



---

# **PATTERN RECOGNITION AND MACHINE LEARNING**

## **CHAPTER 9: MIXTURE MODELS AND EM**

---

# K-means Clustering

---

Suppose we have a data set  $\{x_1, \dots, x_N\}$  consisting of  $N$  observations of a random  $D$ -dimensional Euclidean variable  $x$ .

Our goal is to partition the data set into some number  $K$  of clusters

We can formalize this notion by first introducing a set of  $D$ -dimensional vectors  $\mu_k$ , where  $k = 1, \dots, K$ , in which  $\mu_k$  is a prototype (center of clusters) associated with the  $k$ th cluster.

Our goal is then to find an assignment of data points to clusters, as well as a set of vectors  $\{\mu_k\}$ , such that the sum of the squares of the distances of each data point to its closest vector  $\mu_k$ , is a minimum.

---

# K-means Clustering

---

For each data point  $x_n$ , we introduce a corresponding set of binary indicator variables  $r_{nk} \in \{0, 1\}$ , where  $k = 1, \dots, K$  describing which of the  $K$  clusters the data point  $x_n$  is assigned to, so that if data point  $x_n$  is assigned to cluster  $k$  then  $r_{nk} = 1$ , and  $r_{nj} = 0$  for  $j \neq k$ .

We can then define an objective function, sometimes called a distortion measure, given by

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

Our goal is to find values for the  $\{r_{nk}\}$  and the  $\{\mu_k\}$  to minimize  $J$ .

First, we choose some initial values for the  $\mu_k$  and we minimize  $J$  with respect to the  $r_{nk}$ , keeping the  $\mu_k$  fixed.

In the second phase, we minimize  $J$  with respect to the  $\mu_k$ , keeping  $r_{nk}$  fixed.

---

# K-means Clustering

---

$J$  is a linear function of  $r_{nk}$ , this optimization can be performed easily to give a closed form solution.

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

The objective function  $J$  is a quadratic function of  $\boldsymbol{\mu}_k$ , and it can be minimized by setting its derivative with respect to  $\boldsymbol{\mu}_k$  to zero:

