

Rapport TER

Aymeric VAN MAEL - Julien GHEYSENS

Encadré par Adrien HARDY

Vendredi 26 mai 2017

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Mélange Gaussien</b>	<b>3</b>
2.1	Rappels . . . . .	3
2.2	Simulation/Implémentation . . . . .	5
<b>3</b>	<b>Classification non supervisée</b>	<b>7</b>
3.1	L'algorithme EM . . . . .	7
3.1.1	Définition . . . . .	7
3.1.2	Calcul de l'estimation des paramètres . . . . .	7
3.1.3	Influence de l'initialisation . . . . .	11
3.1.4	Implémentation de l'algorithme EM . . . . .	11
3.2	Région de confiance . . . . .	12
3.3	L'algorithme CEM . . . . .	14
3.3.1	Définition . . . . .	14
3.3.2	Classification avec CEM . . . . .	14
3.4	Comparaison entre EM et CEM . . . . .	15
3.5	Choix du modèle . . . . .	15
3.5.1	Le critère BIC . . . . .	16
3.5.2	Le critère AIC . . . . .	17
3.5.3	Comparaison entre BIC et AIC . . . . .	17
3.5.4	Le critère ICL . . . . .	17
<b>4</b>	<b>Classification supervisée</b>	<b>18</b>
4.1	Modélisation . . . . .	18
4.2	Implémentation des Algorithmes . . . . .	18
4.3	Frontière entre les groupes . . . . .	19
4.4	Les deux problèmes pour classer un individu . . . . .	21
<b>5</b>	<b>Exemple avec le jeu de données <i>Iris</i></b>	<b>22</b>
<b>6</b>	<b>Conclusion</b>	<b>24</b>
<b>A</b>	<b>Annexe</b>	<b>26</b>

# 1 Introduction

La classification en analyse de données *exploratoire* et *décisionnelle* est un grand challenge de nos jours. On cherche en effet de plus en plus à regrouper des personnes/objets selon certaines caractéristiques que ce soit dans le domaine médical, afin de définir une thérapie adaptée à un type particulier de malades, ou en astronomie pour catégoriser les étoiles selon des critères de température ou de luminosité.

L'objectif de l'analyse *exploratoire* est de créer une partition en regroupant les données en fonction de leurs similarités et de leurs dissimilarités. Cette partie de l'analyse de données est la classification dite *non supervisée*. Concernant la partie *décisionnelle*, on s'intéresse plutôt à affecter à une des classes préalablement définie, toute nouvelle observation. Cette partie est représentée par la classification *supervisée*.

Pour cela, considérer que les données sont issues d'une loi de mélanges finis de distributions est devenue très populaire ces dernières années. En effet les modèles de mélanges, notamment gaussiens, bénéficient de propriétés avantageuses permettant de traiter les données efficacement : critères de choix de modèles permettant de définir le nombre de lois du mélange, estimations paramétriques multivariées grâce à des algorithmes de recherche de paramètres des lois du mélange.

Nous ferons tout d'abord des rappels concernant les mélanges gaussiens et nous expliquerons comment simuler algorithmiquement ces mélanges. La section suivante abordera la classification non supervisée. Nous supposerons dans un premier temps connu le nombre  $g$  de classes composant le mélange et nous expliquerons le fonctionnement de deux algorithmes d'estimation des paramètres de la loi. Cependant dans la pratique le nombre de groupes  $g$  n'est pas toujours connu. C'est dans cet objectif que nous présenterons deux points de vue, l'un statisticien et l'autre bayésien, qui permettent d'estimer  $g$  à travers des critères comme AIC, BIC ou encore ICL. La section suivante abordera la classification supervisée où nous détaillerons les algorithmes mis en place afin de classer une nouvelle observation, et nous parlerons de la géométrie des groupes, plus précisément de leurs frontières, dans une troisième partie. Nous terminerons par une illustration de nos algorithmes sur un jeu de données *Iris* de RStudio, les observations étant décrites par deux variables : la durée de l'éruption et le temps d'attente entre deux éruptions.

La réalisation de ce rapport a été basé sur l'article de référence [2].

## 2 Mélange Gaussien

### 2.1 Rappels

Tout d'abord, nous allons faire quelques rappels pour bien définir ce qu'est un vecteur gaussien. [6] [7]

**Définition 1.** Soient  $m \in \mathbb{R}$  et  $\sigma \leq 0$ .  $Z$  une variable aléatoire suivant la loi  $N(0,1)$ . On appelle loi gaussienne de paramètre  $m$  et  $\sigma^2$  la loi de la variable aléatoire  $\sigma Z + \mu$ . La densité est alors :

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

On note  $N(\mu, \sigma)$  cette loi.

**Définition 2.** Un vecteur aléatoire  $\mathbf{X}$  est gaussien si toute combinaison linéaire de ses composantes est une variable aléatoire gaussienne.

Les quelques propriétés suivantes seront notamment utiles pour les algorithmes de génération de vecteurs gaussiens.

**Théorème 1.** Soit  $X_1, \dots, X_d$  des variables réelles gaussiennes. Si elles sont mutuellement indépendantes alors le vecteur  $\mathbf{X} = (X_1, \dots, X_d)$  est gaussien

*Démonstration.* Soient  $X_1, X_2, \dots, X_n$   $n$  variables aléatoires mutuellement indépendantes de loi  $X_i \sim \mathcal{N}(a_i \cdot \mu_i, a_i^2 \cdot \sigma_i^2)$ ,  $i \in \{1, \dots, n\}$ . Montrons que  $\mathbf{X} = (X_1, \dots, X_d)$  est un vecteur gaussien, c'est à dire que  $a_1 \cdot X_1 + a_2 \cdot X_2 + \dots + a_n \cdot X_n$  suit une loi normale.

*Rappel :*

Soit  $X_i$  suit une loi normale de paramètres  $\mu_i$  et  $\Sigma_i$ .

Sa fonction caractéristique est donnée par

$$\phi_x(t) = \exp(it\mu_i - \frac{\sigma^2 t^2}{2})$$

On sait que  $a_i X_i \sim \mathcal{N}(a_i \mu_i, a_i^2 \sigma_i^2)$  et que sa fonction caractéristique est :

$$\phi_{a_i X_i}(t) = \exp(ia_i \mu_i t - \frac{a_i^2 \sigma_i^2 t^2}{2})$$

$$\begin{aligned} \text{Donc : } \phi_{\sum_{i=1}^n a_i X_i}(t) &= \mathbb{E}[e^{it \sum_{i=1}^n a_i X_i}] \\ &= \mathbb{E}\left[\prod_{i=1}^n e^{it a_i X_i}\right] \\ &= \prod_{i=1}^n \phi_{a_i X_i}(t) \\ &= \prod_{i=1}^n \exp(ia_i \mu_i t - \frac{a_i^2 \sigma_i^2 t^2}{2}) \\ &= \exp(it \sum_{i=1}^n a_i \mu_i - \frac{t^2}{2} \sum_{i=1}^n a_i^2 \sigma_i^2) \\ &= \phi_{\mathcal{N}(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2)} \end{aligned}$$



## 2.2 Simulation/Implémentation

Maintenant que nous avons bien défini les mélanges gaussiens, nous avons cherché à implémenter sur Rstudio des fonctions permettant de simuler des mélanges de vecteurs gaussiens de  $\mathbb{R}^d$ .

