

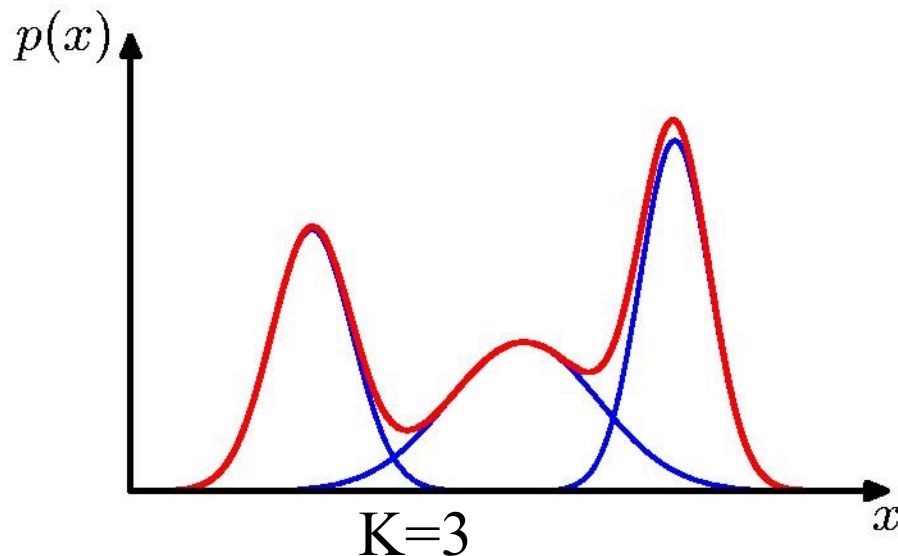
# Mixture of Gaussians

- We combine simple models into a complex model by defining a superposition of  $K$  Gaussian densities of the form:

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

Mixing coefficient

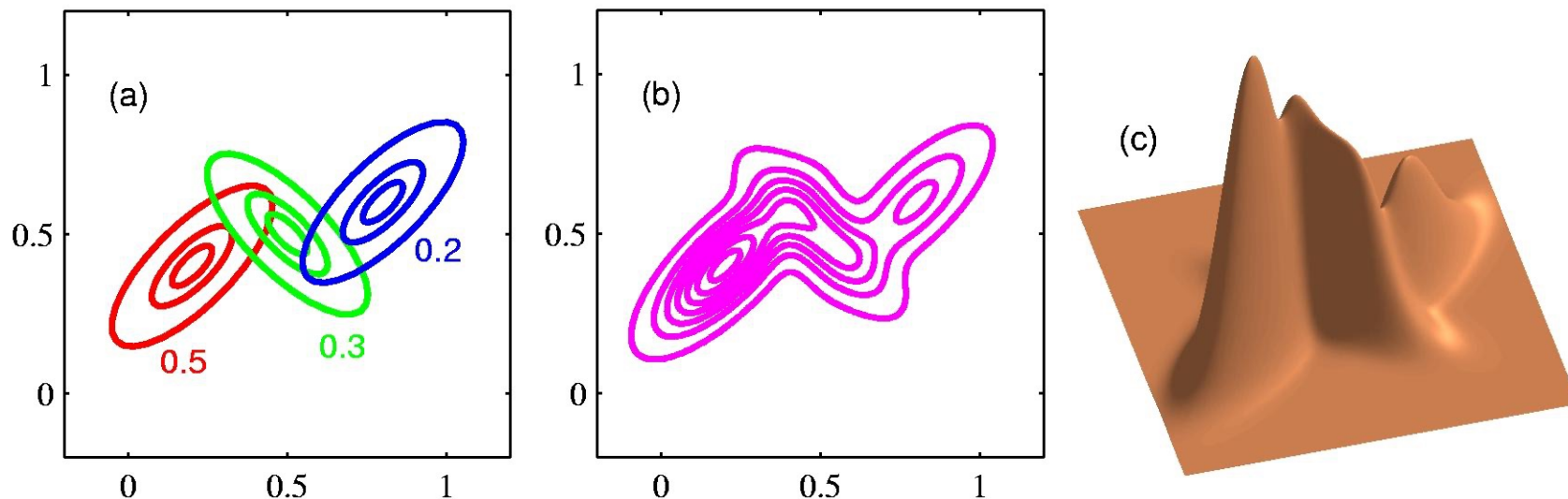
$$\forall k : \pi_k \geq 0 \quad \sum_{k=1}^K \pi_k = 1$$



- Note that each Gaussian component has its own mean and covariance. The parameters  $\pi_k$  are called mixing coefficients.
- More generally, mixture models can comprise linear combinations of other distributions.

# Mixture of Gaussians

- Illustration of a mixture of 3 Gaussians in a 2-dimensional space:



(a) Contours of constant density of each of the mixture components, along with the mixing coefficients

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

(c) A surface plot of the distribution  $p(\mathbf{x})$ .