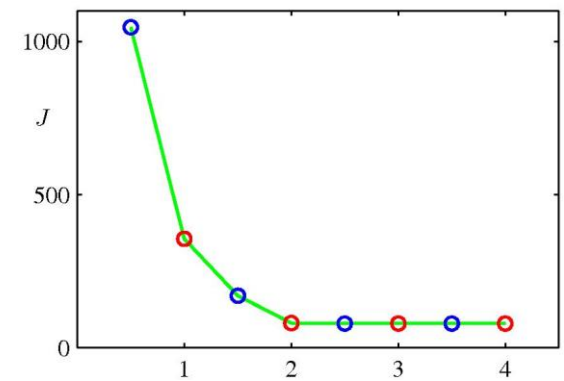
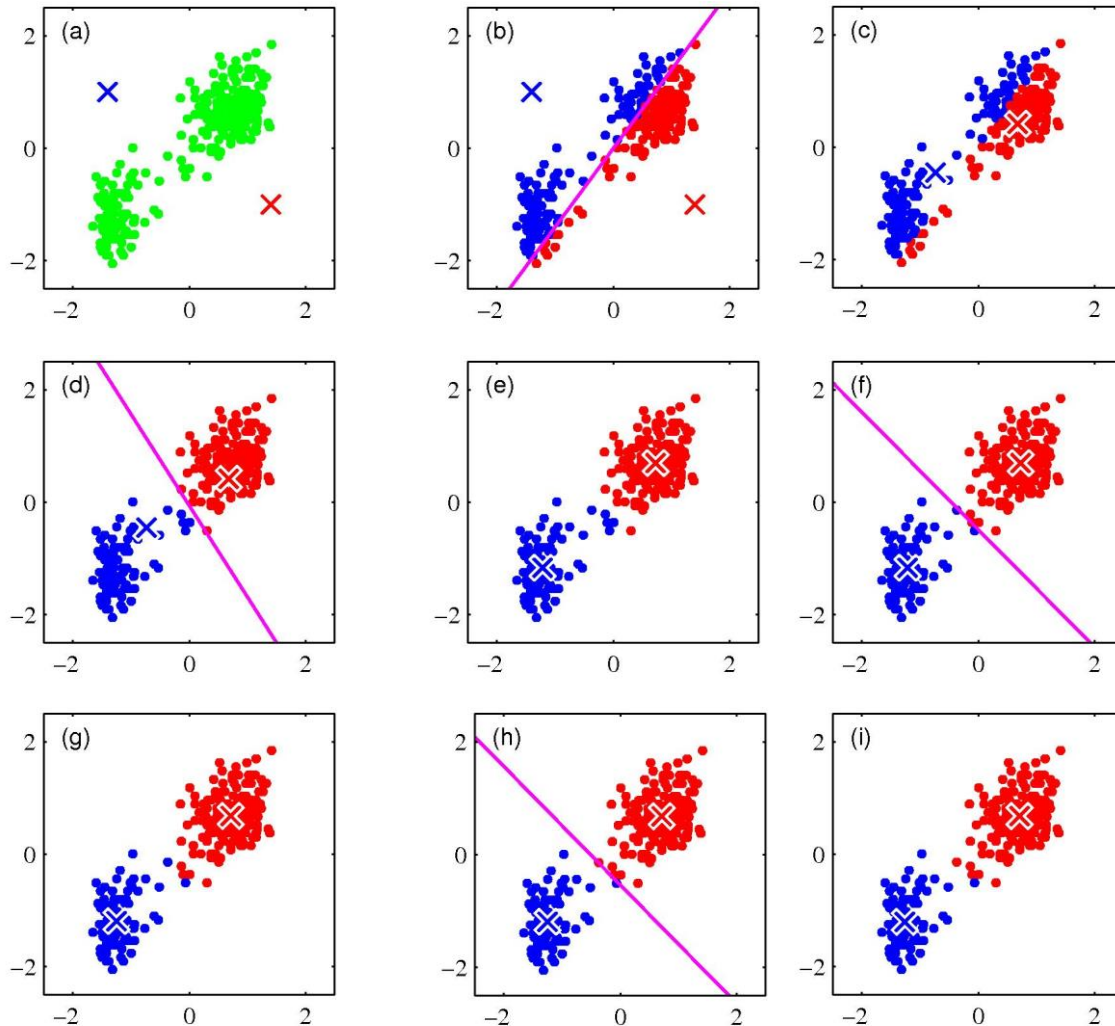
$$2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}.$$

The denominator in this expression is equal to the number of points assigned to cluster  $k$ , and so this result has a simple interpretation, namely set  $\boldsymbol{\mu}_k$  equal to the mean of all of the data points  $\mathbf{x}_n$  assigned to cluster  $k$ . For this reason, the procedure is known as the K-means algorithm.

---

# K-means Clustering





# K-means Clustering

---

$K = 2$



$K = 3$



$K = 10$



Original image



# Mixtures of Gaussians

---

A Gaussian mixture distribution can be written as a linear superposition of Gaussians in the form

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

Let us introduce a K-dimensional binary random variable  $\mathbf{z}$  having a 1-of-K representation in which a particular element  $z_k$  is equal to 1 and all other elements are equal to 0.

The values of  $z_k$  therefore satisfy  $z_k \in \{0, 1\}$  and  $\sum z_k = 1$ , and we see that there are K possible states for the vector  $\mathbf{z}$  according to which element is nonzero.

We shall define the joint distribution  $p(\mathbf{x}, \mathbf{z})$  in terms of a marginal distribution  $p(\mathbf{z})$  and a conditional distribution  $p(\mathbf{x}|\mathbf{z})$ ,

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

---

# Mixtures of Gaussians

---

Another quantity that will play an important role is the conditional probability of  $z$  given  $\mathbf{x}$ .

