missing data such as the latent variables  $\mathbf{z}$  indicating the mixture component.
- Recall the conditions that must be satisfied at a maximum of the likelihood function.

# Expectation-Maximization for GMMs

- **Expectation-Maximization** or EM is an elegant and powerful method for finding MLE solutions in the case of missing data such as the latent variables  $\mathbf{z}$  indicating the mixture component.
- Recall the conditions that must be satisfied at a maximum of the likelihood function.
- For the mean  $\boldsymbol{\mu}_k$ , setting the derivatives of  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  w.r.t.  $\boldsymbol{\mu}_k$  to zero yields

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (17)$$

$$= - \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (18)$$

# Expectation-Maximization for GMMs

- **Expectation-Maximization** or EM is an elegant and powerful method for finding MLE solutions in the case of missing data such as the latent variables  $\mathbf{z}$  indicating the mixture component.
- Recall the conditions that must be satisfied at a maximum of the likelihood function.
- For the mean  $\boldsymbol{\mu}_k$ , setting the derivatives of  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  w.r.t.  $\boldsymbol{\mu}_k$  to zero yields

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (17)$$

$$= - \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (18)$$

- Note the natural appearance of the responsibility terms on the RHS.

- Multiplying by  $\Sigma_k^{-1}$ , which we assume is non-singular, gives

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (19)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (20)$$

- Multiplying by  $\Sigma_k^{-1}$ , which we assume is non-singular, gives

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (19)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (20)$$

- We see the  $k^{\text{th}}$  mean is the weighted mean over all of the points in the dataset.

- Multiplying by  $\Sigma_k^{-1}$ , which we assume is non-singular, gives

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (19)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (20)$$

- We see the  $k^{\text{th}}$  mean is the weighted mean over all of the points in the dataset.
- Interpret  $N_k$  as the number of points assigned to component  $k$ .

- Multiplying by  $\Sigma_k^{-1}$ , which we assume is non-singular, gives

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (19)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (20)$$

- We see the  $k^{\text{th}}$  mean is the weighted mean over all of the points in the dataset.
- Interpret  $N_k$  as the number of points assigned to component  $k$ .
- We find a similar result for the covariance matrix:

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \boldsymbol{\mu}_k)(x_n - \boldsymbol{\mu}_k)^{\top} . \quad (21)$$

- We also need to maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mixing coefficients  $\pi_k$ .

- We also need to maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mixing coefficients  $\pi_k$ .
- Introduce a Lagrange multiplier to enforce the constraint  $\sum_k \pi_k = 1$ .

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) \quad (22)$$

- We also need to maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mixing coefficients  $\pi_k$ .
- Introduce a Lagrange multiplier to enforce the constraint  $\sum_k \pi_k = 1$ .

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) \quad (22)$$

- Maximizing it yields:

$$0 = \frac{1}{N_k} \sum_{n=1} \gamma(z_{nk}) + \lambda \quad (23)$$

- We also need to maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mixing coefficients  $\pi_k$ .
- Introduce a Lagrange multiplier to enforce the constraint  $\sum_k \pi_k = 1$ .

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) \quad (22)$$

- Maximizing it yields:

$$0 = \frac{1}{N_k} \sum_{n=1} \gamma(z_{nk}) + \lambda \quad (23)$$

- After multiplying both sides by  $\pi$  and summing over  $k$ , we get

$$\lambda = -N \quad (24)$$

- We also need to maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mixing coefficients  $\pi_k$ .
- Introduce a Lagrange multiplier to enforce the constraint  $\sum_k \pi_k = 1$ .

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) \quad (22)$$

- Maximizing it yields:

$$0 = \frac{1}{N_k} \sum_{n=1} \gamma(z_{nk}) + \lambda \quad (23)$$

- After multiplying both sides by  $\pi$  and summing over  $k$ , we get

$$\lambda = -N \quad (24)$$

- Eliminate  $\lambda$  and rearrange to obtain:

$$\pi_k = \frac{N_k}{N} \quad (25)$$

# Solved...right?

- So, we're done, right? We've computed the maximum likelihood solutions for each of the unknown parameters.

# Solved...right?

- So, we're done, right? We've computed the maximum likelihood solutions for each of the unknown parameters.
- Wrong!