On souhaite simuler un mélange gaussien de  $g$  lois, chacune de moyenne  $m_j \in \mathbb{R}^d$  et de matrice de covariance  $\Sigma_j$ ,  $j \in \{1, \dots, g\}$ .

Pour cela nous allons d'abord créer une fonction `r_VG` afin de simuler un vecteur gaussien de loi  $\mathcal{N}(\mu, \Sigma)$ .

Comme  $\Sigma_j$  est symétrique définie positive, on peut la diagonaliser  $\Sigma$  pour obtenir la décomposition  $\Sigma = P.D.P^t$  (on pourrait également utiliser la décomposition de Cholesky). Avec

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix} \text{ et } \lambda_k \text{ valeurs propres de } \Sigma_j, k \in \{1, \dots, d\}.$$

Nous utilisons alors la fonction "rnorm" de RStudio pour simuler  $d$  réalisations de  $d$  variables gaussiennes  $Y$  centrées et de variance  $\lambda_k$ ,  $k \in \{1, \dots, d\}$ .

En effet  $D$  étant diagonale, pour simuler un vecteur gaussien il suffit de simuler  $d$  variables gaussiennes indépendantes chacune de variance  $\lambda_k$ ,  $k \in \{1, \dots, d\}$ . Le vecteur  $X = (X_1, \dots, X_d)$  suivra alors une loi normale de variance  $D$ .

Il suffit ensuite d'appliquer la transformation  $P.Y + \mu$ . D'après le corollaire 1 on aura bien  $P.Y + \mu \sim \mathcal{N}(\mu, P.D.P^t)$ .

Il est maintenant facile de créer une fonction servant à simuler un échantillon de taille  $n$  de loi de mélange gaussien  $p$ . Pour cela, on suppose connu les poids  $\pi_k$ , les moyennes  $\mu_k$  et les variances  $\Sigma_k$ . On "tire" alors au hasard une loi en utilisant `rmultinom` puis on stocke dans une matrice  $x \in \mathcal{M}_{n,d}$  la réalisation correspondant la  $k$ -ième loi gaussienne  $p_k$ .

Et enfin on souhaite avoir une fonction qui nous donne la densité d'un vecteur gaussien. Elle nous sera utile dans la suivante. Pour cela il nous faut effectuer le raisonnement inverse par rapport à `r_VG`. Ici notre fonction `d_VG` part d'une réalisation de  $X$  de moyenne  $\mu$  et de variance  $\Sigma$  pour introduire une nouvelle variable  $Y$  centrée de variance  $D$  grâce à la transformation :  $Y = P^t.(X - \mu)$ .  $D$  étant diagonale, la valeur de la densité de  $Y$  sera alors le produit des densités de ces marginales, que l'on calcule avec la fonction "dnorm".

Pour illustrer nos algorithmes, voici quelques exemples d'échantillons avec des mélanges à deux et trois lois gaussiennes. La figure 1 représente un mélange de moyennes  $\mu_1 = (-5, 0)$  et  $\mu_2 = (1, 1)$ , et de variance  $\Sigma_1 = \Sigma_2 = Id$ . Et la figure 2 montre un échantillon ayant comme paramètres  $\mu_1 = (-5, 5)$ ,  $\mu_2 = (4, 4)$  et  $\mu_3 = (3, -3)$ , et de variance  $\Sigma_1 = \begin{pmatrix} 2 & 1.5 \\ 1.5 & 2 \end{pmatrix}$ ,  $\Sigma_2 = Id$  et  $\Sigma_3 = \begin{pmatrix} 4 & -0.5 \\ -0.5 & 1 \end{pmatrix}$ .

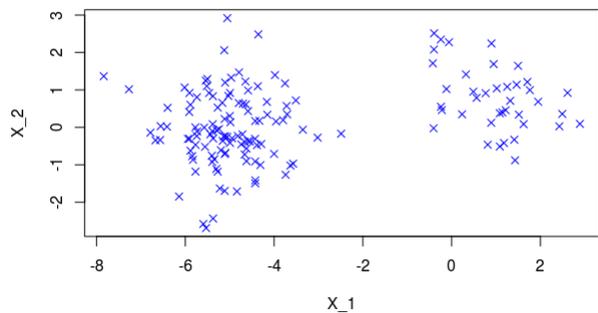


FIGURE 1 – Exemple d'un mélange avec deux lois gaussiennes

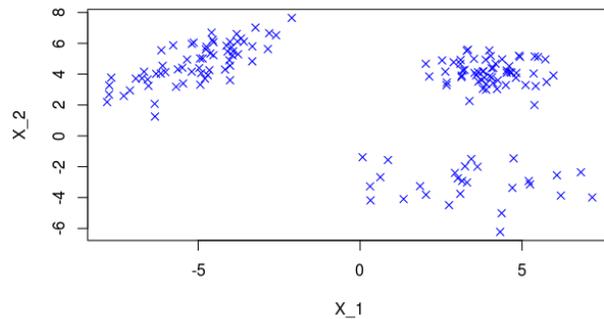


FIGURE 2 – Exemple d'un mélange avec trois lois gaussiennes

### 3 Classification non supervisée

Appelée aussi classification automatique ou regroupement, le but de la classification non supervisée est de partitionner l'ensemble de nos individus en  $g$  groupes  $G_1, \dots, G_g$ . On représente cette partition par un tableau binaire  $\mathbf{z}$  comportant  $n$  lignes et  $g$  colonnes, où  $n$  est le nombre d'individu de notre ensemble. Ce tableau est rempli selon la formule suivante :

$$z_{ik} = \begin{cases} 1 & \text{si } x_i \in G_k \\ 0 & \text{sinon} \end{cases}$$

En outre, la classification non supervisée nous permet de mieux comprendre notre échantillon en regroupant des comportements similaires de différents individus. Pour certains échantillons, il peut être assez facile de trouver les différents groupes à l'oeil nu s'ils sont bien "séparés" (par rapport à une certaine norme) et si les individus sont définis avec deux variables. Mais si les groupes sont "rapprochés" ou que le nombre de variable est plus important, l'analyse peut s'avérer plus difficile.

