

Classification

Dhafer Malouche

ESSAI

dhafer.malouche@me.com

2011-2012

Plan

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

Objectif :

- $E = \{w_1, w_2, \dots, w_n\}$ n individus.
- x^1, \dots, x^d d variables observées, quantitatives et/ou qualitatives
- $x_i^j = x^j(w_i)$ l'observation de w_i % à x^j .
- **Question** : Existe-t-il des groupes homogènes d'individus naturellement constitués dans E , et combien ?

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

Définition d'une distance

On définit sur E l'application suivante

$$d : E \times E \rightarrow \mathbb{R}$$

qui peut vérifier

- 1 Symétrie : $d(w_i, w_{i'}) = d(w_{i'}, w_i)$ pour tout $i, i' = 1, \dots, n$.
- 2 Positivité : $d(w_i, w_{i'}) \geq 0$, pour tout $i, i' = 1, \dots, n$.
- 3 Inégalité triangulaire

$$d(w_i, w_{i'}) \leq d(w_i, w_{i''}) + d(w_{i''}, w_{i'}), \quad \forall w_i, w_{i'}, w_{i''}.$$

- 4 Réflexivité : $d(w_i, w_{i'}) = 0 \iff w_i = w_{i'}$.

d est une distance si 1 et 2 sont satisfaites, et d est une distance métrique si de plus 3 et 4 sont satisfaites.

Distances pour des observations quantitatives

- On suppose que nos données sont exclusivement composées d'observations de variables quantitatives.
- La distance euclidienne (norme L_2) : si $x_i = x(w_i) = (x_i^1, \dots, x_i^d)'$ et $x_{i'} = x(w_{i'}) = (x_{i'}^1, \dots, x_{i'}^d)'$, alors

$$d(w_i, w_{i'}) = \sqrt{\sum_{j=1}^d (x_i^j - x_{i'}^j)^2}.$$

- Propriétés : Clusters hypersphériques, clusters composés avec la d.e. sont invariants par rotation et translation, transformations linéaires peuvent causer des distortions dans les clusters, variables avec une grande variance ou des grandes valeurs peuvent dominer les autres variables.
- Nécessité d'une normalisation des données.

D'autres types de distances :

- Minkowski :

$$d_p(w_i, w_{i'}) = \left(\sum_{j=1}^d (x_i^j - x_{i'}^j)^p \right)^{1/p}.$$

$p = 2$ on obtient la distance euclidienne.

- Distance L_∞ :

$$d_\infty(w_i, w_{i'}) = \max_{1 \leq j \leq d} |x_i^j - x_{i'}^j|.$$

- Distance de Mahalanobis

$$d(w_i, w_{i'}) = (x_i - x_{i'})' S^{-1} (x_i - x_{i'})$$

où S est matrice variance des données. Les clusters ont une forme hyperellipsoïdale et invariants par une transformation linéaire. Si $S = I$, alors elle est équivalente à une distance euclidienne.

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique**
- 3 Partitionnement
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

Algorithme ascendant (agrégatif) hiérarchique

- 1 $k = n$ le nombre de classes
- 2 Chaque classes contient une observation : n classes C^1, \dots, C^n .
- 3 Calculer la matrice de distances (entre classes).
- 4 Considérer les deux classes C^i et C^j telles que

$$d(C^i, C^j) = \min d(C^{i'}, C^{j'})$$

- .
- 5 Agréger C^i et C^j dans une seule classe C^{ij} . Donc $k = k - 1$.
 - 6 $k = 1$? Si oui STOP, sinon, aller à 3.

Méthodes d'agrégations

$$d(C^l, C^{ij}) = \alpha_i d(C^l, C^i) + \alpha_j d(C^l, C^j) + \beta d(C^i, C^j) + \gamma |d(C^l, C^i) - d(C^l, C^j)|$$

Algorithme de classification	α_i	α_j	β	γ
<i>Single Linkage</i>	1/2	1/2	0	-1/2
<i>Complete Linkage</i>	1/2	1/2	0	1/2
<i>Centroid Linkage</i>	$\frac{n_j}{n_i+n_j}$	$\frac{n_i}{n_i+n_j}$	$-\frac{n_i n_j}{n_i+n_j}$	0
Méthode de Ward	$\frac{n_i+n_l}{n_i+n_j+n_l}$	$\frac{n_j+n_l}{n_i+n_j+n_l}$	$-\frac{n_l}{n_i+n_j+n_l}$	0