# Solved...right?

- So, we're done, right? We've computed the maximum likelihood solutions for each of the unknown parameters.
- Wrong!
- The responsibility terms depend on these parameters in an intricate way:

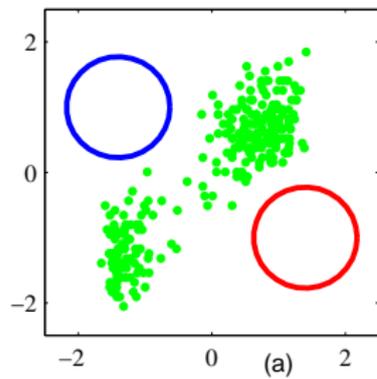
$$\gamma(z_k) \doteq p(z_k = 1 | \mathbf{x}) = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

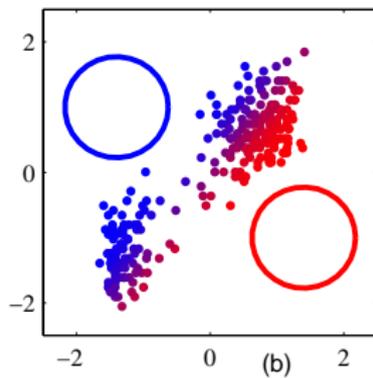
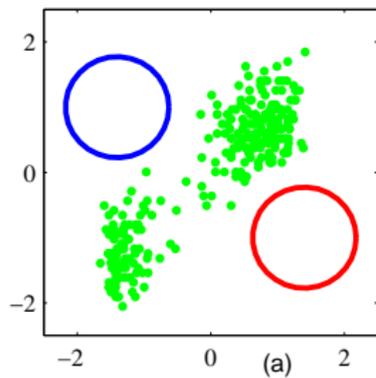
# Solved...right?

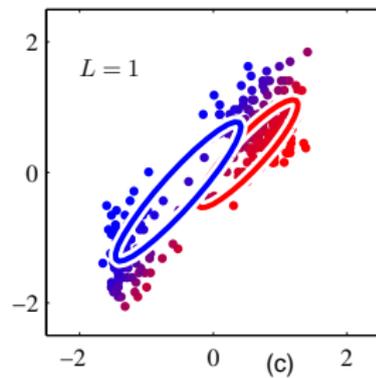
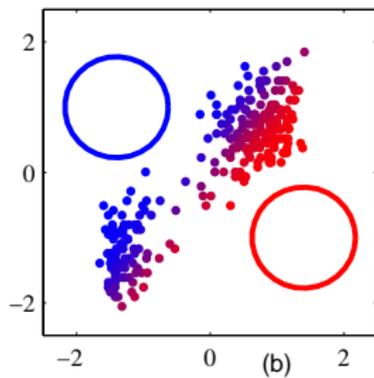
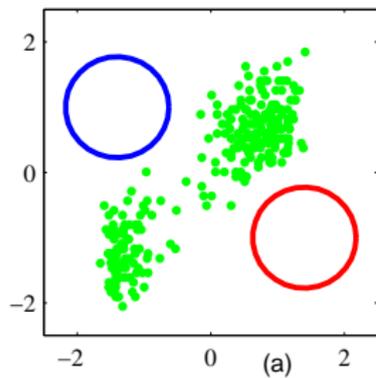
- So, we're done, right? We've computed the maximum likelihood solutions for each of the unknown parameters.
- Wrong!
- The responsibility terms depend on these parameters in an intricate way:

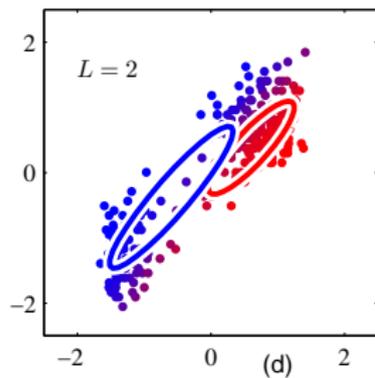
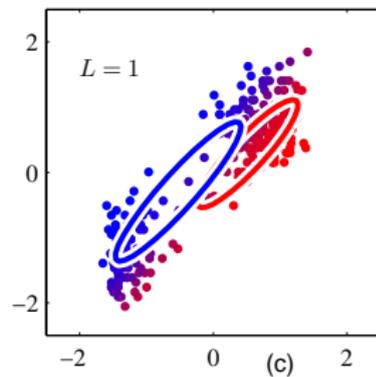
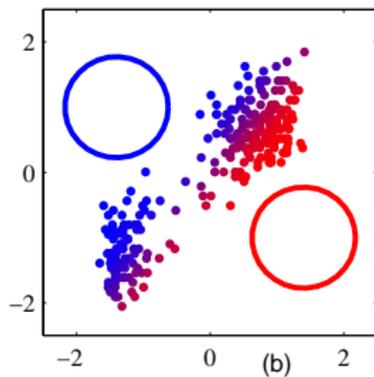
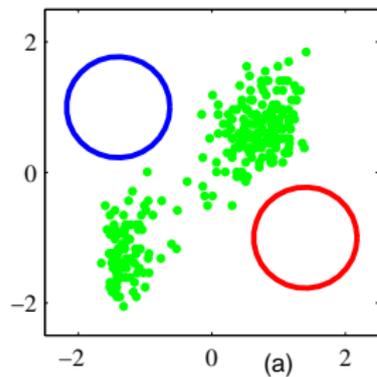
$$\gamma(z_k) \doteq p(z_k = 1 | \mathbf{x}) = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