#### 3.1 L'algorithme EM

Dans cette section, nous allons développer le fonctionnement de l'algorithme *EM* [1]. Notre but est toujours d'estimer les paramètres des différentes lois de notre mélange gaussien, notons cet ensemble  $\theta = \{\pi_1, \dots, \pi_g, \mu_1, \dots, \mu_g, \Sigma_1, \dots, \Sigma_g\}$ . Une représentation graphique sera possible lorsque nous disposons d'un jeu de données caractérisé par deux variables avec la fonction *graph\_EM* en traçant des régions de confiance autour de chaque groupe.

##### 3.1.1 Définition

L'algorithme EM(espérance-maximisation) est un algorithme itératif permettant de trouver les paramètres du maximum de vraisemblance d'un modèle probabiliste. Cette algorithme repose sur deux grandes étapes que l'on réitère un certain nombre de fois.

- Etape (E) : une estimation de l'espérance, où l'on calcule la probabilité conditionnelle qu'un individu  $x_i$  soit issu du groupe  $G_k$ .
- Etape (M) : une estimation du maximum, où l'on maximise l'espérance de la log vraisemblance observée des paramètres en remplaçant la partition  $z$  par l'estimation de l'espérance calculée à l'étape précédente.

##### 3.1.2 Calcul de l'estimation des paramètres

On cherche maintenant à estimer notre paramètre  $\theta = \{\pi_1, \dots, \pi_g, \mu_1, \dots, \mu_g, \Sigma_1, \dots, \Sigma_g\}$  avec l'algorithme EM. On note  $\hat{\theta}_k$  estimation de  $\theta$  à l'étape  $k$  par cet algorithme.

**Proposition 1.** *L'étape (M) de l'algorithme EM consiste à choisir  $\theta_{k+1}$  de la façon suivante :*

$$\begin{aligned}\mu_j^{(k+1)} &= \frac{\sum_{i=1}^n x_i \cdot H_{i,j}(\theta_k)}{\sum_{i=1}^n H_{i,j}(\theta_k)} & 1 \leq j \leq g \\ \Sigma_j^{(k+1)} &= \frac{\sum_{i=1}^n (x_i - \mu_j) \cdot (x_i - \mu_j)^T \cdot H_{i,j}(\theta_k)}{\sum_{i=1}^n H_{i,j}(\theta_k)} & 1 \leq j \leq g \\ \pi_j^{(k+1)} &= \frac{1}{n} \sum_{i=1}^n H_{i,j}(\theta_k) & 1 \leq j \leq g\end{aligned}$$

avec  $H_{i,j}$  est la probabilité conditionnelle que l'individu  $i$  soit issu du groupe  $j$

*Démonstration.* [3] La vraisemblance pour une mélange gaussien s'écrit :

$$V(\theta, X) = \prod_{i=1}^n \left( \sum_{j=1}^g \pi_j f_{\mathcal{N}(\mu_j, \Sigma_j)}(x_i) \right)$$

Et donc sa log vraisemblance vaut :

$$L(\theta, X) = \sum_{i=1}^n \log \left( \sum_{j=1}^g \pi_j f_{\mathcal{N}(\mu_j, \Sigma_j)}(x_i) \right)$$

Nous considérons maintenant  $\mathbf{z}$  défini comme dans le début de cette partie. Grâce à cet élément, nous pouvons définir la log vraisemblance complétée de la façon suivante :

$$\begin{aligned}L(\theta, X, Z) &= \sum_{i=1}^n \sum_{j=1}^g z_{ij} \log (\pi_j f_{\mathcal{N}(\mu_j, \Sigma_j)}(x_i)) \\ &= \sum_{i=1}^n \sum_{j=1}^g z_{ij} \log \left( \pi_j \frac{1}{(2\pi)^{\frac{d}{2}}} |\Sigma_j|^{-\frac{1}{2}} e^{-\frac{(x_i - \mu_j)^t \Sigma_j^{-1} (x_i - \mu_j)}{2}} \right) \\ &= \sum_{i=1}^n \sum_{j=1}^g z_{ij} \left( \log(\pi_j) - \log((2\pi)^{\frac{d}{2}}) + \log(|\Sigma_j|^{-\frac{1}{2}}) - \frac{(x_i - \mu_j)^t \Sigma_j^{-1} (x_i - \mu_j)}{2} \right)\end{aligned}$$

Nous introduisons  $Q(\theta, \theta_k)$  l'espérance conditionnelle de la log vraisemblance complétée sachant  $X$  sous  $\theta_k$ .

$$\begin{aligned}Q(\theta, \theta_k) &= \mathbb{E}_{\theta_k}[L(\theta, X, Z)|X] \\ &= \mathbb{E}_{\theta_k} \left[ \sum_{i=1}^n \sum_{j=1}^g z_{ij} \log (\pi_j f_{\mathcal{N}(\mu_j, \Sigma_j)}(x_i)) |X \right] \\ &= \sum_{i=1}^n \sum_{j=1}^g \mathbb{E}_{\theta_k}[z_{ij}|X] \log (\pi_j f_{\mathcal{N}(\mu_j, \Sigma_j)}(x_i))\end{aligned}$$

Nous cherchons à maximiser  $Q(\theta, \theta_k)$  c'est-à-dire on cherche :

—  $\mu_j^{(k+1)}$  tel que  $\frac{\partial Q(\theta, \theta_k)}{\partial \mu_j} = 0$

—  $\Sigma_j^{(k+1)}$  tel que  $\frac{\partial Q(\theta, \theta_k)}{\partial \Sigma_j} = 0$

—  $\pi_j^{(k+1)}$  tel que  $\frac{\partial Q(\theta, \theta_k)}{\partial \pi_j} = 0$

Nous faisons la notation suivante :  $H_{i,j}(\theta_k) = \mathbb{E}_{\theta_k}[z_{ij}|X]$

Calcul de  $\mu_j^{(k+1)}$  :

$$\begin{aligned} \frac{\partial Q(\theta, \theta_k)}{\partial \mu_j} = 0 &\Leftrightarrow \sum_{i=1}^n \frac{\partial}{\partial \mu_j} (x_i - \mu_j)^T \cdot \Sigma_j^{-1} (x_i - \mu_j) \cdot H_{i,j}(\theta_k) = 0 \\ &\Leftrightarrow \sum_{i=1}^n (x_i - \mu_j)^T \cdot \Sigma_j^{-1} \cdot H_{i,j}(\theta_k) = 0 \\ &\Leftrightarrow \sum_{i=1}^n \mu_j^T \cdot \Sigma_j^{-1} \cdot H_{i,j}(\theta_k) = \sum_{i=1}^n x_i^T \cdot \Sigma_j^{-1} \cdot H_{i,j}(\theta_k) \\ &\Leftrightarrow \sum_{i=1}^n \mu_j \cdot H_{i,j}(\theta_k) = \sum_{i=1}^n x_i \cdot H_{i,j}(\theta_k) \\ \text{donc } \mu_j^{(k+1)} &= \frac{\sum_{i=1}^n x_i \cdot H_{i,j}(\theta_k)}{\sum_{i=1}^n H_{i,j}(\theta_k)} \end{aligned}$$