We shall use  $\gamma(z_k)$  to denote  $p(z_k = 1|\mathbf{x})$ , whose value can be found using Bayes' theorem

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



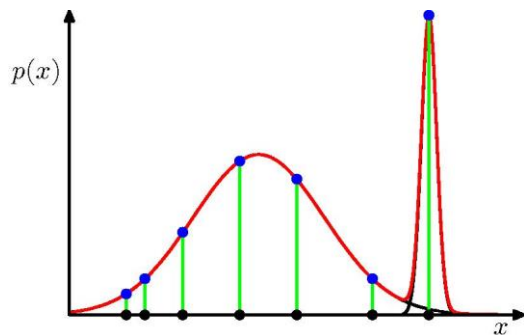
# Mixtures of Gaussians

---

Suppose we have a data set of observations  $\{x_1, \dots, x_N\}$ , and we wish to model this data using a mixture of Gaussians.

If we assume that the data points are drawn independently from the distribution, then we can express the Gaussian mixture model for this i.i.d. data set and the log of the likelihood function is given by

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



A further issue in finding maximum likelihood solutions arises from the fact that for any given maximum likelihood solution, a K-component mixture will have a total of  $K!$  equivalent solutions corresponding to the  $K!$  ways of assigning K sets of parameters to K components.

---

# Mixtures of Gaussians

An elegant and powerful method for finding maximum likelihood solutions for models with latent variables is called the expectation-maximization algorithm, or EM algorithm

## EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means  $\mu_k$ , covariances  $\Sigma_k$  and mixing coefficients  $\pi_k$ , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}. \quad (9.23)$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.24)$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}})(\mathbf{x}_n - \mu_k^{\text{new}})^T \quad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \quad (9.27)$$

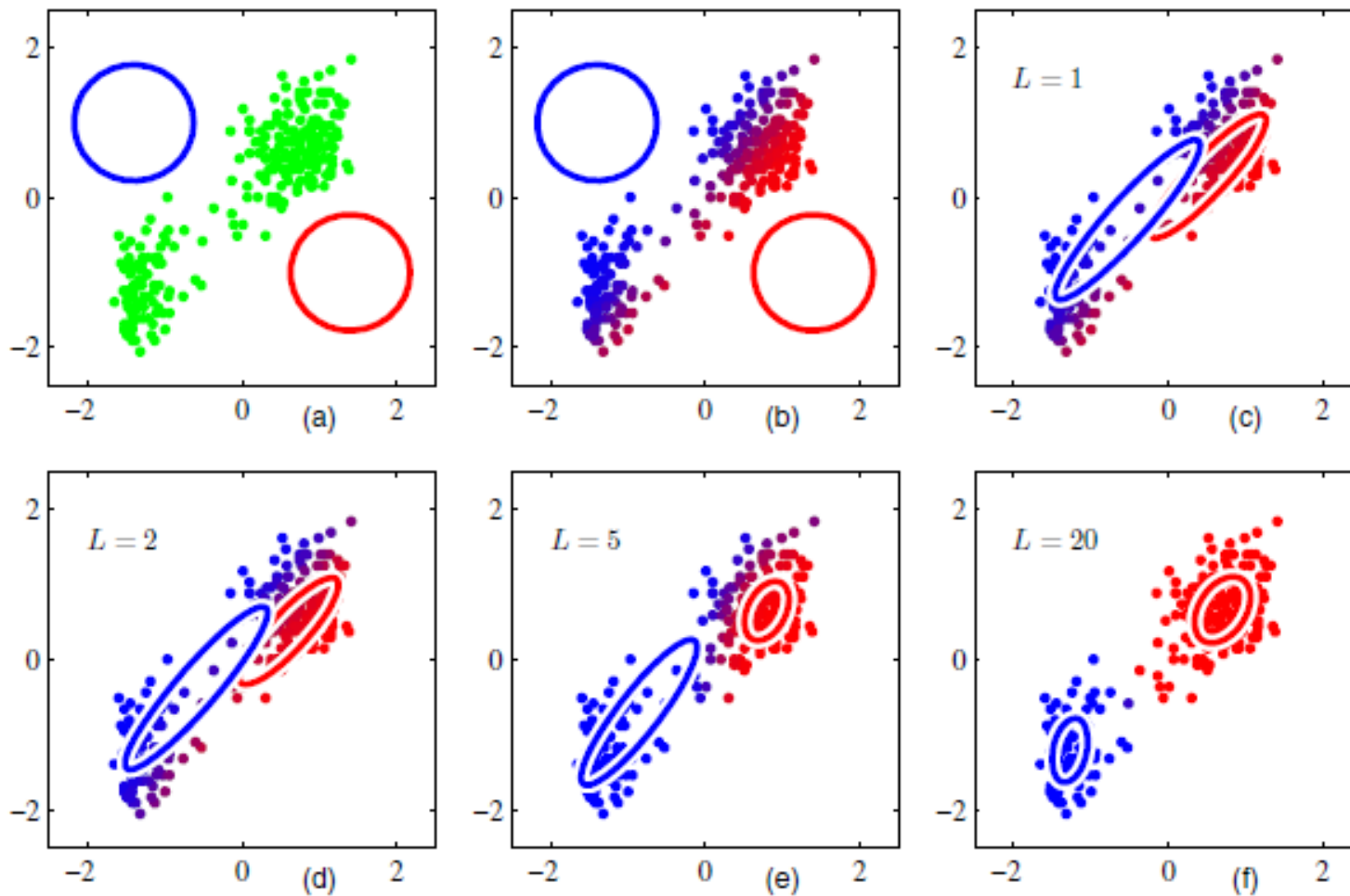
4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\} \quad (9.28)$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