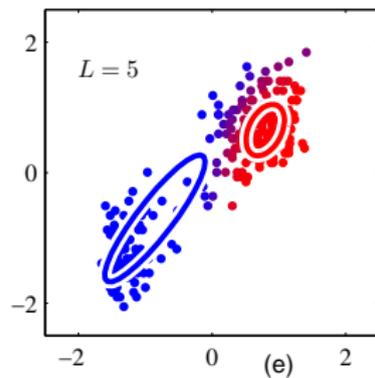
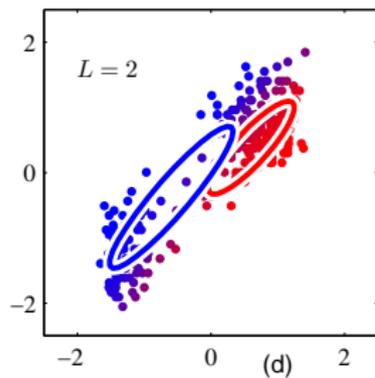
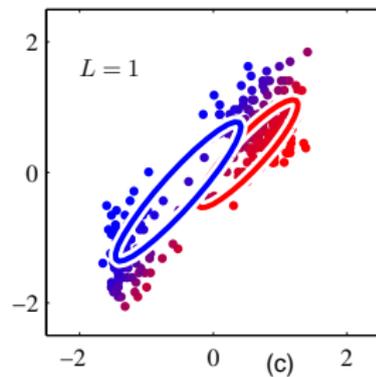
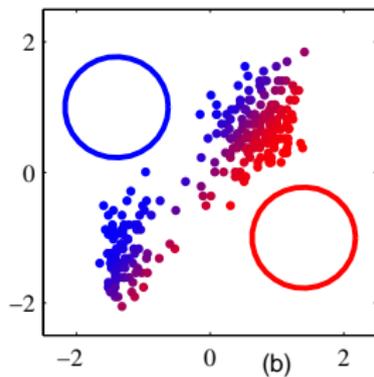
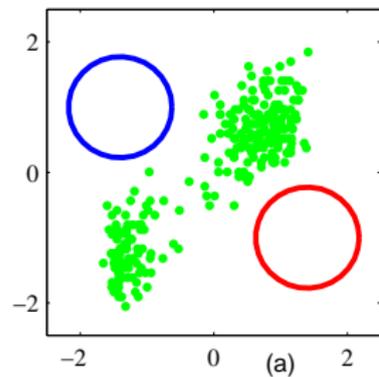
- But, these results do suggest an iterative scheme for finding a solution to the maximum likelihood problem.
  - 1 Choose some initial values for the parameters,  $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ .
  - 2 Use the current parameters estimates to compute the posteriors on the latent terms, i.e., the responsibilities.
  - 3 Use the responsibilities to update the estimates of the parameters.
  - 4 Repeat 2 and 3 until convergence.

