



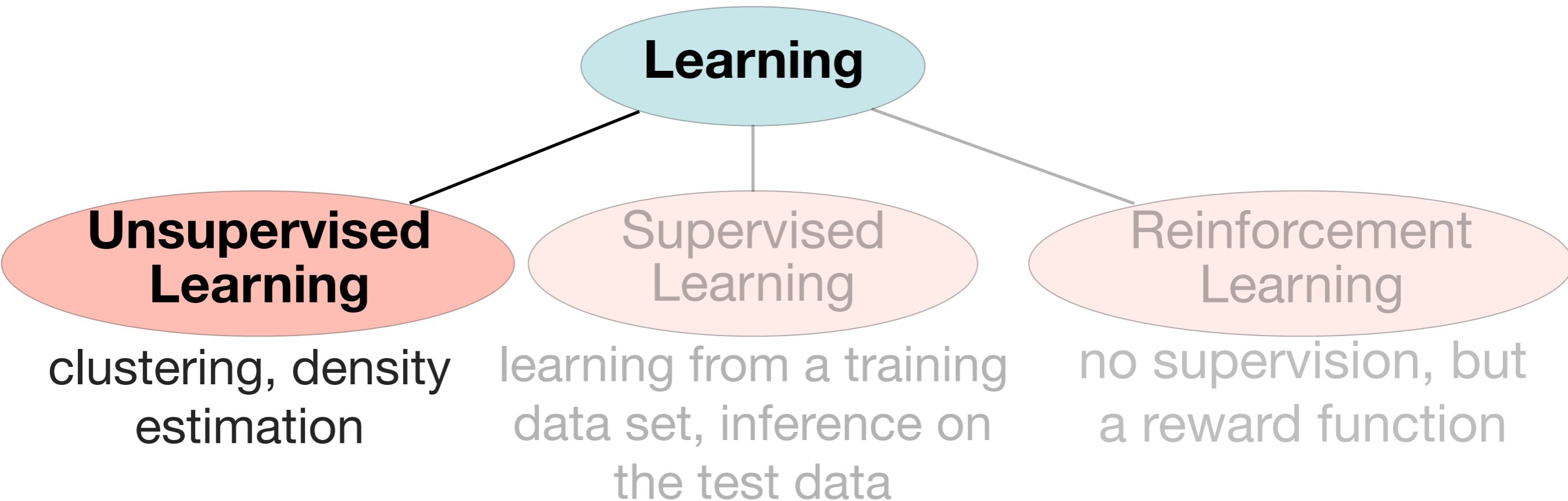
# 10. Clustering

# Motivation

- Supervised learning is good for interaction with humans, but labels from a supervisor are sometimes hard to obtain
- Clustering is **unsupervised** learning, i.e. it tries to learn only from the data
- Main idea: find a similarity measure and group similar data objects together
- Clustering is a very old research field, many approaches have been suggested
- Main problem in most methods: how to find a good number of clusters



# Categories of Learning



In unsupervised learning, there is no **ground truth** information given.

Most Unsupervised Learning methods are based on **Clustering**.



# K-means Clustering

- Given: data set  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , number of clusters  $K$
- Goal: find cluster centers  $\{\mu_1, \dots, \mu_K\}$  so that

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

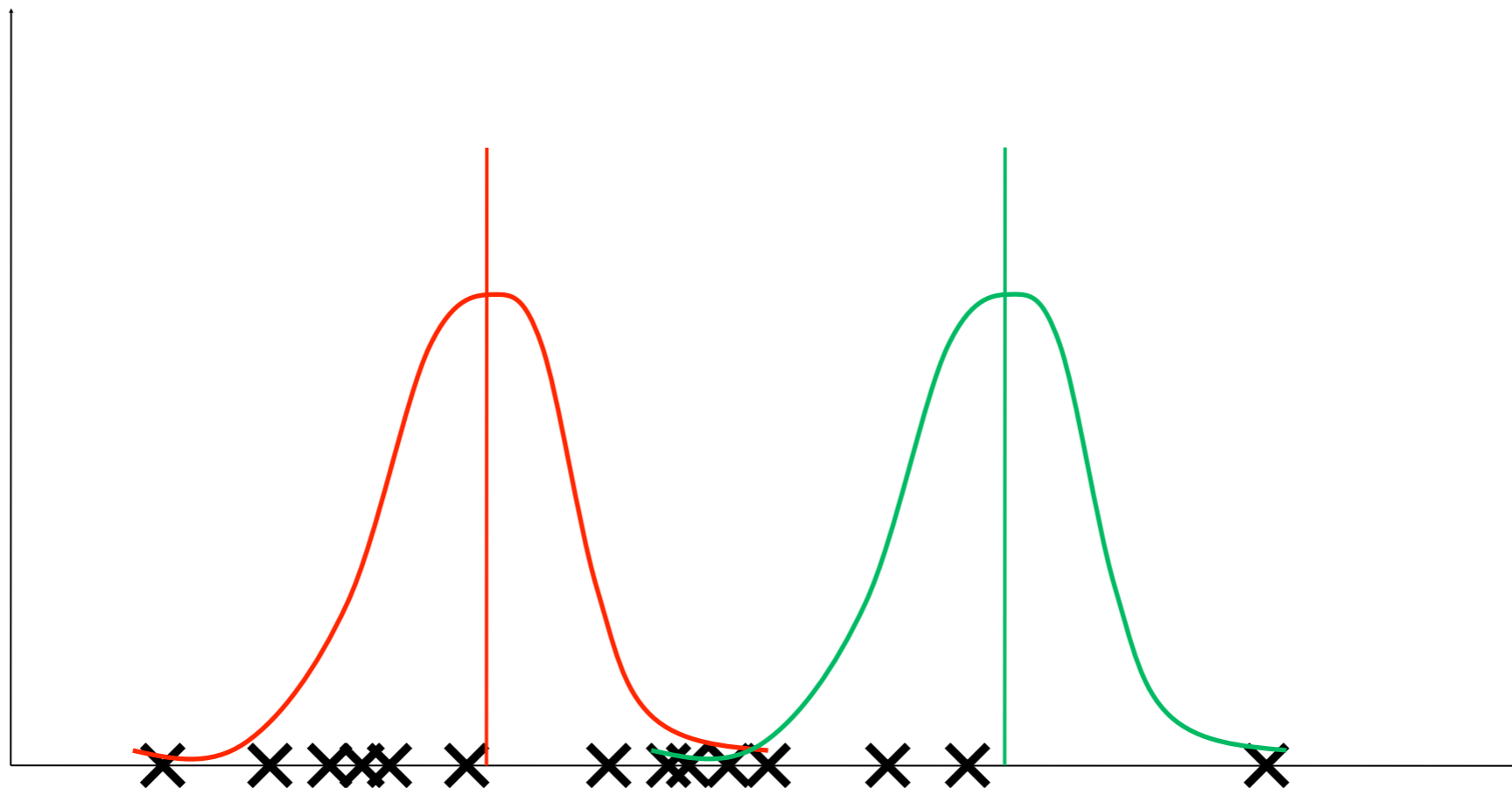
is minimal, where  $r_{nk} = 1$  if  $\mathbf{x}_n$  is assigned to  $\mu_k$

- Idea: compute  $r_{nk}$  and  $\mu_k$  iteratively
- Start with some values for the cluster centers
- Find optimal assignments  $r_{nk}$
- Update cluster centers using these assignments
- Repeat until assignments or centers don't change