Calcul de  $\Sigma_j^{(k+1)}$  (grâce à [8]) :

$$\begin{aligned} \frac{\partial Q(\theta, \theta_k)}{\partial \Sigma_j} = 0 &\Leftrightarrow -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \Sigma_j} \ln |\Sigma_j| \cdot H_{i,j}(\theta_k) - \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \Sigma_j} (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j) \cdot H_{i,j}(\theta_k) = 0 \\ &\Leftrightarrow \sum_{i=1}^n \Sigma_j^{-1} \cdot H_{i,j}(\theta_k) - \sum_{i=1}^n (\Sigma_j^{-1})^T \cdot (x_i - \mu_j) \cdot (x_i - \mu_j)^T (\Sigma_j^{-1})^T \cdot H_{i,j}(\theta_k) = 0 \\ &\Leftrightarrow \sum_{i=1}^n H_{i,j}(\theta_k) - \sum_{i=1}^n \Sigma_j^{-1} \cdot (x_i - \mu_j) \cdot (x_i - \mu_j)^T \cdot H_{i,j}(\theta_k) = 0 \\ &\Leftrightarrow \Sigma_j \cdot \sum_{i=1}^n H_{i,j}(\theta_k) = \sum_{i=1}^n (x_i - \mu_j) \cdot (x_i - \mu_j)^T \cdot H_{i,j}(\theta_k) \\ \text{donc } \Sigma_j^{(k+1)} &= \frac{\sum_{i=1}^n (x_i - \mu_j) \cdot (x_i - \mu_j)^T \cdot H_{i,j}(\theta_k)}{\sum_{i=1}^n H_{i,j}(\theta_k)} \end{aligned}$$

Calcul de  $\pi_j^{(k+1)}$  :

On a comme contrainte d'égalité :  $\sum_{j=1}^g \pi_j = 1$

On pose  $u(\pi_1, \dots, \pi_g) = \sum_{j=1}^g -1$

Le Lagrangien devient alors :  $L(\theta, \theta_k) = Q(\theta, \theta_k) - \lambda \cdot u(\pi_1, \dots, \pi_g)$  avec  $\lambda \in \mathbb{R}$

$$\begin{aligned} \frac{\partial Q(\theta, \theta_k)}{\partial \pi_j} - \frac{\partial u(\theta, \theta_k)}{\partial \pi_j} = 0 &\Leftrightarrow \frac{\sum_{i=1}^n H_{i,j}(\theta_k)}{\pi_j} - \lambda = 0 \\ &\Leftrightarrow \pi_j^{(k+1)} = \frac{\sum_{i=1}^n H_{i,j}(\theta_k)}{\lambda} \end{aligned}$$

or les  $\pi_j^{(k+1)}$  doivent satisfaire  $\sum_{j=1}^g \pi_j = 1$

$$\begin{aligned} \sum_{j=1}^g \pi_j = 1 &\Leftrightarrow \frac{1}{\lambda} \cdot \sum_{i=1}^n \sum_{j=1}^g H_{i,j}(\theta_k) = 1 \\ &\Leftrightarrow \frac{1}{\lambda} \cdot \sum_{i=1}^n 1 = 1 \\ &\Leftrightarrow n = \lambda \end{aligned}$$

Donc,  $\pi_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n H_{i,j}(\theta_k)$

□

### 3.1.3 Influence de l'initialisation

Après plusieurs exécutions de notre algorithme EM, on a pu remarquer qu'il nous retourne quelques fois une mauvaise solution sur un même exmple, surtout si deux groupes sont plutôt rapprochés dans notre exmple. Ceci s'explique par le fait que notre convergence est en fait une convergence locale. Donc l'initialisation de notre algorithme joue un rôle sur la convergence de notre algorithme. Pour remédier à ce problème, nous allons faire un algorithme (*MULTI\_EM*) qui va exécuter plusieurs fois notre algorithme EM avec une initalisation différente. Si on s'intéresse à notre choix d'initialisation, on peut faire différents choix possibles. Le choix que nous avons développer est le suivant.

- Pour chaque moyenne  $\mu_j$ , on choisit aléatoirement un individu de notre ensemble (avec la fonction "sample") et on l'affecte à  $\mu_j$ . Ainsi, on aura des moyennes de départ qui ont tendance à être plus proches d'un certain groupe par rapport aux autres moyennes.
- Pour les variances  $\Sigma_j$ , on les choisit toutes égales à une matrice diagonale qui a pour coefficients la variance empirique de chaque variable.
- Pour les  $\pi_j$ , on les rend tous identiques entre eux dans un premier temps, puis on affecte à tour de rôle l'un des poids d'une valeur plus importante.

C'est donc le choix de  $\pi_j$  qui nous permet d'avoir plusieurs tests de l'algorithme EM. Maintenant que nous avons fait plusieurs fois l'algorithme, il nous faut trouver quelle solution est la meilleure. Comme nous cherchons à faire des groupes, notre idée est de choisir la solution qui possède une variance intra minimale. La variance intra d'un groupe est un indice qui permet de connaître l'effet de dispersion des individus de ce même groupe. Si cette variance est minimale on a donc des groupes de plus petites tailles possibles. Et donc nous n'aurons pas un groupe regroupant potentiellement plusieurs ensembles d'individus et des autres groupes regoupant que quelques individus. C'est dans cet objectif que vous avons implémenter l'algorithme Multi EM présent dans l'annexe. Nous avons utiliser la définition de la variance intra suivante.