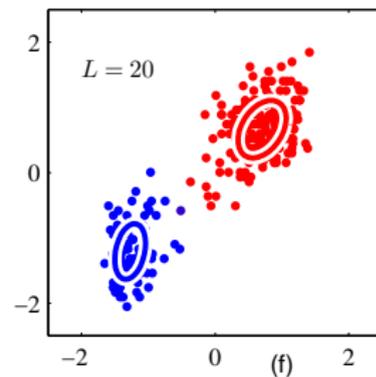
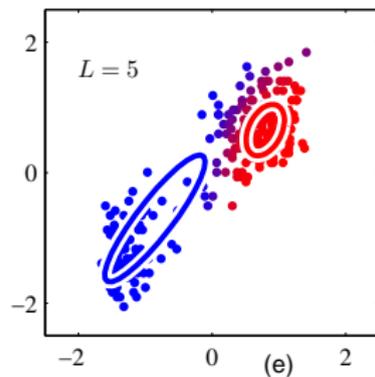
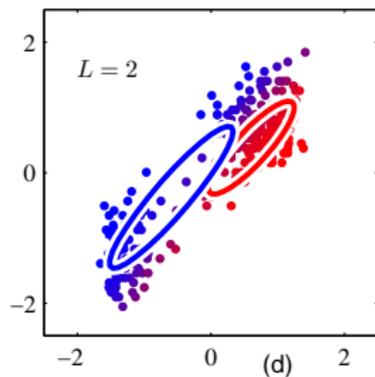
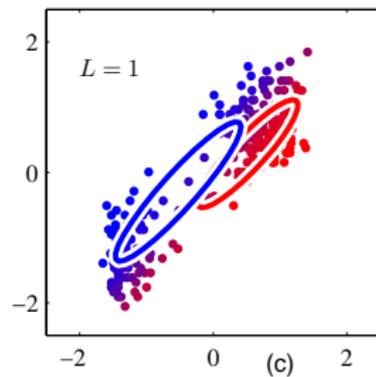
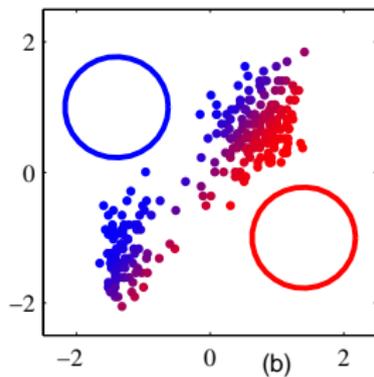
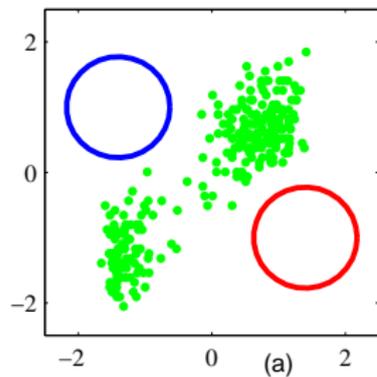












# Some Quick, Early Notes on EM

- EM generally tends to take more steps than the K-Means clustering algorithm.

# Some Quick, Early Notes on EM

- EM generally tends to take more steps than the K-Means clustering algorithm.
- Each step is more computationally intense than with K-Means too.

# Some Quick, Early Notes on EM

- EM generally tends to take more steps than the K-Means clustering algorithm.
- Each step is more computationally intense than with K-Means too.
- So, one commonly computes K-Means first and then initializes EM from the resulting clusters.

# Some Quick, Early Notes on EM

- EM generally tends to take more steps than the K-Means clustering algorithm.
- Each step is more computationally intense than with K-Means too.
- So, one commonly computes K-Means first and then initializes EM from the resulting clusters.
- Care must be taken to avoid singularities in the MLE solution.

# Some Quick, Early Notes on EM

- EM generally tends to take more steps than the K-Means clustering algorithm.
- Each step is more computationally intense than with K-Means too.
- So, one commonly computes K-Means first and then initializes EM from the resulting clusters.
- Care must be taken to avoid singularities in the MLE solution.
- There will generally be multiple local maxima of the likelihood function and EM is not guaranteed to find the largest of these.

Given a GMM, the goal is to maximize the likelihood function with respect to the parameters (the means, the covariances, and the mixing coefficients).

- 1 Initialize the means,  $\mu_k$ , the covariances,  $\Sigma_k$ , and mixing coefficients,  $\pi_k$ . Evaluate the initial value of the log-likelihood.

Given a GMM, the goal is to maximize the likelihood function with respect to the parameters (the means, the covariances, and the mixing coefficients).

- 1 Initialize the means,  $\boldsymbol{\mu}_k$ , the covariances,  $\boldsymbol{\Sigma}_k$ , and mixing coefficients,  $\pi_k$ . Evaluate the initial value of the log-likelihood.
- 2 **E-Step** Evaluate the responsibilities using the current parameter values:

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

Given a GMM, the goal is to maximize the likelihood function with respect to the parameters (the means, the covariances, and the mixing coefficients).

