

# Bayesian inference for mixture models

M. Concepción Ausín  
Universidad Carlos III de Madrid

Master in Business Administration and Quantitative Methods  
Master in Mathematical Engineering

# Contents

- 1. Finite mixtures
  - 1.1 Bayesian inference for finite mixtures
- 2. Infinite mixtures
  - 2.1 Bayesian inference for infinite mixtures

# Finite mixtures

A finite mixture of  $k$  densities of the same distribution is a convex combination,

$$f(x | k, \boldsymbol{\rho}, \boldsymbol{\theta}) = \sum_{i=1}^k \rho_i f(x | \boldsymbol{\theta}_i),$$

of densities  $f(x | \boldsymbol{\theta}_i)$ , where  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_k)$  such that  $\sum_{i=1}^k \rho_i = 1$ .

## Remarks

- Mixture models are frequently referred as semi-parametric models as their flexibility allow to approximate non-parametric problems.
- Mixture component do not always have a physical meaning, they can describe complex behaviour of data in different research areas: biology, astronomy, engineering...

# Finite mixtures

- Note that  $E[X^r] = \sum_{i=1}^k \rho_i E[X^r | \theta_i]$
- Computationally intensive methods must be considered for inference in mixture models: MCMC methods, EM algorithm, . . .
- The Bayesian approach using MCMC methods allows us to transform the complex structure of a mixture model in a set of simple structures using latent variables.

# Finite mixtures

## Example (Gaussian mixtures)

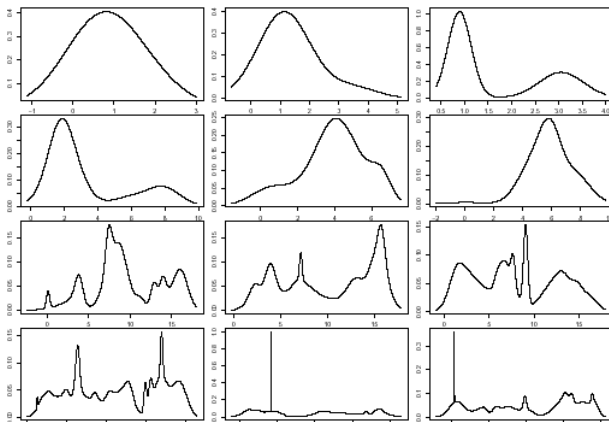
A finite Gaussian mixture of size  $k$  has the following density:

$$f(x | k, \rho, \theta) = \sum_{i=1}^k \rho_i f_N(x | \mu_i, \phi_i),$$

where  $\theta = (\mu_1, \phi_1, \dots, \mu_k, \phi_k)$  and  $f_N(x | \mu_i, \phi_i)$  is the density of a Gaussian distribution with mean  $\mu_i$  and precision  $\phi_i$ .

## Finite mixtures

The following figure shows various density functions of Gaussian mixtures with  $k = 2$  components (first row),  $k = 5$  components (second row),  $k = 25$  components (third row) and  $k = 50$  components (fourth row):



# Bayesian inference for finite mixtures

Assume we have  $n$  observations  $\mathbf{x} = (x_1, \dots, x_n)$  sampled i.i.d. from a finite mixture distribution with density,

$$f(x | \boldsymbol{\rho}, \boldsymbol{\theta}) = \sum_{i=1}^k \rho_i f(x | \boldsymbol{\theta}_i),$$

where  $k$  is finite and known.

We wish to make Bayesian inference for the model parameters  $(\boldsymbol{\rho}, \boldsymbol{\theta})$ . The likelihood is,

$$l(\boldsymbol{\rho}, \boldsymbol{\theta} | \mathbf{x}) = \prod_{j=1}^n \sum_{i=1}^k \rho_i f(x_j | \boldsymbol{\theta}_i),$$

which is given by  $k^n$  terms, which implies a large computational cost for a not very large sample size,  $n$ .

# Bayesian inference for finite mixtures

In order to simplify the likelihood, we can introduce latent variables  $Z_j$  such that:

$$X_j \mid Z_j = i \sim f(x \mid \theta_i) \quad \text{and} \quad P(Z_j = i) = \rho_i.$$

These auxiliary variables allows us to identify the mixture component each observation has been generated from.