# Mixture of Gaussians

- We will look at mixture of Gaussians in terms of **discrete latent variables**.
- The Gaussian mixture:

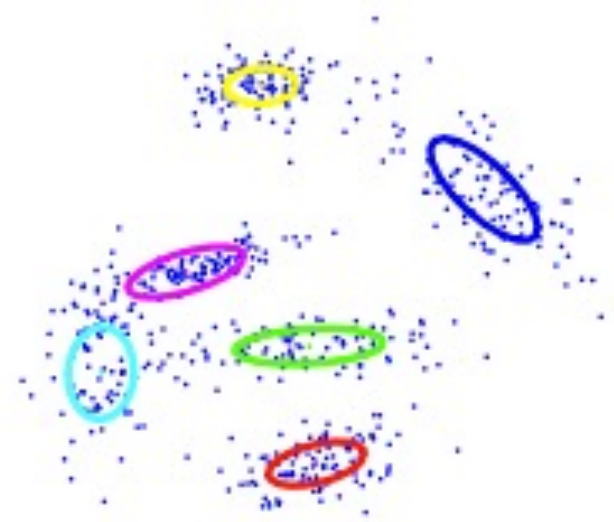
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

- K-dimensional **binary random variable**  $\mathbf{z}$  having a 1-of-K representation is the latent variable:

$$\mathbf{z} = [0, 0, \dots, 1, 0]^T \quad z_k \in \{0, 1\}, \quad \sum_k z_k = 1.$$

- We will specify the distribution over  $\mathbf{z}$  in terms of mixing coefficients:

$$p(z_k = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_k \pi_k = 1.$$



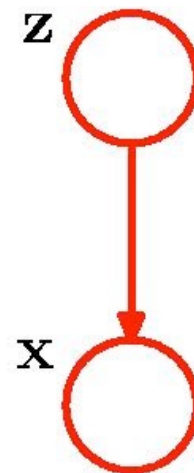
# Mixture of Gaussians

- Because  $\mathbf{z}$  uses **1-of-K encoding**, we have:

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

- We can now specify the conditional distribution:

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



# Mixture of Gaussians

- Because  $\mathbf{z}$  uses **1-of-K encoding**, we have:

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

- We can now specify the conditional distribution:

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

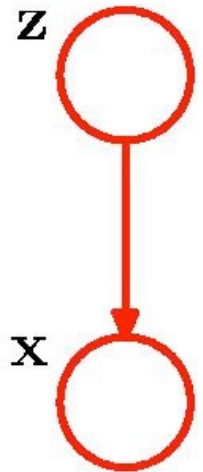
- We have therefore specified the joint distribution:

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

- The **marginal distribution** over  $\mathbf{x}$  is given by:

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

- The marginal distribution over  $\mathbf{x}$  is given by a **Gaussian mixture**.

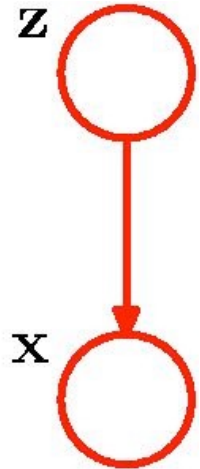


# Mixture of Gaussians

- The marginal distribution:

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

- If we have several observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , it follows that for every observed data point  $\mathbf{x}_n$ , there is a corresponding latent variable  $\mathbf{z}_n$ .



# Mixture of Gaussians

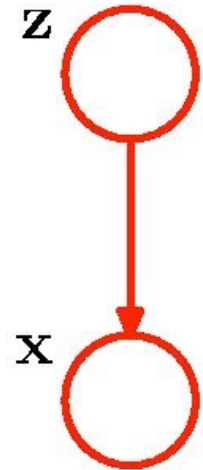
- The marginal distribution:

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

- If we have several observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , it follows that for **every observed data point**  $\mathbf{x}_n$ , there is a corresponding **latent variable**  $\mathbf{z}_n$ .
- Let us look at the conditional  $p(\mathbf{z} | \mathbf{x})$ , “responsibilities”, which we will need for doing inference:

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

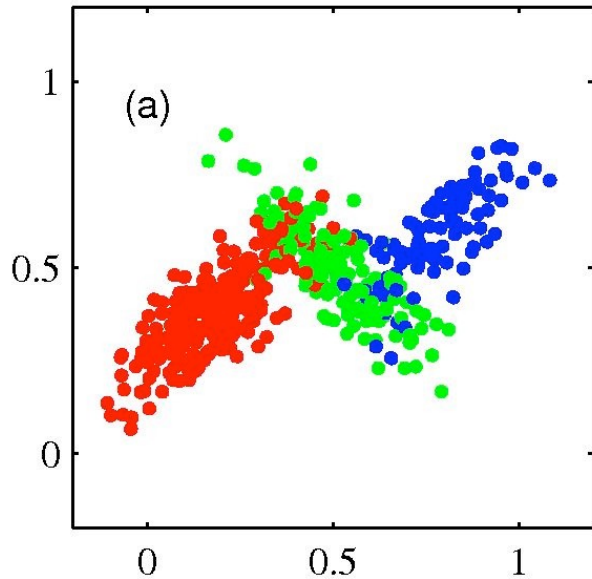
responsibility that component k takes for explaining the data x