# K-means Clustering

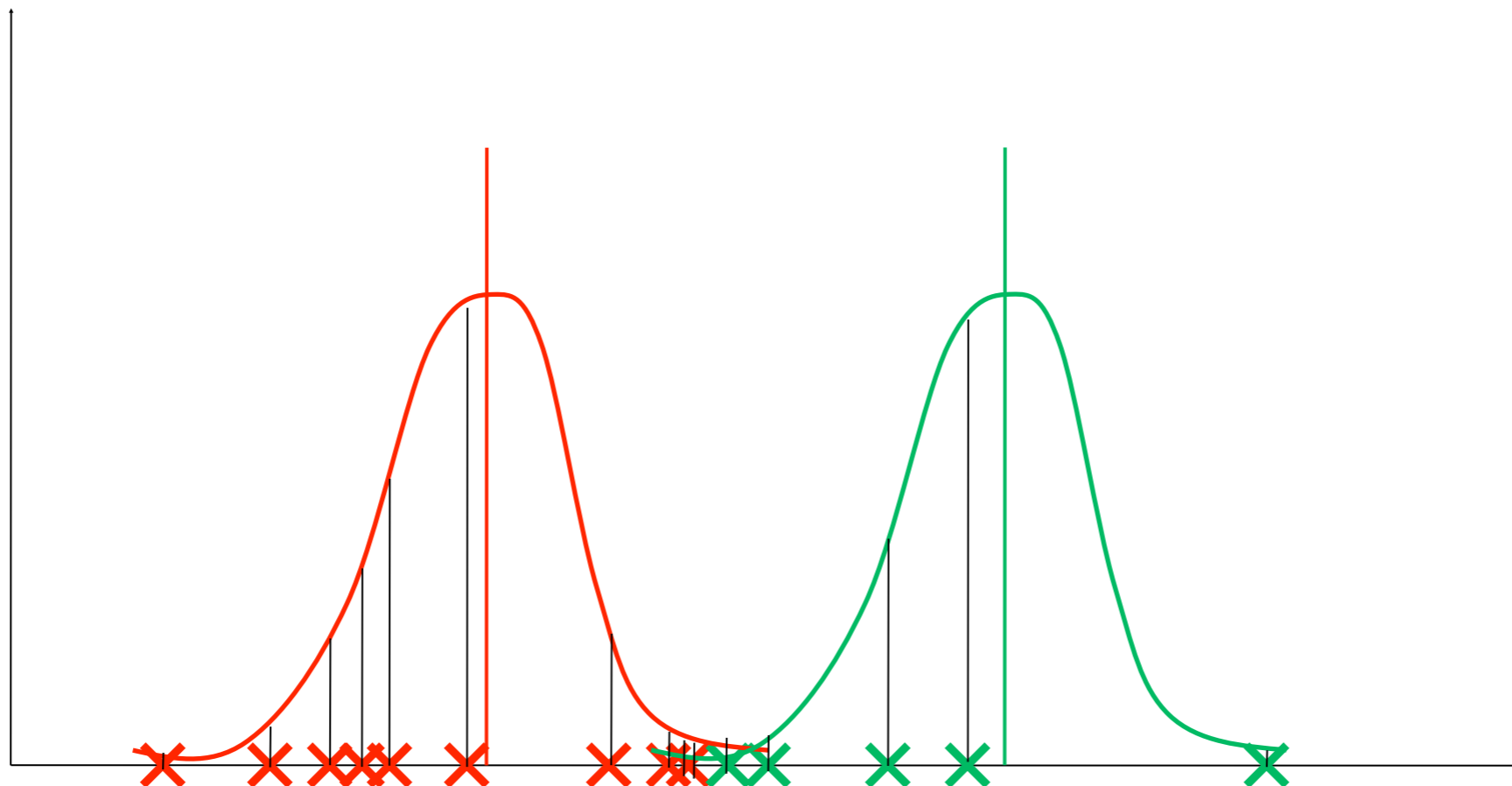
Initialize cluster means:  $\{\mu_1, \dots, \mu_K\}$



# K-means Clustering

Find optimal assignments:

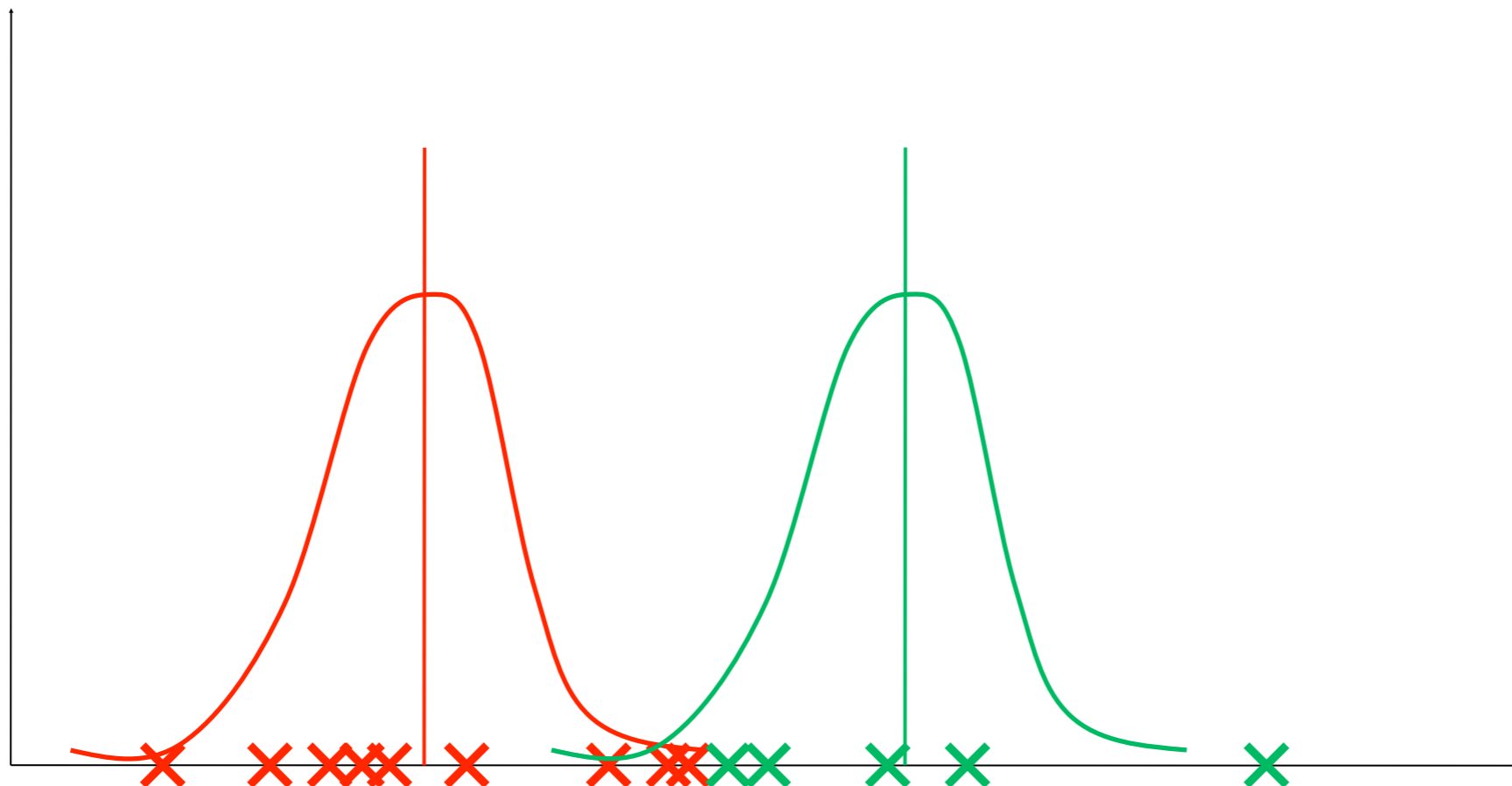
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\| \\ 0 & \text{otherwise} \end{cases}$$



# K-means Clustering

Find new optimal means:  $\frac{\partial J}{\partial \mu_k} = 2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \mu_k) \stackrel{!}{=} 0$

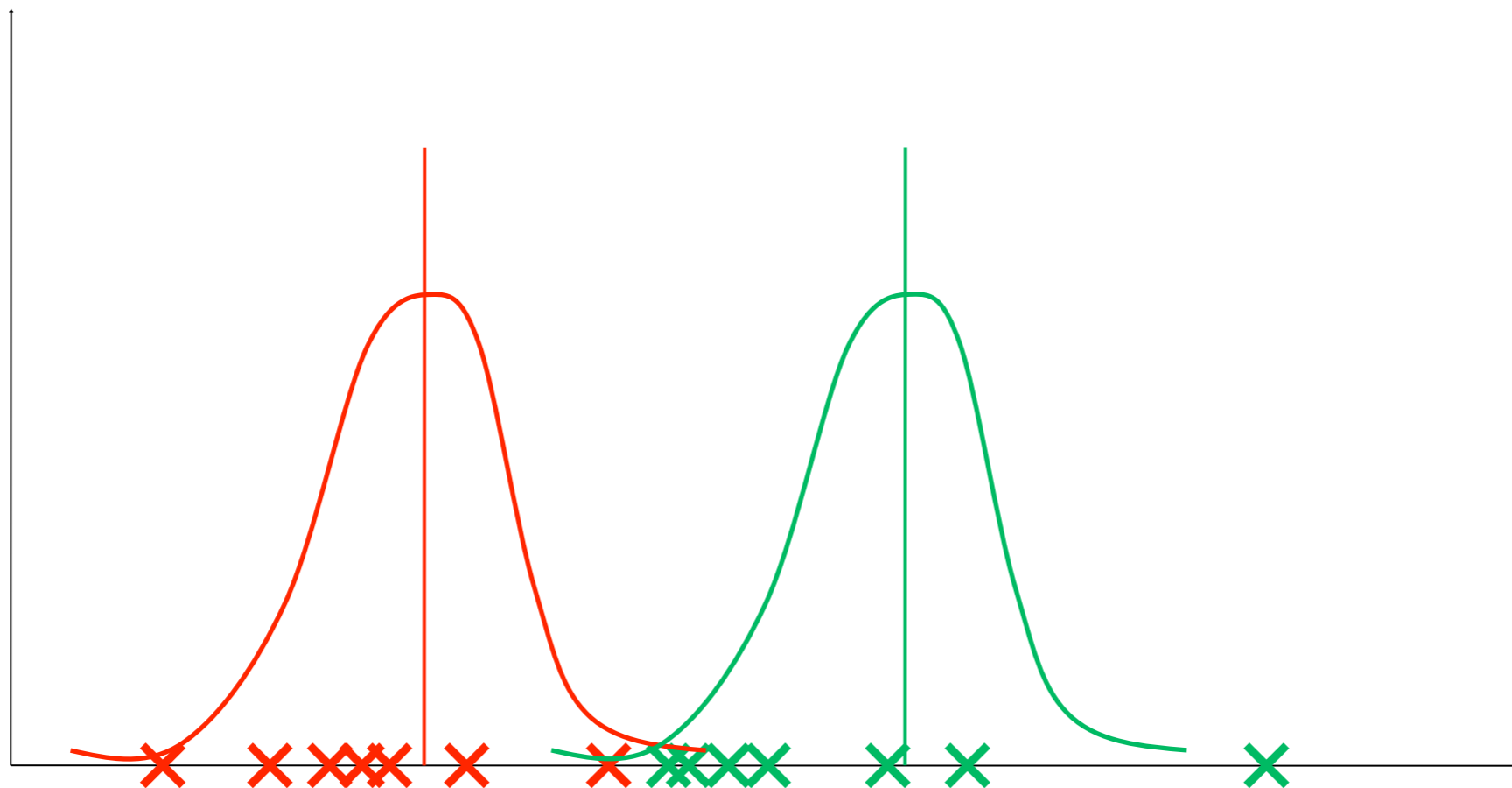
$$\Rightarrow \mu_k = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}}$$