Méthodes d'agrégations

- *Single Linkage* : $d(C^l, C^{ij}) = \min(d(C^l, C^i), d(C^l, C^j))$
- *Complete Linkage* : $d(C^l, C^{ij}) = \max(d(C^l, C^i), d(C^l, C^j))$
- *Centroid Linkage* : $d(C^l, C^{ij}) = d(m_l, m_{ij})$ où m_l et m_{ij} sont resp. les centres de gravité de C^l et C^{ij} .
- Méthode de Ward : elle a pour objectif de minimiser l'accroissement dans la variance intra-classe :

$$E = \sum_{l=1}^k \sum_{x_i \in C^l} d(x_i, m_l)^2.$$

Algorithme descendant (divisif) hiérarchique

- ① $k=0$, Toutes sont dans la même classe C_0
- ② Considérer C_i la classe de diamètre maximale ($d(C_i) = \max_{w, w' \in C_i} d(w, w')$). Toutes sont dans la même classe C_i ,
- ③ Introduire $C_j = \emptyset$.
 - ① Pour chaque $w \in C_i$, calculer

$$d(w, C_i \setminus \{w\}) = \frac{1}{|C_i| - 1} \sum_{w' \in C_i \setminus \{w\}} d(w, w').$$

- ② Calculer la différence $\Delta(w, C_i, C_j) = d(w, C_i \setminus \{w\}) - d(w, C_j)$
- ④ Si $w_m = \operatorname{argmax}_{w \in C_i} d(w, C_i \setminus \{w\})$ et $\Delta(w_m, C_i, C_j) > 0$ alors $w_m \ni C_j$
- ⑤ si $k < n$ alors $k = k + 1$ revient à 1/ sinon STOP.

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement**
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

- **Données** Soient w_1, \dots, w_n n individus pour lesquels on a observé d variables quantitatives x^1, \dots, x^d . Donc chaque w_i est caractérisé par un vecteur d'observations $x'_i = (x_i^1, \dots, x_i^d) \in \mathbb{R}^d$.
- **Objectif** A partir d'une fonction critère J donnée, partitionner w^1, \dots, w_n en K classes (K est fixé d'avance) minimisant ou maximisant la fonction J
- Méthode exhaustive : calculer pour toute les partitions possible C , le critère $J(C)$ et la solution est alors

$$C_{sol} = \operatorname{argmax} J(C).$$

- Solution impossible si n devient grand. En effet le nombre de partitions possible est égal à

$$P(n, K) = \frac{1}{K!} \sum_{m=0}^K \binom{n}{m} (-1)^{K-m} C_K^m m^n.$$

Pour $n = 30$, et $K = 3$, $P(n, K)$ est de l'ordre de 2×10^{14} .

K -means

- Critère à minimiser : la somme des carrés des erreurs (SSE) :

$$J(C) = \sum_{c \in C} \sum_{w \in c} d(w, m(c))$$

où $m(c)$ est le centre de gravité de c .

- Algorithme :

- 1 Initialise K -centres m_1, \dots, m_K et
- 2 Affecter chaque individu vers la classe la plus proche $C = \{c_1, \dots, c_K\}$:

$$w \in c_l \text{ si } d(w, m_l) = \min_{l'=1, \dots, K} d(w, m_{l'}).$$

- 3 Re-calculer les centres de gravité des nouvelles classes,
- 4 Répéter 2 et 3 jusqu'à convergence.

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement
- 4 Indices de validité**
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

Mesure de la validité d'une classification

Deux types de mesures :

- Indices de validité interne : à partir d'un résultat de classification et de l'information intrinsèque dans la base de données on mesure la qualité du résultat :
 - indice de connectivité
 - indice de silhouette
 - indice de *Dunn*
- Mesure de la stabilité de la classification : elle mesure la consistance du résultat d'une classification en la comparant avec une classification obtenue en supprimant chaque fois une colonne (une variable) :
 - proportion moyenne de chevauchement
 - distance moyenne
 - distance moyenne entre les moyennes

Notations

Soit

- n le nombre des observations,
- d le nombre des variables,
- x_1, \dots, x_n les vecteurs d'observations à classer,
- $\mathcal{C} = \{C_1, \dots, C_k\}$ une classification à tester,
- $C(i)$ est la classe contenant x_i , i.e., $C(i) \ni x_i$.