**Définition 4.** Soit un groupe d'individu  $\Gamma = \{x_i; i = 1, \dots, n\}$  et  $\mu$  son individu moyen. On note par  $\|\cdot\|_2$  la norme euclidienne. L'inertie  $I(\Gamma)$  de ce groupe est :

$$I(\Gamma) = \sum_{i=1}^n \frac{1}{n} \|x_i - \mu\|_2^2$$

Si l'on dispose d'un ensemble disposant de  $g$  groupes  $\Gamma_1, \dots, \Gamma_g$ , la somme des inerties de chaque groupe est appelée inertie intraclasse.

$$I_{intra} = I(\Gamma_1) + \dots + I(\Gamma_g)$$

### 3.1.4 Implémentation de l'algorithme EM

Pour l'implémentation de l'algorithme *EM*, nous avons utilisé les techniques et les outils développés précédemment [A : Outils pour les algorithmes EM et CEM]. Nous allons maintenant discuté comment fonctionne notre algorithme EM, celui-ci est présent en annexe page [A : L'algorithme EM]. Après l'initialisation de nos paramètres, nous avons commencé par définir les valeurs de départ de nos paramètres ( $\pi_j, \mu_j, \Sigma_j$ ) présent dans l'ensemble  $\theta_0$ . Dans la boucle de l'algorithme EM, nous calculons tout d'abord les probabilités conditionnelles que l'individu  $i$  appartienne au groupe  $j$  à partir des valeurs de ( $\pi_j, \mu_j, \Sigma_j$ ) calculées précédemment. Ensuite, nous redéfinissons nos paramètres ( $\pi_j, \mu_j, \Sigma_j$ ) avec les formules de la proposition 1. Nous réalisons cette boucle tant

que la précision  $eps$  de nos paramètres n'est pas atteinte. A la fin de la boucle, nous calculons l'inertie intra de notre jeu de données à partir des paramètres  $(\pi_j, \mu_j, \Sigma_j)$  finaux pour l'algorithme  $Multi\_EM$ . Finalement,  $EM$  nous renvoie une liste comportant toutes les valeurs  $(\pi_j, \mu_j, \Sigma_j)$  stockées dans un vecteur pour le  $\pi_j$  et dans deux matrices différentes pour les  $\mu_j$  et les  $\Sigma_j$ , le nombre d'itération et la valeur de l'inertie intra.

Pour l'implémentation de l'algorithme  $Multi\_EM$ , nous avons comparé les valeurs d'inertie intra de plusieurs réalisations de l'algorithme  $EM$ . Cet algorithme est présent en annexe [A : MULTI\_EM et MULTI\_CEM]. Tout d'abord, nous réalisons une première itération avec un poids identique pour chaque classe de notre ensemble. Puis nous réalisons  $g$  itérations (le nombre de la classe) où pour chaque itération, un groupe a un poids plus important que les autres. A chaque fois, nous comparons les valeurs d'inertie intra pour sélectionner le résultat où l'inertie intra est la plus faible. Ainsi,  $Multi\_EM$  nous renvoie une liste comme l'algorithme  $EM$  dont l'inertie intra est la plus faible parmi les différents résultats des algorithmes réalisés.

### 3.2 Région de confiance

Maintenant que nous avons nos résultats, nous souhaitons faire une représentation graphique. Nous avons donc décidé de tracer des *régions de confiance* autour des groupes que nous avons identifiés. Cette représentation nécessite les valeurs des moyennes et des matrices de variance de chaque groupe, ce qui nous apporte un autre moyen de vérifier la cohérence de nos résultats lorsque l'on dispose de deux variables. Au delà, une représentation graphique peut s'avérer plus difficile.

Tout d'abord, on rappelle qu'avec une seule variance  $X \sim N(m, \sigma)$  le plus petit intervalle tel que  $\mathbb{P}(X \in [a; b]) = \alpha$  est un intervalle centré autour de la moyenne (la preuve se situe en "ANNEXE" page 26).

Avec un vecteur gaussien  $X = (X_1, X_2) \sim N(\mu, \Sigma)$ , notre région de confiance est toujours une courbe de niveau qui a la forme d'une ellipse (la preuve se situe en "ANNEXE" page 27). En effet, pour tracer la région de confiance d'un groupe, on diagonalise la matrice de variance de ce groupe. On obtient alors deux variables indépendantes (corollaire 2). On peut donc appliquer la formule pour obtenir un intervalle de confiance pour chacune des deux variables. Ce qui nous donne les distances du grand et du petit axe d'une ellipse centrée. Ensuite, on applique la matrice de permutation inverse sur notre ellipse centrée et on obtient une région de confiance de la forme d'une ellipse.

Sachant maintenant que notre région de confiance pour un vecteur gaussien est une ellipse, nous aimerions avoir une équation définissant cette région en fonction d'un paramètre  $\alpha$ .

**Proposition 2.** Soit  $X = (X_1, X_2) \sim N(\mu, \Sigma)$  où  $\mu = (\mu_1, \mu_2)$  et  $\Sigma = \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}$ .

Soit  $\epsilon_k$  notre ellipse de confiance d'équation  $(\frac{X_1 - \mu_1}{\sigma_1})^2 + (\frac{X_2 - \mu_2}{\sigma_2})^2 = k^2$ . L'équation de l'ellipse de confiance de niveau  $\alpha$  est :

$$\frac{1}{-2\ln(1 - \alpha)} \left( \left( \frac{X_1 - \mu_1}{\sigma_1} \right)^2 + \left( \frac{X_2 - \mu_2}{\sigma_2} \right)^2 \right) = 1$$

*Démonstration.* En remarquant que  $(\frac{X_1 - \mu_1}{\sigma_1})^2 + (\frac{X_2 - \mu_2}{\sigma_2})^2$  suit une loi du Khi-deux de degré de liberté 2 car c'est la somme de deux carrés de loi normale centrée réduite, on a :

$$\begin{aligned}
 \mathbb{P}(X \in \epsilon_k) &= \mathbb{P}\left(\left(\frac{X_1 - \mu_1}{\sigma_1}\right)^2 + \left(\frac{X_2 - \mu_2}{\sigma_2}\right)^2 \leq k^2\right) \\
 &= \int_0^{k^2} \frac{1}{2^{\frac{1}{2}} \Gamma(\frac{2}{2})} t^0 e^{-\frac{t}{2}} dx \\
 &= \int_0^{k^2} \frac{1}{2} e^{-\frac{t}{2}} dx \\
 &= [-e^{-\frac{t}{2}}]_0^{k^2} \\
 &= \boxed{1 - e^{-\frac{k^2}{2}}}
 \end{aligned}$$

On cherche alors à exprimer le coefficient  $k$  en fonction d'un seuil de confiance  $\alpha$ .

$$\begin{aligned}
 \mathbb{P}(X \in \epsilon_k) &= 1 - e^{-\frac{k^2}{2}} = \alpha \\
 \Leftrightarrow 1 - \alpha &= e^{-\frac{k^2}{2}} \\
 \Leftrightarrow -\ln(1 - \alpha) &= \frac{k^2}{2} \\
 \Leftrightarrow k^2 &= -2\ln(1 - \alpha)
 \end{aligned}$$

On obtient alors l'équation de l'ellipse de confiance suivante :

$$\boxed{\frac{1}{-2\ln(1 - \alpha)} \left( \left( \frac{X_1 - \mu_1}{\sigma_1} \right)^2 + \left( \frac{X_2 - \mu_2}{\sigma_2} \right)^2 \right) = 1}$$

□

Nous savons maintenant tracer une ellipse lorsque nous avons un vecteur gaussien dont les composantes sont indépendantes. Pour pouvoir généraliser notre cas à un vecteur gaussien quelconque de variance  $\Sigma$ , nous devons juste diagonaliser cette matrice  $\Sigma = PDP^{-1}$ . Ainsi, si l'on trace d'abord l'ellipse de confiance pour un vecteur gaussien de variance  $D$  puis nous faisons un changement de base avec la matrice  $P$ , nous aurons donc une ellipse de confiance pour notre vecteur gaussien de variance  $\Sigma$ . C'est dans cet objectif que nous avons implémenter la fonction `ellipse_dim2`.