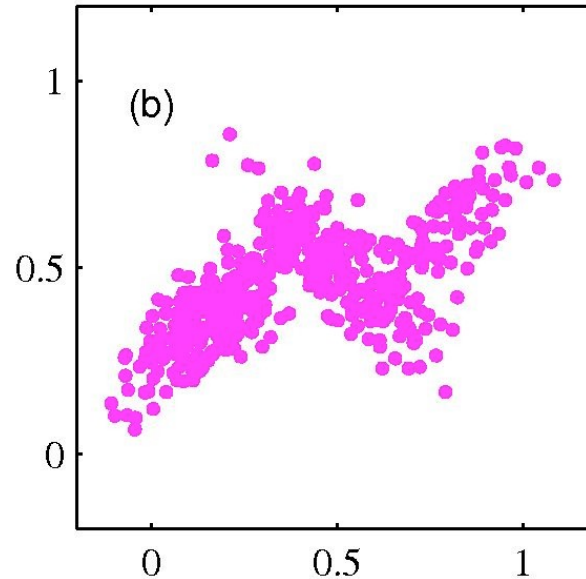
- We will view  $\pi_k$  as **prior probability** that  $z_k = 1$ , and  $\gamma(z_k)$  is the **corresponding posterior** once we have observed the data.

# Example

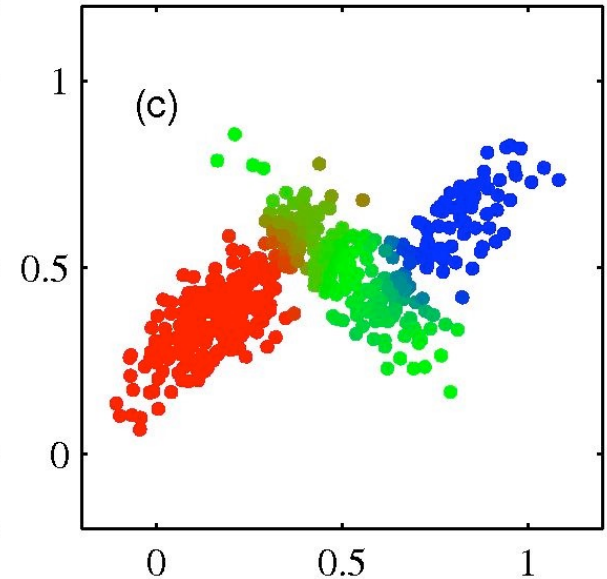
- 500 points drawn from a mixture of 3 Gaussians.



Samples from the **joint**  
**distribution**  $p(x,z)$ .



Samples from the **marginal**  
**distribution**  $p(x)$ .



Same samples where colors  
represent the value of  
responsibilities.



# Maximum Likelihood

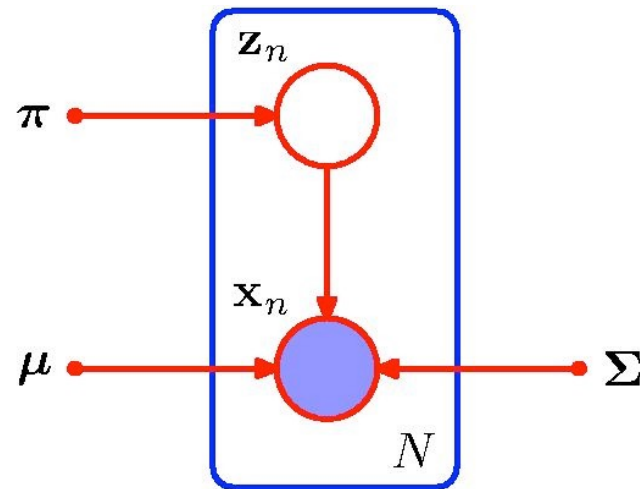
- Suppose we observe a dataset  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and we model the data using mixture of Gaussians.
- We represent the dataset as an  $N$  by  $D$  matrix  $\mathbf{X}$ .
- The corresponding **latent variables** will be represented as an  $N$  by  $K$  matrix  $\mathbf{Z}$ .

- The log-likelihood takes form:

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



Model parameters



Graphical model for a Gaussian mixture model for a set of i.i.d. data point  $\{\mathbf{x}_n\}$ , and corresponding latent variables  $\{\mathbf{z}_n\}$ .

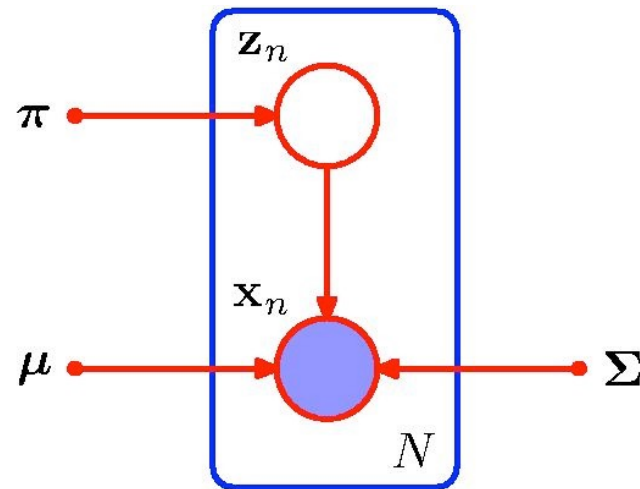
# Maximum Likelihood

- The log-likelihood takes form:

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

- Differentiating with respect to  $\boldsymbol{\mu}_k$  and setting to zero:

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



# Maximum Likelihood

- The log-likelihood takes form:

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

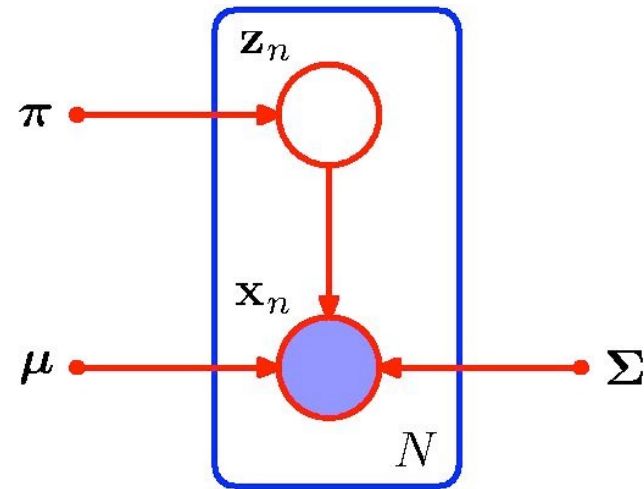
- Differentiating with respect to  $\boldsymbol{\mu}_k$  and setting to zero:

$$0 = \sum_n \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k).$$

$\gamma(z_{nk})$

Soft assignment

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_n \gamma(z_{nk}) \mathbf{x}_n, \quad N_k = \sum_n \gamma(z_{nk}).$$



- We can interpret  $N_k$  as **effective number of points** assigned to cluster  $k$ .
- The mean  $\boldsymbol{\mu}_k$  is given by the mean of all the data points **weighted by the posterior**  $\gamma(z_{nk})$  that component  $k$  was responsible for generating  $\mathbf{x}_n$ .

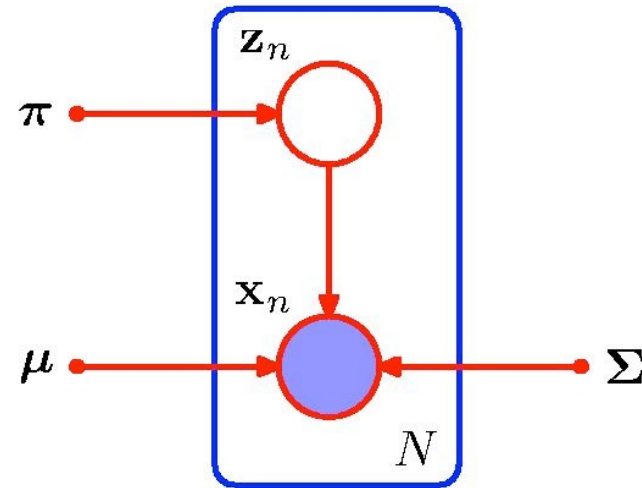
# Maximum Likelihood

- The log-likelihood takes form:

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

- Differentiating with respect to  $\boldsymbol{\Sigma}_k$  and setting to zero:

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



# Maximum Likelihood

- The log-likelihood takes form:

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

- Differentiating with respect to  $\boldsymbol{\Sigma}_k$  and setting to zero:

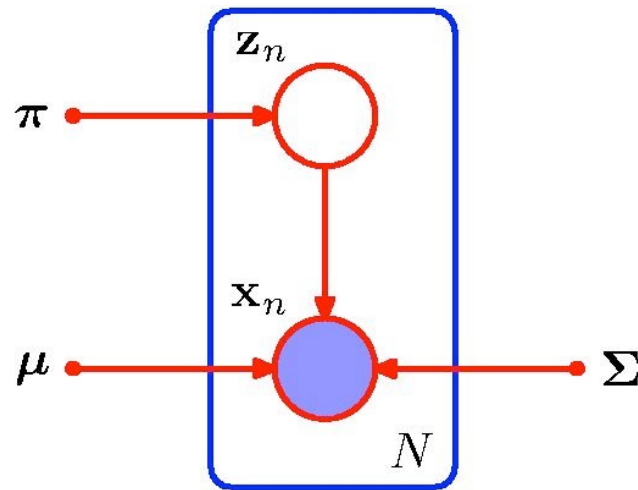
$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T.$$

- Note that the data points are **weighted by the posterior probabilities**.

- Maximizing log-likelihood with respect to mixing proportions:

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

- Mixing proportion for the  $k^{\text{th}}$  component is given by the **average responsibility which that component takes for explaining the data**.



