

## **BSM5** - Bayesian clustering

Lecturer: Riccardo Corradin

### Introduction

Welcome back clustering! Clustering is one of the fundamental techniques in statistical analysis, and any statistician should know the fundamental approaches.

Suppose we have a sequence  $Y_1, \ldots, Y_n$  of observations, where the generic  $Y_i \in \mathbb{Y}$ . Ideally, in cluster analysis, we aim to identify subsets of observations which result to be similar in their observed values.

- $\rightarrow$  Quite a general definition, what does it mean similar?
- $\rightarrow$  Unsupervised learning.

Here we assume that  $\mathbb{Y} \subseteq \mathbb{R}^{p}$ , i.e. we aim to produce homogeneous clusters of real-valued quantities. However, many generalizations can be considered.



### Introduction

We mainly distinguish between two different families of clustering approaches.

- Model-based clustering, where we assume the observations to be distributed according to cluster-specific distributions. Hence, the homogeneity among observations is driven by a probabilistic models.
  - $\rightarrow\,$  Quite flexible approach, the group-specific distribution can be any kind of model and it plays the role of likelihood.
  - ightarrow The group-specific distribution can be useful also to interpret the data behavior.
  - $\rightarrow\,$  In this case, we call the joint model mixture model (different from the mixed model of two slide blocks ago).
- Distance-based clustering, where we use a distance to measure dissimilarities among observations, and clusters are defined by observation closed to each others.
  - ightarrow There are several distances that we can use, depending on specific problems.
  - $\rightarrow~$  There is no interpretation in terms of model structure.
  - $\rightarrow~$  We do not have an usual likelihood function.

In the following, we will study how to perform cluster analysis with a Bayesian approach in the first case, where we assume a cluster-specific distribution.

The key objects to perform model-based clustering are the so-called mixture models.

A mixture model over  $\mathbb Y$  is nothing but a weighted average of different density functions or probability mass functions, of the form

$$f(\mathbf{y}) = \sum_{j=1}^{k} w_j k(\mathbf{y}, \boldsymbol{\theta}_j^*),$$

where  $\theta_i^* \in \Theta$ ,  $j = 1, \ldots, k$ , and  $k \in \mathbb{N} \cup \{\infty\}$ .

• The sequence  $w_1, \ldots, w_k$  is a sequence of non-negative weights taking values in the (k - 1)-dimensional simplex space, i.e.

$$\triangle_{k-1}^{1} = \Big\{ w_{1}, \ldots, w_{k} : 0 \leq w_{j} \leq 1, \ j = 1 \ldots, k, \ \sum_{j=1}^{k} w_{j} = 1 \Big\}.$$

- $k(y, \theta)$  is a kernel function satisfying the followings:
  - $k(, \theta)$  is a density function or a probability mass function for any value of  $\theta \in \Theta$ ;
  - k(y, ) is measurable for any value of  $y \in \mathbb{Y}$ .





#### Pros

- Mixture models are flexible and can capture many complex distributional behavior, such as multimodality and skewness.
- They are naturally tailored to perform model based clustering, as each component can represent a single cluster.
- Despite their simple specification, they can accommodate many different type of data by suitably specifying the kernel function.
  - $\rightarrow k(\mathbf{y}, \boldsymbol{\theta}) \stackrel{d}{=} \textit{N}(\mu, \sigma)$  to cluster univariate data defined on  $\mathbb{R}$ .
  - $\rightarrow k(\mathbf{y}, \mathbf{\theta}) \stackrel{d}{=} N(\mathbf{\mu}, \mathbf{\Sigma})$  to cluster univariate data defined on  $\mathbb{R}^d$ .
  - $\rightarrow k(\mathbf{y}, \theta) \stackrel{d}{=} GP(\mu(t), R(t, t'))$  to cluster functional data taking values on  $C_{\mathbb{X}}$ , space of continuous function with support  $\mathbb{X}$ .

$$\rightarrow k(\mathbf{y}, \mathbf{\theta}) \stackrel{d}{=} ERGM(\mathbf{\theta})$$
 to cluster graph data.

 $\rightarrow \ldots$ 

#### Cons

- The estimation is usually computationally intensive.
- They might be prohibitive with large sample sizes.
- Finite mixture models can be sensible to their specification.

