

# Mixture models

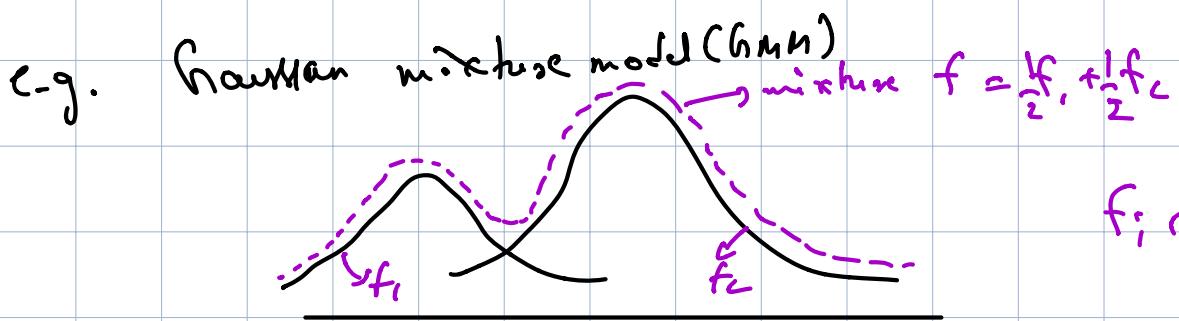
Ref: Bishop's book  
Chapter 9

In practice, class conditional densities may not be in a standard form. E.g. multimodal densities.

Mixture model:  $f(x) = \sum_{j=1}^K \lambda_j f_j(x),$

 $\lambda_j \geq 0, \sum_{j=1}^K \lambda_j = 1$

$f_j, j=1 \dots K$  is a valid density



GMMs are useful to model multimodal densities

$$f(x|\theta) = \sum_{j=1}^K \lambda_j f_j(x), \quad f_j \sim N(\mu_j, \sigma_j^2)$$

$$\Theta = (\lambda_j, j=1 \dots K, \mu_j, \sigma_j^2, j=1 \dots K)$$

These need to be estimated from i.i.d samples

Max-likelihood approach:

Dataset  $D = \{x_1, \dots, x_n\}$  i.i.d from  $f$

$$L(\theta) = \prod_{i=1}^n \left( \sum_{j=1}^K \lambda_j f_j(x_i) \right)$$

$$l(\theta) = \sum_{i=1}^n \log \left( \underbrace{\sum_{j=1}^k \lambda_j f_j(x_i)}_{\text{Sum inside a log}} \right)$$

Sum inside a log  
 $\Rightarrow$  hard to optimize  $l(\theta)$   
 even if  $f_j$ 's are Gaussian

For ease of presentation, we focus on the case  $k=2$

$$\Theta_1 = (\mu_1, \sigma_1^2), \Theta_2 = (\mu_2, \sigma_2^2)$$

$$\phi(x|\Theta_j) = \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{(x-\mu_j)^2}{2\sigma_j^2}\right) \quad j=1, 2$$

$$f(x|\theta) = \lambda_1 \phi(x|\Theta_1) + \lambda_2 \phi(x|\Theta_2)$$

$$\theta = (\underbrace{\Theta_1, \Theta_2}_{\text{component density parameters}}, \underbrace{\lambda_1, \lambda_2}_{\text{mixture co-efficients}})$$

$$l(\theta) = \sum_{i=1}^n \log \left( \lambda_1 \phi(x_i|\Theta_1) + \lambda_2 \phi(x_i|\Theta_2) \right)$$

$$\text{Goal: } \max_{\theta} l(\theta)$$

$$\frac{\partial \phi(x|\Theta_j)}{\partial \mu_j} = \phi(x|\Theta_j) \frac{(x-\mu_j)}{\sigma_j^2}$$

$$\frac{\partial \phi(x|\theta_j)}{\partial \sigma_j} = \phi(x|\theta_j) \left( \frac{(x - \mu_j)^2}{\sigma_j^3} - \frac{1}{\sigma_j} \right)$$

$$\frac{\partial l(\theta)}{\partial \mu_j} = \sum_{i=1}^n \left( \frac{\lambda_j \phi(x_i|\theta_j)}{\lambda_1 \phi(x_i|\theta_1) + \lambda_2 \phi(x_i|\theta_2)} \right) \left( \frac{x_i - \mu_j}{\sigma_j^2} \right)$$