Nous pouvons donc maintenant tracer des ellipses de confiance de niveau  $\alpha$  autour de nos moyennes estimées de chaque groupe. Nous avons aussi implémenté la fonction `graph_EM` [A : Graphique EM et CEM] pour avoir directement une représentation graphique de notre mélange gaussien avec les ellipses de confiance de chaque groupe. Voici le résultat obtenue sur le mélange à deux lois et un niveau de confiance de 95% :

Avec les ellipses de confiance, nous pouvons avoir un idée sur l'appartenance d'un individu à un certain groupe. Mais nous aimerons savoir exactement à quel groupe un individu est le plus susceptible d'appartenir. C'est dans cet objectif que nous avons programmé l'algorithme `CEM`.

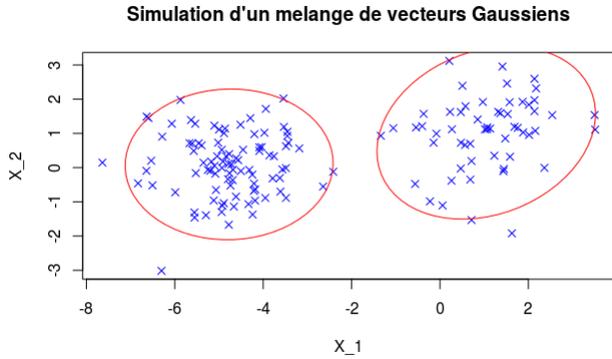


FIGURE 3 – Exemple d’un mélange avec deux lois gaussiennes

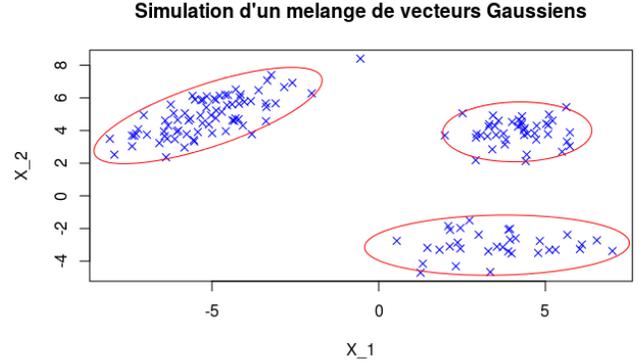


FIGURE 4 – Exemple d’un mélange avec trois lois gaussiennes

### 3.3 L’algorithme CEM

Dans cette section, nous allons développer le fonctionnement de l’algorithme *CEM*. Notre but est toujours d’estimer notre  $\theta = \{\pi_1, \dots, \pi_g, \mu_1, \dots, \mu_g, \Sigma_1, \dots, \Sigma_g\}$  et aussi d’estimer la partition  $z$  qui est une matrice de taille  $n \times g$  où  $n$  est le nombre d’individus de notre jeu de données et  $g$  le nombre de groupe. Ainsi, nous pourrons facilement savoir qui composent nos différents groupes. Une représentation graphique sera possible lorsque nous disposons d’un jeu de données caractérisé par deux variables avec la fonction *graph\_CEM*.

#### 3.3.1 Définition

L’algorithme CEM(classification-espérance-maximisation) fonctionne d’une manière similaire à l’algorithme EM. Cet algorithme repose sur trois grandes étapes que l’on réitère un certain nombre de fois.

- Etape (E) : une estimation de l’espérance, où l’on calcule la probabilité conditionnelle qu’un individu  $x_i$  soit issu du groupe  $G_k$ .
- Etape (C) : une estimation de l’appartenance des individus à l’une des classes en calculant  $\hat{z}$ .
- Etape (M) : une estimation du maximum, où l’on maximise l’espérance de la log vraisemblance observée des paramètres en remplaçant la partition  $z$  par l’estimation de la partition  $\hat{z}$ .

Comme l’initialisation joue un rôle sur notre résultat final, nous avons utilisé le même procédé que pour l’algorithme *EM* en implémentant *Multi\_CEM*. Nous réalisons plusieurs fois l’algorithme CEM avec une initialisation de nos  $\pi_j$  différentes et on sélectionne le résultat possédant une inertie intra la plus faible.

#### 3.3.2 Classification avec CEM

Nous avons donc grâce à l’algorithme CEM une partition  $\hat{z}$  de notre échantillon. Cette partition repose sur le principe MAP (Maximim A Posteriori) affectant chaque individus  $x_i$  à la classe de plus grande probabilité conditionnelle estimée :

$$\hat{z}_{ik} = \begin{cases} 1 & \text{si } k = \operatorname{argmax}_{k'=1, \dots, g} (\mathbb{P}_{\hat{\theta}}(z_{ik'} | X_i = x_i)) \\ 0 & \text{sinon} \end{cases}$$

Nous pouvons utiliser cette partition pour notre représentation graphique. La fonction `graph_CEM` récupère la partition  $\hat{z}$  et associe un symbole à chaque individu  $x_i$  en fonction de son appartenance à une classe, c'est-à-dire en fonction de la position du 1 sur la  $i^{me}$  lignes de  $\hat{z}$ . Elle représente aussi les ellipse de confiance avec les paramètres estimés par l'algorithme `CEM`. L'implémentation de l'algorithme `CEM` se fait de la même manière que pour `EM`. Cet algorithme est présent en annexe A : L'algorithme `CEM`. Dans `CEM`, les paramètres  $(\pi_i, \mu_i, \Sigma_i)$  sont calculés ne sont plus avec les valeurs des probabilités conditionnelles mais avec la partition  $\hat{z}$ . Le reste de l'algorithme `CEM` est identique à l'algorithme `EM`. De plus, notre l'algorithme `MULTI_CEM` [A : `MULTI_EM` et `MULTI_CEM`], fonctionne aussi de la même façon que `MULTI_EM` sauf qu'il utilise `CEM` à la place de `EM`.