Pour tout x_i , et $j \in \{1, \dots, n\} \setminus \{i\}$: $x_i^{(j)}$ est le $j^{\text{ième}}$ voisin le plus proche de x_i :

$$d(x_i, x_{(1)}) \leq d(x_i, x_{(2)}) \leq \dots d(x_i, x_{(j)}) \leq \dots \leq d(x_i, x_{(n-1)})$$

On pose pour tout $i, j \in \{1, \dots, n\}$

$$\delta(x_i, x_i^{(j)}) = \begin{cases} 0 & \text{si } x_i, x_i^{(j)} \in C(i) \\ 1/j & \text{sinon} \end{cases}$$

```
> data(ruspini)
> X=ruspini[1:5,]
> dist(X)
```

	1	2	3	4
2	10.049876			
3	8.485281	6.403124		
4	24.515301	14.560220	18.027756	
5	9.848858	16.124515	10.440307	28.284271

Donc $x_1^{(1)} = x_3$, $x_1^{(2)} = x_5$, $x_1^{(3)} = x_2$, $x_1^{(4)} = x_4$.

Si $x_1, x_3 \in C_1$ et $x_2, x_4, x_5 \in C_2$ alors

$\delta(x_1, x_3) = \delta(x_1, x_1^{(1)}) = 0$ et $\delta(x_1, x_1^{(4)}) = 1/4$.

Exercice

Comparer $\delta(x_i, x_i^{(j)})$ et $\delta(x_j, x_j^{(i)})$.

Indices de Connectivité

Pour une classification \mathcal{C} ,

$$\text{Conn}(\mathcal{C}) = \sum_{i,j=1,\dots,n} \delta(x_i, x_i^{(j)})$$

$\text{Conn}(\mathcal{C})$ prend ses valeurs dans $[0, +\infty[$ et c'est un coefficient qui doit être minimisée :

$\text{Conn}(\mathcal{C})$ proche de zéro \iff une bonne connectivité à l'intérieure des classes.

R : clValid package, connectivity

Largeur de la silhouette

- C'est la moyenne de toutes les valeurs de la silhouette.
- Valeur de la silhouette pour une observation x_i :

$$S(i) = \frac{b_i - a_i}{\max(b_i, a_i)}$$

où

- a_i est la distance moyenne entre x_i et tous les autres $x_j \in C(i)$.

$$a_i = \frac{1}{n(C(i))} \sum_{x \in C(i)} d(x_i, x).$$

- b_i est la distance moyenne entre x_i et les observations dans la classe la plus proche de x_i :

$$b_i = \min_{C \in \mathcal{C} \setminus C(i)} \frac{1}{n(C)} \sum_{x \in C} d(x_i, x).$$

Largeur de la silhouette

- Définie par

$$S(C) = \frac{1}{n} \sum_i S(i)$$

- La largeur de la silhouette est un nombre qui est toujours compris entre -1 et $+1$.
- S doit être maximisée.
- une valeur de $S(i)$ négative indique que cet individu n'est pas dans "la bonne classe" et il pourrait être déplacé vers la classe la plus proche.
- R, package cluster, silhouette

Indice de *Dunn*

L'indice de *Dunn* est le ratio entre la plus petite distance entre les individus n'appartenant pas aux mêmes classes par rapport à la plus grande distance intra-classe :

$$D(C) = \frac{\min_{C, C' \in \mathcal{C}, C \neq C'} (\min_{x \in C, x' \in C'} d(x, x'))}{\max_{C \in \mathcal{C}} \text{diam}(C)}$$

où

$$\text{diam}(C) = \max_{x, x' \in C} d(x, x')$$

L'indice de *Dunn* prend que des valeurs strictement positives et il doit être maximisé.

R, package `clValid`, `dunn`

Travaux pratiques

- Installer les packages, `clValid`, `cluster` et `clusterGeneration` et charger les dans votre environnement R
- A l'aide de la commande `simClustDesign` générer un jeu de données ayant déjà la répartition réelle des individus. On commencera par un jeu de données de dimension 2, on fera varier l'argument `sepVal` mesurant le degré de "séparabilité" entre les classes.
- Explorer les commandes `plot2DProjection` et `nearestNeighborSepVal`
- Créer des différents jeux de données, calculer les indices internes connectivité, silhouette et Dunn et conclure.