Letting  $\alpha_{ij} = \left( \frac{\lambda_j \phi(x_i|\theta_j)}{\lambda_1 \phi(x_i|\theta_1) + \lambda_2 \phi(x_i|\theta_2)} \right)$ , we have

$$\frac{\partial l(\theta)}{\partial \mu_j} = \sum_{i=1}^n \alpha_{ij} \left( \frac{x_i - \mu_j}{\sigma_j^2} \right)$$

$$\frac{\partial l(\theta)}{\partial \sigma_j} = \sum_{i=1}^n \alpha_{ij} \left( \frac{x_i - \mu_j}{\sigma_j^3} - \frac{1}{\sigma_j} \right)$$

ML-estimates satisfy

$$\hat{\mu}_j = \frac{\sum_{i=1}^n \alpha_{ij} x_i}{\sum_{i=1}^n \alpha_{ij}}, \quad \hat{\sigma}_j^2 = \frac{\sum_{i=1}^n \alpha_{ij} (x_i - \hat{\mu}_j)^2}{\sum_{i=1}^n \alpha_{ij}}, \quad j = 1, 2$$

Maximizing  $l(\theta)$  w.r.t.  $\lambda_1$  &  $\lambda_2$ :

$$\max_{\theta} l(\theta) \quad \text{s.t. } \lambda_1 + \lambda_2 = 1$$

Lagrangian:  $g(\theta, \eta) = l(\theta) + \eta(\lambda_1 + \lambda_2 - 1)$

$\nearrow$   
Lagrange multiplier

$$\nabla_{\lambda_1} g(\theta, \eta) = 0$$

∴  $\sum_{i=1}^n \frac{\phi(x_i | \theta_1)}{\lambda_1 \phi(x_i | \theta_1) + \lambda_2 \phi(x_i | \theta_2)} + \eta = 0$

or, equivalently,  $\sum_{i=1}^n \frac{\lambda_{i1}}{\lambda_1} + \eta = 0 \quad \text{--- (1)}$

Similarly,  $\sum_{i=1}^n \frac{\lambda_{i2}}{\lambda_2} + \eta = 0 \quad \text{--- (2)}$

Summing (1) & (2), we have

$$\sum_{i=1}^n \lambda_{i1} + \sum_{i=1}^n \lambda_{i2} + \eta(\lambda_1 + \lambda_2) = 0$$

$$\begin{aligned} \eta &= - \left( \sum_{i=1}^n \lambda_{i1} + \sum_{i=1}^n \lambda_{i2} \right) \\ &= -n \end{aligned}$$

So, ML estimates for  $\lambda_j$  satisfy

$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n \lambda_{ij}$$

So, ML estimation for  $\theta$  involves solving the following system of equations:

$$\hat{\mu}_j = \frac{\sum_{i=1}^n \alpha_{ij} x_i}{\sum_{i=1}^n \alpha_{ij}}, \quad \hat{\sigma}_j^2 = \frac{\sum_{i=1}^n \alpha_{ij} (x_i - \mu_j)^2}{\sum_{i=1}^n \alpha_{ij}}$$

$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n \alpha_{ij}$$

An iterative scheme for solving the above set of equations:-

$$\hat{\mu}_j^{(k+1)} = \frac{\sum_{i=1}^n \alpha_{ij}^{(k)} x_i}{\sum_{i=1}^n \alpha_{ij}^{(k)}}, \quad \hat{\sigma}_j^2 = \frac{\sum_{i=1}^n \alpha_{ij}^{(k)} (x_i - \mu_j^{(k)})^2}{\sum_{i=1}^n \alpha_{ij}^{(k)}}$$

$$\hat{\lambda}_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \alpha_{ij}^{(k)}$$

$$\alpha_{ij}^{(k+1)} = \frac{\lambda_j^{(k+1)} \phi(x_i | \theta_j^{(k+1)})}{\lambda_1^{(k+1)} \phi(x_i | \theta_1^{(k+1)}) + \lambda_2^{(k+1)} \phi(x_i | \theta_2^{(k+1)})}$$