# K-means Clustering

Find new optimal assignments:

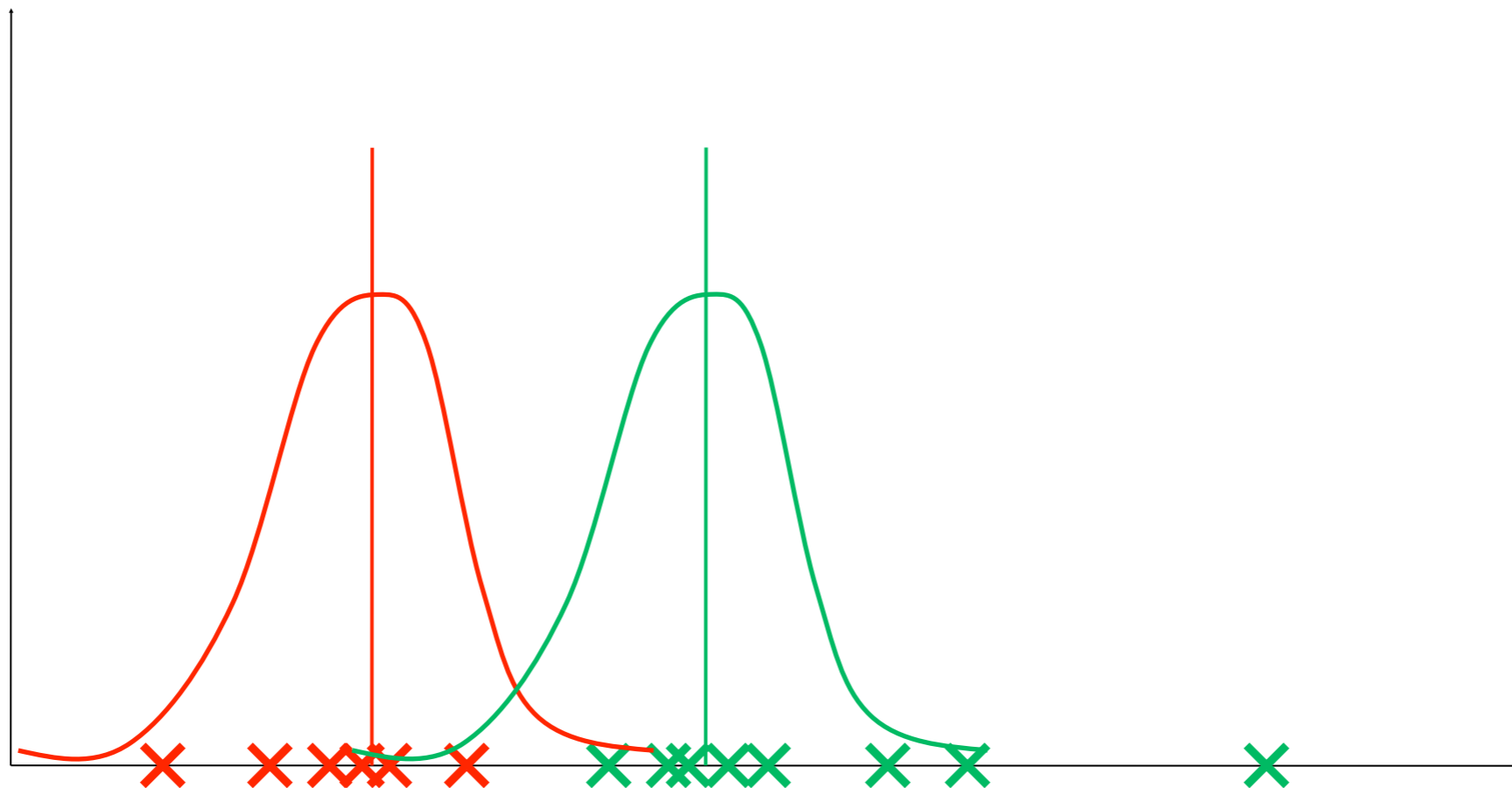
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\| \\ 0 & \text{otherwise} \end{cases}$$