Therefore, for each sample of data  $\mathbf{x} = (x_1, \dots, x_n)$ , we assume a missing data set  $\mathbf{z} = (z_1, \dots, z_n)$ , which provide the labels indicating the mixture components from which the observations have been generated.



# Bayesian inference for finite mixtures

Using this missing data set, the likelihood simplifies to:

$$\begin{aligned}l(\boldsymbol{\rho}, \boldsymbol{\theta} \mid \mathbf{x}, \mathbf{z}) &= \prod_{j=1}^n \rho_{z_j} f(x_j \mid \boldsymbol{\theta}_{z_j}) \\ &= \prod_{i=1}^k \rho_i^{n_i} \left[ \prod_{j:z_j=i} f(x_j \mid \boldsymbol{\theta}_i) \right],\end{aligned}$$

where  $n_i = \#\{z_j = i\}$  and  $\sum n_i = n$ .

Then, the posterior probability that the observation  $x_j$  has been generated from the  $i$ -th component is:

$$P(z_j = i \mid x_j, \boldsymbol{\rho}, \boldsymbol{\theta}) = \frac{\rho_i f(x_j \mid \boldsymbol{\theta}_i)}{\sum_{i=1}^k \rho_i f(x_j \mid \boldsymbol{\theta}_i)}.$$

# Bayesian inference for finite mixtures

## Example (Bayesian inference for Gaussian mixtures)

Using the missing data,  $\mathbf{z} = (z_1, \dots, z_n)$  the likelihood simplifies to:

$$l(\boldsymbol{\rho}, \boldsymbol{\mu}, \boldsymbol{\phi} \mid \mathbf{x}, \mathbf{z}) \propto \prod_{i=1}^k (\rho_i \phi_i)^{n_i} \exp\left(-\frac{\phi_i}{2} \sum_{j:z_j=i} (x_j - \mu_i)^2\right),$$

where  $n_i = \#\{z_j = i\}$ .

And we have that:

$$P(z_j = i \mid x_j, \boldsymbol{\rho}, \boldsymbol{\mu}, \boldsymbol{\phi}) = \frac{\rho_i \phi_i \exp\left\{-\frac{\phi_i}{2} (x_j - \mu_i)^2\right\}}{\sum_{i=1}^k \rho_i \phi_i \exp\left\{-\frac{\phi_i}{2} (x_j - \mu_i)^2\right\}}$$

## Bayesian inference for finite mixtures

For the model parameters,  $\rho$ ,  $\mu$  and  $\phi$ , we assume conjugate priors:

Prior	Posterior
$\rho \sim D(\delta_1, \dots, \delta_k)$	$\rho \mid \mathbf{x}, \mathbf{z} \sim D(\delta_1^*, \dots, \delta_k^*)$
$\phi_i \sim G(a/2, b/2)$	$\phi_i \mid \mathbf{x}, \mathbf{z} \sim G(a_i^*/2, b_i^*/2)$
$\mu_i \mid \phi_i \sim N\left(m_i, \frac{1}{\alpha_i \phi_i}\right)$	$\mu_i \mid \mathbf{x}, \mathbf{z}, \phi_i \sim N\left(m_i^*, \frac{1}{\alpha_i^* \phi_i}\right)$

where

$$\delta_i^* = \delta_i + n_i, \quad a_i^* = a + n_i,$$

$$b_i^* = b + \sum_{j:z_j=i} (x_j - \mu_i)^2, \quad \alpha_i^* = \alpha_i + n_i,$$

$$m_i^* = \frac{\alpha_i m_i + n_i \bar{x}_i}{\alpha_i + n_i}, \quad \text{where } \bar{x}_i = \frac{1}{n_i} \sum_{j:z_j=i} x_j.$$

For identifiability reasons, we assume that  $\mu_1 < \dots < \mu_k$ .

# Bayesian inference for finite mixtures

Note that  $D(\delta_1, \dots, \delta_k)$  a Dirichlet distribution with density:

$$f(\rho_1, \dots, \rho_k) \propto \prod_{i=1}^k \rho_i^{\delta_i - 1}.$$

The usual prior choice is to take  $(\delta_1, \dots, \delta_k) = (1, \dots, 1)$  to impose a uniform prior over the mixture weights.