Prelude to EM algorithm:

$$f(x | \theta) = \lambda_1 \phi(x | \theta_1) + \lambda_2 \phi(x | \theta_2)$$

$x_i \sim f(x | \theta)$  can be generated as follows:

$$(i) \quad Z = \begin{cases} 1 & \text{w.p. } \lambda_1 \\ 2 & \text{w.p. } \lambda_2 \end{cases}$$

(i.) Generate  $x_i$  from  $\phi(x|\theta_2)$

Dataset  $D = \{x_i, i=1 \dots n\}$

If we are given component information for each  $x_i$ ,  
then the "parameter estimation" task is easy

Let  $Z_{ij}$ ,  $i=1 \dots n$ ,  $j=1, 2$  denote component information

$$Z_{ij} = 1 \text{ if } x_i \sim \phi(\cdot | \theta_j)$$

$$P(Z_{ij}=1) = \lambda_j, \quad i=1 \dots n, j=1, 2$$

$$f(x_i | Z_{ij}=1) = \phi(x_i | \theta_j)$$

$$Z_i = (Z_{i1}, Z_{i2})$$

Complete dataset  $\{(x_i, Z_i), i=1 \dots n\}$

Note:  $Z_i$  is "missing" information

However, given  $Z_i$ , the MLE-estimates are

$$\hat{\mu}_1 = \frac{\sum_{i=1}^n Z_{i1} x_i}{\sum_{i=1}^n Z_{i1}}, \quad \hat{\mu}_2 = \frac{\sum_{i=1}^n Z_{i2} x_i}{\sum_{i=1}^n Z_{i2}}$$

The ML-estimates for  $\hat{\gamma}_j^2$  can be worked out in a similar fashion.

### EM algorithm:

Iterative procedure for ML estimation  
works with a regular dataset, i.e.,  $Z_i$  (missing information) is not provided

$$D = \{x_i, i=1-n\} \quad \bar{D} = \{x_i, z_i, i=1-n\}$$

$f(x, z | \theta)$  → density for  $\bar{D}$

$$\bar{l}(\theta) = \log \left( \prod_{i=1}^n f(x_i, z_i | \theta) \right)$$

Let  $\log(f(x, z | \theta))$

$$\text{denote } \log \left( \prod_{i=1}^n f(x_i, z_i | \theta) \right)$$

$$x = (x_1, \dots, x_n), z = (z_1, \dots, z_n)$$

### EM algorithm:-

Expectation step (E-step):

$$Q(\theta, \theta^{(k)}) = E_{z|x, \theta^{(k)}} \log(f(x, z | \theta))$$

↙

Expectation of complete data log likelihood  
w.r.t. conditional distribution of missing data  
given incomplete data  $x$  & iterate  $\theta^{(k)}$ .

Maximization step (M-step):

$$\theta^{(k+1)} = \arg \max_{\theta} Q(\theta, \theta^{(k)})$$

Compute next value of  $\theta$  i.e.,  $\theta^{(k+1)}$  by maxing  
the  $Q(\cdot, \cdot)$  found in E-step.

---

Example:

GMMs with two components

$$f(x|\theta) = \lambda_1 \phi(x|\theta_1) + \lambda_2 \phi(x|\theta_2)$$

Recall  $P(2_{ij}=1) = \lambda_j$ ,  $f(x_i|2_{ij}=1) = \phi(x_i|\theta_j)$

$$2_i = (2_{i1}, 2_{i2})$$

$$f(2_i|\theta) = (\lambda_1)^{2_{i1}} (\lambda_2)^{2_{i2}} = \prod_{j=1}^2 (\lambda_j)^{2_{ij}}$$

$$f(x_i|2_i, \theta) = [\phi(x_i|\theta_1)]^{2_{i1}} [\phi(x_i|\theta_2)]^{2_{i2}}$$