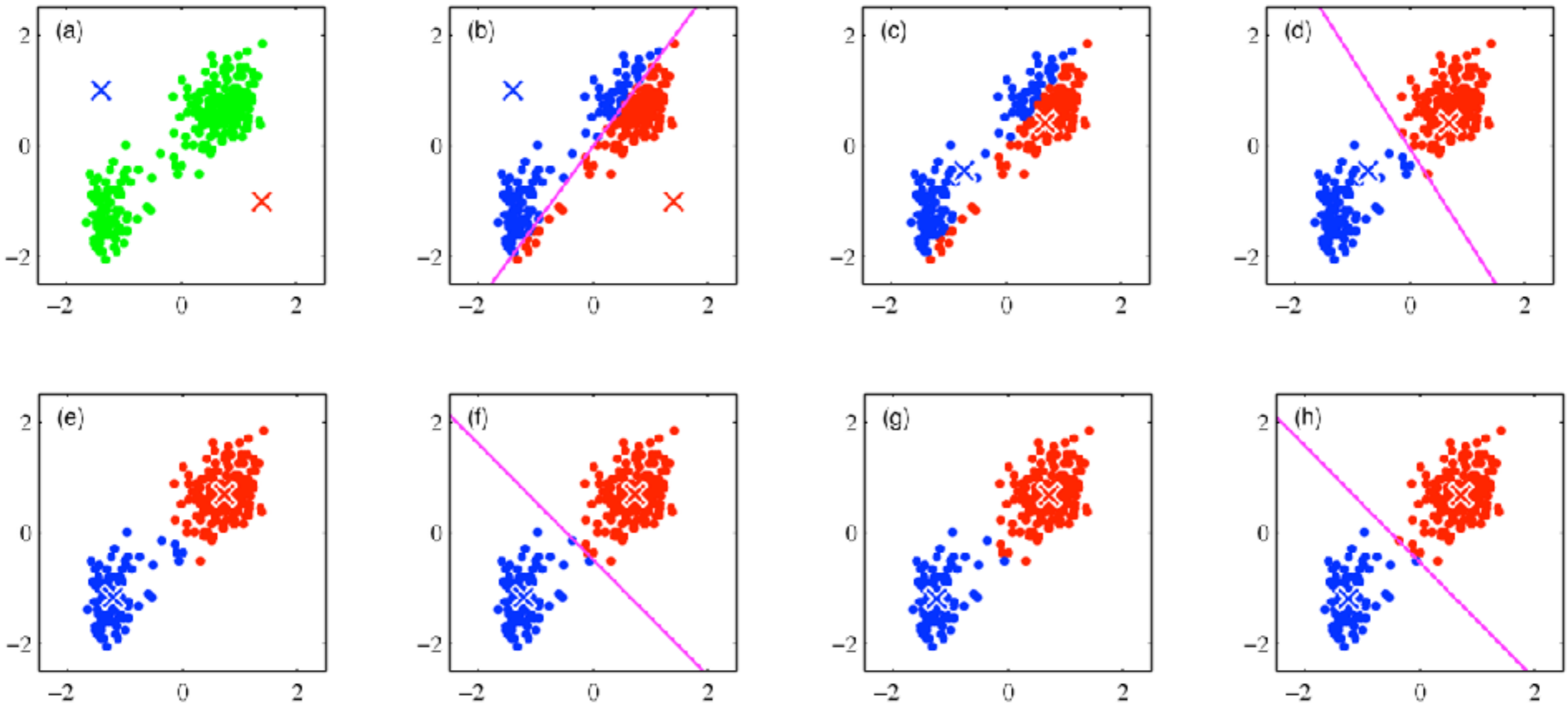


# K-means Clustering

Iterate these steps until means and assignments do not change any more



# 2D Example

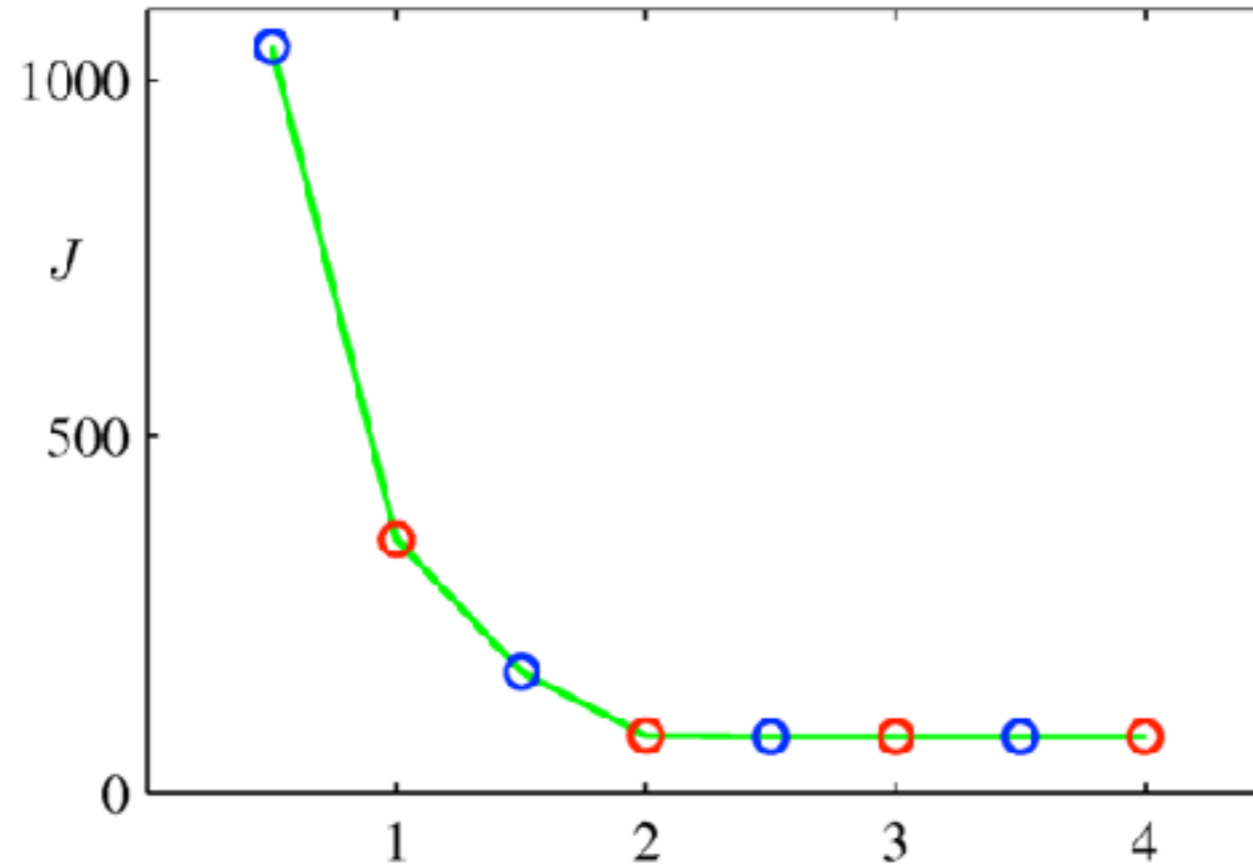


- Real data set
- Random initialization

- Magenta line is “decision boundary”



# The Cost Function



- After every step the cost function  $J$  is minimized
- Blue steps: update assignments
- Red steps: update means
- Convergence after 4 rounds



# K-means for Segmentation

$K = 2$



$K = 3$



$K = 10$



Original image



# K-Means: Additional Remarks

- K-means converges always, but the minimum is not guaranteed to be a global one
- There is an **online** version of  $K$ -means
  - After each addition of  $\mathbf{x}_n$ , the nearest center  $\boldsymbol{\mu}_k$  is updated:
- The  **$K$ -medoid** variant:
  - Replace the Euclidean distance by a general measure  $V$ .

$$\boldsymbol{\mu}_k^{\text{new}} = \boldsymbol{\mu}_k^{\text{old}} + \eta_n (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{old}})$$

$$\tilde{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$



# Mixtures of Gaussians

- Assume that the data consists of  $K$  clusters
- The data within each cluster is Gaussian
- For any data point  $\mathbf{x}$  we introduce a  $K$ -dimensional binary random variable  $\mathbf{z}$  so that:

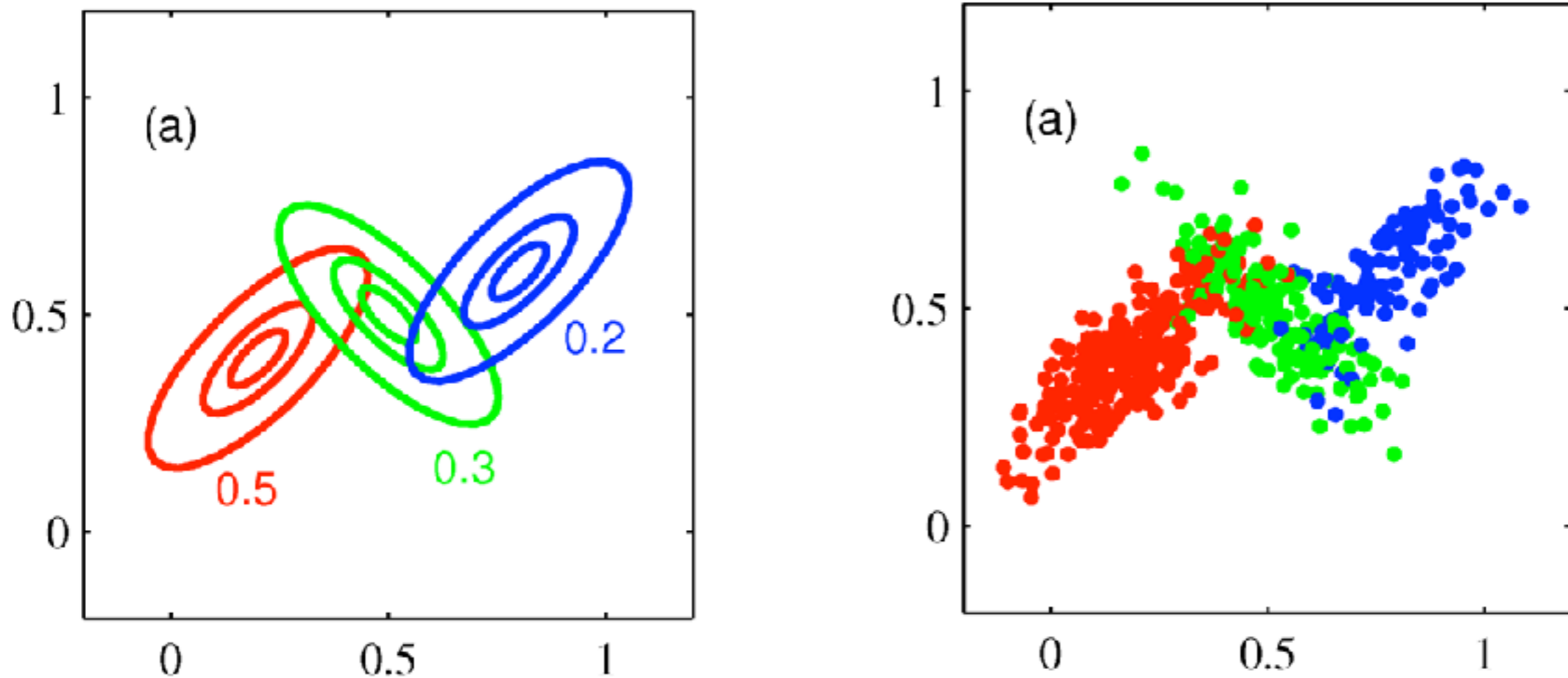
$$p(\mathbf{x}) = \sum_{k=1}^K \underbrace{p(z_k = 1)}_{=: \pi_k} \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where

$$z_k \in \{0, 1\}, \quad \sum_{k=1}^K z_k = 1$$



# A Simple Example



- Mixture of three Gaussians with mixing coefficients
- Left: all three Gaussians as contour plot
- Right: samples from the mixture model, the red component has the most samples



# Parameter Estimation

- From a given set of training data  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  we want to find parameters  $(\pi_{1,\dots,K}, \boldsymbol{\mu}_{1,\dots,K}, \Sigma_{1,\dots,K})$  so that the likelihood is maximized (MLE):

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N \mid \pi_{1,\dots,K}, \boldsymbol{\mu}_{1,\dots,K}, \Sigma_{1,\dots,K}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \Sigma_k)$$

or, applying the logarithm:

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

- However: this is not as easy as maximum-likelihood for single Gaussians!