Note that this prior choice is equivalent to use the following reparameterization:

$$\begin{aligned}\rho_1 &= \eta_1, \\ \rho_i &= (1 - \eta_1) \dots (1 - \eta_{i-1}) \eta_i\end{aligned}$$

assuming that  $\eta_i \sim \mathcal{B}(1, k - i + 1)$ .

# Bayesian inference for finite mixtures

## MCMC algorithm

1. Set initial values  $\eta^{(0)}, \mu^{(0)}$  and  $\phi^{(0)}$ .
2. Update  $\mathbf{z}$  sampling from  $\mathbf{z}^{(j+1)} \sim \mathbf{z} | \mathbf{x}, \rho^{(j)}, \mu^{(j)}, \phi^{(j)}$ .
3. Update  $\eta$  sampling from  $\eta^{(j+1)} \sim \eta | \mathbf{x}, \mathbf{z}^{(j+1)}$ .
4. Update  $\phi_i$  sampling from  $\phi_i^{(j+1)} \sim \phi_i | \mathbf{x}, \mathbf{z}^{(j+1)}$ .
5. Update  $\mu_i$  sampling from  $\mu_i^{(j+1)} \sim \mu_i | \mathbf{x}, \mathbf{z}^{(j+1)}, \phi_i^{(j+1)}$ .
6. Order  $\mu^{(j+1)}$  and arrange  $\rho^{(j+1)}$  y  $\phi^{(j+1)}$  with this order.
7.  $j = j + 1$ . Go to 2.

## Infinite mixtures

Now, consider an infinite mixture of densities of the same distribution,

$$f(x | \boldsymbol{\rho}, \boldsymbol{\theta}) = \sum_{i=1}^{\infty} \rho_i f(x | \boldsymbol{\theta}_i),$$

of densities  $f(x | \boldsymbol{\theta}_i)$ , where  $\boldsymbol{\rho} = (\rho_1, \rho_2, \dots)$  such that  $\sum_{i=1}^{\infty} \rho_i = 1$ .

Suppose that we reparametrize the weights such that:

$$\begin{aligned}\rho_1 &= \eta_1, \\ \rho_i &= (1 - \eta_1) \dots (1 - \eta_{i-1}) \eta_i\end{aligned}$$

and assume a priori that:

$$\begin{aligned}\eta_i &\sim \mathcal{B}(1, \alpha), \\ \boldsymbol{\theta}_i &\sim P_0,\end{aligned}$$

for  $i = 1, 2, \dots$

## Infinite mixtures

Note that this infinite mixture model with the considered prior choice corresponds to a Dirichlet process mixture model (DPM model) given by,

$$\begin{aligned}X_i | \theta_i &\sim f(x | \theta_i), \\ \theta_i | P &\sim P(\theta) \\ P &\sim DP(\alpha, P_0)\end{aligned}$$

or equivalently, using the stick-breaking representation,

$$\begin{aligned}x_j | z_j &\sim f(x | \theta_{z_j}) \\ \Pr(z_j = i) &= \rho_i, \\ \theta &\sim P_0 \\ \rho_1 &= \eta_1, \quad \rho_i = (1 - \eta_1) \dots (1 - \eta_{i-1}) \eta_s \\ \eta_i &\sim \mathcal{B}(1, \alpha)\end{aligned}$$

# Bayesian inference for infinite mixtures

Observe that even using the latent variables,  $z_j$ , the likelihood is complicated:

$$l(\boldsymbol{\rho}, \boldsymbol{\theta} \mid \mathbf{x}, \mathbf{z}) = \prod_{i=1}^{\infty} \rho_i^{n_i} \left[ \prod_{j: z_j=i} f(x_j \mid \boldsymbol{\theta}_i) \right],$$

where  $n_i = \#\{z_j = i\}$  and  $\sum n_i = n$ .