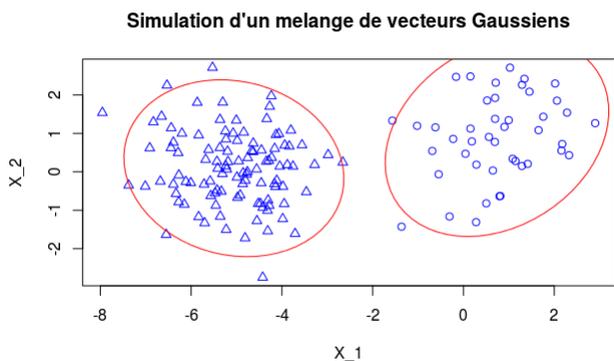


FIGURE 5 – Exemple d'un mélange avec deux lois gaussiennes

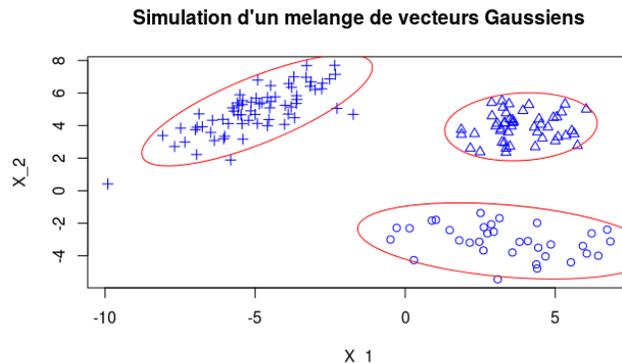


FIGURE 6 – Exemple d'un mélange avec trois lois gaussiennes

### 3.4 Comparaison entre EM et CEM

Outre le fait que l'algorithme `CEM` nous donne une sortie supplémentaire par rapport à l'algorithme `EM` (la partition  $\hat{z}$ ), on peut remarquer quelques différences en pratique entre ces deux algorithmes. En effet, `CEM` a l'avantage d'être un algorithme stationnaire par rapport à `EM` qui a un paramètre de précision pour arrêter la boucle. De plus, nous remarquons que ces deux algorithmes convergent localement vers une solution, d'où l'importance des fonctions `Multi_EM` et `Multi_CEM`. En pratique, nous remarquons que `CEM` a tendance à converger plus rapidement que `EM`. Mais le résultat peut comporter un biais qui est d'autant plus important si les groupes sont fortement imbriqués. Ceci s'explique par le fait que l'algorithme `CEM` remplace les probabilités conditionnelles à l'étape (E) par notre  $\hat{z}$  calculé grâce au principe *MAP*. Par exemple, avec un jeu de données composé de deux classes, si nous avons que la probabilité conditionnelle d'appartenir au premier groupe est de 0.6 et celle du second de 0.4, `EM` utilise ces probabilités comme des poids alors que `CEM` remplace ces probabilités par 1 et 0 pour les premiers et seconds groupes respectivement. Nous pouvons aussi remarquer que `CEM` apporte une meilleure estimation des paramètres si l'échantillon est petit.

### 3.5 Choix du modèle

Dans cette section, nous allons chercher le meilleur modèle pour notre jeu de données. Nous cherchons donc à définir le nombre de classes optimal  $g$  de notre jeu de données. C'est-à-dire nous

cherchons à définir le nombre de classes sans pour autant qu'il ne soit trop excessif. Si nous avons une classe pour chaque individu, notre analyse peut s'avérer inutile. En effet, augmenter le nombre de paramètre améliore nécessairement la qualité de notre analyse. Il nous faut donc un nombre de classe minimale pouvant expliquer au maximum les différentes classes de notre jeu de donnée. Nous appelons ce modèle un modèle parcimonieux. Pour choisir le meilleur modèle répondant à nos attentes, il existe plusieurs critères classiques en statistiques qui se définissent comme une pénalisation de la log-vraisemblance maximale.

$$\star IC_m = L(\hat{\theta}_m; x) - v_m f(n)$$

où  $\hat{\theta}_m$  désigne l'estimation du maximum de la vraisemblance sous le modèle  $m$ ,  $v_m$  le nombre de paramètre à estimer et  $f$  est une fonction de  $n$ . On retient alors modèle  $m$  ayant la plus grande valeur du critère. Il existe plusieurs choix pour la fonction  $f$ . Nous allons développer deux entre elles.

### 3.5.1 Le critère BIC

Le critère BIC [A : Critere AIC ou critere BIC] (Bayesian Information Criterion) propose le choix de  $f$  suivant [5].

$$f(n) = \frac{1}{2} \log(n)$$

BIC se place dans un contexte bayésien (nos paramètres  $\theta_i$  et nos différents modèles à tester  $m_i$  sont vu comme des variables aléatoires). Ce critère cherche à sélectionner le modèle  $m_i$  qui maximise la probabilité *a posteriori*  $\mathbb{P}(m_i|X)$  :

$$m_{BIC} = \arg \max_{m_i} \mathbb{P}(m_i|X)$$

D'après la formule de Bayes, nous pouvons réécrire  $\mathbb{P}(m_i|X)$  comme étant :

$$\mathbb{P}(m_i|X) = \frac{\mathbb{P}(X|m_i)\mathbb{P}(m_i)}{\mathbb{P}(X)}$$

Si l'on suppose que chaque  $\mathbb{P}(m_i)$  ait la même valeur quelque soit le modèle, la recherche du meilleur modèle selon BIC dépend seulement du calcul de  $\mathbb{P}(X|m_i)$ . De plus, avec la méthode d'approximation de Laplace, nous obtenons une valeur approchée de  $\mathbb{P}(X|m_i)$ .

$$\mathbb{P}(X|m_i) = e^{V(m_i, \theta_i^*)} \left( \frac{2\pi}{n} \right)^{v_{m_i}/2} |A_{\theta_i^*}|^{-1/2} + \mathcal{O}(n^{-1})$$

où  $V$  désigne la vraisemblance,  $\theta_i^*$  est l'arg max  $\frac{1}{n} \sum_{k=1}^n \log(V(X_k, \theta_i))$  et  $A_{\theta_i^*}$  l'opposé de la matrice hessienne des dérivées secondes partielles de la fonction  $\frac{1}{n} \sum_{k=1}^n \log(V(X_k, \theta_i))$  en  $\theta_i$ . En remarquant que  $\log(\mathbb{P}(X|m_i)) \approx \log(V(m_i, \hat{\theta}_i)) - \frac{v_{m_i}}{2} \log(n)$ , nous pouvons définir  $BIC_i$  comme :

$$BIC_i = \log(V(m_i, \hat{\theta}_i)) - \frac{v_{m_i}}{2} \log(n)$$

Le modèle sélectionné par ce critère est alors :

$$m_{BIC} = \arg \max_{m_i} BIC_i$$

### 3.5.2 Le critère AIC

Le critère AIC [A : Critere AIC ou critere BIC] (Akaike Information Criterion) propose le choix de  $f$  suivant [4].