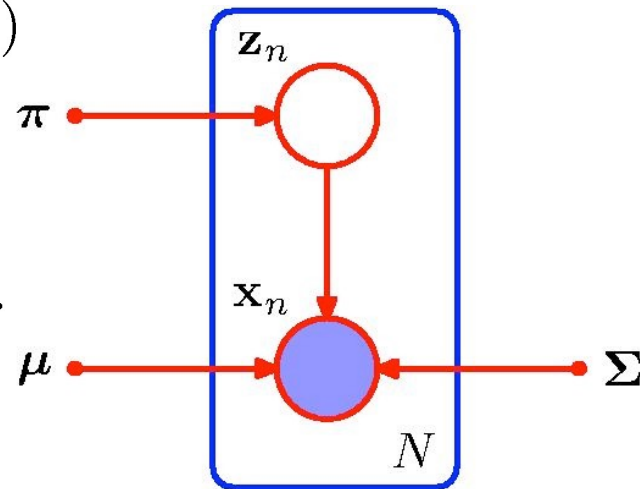
# Maximum Likelihood

- The log-likelihood takes form:

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

- Note that the maximum likelihood **does not have a closed form solution**.
- Parameter updates **depend on responsibilities**  $\gamma(z_{nk})$  which themselves depend on those parameters:

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



# Maximum Likelihood

- The log-likelihood:

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

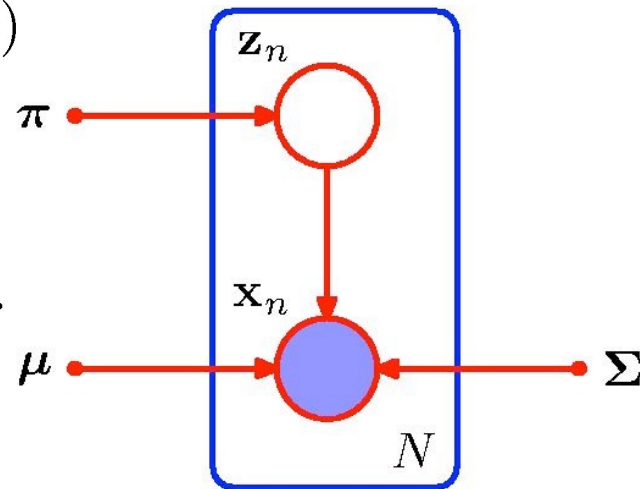
- Note that the maximum likelihood **does not have a closed form solution**.
- Parameter updates **depend on responsibilities**  $\gamma(z_{nk})$  which themselves depend on those parameters:

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

- Iterative Solution:

**E-step:** Update responsibilities  $\gamma(z_{nk})$

**M-step:** Update model parameters:  $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ , for  $k = 1, \dots, K$ .

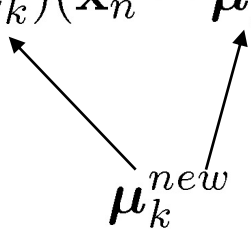


# EM algorithm

- Initialize the means  $\mu_k$ , covariances  $\Sigma_k$ , and mixing proportions  $\pi_k$
- **E-step**: Evaluate responsibilities using current parameter values:

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

- **M-step**: Re-estimate model parameters using the current responsibilities:

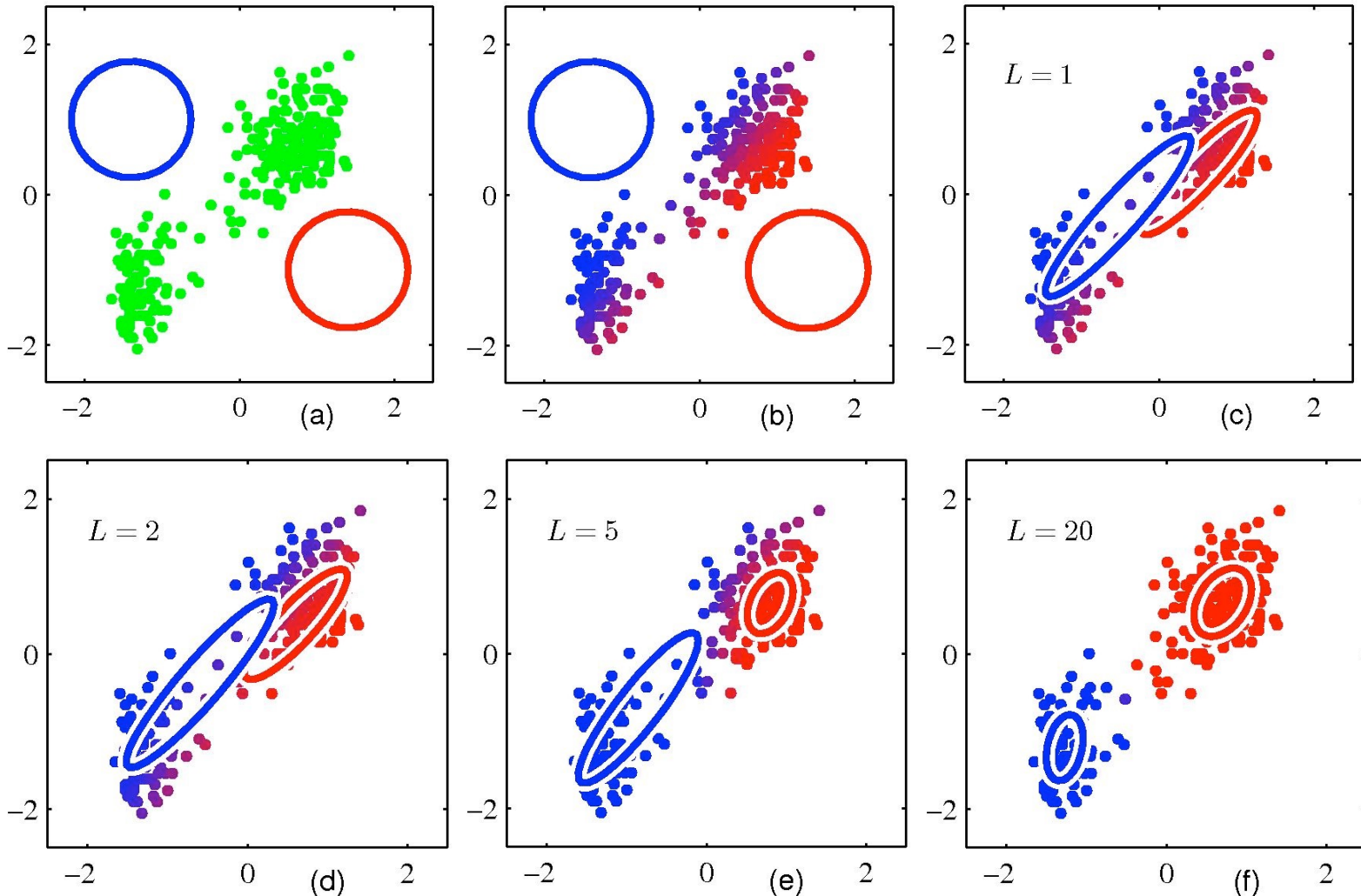
$$\begin{aligned}\mu_k^{new} &= \frac{1}{N_k} \sum_n \gamma(z_{nk}) \mathbf{x}_n, & N_k &= \sum_n \gamma(z_{nk}), \\ \Sigma_k^{new} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(y_{nk}) (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T, \\ \pi_k^{new} &= \frac{N_k}{N}.\end{aligned}$$