To complete the model specification, we need to assume a specific weights distribution  $\boldsymbol{w} \sim \pi(\boldsymbol{w})$ , the kernel function  $k(\boldsymbol{y}, \boldsymbol{\theta})$ , and a distribution for the cluster-specific parameters  $\theta_i^*$ .

• One common choice for the weights distribution is the Dirichlet distribution, i.e.  $w \sim Dirichlet(\alpha)$  with

$$\pi(\boldsymbol{w}) = \frac{\Gamma(\alpha_1 + \cdots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \prod_{j=1}^k w_j^{\alpha_j - 1},$$

having support  $\triangle_{k-1}^1$ , the k-1-dimensional simplex with total mass 1.

- $\rightarrow$  Recall that  $E[W_j] = \frac{\alpha_j}{\alpha^+}$  and  $var(W_j) = \frac{\alpha_j(\alpha^+ \alpha_j)}{\alpha^+(\alpha^+ + 1)}$ , where  $\alpha^+ = \sum_{j=1}^k \alpha_j$ .
- $\rightarrow$  Different choices of  $\alpha_1, \ldots, \alpha_k$  lead to different specification.
- $\rightarrow$  The case  $\alpha_1 = \cdots = \alpha_k = \alpha$  is called symmetric.
- $\rightarrow$  In the symmetric case, common settings are  $\alpha = 1$  or  $\alpha = \frac{1}{k}$ .
- Specific kernel function choices depend on the type of data we are analyzing.
  - $\rightarrow~$  If the data are discrete, continuous, functions, etc.
  - $\rightarrow\,$  If we want a particular feature from the kernel function, such as skewness, heavy tails, etc.
- The distribution of  $\theta_i^*$ , denoted by  $\pi(\theta_i^*)$ , depends on the kernel assumption.

Ideally, we are in a situation where **our observed data** are distributed according to a **mixture model** 

$$\begin{aligned} \mathbf{Y}_1, \dots, \mathbf{Y}_n \mid \mathbf{w}, \boldsymbol{\theta}_{1:k}^* \stackrel{\text{iid}}{\sim} f(y) &= \sum_{j=1}^k w_j k(\mathbf{y}, \boldsymbol{\theta}_j^*) \\ \mathbf{w} &\sim Dir(\boldsymbol{\alpha}), \\ \boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_k^* \stackrel{\text{iid}}{\sim} \pi(\boldsymbol{\theta}), \end{aligned}$$

where  $\pi(\theta)$  is a suitable prior for the component-specific parameter, matching its support.

#### A first augmentation

Equivalently, we can consider a sequence of latent parameters  $\theta_1, \ldots, \theta_n$ , each associated to a specific observations, having

We can easily see from the previous that marginalizing with respect to a generic  $\theta_i$  the joint distribution of  $\mathbf{Y}_i, \theta_i \mid \mathbf{w}, \theta_{1:k}^*$ , we recover the distribution of  $\mathbf{Y}_i \mid \mathbf{w}, \theta_{1:k}^*$ .

In practice, we are matching each observation with a latent parameter, with the distributional assumption

$$\boldsymbol{ heta}_1,\ldots,\boldsymbol{ heta}_n\mid \boldsymbol{w}, \boldsymbol{ heta}_{1:k}^* \stackrel{\mathrm{iid}}{\sim} \sum_{j=1}^k w_j \delta_{\boldsymbol{ heta}_j}^*(\boldsymbol{ heta}).$$

- The previous is a case of random probability measure, since
  - $\rightarrow$  it is a probability measure;
  - $\rightarrow$  both weights **w** and atoms  $\boldsymbol{\theta}^*_{1:k}$  are random quantities.
- Further, the previous is a discrete random probability measure.
  - ightarrow The sequence of latent parameters  $oldsymbol{ heta}_1,\ldots,oldsymbol{ heta}_n$  can then have ties, meaning that

$$P(\boldsymbol{\theta}_i = \boldsymbol{\theta}_j) > 0, \qquad i \neq j.$$

- $\rightarrow$  Ties implies that the observations can be grouped together when the corresponding latent parameters take the same unique values  $\theta_j^*$ , hence they are generated from the same kernel function.
- $\rightarrow\,$  This is naturally inducing a clustering within the data, where clusters mean that observations share the same distribution in a model-based case.