Proportion moyenne de chevauchement

- *Average Proportion of Non-overlap* (APN)
- l'APN mesure la proportion moyenne des individus qui changent de classes quand on effectue une classification en supprimant une variable de la base des données.
- *Calcul* :
 - Soit $\mathcal{C}^0 = \{C_1^0, \dots, C_K^0\}$ une classification obtenue en utilisant toute la base de données.
Pour un individu x_i , on note $C(i)^0$ la classe contenant x_i .
 - L'APN se calcule de la façon suivante

$$APN(K) = \frac{1}{nd} \sum_{i=1}^n \sum_{l=1}^d \left(1 - \frac{|C(i)^l \cap C(i)^0|}{|C(i)^0|} \right)$$

- $APN(K) \in [0, 1]$, des valeurs proches de zéro indiquent une forte consistance ou stabilité dans la classification originale \mathcal{C}^0 .

Proportion moyenne de chevauchement

- *Average Proportion of Non-overlap* (APN)
- l'APN mesure la proportion moyenne des individus qui changent de classes quand on effectue une classification en supprimant une variable de la base des données.
- *Calcul* :
 - Soit $\mathcal{C}^l = \{C_1^l, \dots, C_K^l\}$ une classification obtenue en utilisant la base privée de la l - ième variable.
Pour un individu x_i , on note $C(i)^l$ la classe contenant x_i .
 - L'APN se calcule de la façon suivante

$$APN(K) = \frac{1}{nd} \sum_{i=1}^n \sum_{l=1}^d \left(1 - \frac{|C(i)^l \cap C(i)^0|}{|C(i)^0|} \right)$$

- $APN(K) \in [0, 1]$, des valeurs proches de zéro indiquent une forte consistance ou stabilité dans la classification originale \mathcal{C}^0 .

Trois autres indices de stabilité

- Distance moyenne, *Average Distance* (AD)

$$AD(K) = \frac{1}{nd} \sum_{i=1}^n \sum_{l=1}^d \frac{1}{|C(i)^l| |C(i)^0|} \left(\sum_{x \in C(i)^0, x' \in C(i)^l} \text{dist}(x, x') \right)$$

$AD(K) \in [0, +\infty[$ et des valeurs de $AD(K)$ proches de zéro indique une forte consistance.

- Distance moyenne entre les moyennes, *Average Distance between means* (ADM)

$$ADM(K) = \frac{1}{nd} \sum_{i=1}^n \sum_{l=1}^d \text{dist} \left(\bar{x}_{C(i)^l}, \bar{x}_{C(i)^0} \right)$$

$ADM(K) \in [0, +\infty[$ et des valeurs de $ADM(K)$ proches de zéro indique une forte consistance.

Trois autres indices de stabilité

- Facteur de mérite, *Figure of Merit* (FOM) : pour chaque $l \in \{1, \dots, d\}$, on calcule d'abord

$$\text{FOM}(l, K)^2 = \frac{1}{n} \sum_{C^l \in \mathcal{C}^l} \sum_{x \in C^l} \text{dist}(x, \bar{x}_{C^l})$$

et l'indice FOM est alors égal à

$$\text{FOM}(K) = \sqrt{\frac{n}{n-K} \sum_{l=1}^d \text{FOM}(l, K)}$$

$\text{FOM}(K) \in [0, +\infty[$ et des valeurs de $\text{FOM}(K)$ proches de zéro indique une forte consistance.

Travaux pratiques

- 1 Utilisant la commande `simClustDesign` du package `clusterGeneration` un jeu de données avec un nombre de classes fixe et un nombre de variables égal à 5 (dont une est une variable “bruitée”).
- 2 L'objectif de ce TP est d'utiliser les indices de validation définis ci-dessus pour retrouver le nombre de classes après avoir effectué une classification selon les méthodes : hiérarchiques, k-means et diana.
- 3 Explorer la commande `clValid` du package `clValid` qui permet de calculer à la fois les indices internes et de stabilités.
- 4 Représenter le résultat dans un graphique, donner la méthode adéquate pour retrouver la classification réelle des individus.

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes**
- 6 Méthodes probabilistes pour la classification

Indices externes

- \mathcal{P} une partition existante de l'ensemble des individus E ,
 $\mathcal{P} = \{G_1, \dots, G_r\}$.
- $\mathcal{C} = \{C_1, \dots, C_k\}$ une classification obtenue à l'aide d'un algorithme de classification (*hiérarchique, kmeans...*)
- Objectif : comparer la classification \mathcal{C} et la partition \mathcal{P}
- Définition de quelques indices de validité externes : package `clv` sous R