- Evaluate the log-likelihood and check for convergence.



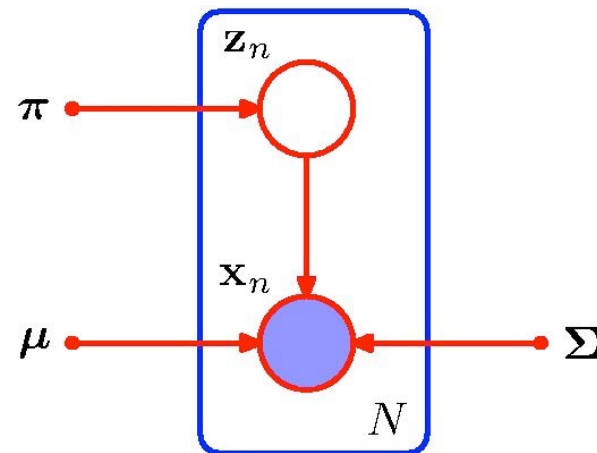
# Mixture of Gaussians: Example

- Illustration of the EM algorithm (much slower convergence compared to K-means)



# General EM formulation

- The goal of EM is to **find maximum likelihood solutions** for models with latent variables.
- We represent the **observed dataset** as an  $N$  by  $D$  matrix  $\mathbf{X}$ .
- **Latent variables** will be represented and an  $N$  by  $K$  matrix  $\mathbf{Z}$ .
- The set of all **model parameters** is denoted by  $\theta$ .



# General EM formulation

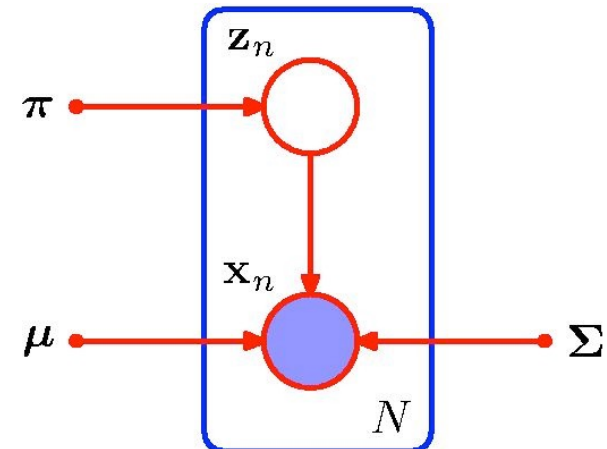
- The goal of EM is to **find maximum likelihood solutions** for models with latent variables.
- We represent the **observed dataset** as an N by D matrix  $\mathbf{X}$ .
- **Latent variables** will be represented and an N by K matrix  $\mathbf{Z}$ .
- The set of all **model parameters** is denoted by  $\theta$ .
- The log-likelihood takes form:

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

- We will call:

$\{\mathbf{X}, \mathbf{Z}\}$  as **complete** dataset.

$\{\mathbf{X}\}$  as **incomplete** dataset.