Ideally, we have a latent partition induced by ties among latent parameters, defining equivalence classes. Unfortunately, the space spanning possible partitions is complex and grows rapidly as far the sample size increases.



- A partition, here denoted by ρ<sub>n</sub>, is then a set of blocks {A<sub>1</sub>,..., A<sub>k</sub>}, where the generic *j*th element A<sub>j</sub> = {*i* ∈ {1,..., n} s.t. S<sub>j</sub> = *j*}, *j* = 1,..., *k*.
  - The sets are disjoint, i.e.  $A_i \cap A_j = \emptyset$ , for  $i \neq j$ . This means that an element can belong only to a single cluster.
  - The union of all the blocks recover the set of observed indices, i.e.  $A_1 \cup A_2 \cup \cdots \cup A_k = \{1, \dots, n\}.$
- We denote by  $n_j$  the size of each block, i.e.  $n_j = |A_j| = \sum_{i=1}^n \mathbf{1}_{[S_i=j]}$ .

We have

$$A_{n,k} = \frac{1}{k!} \sum_{j=0}^{k} (-1)^{j} {\binom{k}{j}} (k-j)^{n}$$

(Stirling number of the second kind) ways to partition n elements in k groups, and

$$B_n = \sum_{k=1}^n A_{n,k}$$

(Bell number) ways to partition *n* elements, which explodes as far *n* grows.

Scan the entire partitions space is unfeasible in a reasonable time, as far as  $n \nearrow$ .

n	1	2	3	4	5	6	7	8	9	10
Bn	1	2	5	15	52	203	877	4140	21147	115975

In practice, we produce a sample of the posterior distribution obtaining a set of partitions which are representative of the latent cluster of the data.

While the previous augmentation has a nice interpretation, in practice is more convenient to work with another one.

#### A second augmentation

We can consider a sequence of latent indicators  $S_1, \ldots, S_n$ , one for each observations and describing which component is the observation associated to, having

$$\begin{split} \mathbf{Y}_i \mid S_i, \boldsymbol{\theta}_{1:k}^* \sim k(\mathbf{y}_i, \boldsymbol{\theta}_{S_i}^*), & i = 1, \dots, n \\ S_1, \dots, S_n \mid \mathbf{w} \stackrel{\text{iid}}{\sim} Cat(w_1, \dots, w_k), \\ & \mathbf{w} \sim Dir(\boldsymbol{\alpha}), \\ \boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_k^* \stackrel{\text{iid}}{\sim} \pi(\boldsymbol{\theta}). \end{split}$$

- Two observations i and j, with i ≠ j, are in the same cluster if they are associated to the same components, i.e. if S<sub>i</sub> = S<sub>j</sub>.
- This different augmentation, even if is not easily interpretable in terms of latent parameter, is quite useful to perform posterior inference with these models.

Also for the second augmentation, we can easily see that marginalizing with respect to a generic  $S_i$  the joint distribution of  $\mathbf{Y}_i, S_i \mid \mathbf{w}, \boldsymbol{\theta}_{1:k}^*$ , we recover the distribution of  $\mathbf{Y}_i \mid \mathbf{w}, \boldsymbol{\theta}_{1:k}^*$ .

With the previous model setting, we are able to derive the full conditional distributions of **the weights** w, the group-specific parameters  $\theta_{1:k}^*$ , and the augmented variables.

Let us assume a mixture model with a generic kernel function  $k(\mathbf{y}, \boldsymbol{\theta})$ , a Dirichlet prior for the weights  $\mathbf{w} \sim Dir(\alpha)$ , and a suitable prior for the group-specific parameters  $\pi(\boldsymbol{\theta})$ . Then we have