- 1 Initialize the means,  $\boldsymbol{\mu}_k$ , the covariances,  $\boldsymbol{\Sigma}_k$ , and mixing coefficients,  $\pi_k$ . Evaluate the initial value of the log-likelihood.
- 2 **E-Step** Evaluate the responsibilities using the current parameter values:

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

- 3 **M-Step** Update the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (26)$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^{\top} \quad (27)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (28)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (29)$$

## ④ Evaluate the log-likelihood

$$\ln p(\mathbf{X} | \boldsymbol{\mu}^{\text{new}}, \boldsymbol{\Sigma}^{\text{new}}, \boldsymbol{\pi}^{\text{new}}) = \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k^{\text{new}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{\text{new}}, \boldsymbol{\Sigma}_k^{\text{new}}) \right] \quad (30)$$

4 Evaluate the log-likelihood

$$\ln p(\mathbf{X} | \boldsymbol{\mu}^{\text{new}}, \boldsymbol{\Sigma}^{\text{new}}, \boldsymbol{\pi}^{\text{new}}) = \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k^{\text{new}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{\text{new}}, \boldsymbol{\Sigma}_k^{\text{new}}) \right] \quad (30)$$

5 Check for convergence of either the parameters of the log-likelihood. If the convergence is not satisfied, set the parameters:

$$\boldsymbol{\mu} = \boldsymbol{\mu}^{\text{new}} \quad (31)$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^{\text{new}} \quad (32)$$

$$\boldsymbol{\pi} = \boldsymbol{\pi}^{\text{new}} \quad (33)$$

and goto step 2.

# A More General View of EM

- The goal of EM is to find maximum likelihood solutions for models having latent variables.

# A More General View of EM

- The goal of EM is to find maximum likelihood solutions for models having latent variables.
- Denote the set of all model parameters as  $\theta$ , and so the log-likelihood function is

$$\ln p(\mathbf{X}|\theta) = \ln \left[ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right] \quad (34)$$

# A More General View of EM

- The goal of EM is to find maximum likelihood solutions for models having latent variables.
- Denote the set of all model parameters as  $\theta$ , and so the log-likelihood function is

$$\ln p(\mathbf{X}|\theta) = \ln \left[ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right] \quad (34)$$

- Note how the summation over the latent variables appears inside of the log.
  - Even if the joint distribution  $p(\mathbf{X}, \mathbf{Z}|\theta)$  belongs to the exponential family, the marginal  $p(\mathbf{X}|\theta)$  typically does not.

# A More General View of EM

- The goal of EM is to find maximum likelihood solutions for models having latent variables.
- Denote the set of all model parameters as  $\theta$ , and so the log-likelihood function is

$$\ln p(\mathbf{X}|\theta) = \ln \left[ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right] \quad (34)$$

- Note how the summation over the latent variables appears inside of the log.
  - Even if the joint distribution  $p(\mathbf{X}, \mathbf{Z}|\theta)$  belongs to the exponential family, the marginal  $p(\mathbf{X}|\theta)$  typically does not.
- If, for each sample  $\mathbf{x}_n$  we were given the value of the latent variable  $\mathbf{z}_n$ , then we would have a **complete** data set,  $\{\mathbf{X}, \mathbf{Z}\}$ , with which maximizing this likelihood term would be straightforward.

- However, in practice, we are not given the latent variables values.

- However, in practice, we are not given the latent variables values.
- So, instead, we focus on the expectation of the log-likelihood under the posterior distribution of the latent variables.

- However, in practice, we are not given the latent variables values.
- So, instead, we focus on the expectation of the log-likelihood under the posterior distribution of the latent variables.
- In the E-Step, we use the current parameter values  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ .

- However, in practice, we are not given the latent variables values.
- So, instead, we focus on the expectation of the log-likelihood under the posterior distribution of the latent variables.
- In the E-Step, we use the current parameter values  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ .
- This posterior is used to define the **expectation of the complete-data log-likelihood**, denoted  $Q(\theta, \theta^{\text{old}})$ , which is given by

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (35)$$

- However, in practice, we are not given the latent variables values.
- So, instead, we focus on the expectation of the log-likelihood under the posterior distribution of the latent variables.
- In the E-Step, we use the current parameter values  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ .
- This posterior is used to define the **expectation of the complete-data log-likelihood**, denoted  $Q(\theta, \theta^{\text{old}})$ , which is given by

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (35)$$

- Then, in the M-step, we revise the parameters to  $\theta^{\text{new}}$  by maximizing this function:

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad (36)$$

- However, in practice, we are not given the latent variables values.
- So, instead, we focus on the expectation of the log-likelihood under the posterior distribution of the latent variables.
- In the E-Step, we use the current parameter values  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ .
- This posterior is used to define the **expectation of the complete-data log-likelihood**, denoted  $Q(\theta, \theta^{\text{old}})$ , which is given by

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (35)$$

- Then, in the M-step, we revise the parameters to  $\theta^{\text{new}}$  by maximizing this function:

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad (36)$$

- Note that the log acts directly on the joint distribution  $p(\mathbf{X}, \mathbf{Z}|\theta)$  and so the M-step maximization will likely be tractable.