# Mixtures of Gaussians

---



# Mixtures of Gaussians

---

## The General EM Algorithm

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

1. Choose an initial setting for the parameters  $\boldsymbol{\theta}^{\text{old}}$ .

2. **E step** Evaluate  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ .

3. **M step** Evaluate  $\boldsymbol{\theta}^{\text{new}}$  given by

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (9.32)$$

where

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.33)$$

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}} \quad (9.34)$$

and return to step 2.

Comparison of the K-means algorithm with the EM algorithm for Gaussian mixtures shows that there is a close similarity.

Whereas the K-means algorithm performs a hard assignment of data points to clusters, in which each data point is associated uniquely with one cluster, the EM algorithm makes a soft assignment based on the posterior probabilities.

# The EM Algorithm in General

---

Consider a probabilistic model in which we collectively denote all of the observed variables by  $\mathbf{X}$  and all of the hidden variables by  $\mathbf{Z}$ . The joint distribution  $p(\mathbf{X}, \mathbf{Z} | \theta)$  is governed by a set of parameters denoted  $\theta$ .

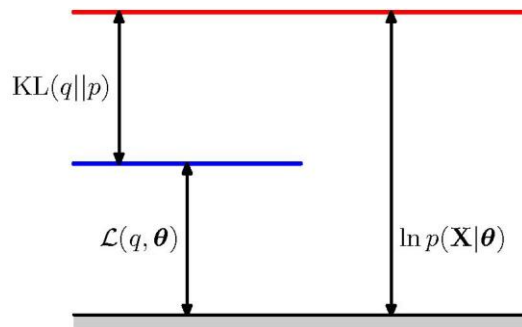
$$p(\mathbf{X} | \theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta).$$

Next, we introduce a distribution  $q(\mathbf{Z})$  defined over the latent variables, and we observe that, for any choice of  $q(\mathbf{Z})$ , the following decomposition holds

$$\ln p(\mathbf{X} | \theta) = \mathcal{L}(q, \theta) + \text{KL}(q \| p)$$

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

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



The Kullback-Leibler divergence  $\geq 0$ .

$\mathcal{L}(q, \theta)$  = a lower bound