$$\begin{split} S_i \mid \mathbf{y}_i, \mathbf{w}, \mathbf{\theta}_{1:k}^* &\stackrel{\text{ind}}{\sim} Cat(w_1k(\mathbf{y}_i, \mathbf{\theta}_1^*), \dots, w_kk(\mathbf{y}_i, \mathbf{\theta}_k^*)), \qquad i = 1, \dots, n, \\ \mathbf{w} \mid S_{1:n} \sim Dir(\alpha_1 + n_1, \dots, \alpha_k + n_k), \\ \mathbf{\theta}_j^* \mid S_{1:n}, \mathbf{y}_{1:n} &\stackrel{\text{iid}}{\sim} \pi(\mathbf{\theta}_j^* \mid S_{1:n}, \mathbf{y}_{1:n}) \propto \pi(\mathbf{\theta}_j^*) \prod_{i: S_j = j} k(\mathbf{y}_i, \mathbf{\theta}_j^*), \qquad j = 1, \dots, k. \end{split}$$

where  $n_j = \sum_{i=1}^n \mathbf{1}_{[S_i=j]}$  is the number of observations assigned to the *j*th component of the mixture.

The previous can guide us in the specification of the prior distributions parameters, in particular

- $\rightarrow$  The Dirichlet distribution should be symmetric, in force of the label switching, whereas each  $\alpha_j$  can be interpreted as the **prior sample sizes** of the generic *j*th cluster.
- $\rightarrow$  The prior distribution  $\pi(\theta_j^*)$  depends on the specific kernel function, but usually is set to be vague or noninformative.

We can exploit the previous full conditionals to derive a sampling strategies to perform posterior inference with those models, by iteratively sampling the quantities we need.

#### Conditional algorithm to perform model-based clustering

```
Input: \alpha_1, \ldots, \alpha_k, parameter of \pi(\theta_i^*). Initial values for \theta_{1:k}^*.
for r = 1 to R do
      for i = 1 to n do
            Sample the cluster allocation of the ith observation, with
                             P(S_i = j \mid \mathbf{y}_{1:n}, \mathbf{w}, \boldsymbol{\theta}_{1:k}^*) \propto w_i k(\mathbf{y}_i, \boldsymbol{\theta}_i^*), \qquad j = 1, \dots, k.
     end
      for i = 1 to k do
           Update the cluster-specific parameters, with
                           \pi(\boldsymbol{\theta}_j^* \mid y_{1:n}, S_{1:n}) \propto \pi(\boldsymbol{\theta}_j^*) \prod k(y_i, \boldsymbol{\theta}_j^*), \qquad j = 1, \dots, k.
     end
```

Update the weights with  $w \sim Dir(\alpha_1 + n_1, \dots, \alpha_k + n_k)$ , where  $n_j$  is the size of the generic *j*th cluster.

#### end

#### Example

Let us consider the following data

set.seed(42)

```
y <- c(rnorm(50, -3, 1), rnorm(25, 3, 1))
```

We want to cluster those observations, assuming a mixture of univariate Gaussian, using a  $N(\mu, \sigma^2)$  as kernel function. Further, we set as prior for the group-specific parameter

$$\mu \mid \sigma^2 \sim N(\mu, \sigma^2/k_0),$$
  
 $\sigma^2 \sim IG(a_0, b_0).$ 

• Show that the previous prior assumption is conjugate for the Gaussian likelihood, with  $\mu, \sigma^2 \mid y_{1:n} \sim NIG(m_n, k_n, a_n, b_n)$ , where

$$k_n = k_0 + n, \qquad m_n = (k_0 m_0 + n \bar{y})/(k_0 + n)$$
  
$$a_n = a_0 + n/2, \qquad b_n = b_0 + \frac{1}{2} \left( \sum_{i=1}^n (y_i - \bar{y})^2 + \frac{k_0 n}{k_n} (\bar{y} - m_0)^2 \right).$$

• implement a Gibbs sampler to sample from the posterior distribution of the latent partition in the data. Check the convergence of the algorithm by looking at the partition entropy at each iteration.

Summarizing our posterior inference

We have a new problem, the so-called label switching.

Once we produce a sample from the posterior distribution of interest, i.e. a sample of latent partitions  $\rho_1, \ldots, \rho_R$ , we want to produce summaries of these sampled values, such as a single point estimate.

Data are assumed to be exchangeable, the model is invariant with respect to permutation of the data or the components.

Two observations can be both assigned to a specific component at a generic iteration r, and both to a different component at iteration  $r + \ell$ .

- From a clustering perspective, they belong to the same cluster, but from a label perspective they have different values.
- In our inferential procedure, we should take in account such a problem.