Calcul des indices externes

Soient deux individus x_i et $x_{i'}$. Quatre cas sont possibles

Cas 1 $\exists C \in \mathcal{C}$ et $\exists G \in \mathcal{P}$ tel que

$$x_i, x_{i'} \in C \text{ et } x_i, x_{i'} \in G$$

Cas 2 $\exists C \in \mathcal{C}$ et $\exists G \neq G' \in \mathcal{P}$ tel que

$$x_i, x_{i'} \in C \text{ et } x_i \in G, x_{i'} \in G'$$

Cas 3 $\exists C \neq C' \in \mathcal{C}$ et $\exists G \in \mathcal{P}$ tel que

$$x_i \in C, x_{i'} \in C' \text{ et } x_i, x_{i'} \in G$$

Cas 4 $\exists C \neq C' \in \mathcal{C}$ et $\exists G \neq G' \in \mathcal{P}$ tel que

$$x_i \in C, x_{i'} \in C' \text{ et } x_i \in G, x_{i'} \in G'$$

Calcul des indices externes

Notons par \mathcal{A}_l l'ensemble des paires $(x_i, x_{i'})$ vérifiant le cas l et par $a_l = |\mathcal{A}_l|$. Trois indices :

- Indice de Rand

$$\text{Rand} = \frac{a_1 + a_4}{a_1 + a_2 + a_3 + a_4}$$

- Indice de Jaccard

$$\text{Jaccard} = \frac{a_1}{a_1 + a_2 + a_3}$$

- Indice de Folkes et Mallows

$$\text{FM} = \sqrt{\frac{a_1}{a_1 + a_3} \frac{a_1}{a_1 + a_2}}$$

Usage sous R

- ① Charger `clv`
- ② Effectuer une classification à l'aide d'une des algorithmes de classification.
- ③ Utiliser la commande `std.ext` pour calculer a_1 , a_2 , a_3 et a_4
- ④ Utiliser les commandes `clv.Rand`, `clv.Jaccard` et `clv.Folkes.Mallows` pour calculer respectivement les indices de Rand, Jaccard et FM. L'argument de ses principales fonctions est l'objet obtenu à l'aide de `std.ext`

Application :

- Appliquer toutes ces commandes sur des données simulées
- Créer une procédure permettant de déterminer les variables minimales construisant les classe.

Travaux pratiques sous R

- 1 Construire une fonction R permettant de calculer tous les indices de validité étudiés.
- 2 Construire une fonction R permettant de permuter la position d'un individu d'une classe à l'autre.
- 3 Construire une fonction R permettant de détecter la stabilité des individus dans une classification.
- 4 Créer des fonctions R permettant d'illustrer tous les résultats des fonctions précédentes à l'aide de graphiques.

- 1 Distances et mesures de proximité
- 2 Classification hiérarchique
- 3 Partitionnement
- 4 Indices de validité
 - Indices de validité interne
 - Indices de stabilité
- 5 Indices Externes
- 6 Méthodes probabilistes pour la classification

l'EM-algorithme

- Soit $\underline{X} = (X_1, \dots, X_n) \sim g(x | \theta)$ un n -échantillon de vecteurs aléatoire dans \mathbb{R}^d . On suppose

$$g(x | \theta) = \int f(x, z | \theta) dz$$

- Objectif : calculer $\hat{\theta} = \operatorname{argmax}_{\theta} L(\underline{x} | \theta) = \operatorname{argmax}_{\theta} \prod_{i=1}^n g(x_i | \theta)$
où $\underline{x} = (x_1, \dots, x_n)$ est une réalisation de \underline{X} .
- On pose $(X, Z) \sim f(x, z | \theta)$, et la densité de $Z | \theta, X = x$ est

$$k(z | \theta, x) = \frac{f(x, z | \theta)}{g(x | \theta)}.$$

l'EM-algorithme

- Donc

$$\log g(x | \theta) = \log f(x, z | \theta) - \log k(z | \theta, x)$$

- Pour un θ_0 fixé,

$$\begin{aligned} \log g(x | \theta)k(z | \theta_0, x) &= \log f(x, z | \theta)k(z | \theta_0, x) \\ &\quad - \log k(z | \theta, x)k(z | \theta_0, x) \quad (1) \end{aligned}$$

- En intégrant (1) par rapport à z , on obtient

$$\log g(x | \theta) = \mathbb{E}_{\theta_0}(\log f(x, Z | \theta)) - \mathbb{E}_{\theta_0}(\log k(Z | x, \theta))$$

l'EM-algorithme

- La fonction vraisemblance s'écrit alors

$$\begin{aligned}\log L(\underline{x} \mid \theta) &= \sum_{i=1}^n \log g(x_i \mid \theta) \\ &= \underbrace{\mathbb{E}_{\theta_0}(\log L^c(\underline{x}, Z \mid \theta))}_{Q(\theta \mid \underline{x}, \theta_0)} - \sum_{i=1}^n \mathbb{E}_{\theta_0}(\log k(Z \mid x_i, \theta))\end{aligned}$$

où L^c est appelée la vraisemblance *complète*.