# General EM formulation

- In practice, we are **not given a complete dataset**  $\{\mathbf{X}, \mathbf{Z}\}$ , but only incomplete dataset  $\{\mathbf{X}\}$ .
- Our knowledge about the latent variables is given only by **the posterior distribution**  $p(\mathbf{Z}|\mathbf{X}, \theta)$ .
- Because we cannot use the complete data log-likelihood, we can consider **expected complete-data log-likelihood**:

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

# General EM formulation

- In practice, we are **not given a complete dataset**  $\{\mathbf{X}, \mathbf{Z}\}$ , but only incomplete dataset  $\{\mathbf{X}\}$ .
- Our knowledge about the latent variables is given only by **the posterior distribution**  $p(\mathbf{Z}|\mathbf{X}, \theta)$ .
- Because we cannot use the complete data log-likelihood, we can consider **expected complete-data log-likelihood**:

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

- In the E-step, we use the current parameters  $\theta^{old}$  to compute **the posterior over the latent variables**  $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ .
- We use this posterior to compute expected complete log-likelihood.
- In the M-step, we find the revised parameter estimate  $\theta^{new}$  by **maximizing the expected complete log-likelihood**:

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

← Tractable

# The General EM algorithm

- Given a joint distribution  $p(\mathbf{Z}, \mathbf{X} | \theta)$  over observed and latent variables, the goal is to maximize the likelihood function  $p(\mathbf{X} | \theta)$  with respect to  $\theta$ .
- Initialize parameters  $\theta^{\text{old}}$ .
- **E-step:** Compute posterior over latent variables:  $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$  and  $Q(\theta, \theta^{\text{old}})$ .
- **M-step:** Find the new estimate of parameters  $\theta^{\text{new}}$ :

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

where

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

# The General EM algorithm

- Given a joint distribution  $p(\mathbf{Z}, \mathbf{X} | \theta)$  over observed and latent variables, the goal is to maximize the likelihood function  $p(\mathbf{X} | \theta)$  with respect to  $\theta$ .
- Initialize parameters  $\theta^{\text{old}}$ .
- **E-step:** Compute posterior over latent variables:  $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$  and  $Q(\theta, \theta^{\text{old}})$ .
- **M-step:** Find the new estimate of parameters  $\theta^{\text{new}}$ :

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

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 the parameter values.

Otherwise:

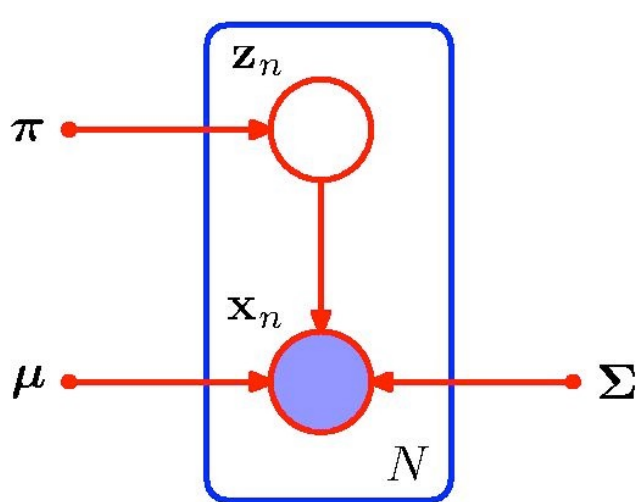
$$\theta^{\text{new}} \leftarrow \theta^{\text{old}}, \quad \text{and iterate.}$$

# Gaussian Mixtures Revisited

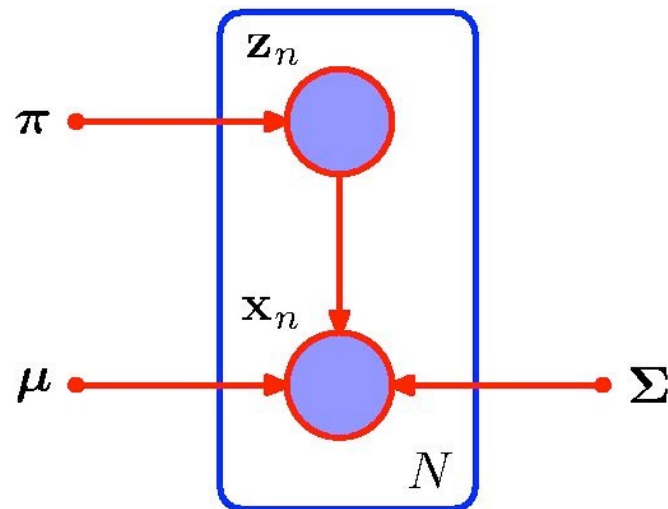
- We now consider the application of the latent variable view of EM to the case of **Gaussian mixture model**.

- Recall:

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



$\{\mathbf{X}\}$  -- incomplete dataset.



$\{\mathbf{X}, \mathbf{Z}\}$  -- complete dataset.



# Maximizing Complete Data

- Consider the problem of maximizing the likelihood for the complete data:

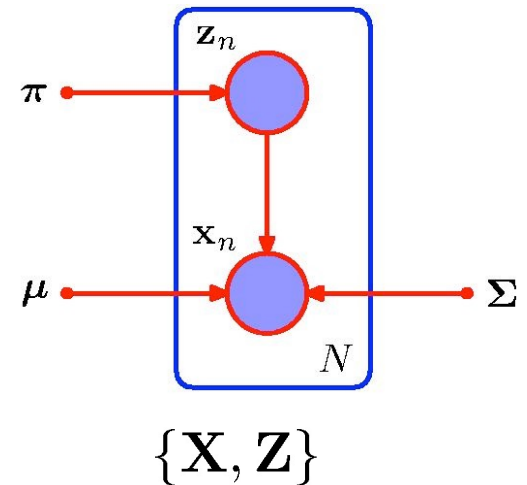
$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \prod_{k=1}^K \left[ \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]^{z_{nk}}.$$

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

Sum of K independent contributions, one for each mixture component.