$$\begin{aligned}
 f(x_i, z_i | \theta) &= f(x_i | z_i, \theta) f(z_i | \theta) \\
 &= \prod_{j=1}^J \left[ \lambda_j \phi(x_i | \theta_j) \right]^{z_{ij}}
 \end{aligned}$$

$$f(x, z | \theta) = \prod_{i=1}^n \prod_{j=1}^J \left[ \lambda_j \phi(x_i | \theta_j) \right]^{z_{ij}}$$

$$\begin{aligned}
 \log f(x, z | \theta) &= \sum_{i=1}^n \sum_{j=1}^J z_{ij} \log (\lambda_j \phi(x_i | \theta_j)) \\
 &= \sum_{i=1}^n \sum_{j=1}^J z_{ij} (\log \lambda_j + \log (\phi(x_i | \theta_j)))
 \end{aligned}$$



Sum of log  $\Rightarrow$  ML estimation is easier

  
we got sum of log because we are given

the missing info  $z_i$  in this

case

E-step:  $Q(\theta, \theta^{(k)}) = E_{z|x, \theta^{(k)}} \log(f(x, z | \theta))$

$$Q(\theta, \theta^{(k)}) = \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij} | \mathbf{x}, \theta^{(k)}] \log(\lambda_j \phi(x_i | \theta_j))$$

$$E(Z_{ij} | \mathbf{x}, \theta) = P(Z_{ij}=1 | \mathbf{x}, \theta)$$

$$= \frac{f(x_i | Z_{ij}=1, \theta) P(Z_{ij}=1)}{\sum_{k=1}^2 f(x_i | Z_{ik}=1, \theta) P(Z_{ik}=1)}$$

$$= \frac{\lambda_j \phi(x_i | \theta_j)}{\sum_{k=1}^2 \lambda_k \phi(x_i | \theta_k)} = \lambda_{ij}(\theta)$$

So,

$$Q(\theta, \theta^{(k)}) = \sum_{i=1}^n \sum_{j=1}^2 \lambda_{ij}(\theta^{(k)}) \log(\lambda_j \phi(x_i | \theta_j))$$

M-step:

$$\theta^{(k+1)} = \arg \max_{\theta} Q(\theta, \theta^{(k)})$$

$$Q(\theta, \theta^{(k)}) = \sum_{i=1}^n \sum_{j=1}^2 \lambda_{ij}(\theta^{(k)}) \left[ \log \lambda_j - \log(\sigma_j \sqrt{2\pi}) - \frac{(x_i - \mu_j)^2}{2\sigma_j^2} \right]$$

Need to maximize the above w.r.t.  $\theta$

$$\frac{\partial Q}{\partial \mu_i} = 0 \Rightarrow \sum_{i=1}^n \lambda_{ii}(\theta^{(k)}) \left( \frac{x_i - \mu_i}{\sigma_i^2} \right) = 0$$

$$S_o, \quad \mu_i^{(k+1)} = \frac{\sum_{i=1}^n \lambda_{ii}(\theta^k) x_i}{\sum_{i=1}^n \lambda_{ii}(\theta^k)} \quad --- (*)$$

$$\frac{\partial Q}{\partial \sigma_i} = 0 \Rightarrow \sum_{i=1}^n \lambda_{ii}(\theta^k) \left( -\frac{1}{\sigma_i} + \frac{(x_i - \mu_i)^2}{\sigma_i^3} \right) = 0$$

$$S_o, \quad (\sigma_i^2)^{(k+1)} = \frac{\sum_{i=1}^n \lambda_{ii}(\theta^k) (x_i - \mu_i^{(k)})^2}{\sum_{i=1}^n \lambda_{ii}(\theta^k)} \quad --- (***)$$

(\*) and (\*\*) are the same as that obtained earlier by a direct maximum likelihood approach.

$$\lambda_{ij}^{(k+1)} = \frac{\lambda_j^{(k+1)} \phi(x_i | \theta_j^{(k+1)})}{\sum_{l=1}^L \lambda_l^{(k+1)} \phi(x_i | \theta_l^{(k+1)})} \approx d_{ij}(\theta^{(k+1)})$$

H.W.