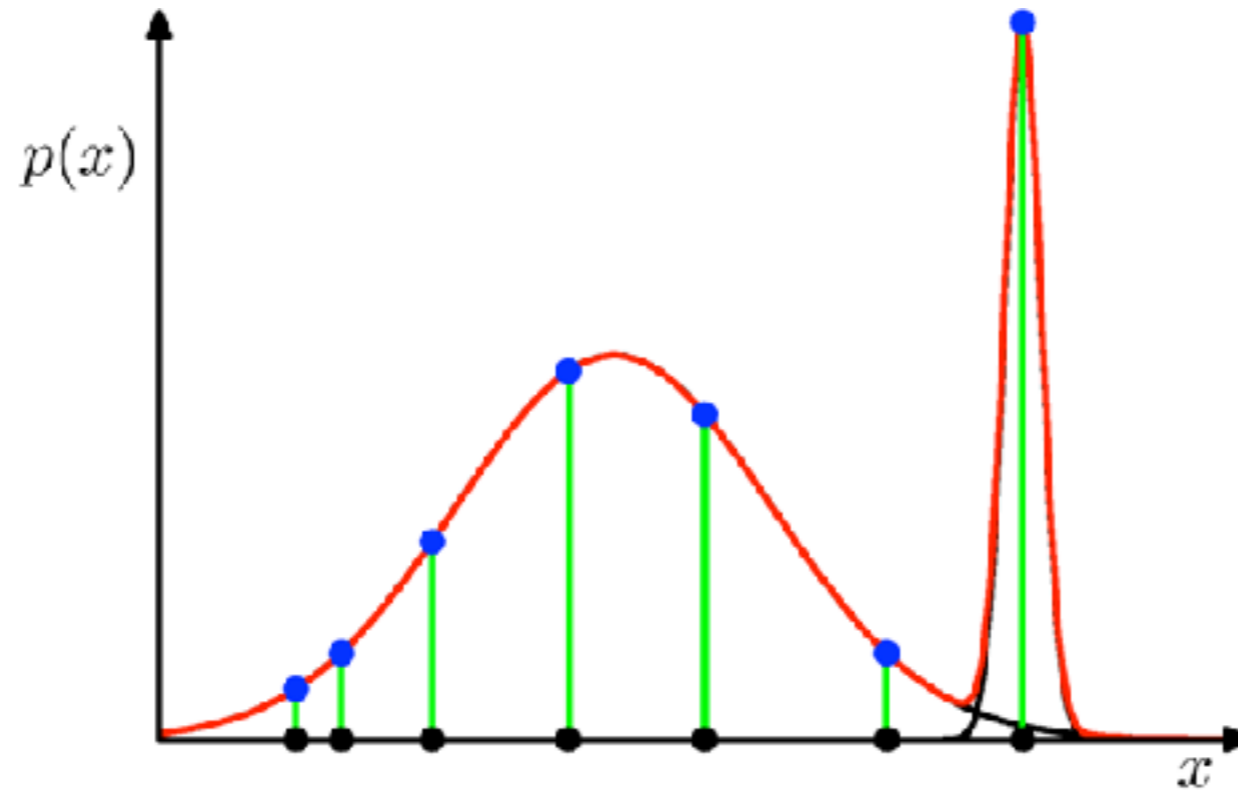


# Problems with MLE for Gaussian Mixtures

- Assume that for one  $k$  the mean  $\mu_k$  is exactly at a data point  $\mathbf{x}_n$ 
  - For simplicity: assume that  $\Sigma_k = \sigma_k^2 I$
  - Then: 
$$\mathcal{N}(\mathbf{x}_n \mid \mathbf{x}_n, \sigma_k^2 I) = \frac{1}{\sqrt{2\pi\sigma_k^D}}$$
  - This means that the overall log-likelihood can be maximized arbitrarily by letting  $\sigma_k \rightarrow 0$  (**overfitting**)
- Another problem is the **identifiability**:
  - The order of the Gaussians is not fixed, therefore:
  - There are  $K!$  equivalent solutions to the MLE problem



# Overfitting with MLE for Gaussian Mixtures



- One Gaussian fits exactly to one data point
- It has a very small variance, i.e. contributes strongly to the overall likelihood
- In standard MLE, there is no way to avoid this!



# Expectation-Maximization

- EM is an elegant and powerful method for MLE problems with latent variables
- Main idea: model parameters and latent variables are estimated iteratively, where average over the latent variables (expectation)
- A typical example application of EM is the Gaussian Mixture model (GMM)
- However, EM has many other applications
- First, we consider EM for GMMs



# Expectation-Maximization for GMM

- First, we define the **responsibilities**:

$$\gamma(z_{nk}) = p(z_{nk} = 1 \mid \mathbf{x}_n) \quad z_{nk} \in \{0, 1\}$$
$$\sum_k z_{nk} = 1$$



# Expectation-Maximization for GMM

- First, we define the **responsibilities**:

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



# Expectation-Maximization for GMM

- First, we define the **responsibilities**:

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

- Next, we derive the log-likelihood wrt. to  $\boldsymbol{\mu}_k$  :

$$\frac{\partial \log p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_k} \stackrel{!}{=} \mathbf{0}$$



# Expectation-Maximization for GMM

- First, we define the **responsibilities**:

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

- Next, we derive the log-likelihood wrt. to  $\boldsymbol{\mu}_k$ :

$$\frac{\partial \log p(X \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_k} \stackrel{!}{=} \mathbf{0}$$

and we obtain:

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



# Expectation-Maximization for GMM

- We can do the same for the covariances:

$$\frac{\partial \log p(X | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}_k} \stackrel{!}{=} \mathbf{0}$$

and we obtain:

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

- Finally, we derive wrt. the mixing coefficients  $\pi_k$  :

$$\frac{\partial \log p(X | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \pi_k} \stackrel{!}{=} \mathbf{0} \quad \text{where:} \quad \sum_{k=1}^K \pi_k = 1$$





# Expectation-Maximization for GMM

- We can do the same for the covariances:

$$\frac{\partial \log p(X | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}_k} \stackrel{!}{=} \mathbf{0}$$

and we obtain:

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

- Finally, we derive wrt. the mixing coefficients  $\pi_k$  :

$$\frac{\partial \log p(X | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \pi_k} \stackrel{!}{=} \mathbf{0} \quad \text{where:} \quad \sum_{k=1}^K \pi_k = 1$$

and the result is:

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



# Algorithm Summary

1. Initialize means  $\mu_k$  covariance matrices  $\Sigma_k$  and mixing coefficients  $\pi_k$
2. Compute the initial log-likelihood  $\log p(X | \pi, \mu, \Sigma)$
- 3. E-Step.** Compute the responsibilities:

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

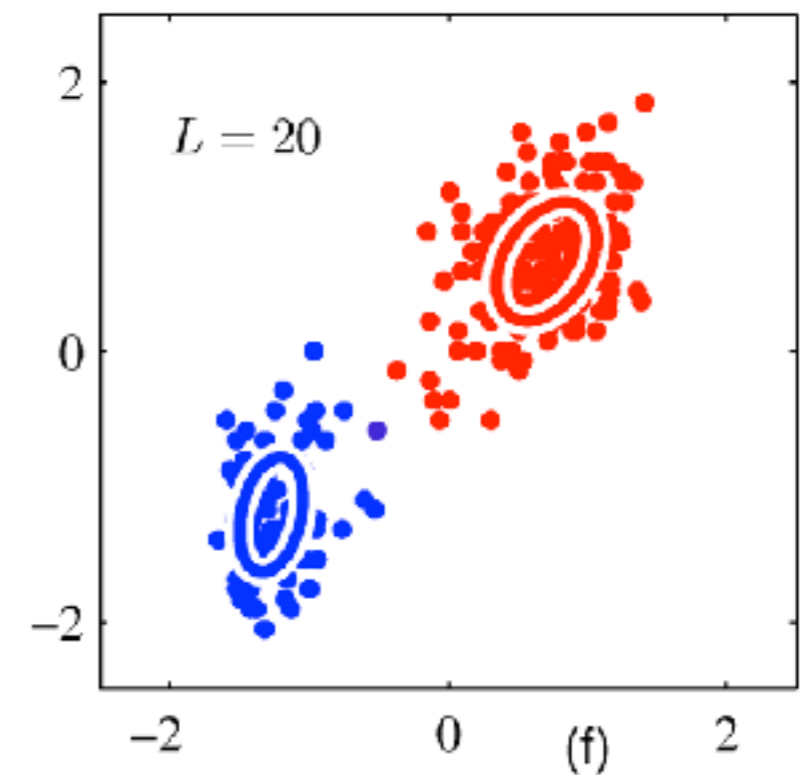
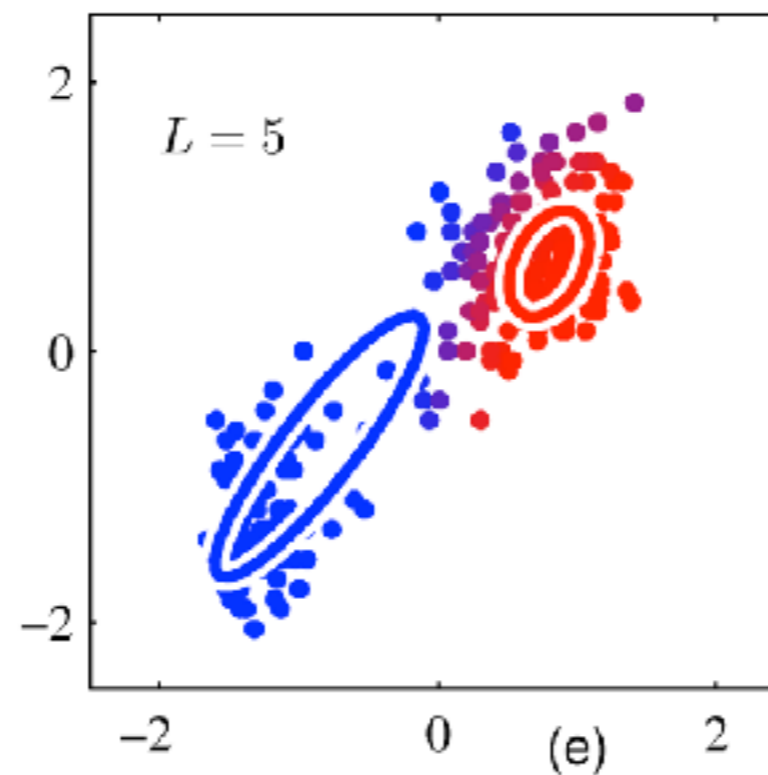
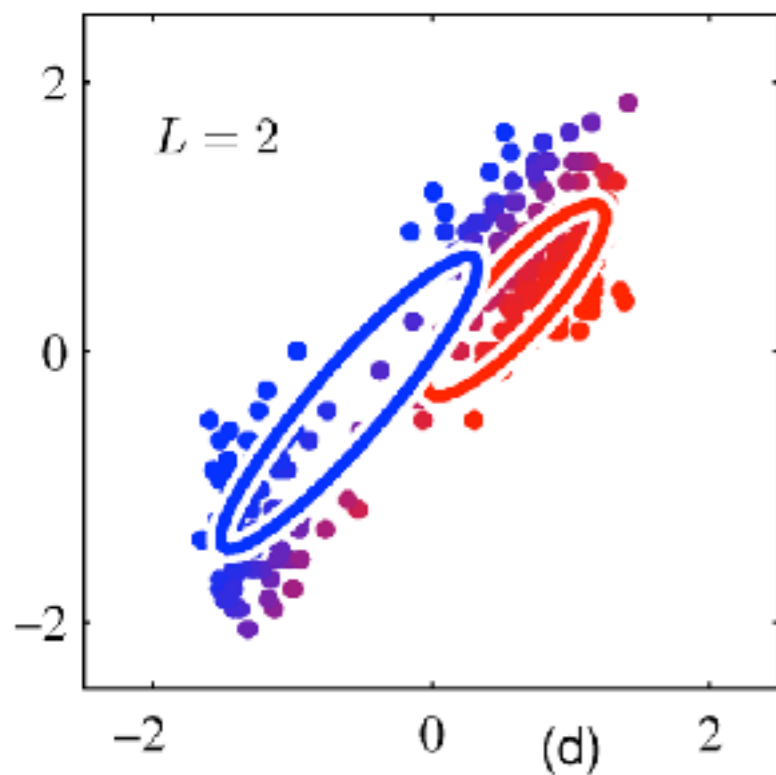
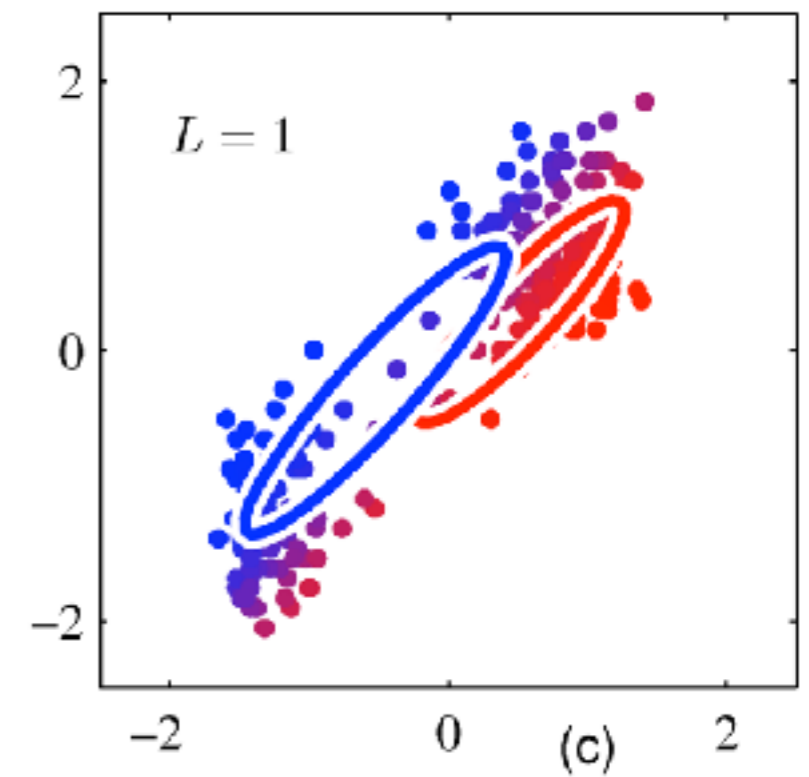
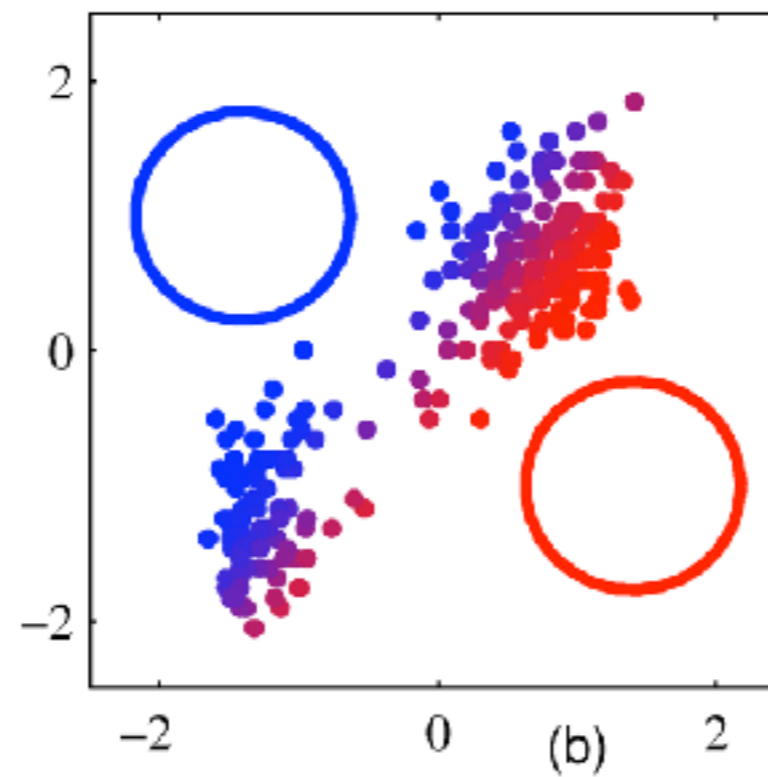
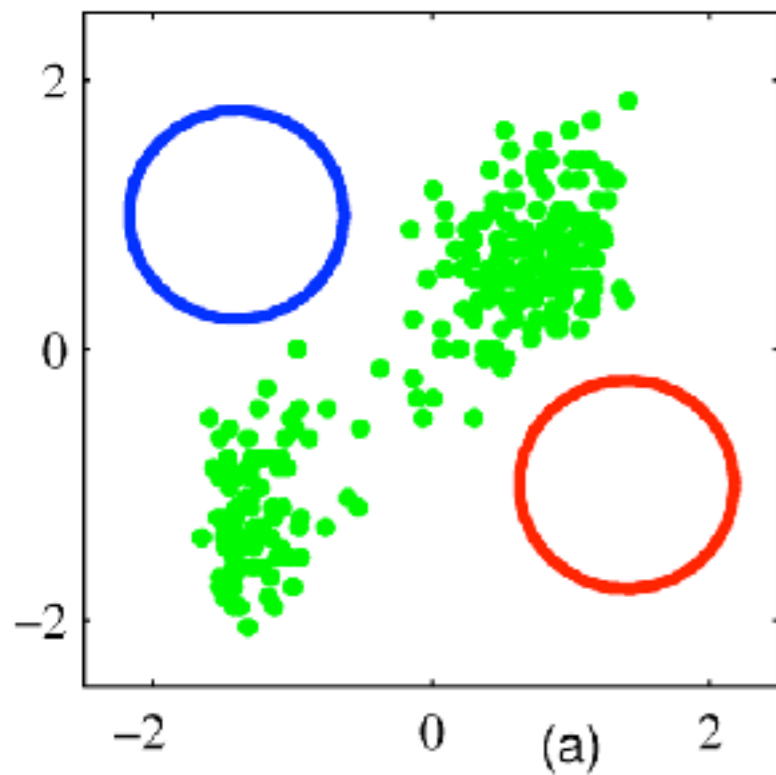
- 4. M-Step.** Update the parameters:

$$\mu_k^{\text{new}} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad \Sigma_k^{\text{new}} = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T}{\sum_{n=1}^N \gamma(z_{nk})} \quad \pi_k^{\text{new}} = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nk})$$