We can approach the problem from a decision theory point of view. Let  $L(\cdot, \cdot)$  denotes a loss function, having two partitions of *n* elements as arguments. The **optimal partition**  $\rho_n^* \in \mathbb{B}$  (or equivalently  $(S_1^*, \ldots, S_n^*) \in \mathbb{N}^n$ ) is the solution of

$$\begin{split} \rho_n^* &= \operatorname*{arg\,min}_{\hat{\rho}_n \in \mathbb{B}} \left\{ \mathrm{E}[L(\rho_n, \hat{\rho}_n) \mid \boldsymbol{y}_{1:n}] \right\} \\ &= \operatorname*{arg\,min}_{\hat{\rho}_n \in \mathbb{B}} \left\{ \sum_{\rho_n \in \mathbb{B}} L(\rho_n, \hat{\rho}_n) \mathrm{P}(\rho_n \mid \boldsymbol{y}_{1:n}) \right\} \end{split}$$

where  $\rho_n^*$  denotes the optimal partition,  $\hat{\rho}_n$  the partition we are considering,  $\rho_n$  the partitions wrt we take the expectation, and  $P(\rho_n | \mathbf{Y})$  is the posterior probability of  $\rho_n$ , given a set of observations  $\mathbf{y}_{1:n}$ .

- The space  ${\mathbb B}$  is too large
- We can choose  $L(\cdot, \cdot)$  in several ways

A first loss function that we can consider is the 0-1 loss function, i.e.

$$L_{0-1}(\rho_n, \hat{\rho}_n) = \mathbf{1}_{[\rho_n \neq \hat{\rho}_n]}$$

- Such loss function is taking value 1 when all the blocks of ρ<sub>n</sub> and ρ̂<sub>n</sub> coincide, and 0 otherwise.
- Is not accommodating possible similarities in the partitions, all the partitions different from  $\hat{\rho}_n$  are penalized in the same way
- the point estimate coincides with

$$\arg\min_{\hat{\rho}_n \in \mathbb{B}} \left\{ \sum_{\rho_n \in \mathbb{B}} \mathbf{1}_{[\rho_n \neq \hat{\rho}_n]} \mathcal{L}(\rho_n = \rho_n \mid \mathbf{Y}) \right\} = \arg\max_{\hat{\rho}_n \in \mathbb{B}} \left\{ \mathcal{L}(\rho_n = \hat{\rho}_n \mid \mathbf{Y}) \right\}$$

the maximum a posteriori of the distribution of  $\rho_n$ .

#### Summarizing our posterior inference

A more relaxed loss function: the Binder loss function.

$$L_B(\rho_n, \hat{\rho}_n) = \sum_{j < i} \left[ C_1 \mathbf{1}_{[S_i = S_j]} \mathbf{1}_{[\hat{S}_i \neq \hat{S}_j]} + C_2 \mathbf{1}_{[S_i \neq S_j]} \mathbf{1}_{[\hat{S}_i = \hat{S}_j]} \right]$$

and by setting  $C_1 = C_2$  it can be written as

$$L_B(
ho_n, \hat{
ho}_n) = rac{1}{2} \left( \sum_{i=1}^{k_n} n_{i \bullet}^2 + \sum_{j=1}^{\hat{k}_n} n_{\bullet j}^2 - 2 \sum_{j=1}^{\hat{k}_n} \sum_{i=1}^{k_n} n_{ij}^2 
ight)$$

It is based on the cluster frequencies, where

- $k_n$  is the number of blocks in  $\rho_n$
- $\hat{k}_n$  is the number of blocks in  $\hat{\rho}_n$
- $n_{ij} = \sum_{\ell=1}^{n} \mathbf{1}_{[S_{\ell}=i]} \mathbf{1}_{[\hat{S}_{\ell}=j]}$  is the number of observations in the *i*-th block of  $\rho_n$  and in the *j*-th block of  $\hat{\rho}_n$
- $n_{i\bullet} = \sum_{j=1}^{\hat{k}_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[S_\ell = i]}$  is the number of observations in the *i*-th block of  $\rho_n$
- $n_{\bullet j} = \sum_{i=1}^{k_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[\hat{S}_{\ell}=j]}$  is the number of observations in the *j*-th block of  $\hat{\rho}_n$