(1) Derive the EM algorithm for Bernoulli mixtures

(2) Work out the EM algorithm for HMMs with

two components, both having the same variance

## K-means clustering

Given:  $\{x_1, \dots, x_n\}$   $x_i \in \mathbb{R}^d$

Goal: Partition these points into K clusters,  
where K is given.

Let  $\mu_j, j=1 \dots K$  be the cluster centers.

To restate the goal:

Find cluster centers  $\mu_j, j=1 \dots K$  such that  
the sum of square distances of each data point  
to closest  $\mu_j$  is minimum

Notation:-  $z_{ij} \in \{0, 1\}$   $i=1 \dots n, j=1 \dots K$

$z_{ij}$  = indicator that takes value 1 if  
 $x_i$  is mapped to cluster  $j$

Objective  $J = \sum_{i=1}^n \sum_{j=1}^K z_{ij} \|x_i - \mu_j\|^2$

Need a clustering scheme that minimizes  $J$ , i.e.,

$$\min_{\substack{z_{ij}, \mu_j \\ i=1 \dots n \\ j=1 \dots K}} J$$

K-means clustering is an iterative scheme with the following steps:

Step I: Fix  $\mu_j, \gamma_{ij}$  and  $\min J$

$\gamma_{ij}$   
 $\gamma_{i,j}$

} fix cluster-head & find the assignment of points to clusters

Step II: Fix  $\gamma_{ij}, \gamma_{i,j}$  and  $\min J$

$\mu_j$   
 $\gamma_j$

} fix the assigned ( $\gamma_{ij}$ ) & find the cluster-heads

Repeat Steps I & II until no change in the assignment;  
(or) after a large # of steps.

How to perform Step I:

$J$  = linear function of  $\gamma_{ij}$

So, each data point can be optimized separately

Since  $J$  is in an additive form.

For point  $x_i$ , minimize  $\sum_{j=1}^K \gamma_{ij} \|x_i - \mu_j\|^2$

Solution:  $\gamma_{ij} = \begin{cases} 1 & \text{if } j = \arg \min_k \|x_i - \mu_k\|^2 \\ 0 & \text{else} \end{cases}$

Now to perform Step II: "Z<sub>ij</sub>'s are fixed"

J: quadratic in  $\mu_j$

$$\nabla_{\mu_j} J = 0 \Rightarrow 2 \sum_{i=1}^n Z_{ij} (x_i - \mu_j) = 0$$

leading to

$$\mu_j = \frac{\sum_{i=1}^n Z_{ij} x_i}{\sum_{i=1}^n Z_{ij}}$$

average of all  
points in cluster j

# of points  
in cluster j

Note! Each round (= step I & II) of k-means clustering  
reduces J & hence k-means converges.

Connection of k-means to EM for GMMs:

GMM with a common covariance matrix  $\epsilon I$  for each component, where  $\epsilon > 0$ , I is identity matrix.  
# components = K.

$$\phi(x | \mu_j, \Sigma_j) = \frac{1}{(2\pi\epsilon)^{d/2}} \exp\left(-\frac{1}{2\epsilon} \|x - \mu_j\|^2\right)$$

$\boxed{\Sigma_j = \epsilon I \ \forall j}$

$$d_{ij} = \frac{\lambda_j \exp(-\|x_i - \mu_j\|^2/2\sigma^2)}{\sum_{k=1}^K \lambda_k \exp(-\|x_i - \mu_k\|^2/2\sigma^2)}$$

As  $\epsilon \rightarrow 0$ ,

$$d_{ij} = \begin{cases} 1 & \text{if } j = \arg\min_k \|x_i - \mu_k\|^2 \\ 0 & \text{else} \end{cases}$$

the assignment that one would obtain from k-means clustering.

Cluster head  $\rightarrow$  averages of points assigned to a particular cluster.

$d_{ij} \rightarrow$  hard assignment.

So, GMMs with common covariance  $\in \Gamma$  would become the K-means clustering algorithm in the limit as  $\epsilon \rightarrow 0$ .

H.W.! Check the log likelihood of complete data becomes the sum of squares ( $\bar{J}$ ) in K-means case i.e., as  $\epsilon \rightarrow 0$