5. Compute log-likelihood; if not converged go to 3.



# The Same Example Again



# Why is it Called “EM”?

Assume for a moment that we observe  $X$  and the binary latent variables  $Z$ . The likelihood is then:

$$p(X, Z \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N p(\mathbf{z}_n \mid \boldsymbol{\pi}) p(\mathbf{x}_n \mid \mathbf{z}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

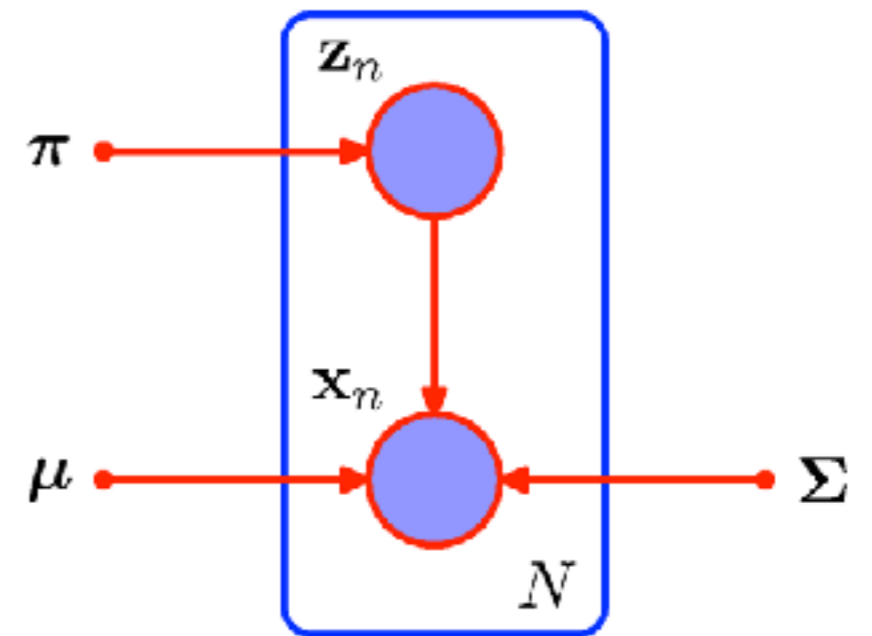
“Complete-data log-likelihood”

where  $p(\mathbf{z}_n \mid \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_{nk}}$  and

$$p(\mathbf{x}_n \mid \mathbf{z}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

which leads to the log-formulation:

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



# Why is it Called “EM”?

Instead of maximizing the joint log-likelihood, we maximize its **expectation** under the latent variable distribution:

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



# Why is it Called “EM”?

Instead of maximizing the joint log-likelihood, we maximize its **expectation** under the latent variable distribution:

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

where the latent variable distribution per point is:

$$\begin{aligned} p(\mathbf{z}_n \mid \mathbf{x}_n, \boldsymbol{\theta}) &= \frac{p(\mathbf{x}_n \mid \mathbf{z}_n, \boldsymbol{\theta}) p(\mathbf{z}_n \mid \boldsymbol{\theta})}{p(\mathbf{x}_n \mid \boldsymbol{\theta})} & \boldsymbol{\theta} &= (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \frac{\prod_{l=1}^K (\pi_l \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l))^{z_{nl}}}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$



# Observations

- Compared to K-means, points can now belong to both clusters (**soft assignment**)
- In addition to the cluster center, a covariance is estimated by EM
- Initialization is the same as used for K-means
- Number of iterations needed for EM is much higher
- Also: each cycle requires much more computation
- Therefore: start with K-means and run EM on the result of K-means (covariances can be initialized to the sample covariances of K-means)
- EM only finds a **local** maximum of the likelihood!



# Clustering using Mixture Models

The full posterior of the Gaussian Mixture Model is

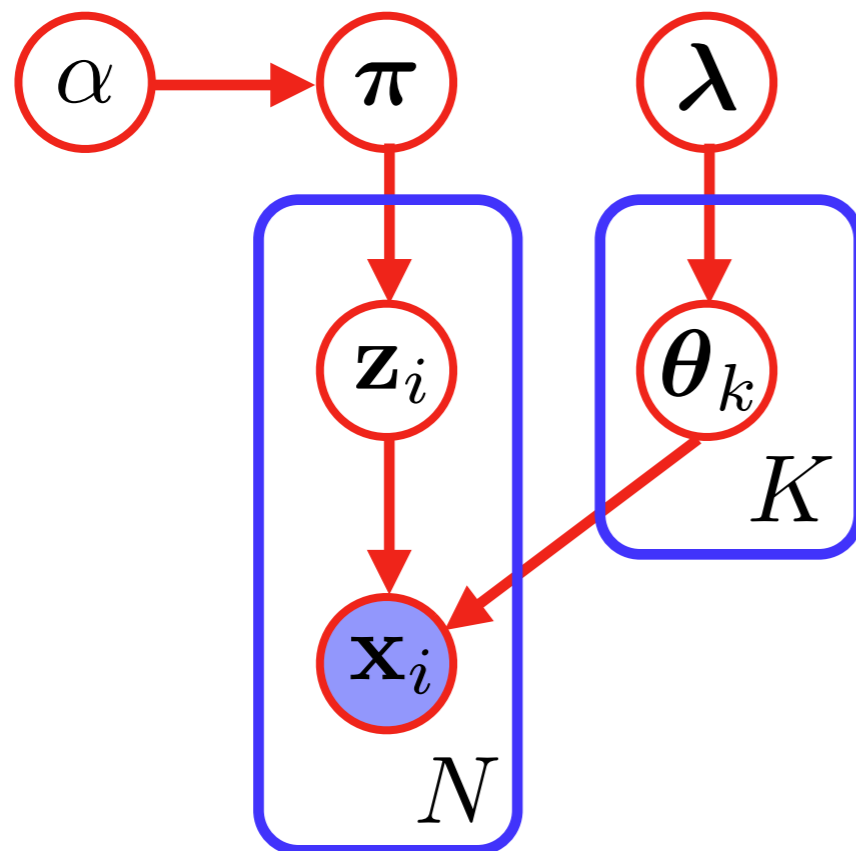
$$p(X, Z, \mu, \Sigma, \pi) = \underbrace{p(X | Z, \mu, \Sigma)}_{\text{data likelihood (Gaussian)}} \underbrace{p(Z | \pi)}_{\text{correspondence prob. (Multinomial)}} \underbrace{p(\pi | \alpha)}_{\text{mixture prior (Dirichlet)}} \underbrace{p(\mu, \Sigma | \lambda)}_{\text{parameter prior (Gauss-IW)}}$$

data likelihood  
(Gaussian)

correspondence  
prob. (Multinomial)

mixture prior  
(Dirichlet)

parameter prior  
(Gauss-IW)



In this model, we use:

- $\mu = (\mu_1, \dots, \mu_K)$
- $\Sigma = (\Sigma_1, \dots, \Sigma_K)$
- $(\mu_k, \Sigma_k) = \theta_k$