And the posterior probability that the observation  $x_j$  has been generated from the  $i$ -th component is difficult to evaluate:

$$P(z_j = i \mid x_j, \boldsymbol{\rho}, \boldsymbol{\theta}) = \frac{\rho_i f(x_j \mid \boldsymbol{\theta}_i)}{\sum_{i=1}^{\infty} \rho_i f(x_j \mid \boldsymbol{\theta}_i)}.$$



# Bayesian inference for infinite mixtures

To solve this problem, Walker (2007) proposes to introduce a new set of latent variables,  $\mathbf{u} = (u_1, \dots, u_n)$  such that,

$$f(x_j, u_j | \boldsymbol{\rho}, \boldsymbol{\theta}) = \sum_{i=1}^{\infty} I(u_j < \rho_i) f(x_j | \theta_i),$$

where  $I$  is the indicator function. Observe that integrating over  $u_j$  the marginal density is  $f(x | \boldsymbol{\rho}, \boldsymbol{\theta})$ . Also note that we can write,

$$f(x_j, u_j | \boldsymbol{\rho}, \boldsymbol{\theta}) = \sum_{i=1}^{\infty} \rho_i f_U(u_j | 0, \rho_i) f(x_j | \theta_i),$$

where  $f_U$  is the density of a uniform  $U(0, \rho_i)$ . Then, with probability  $\rho_i$ , the auxiliary variable  $u_j$  follows a uniform distribution in  $(0, \rho_i)$  and the variable  $x_j$  follows the density  $f(x_j | \theta_i)$ .

# Bayesian inference for infinite mixtures

With this new set of latent variables, the complete likelihood function is,

$$l(\boldsymbol{\rho}, \boldsymbol{\theta} \mid \mathbf{x}, \mathbf{u}, \mathbf{z}) \propto \prod_{j=1}^n l(u_j < \rho_{z_j}) f(x_j \mid \theta_{z_j}).$$

And the posterior probability that the observation  $x_j$  has been generated from the  $i$ -th component is:

$$P(z_j = i \mid x_j, u_j, \boldsymbol{\rho}, \boldsymbol{\theta}) = \frac{f(x_j \mid \theta_i)}{\sum_{i: \rho_i > u_j} f(x_j \mid \theta_i)}.$$

# Bayesian inference for infinite mixtures

Given  $\rho$ , the posterior distribution of  $u_j$  is:

$$u_j \sim U(0, \rho_{z_j}),$$

for  $j = 1, \dots, n$ , where  $\rho_{z_j} = (1 - \eta_1) \dots (1 - \eta_{z_j-1}) \eta_{z_j}$ .

Given  $\mathbf{z}$ , the posterior distribution of  $\boldsymbol{\eta}$  is:

$$\eta_j | \mathbf{z} \sim \text{Beta} \left( n_s + 1, n - \sum_{l=1}^s n_l + \alpha \right)$$

where  $n_i = \sum_{j=1}^n I(z_j = i)$ .

Clearly, assuming a conjugate prior,  $P_0$ , for all  $\theta_i$ , the conditional posterior distribution of  $\theta_i$  given  $\mathbf{z}$  is straightforward to obtain.

# Bayesian inference for infinite mixtures

## MCMC algorithm

1. Set an initial allocation  $\mathbf{z} = \{z_1, \dots, z_n\}$ .
2. Update  $\eta_i$  by simulating from the beta distribution for  $i = 1, \dots, z^*$ , where  $z^* = \max\{z_j\}_{j=1}^n$ .
3. Update  $u_j$  by simulating from  $u_j \sim U(0, \rho_{z_j})$  for  $j = 1, \dots, n$ .
4. Update  $\eta_i$  by simulating from  $\eta_i \sim \text{Beta}(1, \alpha)$  for  $i = z^* + 1, \dots, s^*$ , where  $s^*$  is the smallest value such that:  $\sum_{i=1}^{s^*} \rho_i > 1 - u^*$  where  $u^* = \min\{u_1, \dots, u_n\}$ .
5. Update  $\theta_i$  by simulating from its conditional posterior distribution for  $i = 1, \dots, s^*$ .
6. Update  $z_j$  by simulating from  $Z_j \mid x_j, u_j, \rho, \theta$  for  $j = 1, \dots, n$ .