You can show that the solution of

$$\rho_n^* = \arg\min_{\hat{\rho}_n \in \mathbb{B}} \left\{ \sum_{\rho_n \in \mathbb{B}} L_B(\rho_n, \hat{\rho}_n) \mathrm{P}(\rho_n \mid \boldsymbol{y}_{1:n}) \right\}$$

is given by the partition which, in binary representation, minimize the distance from the posterior similarity matrix  $\mathcal{M}$ , where

$$\mathcal{M}_{ij} = P(S_i = S_j \mid \mathbf{Y})$$

and

$$\boldsymbol{S}^* = \argmin_{\boldsymbol{\hat{S}} \in \mathbb{N}^n} \left\{ \sum_{i < j} |\boldsymbol{1}_{[\hat{S}_i = \hat{S}_j]} - \mathcal{M}_{ij}| \right\}$$

The Binder loss function, differently from the 0-1 loss function, is not penalizing in the same way the partitions different from the optimal one.

## Summarizing our posterior inference

Let's check the previous!

#### Summarizing our posterior inference

A practical example: consider n = 3, and all the possible partitions are

$$\mathbb{B} = \{\{1, 2, 3\}, \{12, 3\}, \{1, 23\}, \{13, 2\}, \{123\}\}$$

Equivalently in a binary representation

$$W_{1} = \begin{bmatrix} 1 \\ 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \qquad W_{2} = \begin{bmatrix} 1 \\ 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \qquad W_{3} = \begin{bmatrix} 1 \\ 0 & 1 \\ 0 & 1 & 1 \end{bmatrix},$$
$$W_{4} = \begin{bmatrix} 1 \\ 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}, \qquad W_{5} = \begin{bmatrix} 1 \\ 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$
and  $\mathcal{M} = \begin{bmatrix} 1 \\ 0.8 & 1 \\ 0.2 & 0.2 & 1 \end{bmatrix}$ , then we have
$$d(W_{1}, \mathcal{M}) = 0.8^{2} + 0.2^{2} + 0.2^{2} = 0.72; \quad d(W_{2}, \mathcal{M}) = 0.08; \quad d(W_{3}, \mathcal{M}) = 0.68;$$
$$d(W_{4}, \mathcal{M}) = 0.68; \quad d(W_{5}, \mathcal{M}) = 1.32,$$

and the point estimate is given by  $R_2$ .

Upon the previous, we can build an algorithm that produce a point estimate of the latent partition, under the Binder loss function, starting from a MCMC sample from the posterior distribution of interest.

#### Point estimate under Binder loss function

**Input:** A sample of partitions from the posterior distribution  $\{\rho^{(1)}, \ldots, \rho^{(R)}\}$ . Produce an estimate of the posterior similarity matrix  $\mathcal{M}$ , where

$$[\mathcal{M}]_{i,j} = \frac{1}{R} \sum_{r=1}^{R} \mathbf{1}_{[S_i^{(r)} = S_j^{(r)}]}$$

for r = 1 to R do

Reconstruct the binary matrix for the rth partition, with

$$[W_r]_{ij} = \begin{cases} 1 : S_i^{(r)} = S_j^{(r)}, \\ 0 : \text{ otherwise} \end{cases}$$

Compute the squared distance  $d_r$  between  $\mathcal{M}$  and  $W_r$ . end

Return the partitions  $R_r$  that minimize the distance.

### Summarizing our posterior inference

#### Example

Let us consider the previous data

set.seed(42)
y <- c(rnorm(50, -3, 1), rnorm(25, 3, 1))</pre>

with the same assumption of the previous example, i.e. using a  $N(\mu, \sigma^2)$  as kernel function. Further, we set as prior for the group-specific parameter

$$egin{aligned} \mu \mid \sigma^2 &\sim \textit{N}(\mu, \sigma^2/k_0), \ \sigma^2 &\sim \textit{IG}(\textit{a}_0, \textit{b}_0), \end{aligned}$$

that a posteriori leads to  $\mu, \sigma^2 \mid y_{1:n} \sim NIG(m_n, k_n, a_n, b_n)$ , where

$$k_n = k_0 + n, \qquad m_n = (k_0 m_0 + n \bar{y})/(k_0 + n)$$
  
$$a_n = a_0 + n/2, \qquad b_n = b_0 + \frac{1}{2} \left( \sum_{i=1}^n (y_i - \bar{y})^2 + \frac{k_0 n}{k_n} (\bar{y} - m_0)^2 \right).$$

- Implement an R function to obtain a point estimate with the Binder loss function, starting from an MCMC output.
- Use the function to find the point estimate with the partitions sampled in the previous example.