- Maximizing with respect to **mixing proportions** yields:

$$\pi_k = \frac{1}{N} \sum_{n=1}^N z_{nk}.$$



-- complete dataset.

# Posterior Over Latent Variables

- Remember:

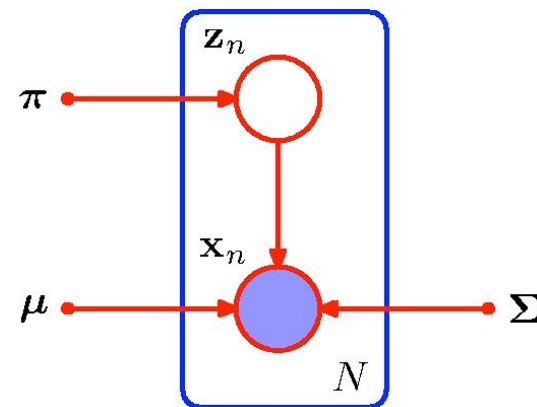
$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}, \quad p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}.$$

- The **posterior over latent variables** takes form:

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


- Note that **the posterior factorizes over n points**, so that under the posterior distribution  $\{\mathbf{z}_n\}$  are independent.

$$p(z_{nk} = 1|\mathbf{X}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\pi_k \mathcal{N}(x_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$



# Expected Complete Log-Likelihood


- The expected value of indicator variable  $z_{nk}$  under the posterior distribution is:

Under posterior 

$$\begin{aligned}\mathbb{E}[z_{nk}] &= \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)]^{z_{nj}}} \\ &= \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \gamma(z_{nk}).\end{aligned}$$

# Expected Complete Log-Likelihood

- The expected value of indicator variable  $z_{nk}$  under the posterior distribution is:

Under posterior 

$$\begin{aligned}\mathbb{E}[z_{nk}] &= \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)]^{z_{nj}}} \\ &= \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \gamma(z_{nk}).\end{aligned}$$

- This represent **the responsibility** of component  $k$  for data point  $\mathbf{x}_n$ .
- The **complete-data log-likelihood**:

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

- The **expected complete data log-likelihood** is:

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

# Expected Complete Log-Likelihood

- The expected complete data log-likelihood is:

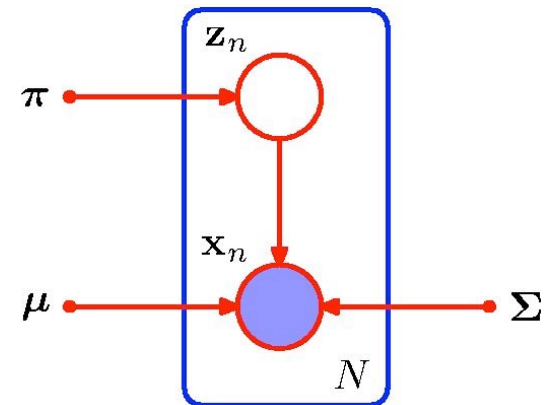
$$\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right].$$

- Maximizing with respect to model parameters we obtain:

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_n \gamma(z_{nk}) \mathbf{x}_n, \quad N_k = \sum_n \gamma(z_{nk}),$$

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

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



# Relationship to K-Means

- Consider a Gaussian mixture model in which **covariances are shared** and are given by  $\epsilon I$

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

- Consider EM algorithm for a mixture of K Gaussians, in which **we treat  $\boldsymbol{\Sigma}_k$  as a fixed constant**. The **posterior responsibilities** take form:

$$\gamma(z_{nk}) = \frac{\pi_k \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2/2\epsilon)}{\sum_{j=1}^K \pi_j \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2/2\epsilon)}.$$

# Relationship to K-Means

- Consider a Gaussian mixture model in which **covariances are shared** and are given by  $\epsilon I$

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

- Consider EM algorithm for a mixture of K Gaussians, in which **we treat  $\boldsymbol{\Sigma}_k$  as a fixed constant**. The **posterior responsibilities** take form:

$$\gamma(z_{nk}) = \frac{\pi_k \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2/2\epsilon)}{\sum_{j=1}^K \pi_j \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2/2\epsilon)}.$$

- Consider the limit  $\epsilon$  goes to 0.
- In the denominator, the term for which  $\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2$  is smallest will go to zero **most slowly**. Hence  $\gamma(z_{nk})$  goes to  $r_{nk}$  where

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

# Relationship to K-Means

- Consider EM algorithm for a mixture of  $K$  Gaussians, in which we treat  $\epsilon$  as a fixed constant. The posterior responsibilities take form:

$$\gamma(z_{nk}) = \frac{\pi_k \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 / 2\epsilon)}{\sum_{j=1}^K \pi_j \exp(-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon)}.$$

- Finally, in the limit  $\epsilon$  goes to 0, the expected complete log-likelihood becomes:

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

- Hence in the limit, maximizing the expected complete log-likelihood is equivalent to minimizing the distortion measure in the K-means algorithm.
- This is why assignment step in k-Means is called E-step, and finding the cluster centers is called the M-step.