- Ainsi, dans la recherche du maximum de $\log L(\underline{x} \mid \theta)$, dans l'EM-algorithme on maximize, d'une façon itérative, $Q(\theta \mid \underline{x}, \theta_0)$.

l'EM-algorithme

- 1 Considérons une valeur initiale $\hat{\theta}_{(0)}$
- 2 Calculer l'espérance (*E-step*)

$$Q(\theta \mid \underline{x}, \hat{\theta}_{(m)}) = \mathbb{E}_{\hat{\theta}_{(m)}}(\log L^c(\underline{x}, Z \mid \theta))$$

où la moyenne a été calculée par rapport à $k(z \mid \hat{\theta}_{(m)}, x)$ et $m = 0$

- 3 Maximiser la fonction $Q(\theta \mid \underline{x}, \hat{\theta}_{(m)})$ en θ (*M-step*) :

$$\hat{\theta}_{(m+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \underline{x}, \hat{\theta}_{(m)})$$

et $m = m + 1$.

- 4 Répéter les étapes 2-3 jusqu'à le point fix soit atteint :
 $\hat{\theta}_{(m+1)} = \hat{\theta}_{(m)}$

l'EM-algorithme, dans la classification

- On suppose que les données $\underline{x} = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ est un n -échantillon qu'on suppose issu d'une loi de probabilité

$$g(\underline{x} \mid \theta) = \sum_{l=1}^k \alpha_l g_l(\underline{x} \mid \theta_l)$$

où $\alpha_l \in]0, 1[$, $l = 1, \dots, k$ et $\sum_{l=1}^k \alpha_l = 1$ et $g_l(\underline{x} \mid \theta_l)$ est la densité d'une loi Gaussienne multivariée :

$$f_l(\underline{x} \mid \theta_l) = \frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \exp\left(-\frac{1}{2}(\underline{x} - \mu_l)' \Sigma_l^{-1} (\underline{x} - \mu_l)\right)$$

- Donc $X \sim g(\underline{x} \mid \theta)$ et Z est la variable à valeurs dans $\{1, \dots, k\}$, où k est le nombre de classes et

$$\mathbb{P}(Z = l) = \alpha_l, \quad \forall l = 1, \dots, k.$$

Mise en œuvre sous R : le package `mclust`

- Description détaillé du package dans
 - C. Fraley and A. E. Raftery (2006). MCLUST Version 3 for R : Normal Mixture Modeling and Model-Based Clustering, Technical Report no. 504, Department of Statistics, University of Washington
 - C. Fraley and A. E. Raftery (2002). Model-based clustering, discriminant analysis, and density estimation. Journal of the American Statistical Association 97 :611-631
- Utiliser la commande `Mclust`,
- `Mclust` applique l'EM-algorithme (avec des différentes paramétrisations), et utilise le BIC comme critère de sélection du modèle : chercher le modèle M_k correspondant au BIC maximal.

$$\text{BIC} = 2 \log \text{Vraisemblance} - \text{Nbre paramètres} \times \log(n)$$

Paramétrisations...

identifieur	Model	HC	EM	Distribution	Volume	Shape	Orientation
E		•	•	(univariate)	equal		
V		•	•	(univariate)	variable		
EII	λI	•	•	Spherical	equal	equal	NA
VII	$\lambda_k I$	•	•	Spherical	variable	equal	NA
EEI	λA		•	Diagonal	equal	equal	coordinate axes
VEI	$\lambda_k A$		•	Diagonal	variable	equal	coordinate axes
EVI	λA_k		•	Diagonal	equal	variable	coordinate axes
VVI	$\lambda_k A_k$		•	Diagonal	variable	variable	coordinate axes
EEE	$\lambda D A D^T$	•	•	Ellipsoidal	equal	equal	equal
EEV	$\lambda D_k A D_k^T$		•	Ellipsoidal	equal	equal	variable
VEV	$\lambda_k D_k A D_k^T$		•	Ellipsoidal	variable	equal	variable
VVV	$\lambda_k D_k A_k D_k^T$	•	•	Ellipsoidal	variable	variable	variable