# Clustering using Mixture Models

The full posterior of the Gaussian Mixture Model is

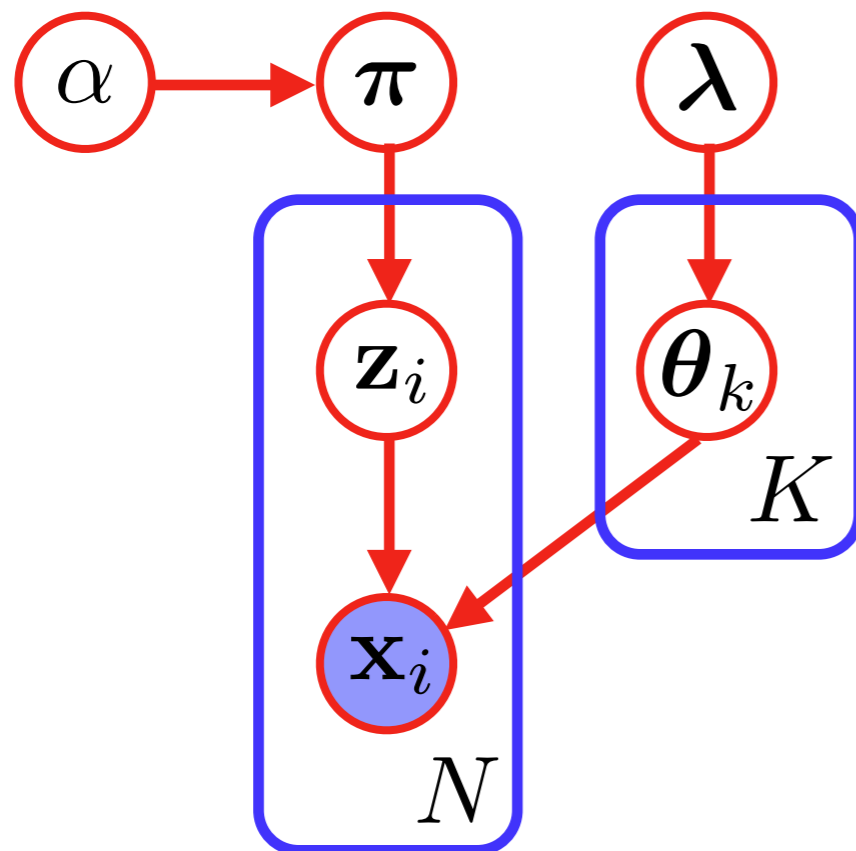
$$p(X, Z, \mu, \Sigma, \pi) = \underbrace{p(X | Z, \mu, \Sigma)}_{\text{data likelihood (Gaussian)}} \underbrace{p(Z | \pi)}_{\text{correspondence prob. (Multinomial)}} \underbrace{p(\pi | \alpha)}_{\text{mixture prior (Dirichlet)}} \underbrace{p(\mu, \Sigma | \lambda)}_{\text{parameter prior (Gauss-IW)}}$$

data likelihood  
(Gaussian)

correspondence  
prob. (Multinomial)

mixture prior  
(Dirichlet)

parameter prior  
(Gauss-IW)



Given this model, we can create new samples:

1. Sample  $\pi, \theta_k$  from priors
2. Sample corresp.  $z_i$
3. Sample data point  $x_i$



# Clustering using Mixture Models

The full posterior of the Gaussian Mixture Model is

$$p(X, Z, \mu, \Sigma, \pi) = \underbrace{p(X | Z, \mu, \Sigma)}_{\text{data likelihood (Gaussian)}} \underbrace{p(Z | \pi)}_{\text{correspondence prob. (Multinomial)}} \underbrace{p(\pi | \alpha)}_{\text{mixture prior (Dirichlet)}} \underbrace{p(\mu, \Sigma | \lambda)}_{\text{parameter prior (Gauss-IW)}}$$

data likelihood  
(Gaussian)

correspondence  
prob. (Multinomial)

mixture prior  
(Dirichlet)

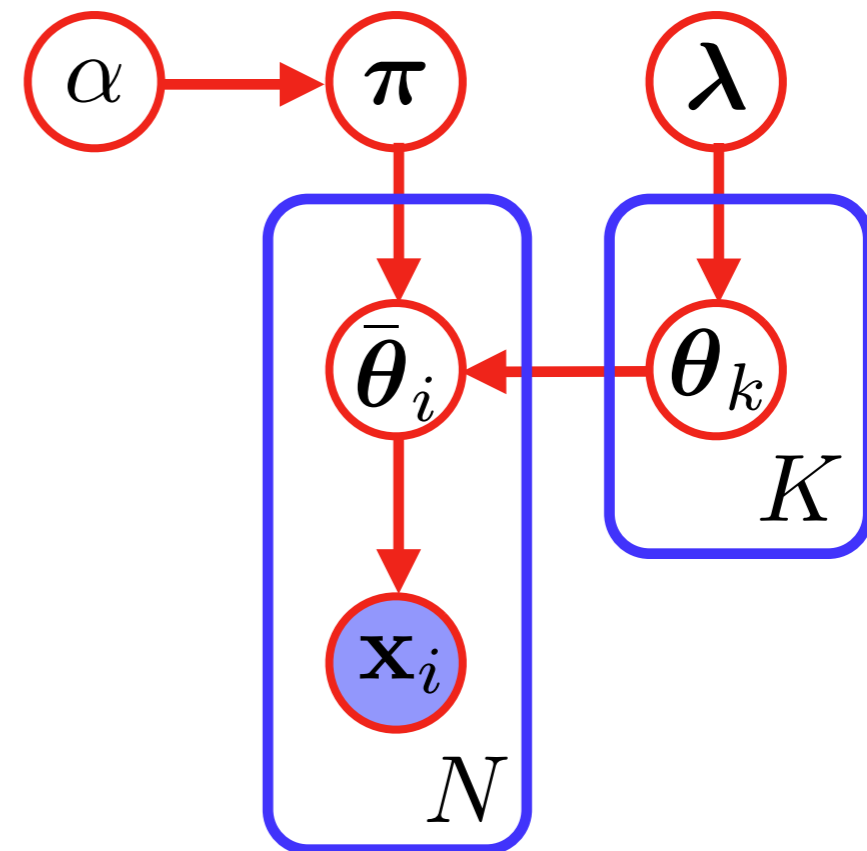
parameter prior  
(Gauss-IW)

An equivalent formulation of this model is this:

1. Sample  $\pi, \theta_k$  from priors
2. Sample params  $\bar{\theta}_i$  from:

$$p(\bar{\theta}_i | \pi, \theta_k) = \sum_{k=1}^K \pi_k \delta(\theta_k, \bar{\theta}_i)$$

3. Sample data point  $\mathbf{x}_i$



# Clustering using Mixture Models

What is the difference in that model?

- there is one parameter  $\bar{\theta}_i$  for each observation  $x_i$
- intuitively: we first sample the location of the cluster and then the data that corresponds to it

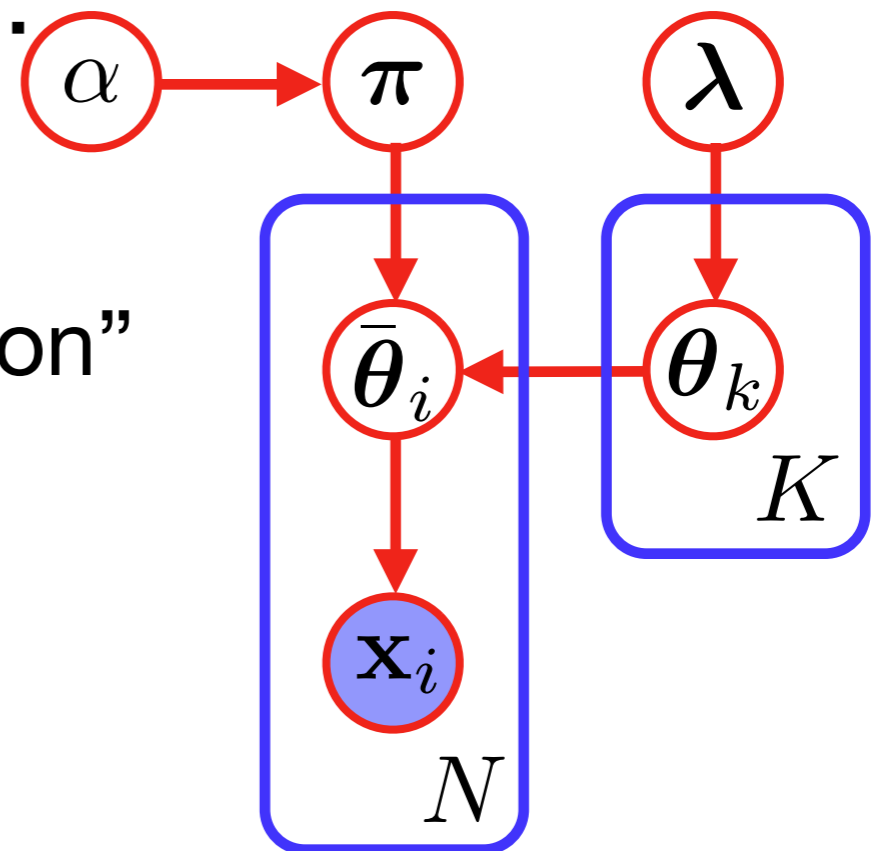
In general, we use the notation:

$$\pi \sim \text{Dir}\left(\frac{\alpha}{K} \mathbf{1}\right)$$

$$\theta_k \sim H(\lambda) \quad \text{“Base distribution”}$$

$$\bar{\theta}_i \sim G(\pi, \theta_k) \quad \text{where}$$

$$G(\pi, \theta_k) = \sum_{k=1}^K \pi_k \delta(\theta_k, \bar{\theta}_i)$$



**However: We need to know K**