The multivariate Gaussian case

A common case while dealing with multivariate real-valued data is to consider a **mixture of multivariate Gaussian distribution**. Suppose we observe a sample  $y_1, \ldots, y_n$ . Here, we want to identify clusters of observations where the generic  $y_i \in \mathbb{R}^p$ ,  $i = 1, \ldots, n$ .

In this framework, each term describing a cluster of data is given by a multivariate Gaussian component, with its own location and covariance matrix.

The kernel function is then given by the density function of a multivariate Gaussian distribution, which is

$$f(\mathbf{y}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{y}_i - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})\right\}$$

 $\bullet\,$  Observations close to the centroid  $\mu$  belong to the cluster with high probability.

- Implicitly, we are looking at the Mahalanobis distance among observations belonging to the same cluster.
- The covariance matrix  $\Sigma$  drives the cluster dispersion and shape.

#### The multivariate Gaussian case

To implement the algorithm we saw earlier in this slide block, we need first to set a distribution on the cluster-specific parameters.

- The location parameter  $\mu$  is a *p*-dimensional real-valued vector, hence we can set a multivariate Gaussian distribution.
- The covariance matrix is a *p* × *p* positive definite real-valued matrix. A possible distributional assumption for this parameter is given by the inverse-Wishart distribution.

Let X be a  $p \times p$  positive definite real-valued matrix. We say that X is distributed as an inverse-Wishart distribution, with  $\nu_0$  degrees of freedom and scale parameter  $\Lambda_0$ , if its density function corresponds to

$$f(\Sigma \mid \nu_0, \Lambda_0) = \frac{|\Lambda_0|^{\nu_0/2}}{2^{\nu_0 \rho/2} \Gamma_{\rho}(\nu_0/2)} |\Sigma|^{-(\nu_0 + \rho + 1)/2} \mathrm{e}^{-\frac{1}{2} \mathrm{tr}(\Lambda_0 \Sigma^{-1})},$$

where  $\Gamma_p(\nu_0/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma(\nu_0/2 + (1-j)/2)$  is the mutivariate gamma function.

Note that, with the previous distributional assumption,

$$\mathbf{E}[\boldsymbol{\Sigma}] = \Lambda_0/(\nu_0 - \boldsymbol{p} - 1),$$

defined for  $\nu_0 > p + 1$ . This might help us to specify the parameters  $\nu_0$  and  $\Lambda_0$ .

#### The multivariate Gaussian case

Specifically, as prior assumption for the main parameter of the model, we consider an hierarchical specification of the form

$$egin{aligned} \mu \mid \Sigma \sim \mathit{N}(\emph{\emph{m}}_{0}, \Sigma/k_{0}), \ \Sigma \sim \mathit{IW}(
u_{0}, \Lambda_{0}). \end{aligned}$$

which is called normal-inverse-Wishart distribution.

Let us assume the data Gaussian distributed. Under the previous prior assumption, we have

 $\mu \mid \Sigma \sim N(\boldsymbol{m}_n, \Sigma/k_n),$  $\Sigma \sim IW(\nu_n, \Lambda_n),$ 

with

$$k_n = k_0 + n, \qquad \mathbf{m}_n = (k_0 \mathbf{m}_0 + n \bar{\mathbf{y}})/k_n$$
  

$$\nu_n = \nu_0 + n, \qquad \Lambda_n = \Lambda_0 + \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^{\mathsf{T}} + \frac{k_0 n}{k_n} (\bar{\mathbf{y}} - \mathbf{m}_0)(\bar{\mathbf{y}} - \mathbf{m}_0)^{\mathsf{T}}.$$

Hence, the normal-inverse-Wishart prior assumption is conjugate to the multivariate Gaussian likelihood.