$$f(n) = 1$$

L'AIC est basé sur la théorie de l'information. L'objectif de ce critère est de choisir le modèle vérifiant l'expression suivante.

$$m_{AIC} = \arg \max_{m_i} \mathbb{E} \left[ \int \log \left( \frac{f(x)}{V(m_i, \theta_i)} \right) f(x) dx \right]$$

Nous pouvons retrouver ce modèle en maximisant la valeur suivante :

$$m_{AIC} = \arg \max_{m_i} \log(V(m_i, \hat{\theta}_i)) - v_{m_i}$$

### 3.5.3 Comparaison entre BIC et AIC

Le choix d'un des deux critères se porte sur le point de vue que nous adoptons sur le problème. Si nous voulons savoir quelles méthodes nous donnent le meilleur modèle en simulant des données à partir d'un modèle  $m_t$ , le critère BIC a plus tendance à nous retrouver notre modèle de départ  $m_t$ . Mais si nous voulons choisir le critère en fonction de la qualité de la prédiction, le critère AIC s'avère plus performant. En conclusion, nous pouvons dire que le choix d'un critère de sélection de modèles dépend de l'objectif de l'analyse et de la connaissance des données. C'est dans ce but que nous avons décidé de nous intéresser au critère de sélection ICL.

### 3.5.4 Le critère ICL

Le critère ICL [A : Critere ICL] (Integrated Complete Likelihood) a pour philosophie de reporter l'objectif de la classification de la méthode d'estimation vers la méthode de sélection du modèle. En effet, au lieu de pénaliser la log-vraisemblance comme dans AIC ou BIC, le critère ICL pénalise la log-vraisemblance complétée du même terme que le critère BIC. Ainsi, nous définissons ce critère de la manière suivante :

$$\begin{aligned} ICL_{m_i} &= L(\hat{\theta}_{m_i}, x, \hat{z}_{m_i}) - \frac{v_{m_i}}{2} \log(n) \\ &= L(\hat{\theta}_{m_i}, x) - \mathbb{E}(\hat{t}_{m_i}, \hat{z}_{m_i}) - \frac{v_{m_i}}{2} \log(n) \\ &= BIC_i - \mathbb{E}(\hat{t}_{m_i}, \hat{z}_{m_i}) \end{aligned}$$

où  $\hat{t}_{m_i}$  désigne l'estimation de la probabilité conditionnelle qu'un individu appartienne à un certain groupe avec le modèle  $m_i$ .

## 4 Classification supervisée

### 4.1 Modélisation

En classification supervisée, nous disposons en théorie du couple  $(x, z)$  pour notre échantillon, c'est-à-dire que nous connaissons quel individu appartient à quel groupe. Cet échantillon est appelé *ensemble d'apprentissage*. Et nous souhaitons, à l'instar de la régression linéaire, prévoir pour toute nouvelle observation le groupe auquel elle appartient. En définitive, nous cherchons à estimer une règle de classement  $r$  définie par :

$$\begin{aligned} r : \mathcal{X} &\rightarrow \{1, \dots, g\} \\ x_{n+1} &\mapsto r(x_{n+1}) \end{aligned}$$

D'autre part, nous souhaitons également avoir une description des groupes et plus précisément de la frontière qui les sépare.

Pour résoudre ce problème de classification, nous nous appuyons sur les probabilités conditionnelles pour obtenir une règle de classement  $r$ . Cette règle correspond au critère MAP que l'on a vu précédemment :

$$\forall x_{n+1} \in \mathcal{X}, r(x_{n+1}) = \underset{k \in \{1, \dots, g\}}{\operatorname{Argmax}} \mathbb{P}(Z_{n+1, k} = 1 | X_{n+1} = x_{n+1})$$

### 4.2 Implémentation des Algorithmes

Afin de visualiser les résultats de ces règles de classement, nous avons créé des algorithmes qui nous trace les frontières de groupe (figure 8) pour notre exemple avec deux lois.

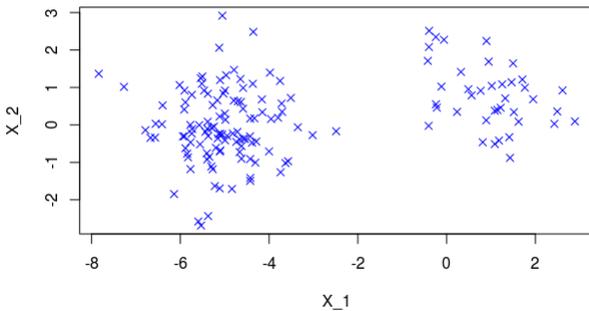


FIGURE 7 – Mélange Gaussien avec deux lois

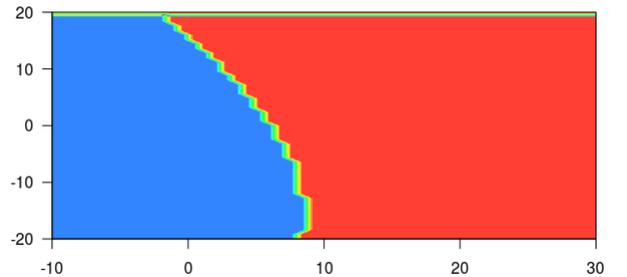


FIGURE 8 – Graphique des frontières

Pour obtenir ce graphique nous avons tracé un maillage de  $2(N + 1)$  points du plan. On souhaite stocker chaque point du maillage dans une matrice  $A \in \mathcal{M}_{N+1, N+1}$  la valeur du groupe auquel appartient l'individu correspondant. Pour cela nous calculons le critère MAP en chaque  $A_{i, j}, \forall i \in \{1, \dots, N + 1\}$ .

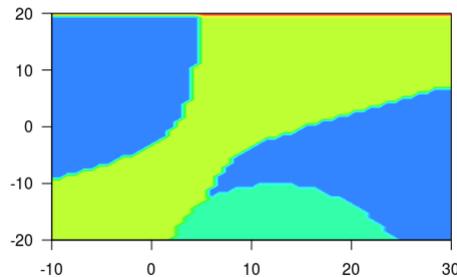
Nous avons donc besoin de deux fonctions, une première nommée *appartenance* qui implémente MAP et une seconde *graph\_coloree* qui crée le maillage et affiche le graphe [A : Graphique colore].

*Appartenance* va prendre en entrée un individu  $x_i$ , le nombre de groupe  $g$  ainsi que notre paramètre  $\theta$ . Ce programme calcule les probabilités conditionnelles  $\mathbb{P}_\theta(Z_{i,k} = 1 | X_i = x_i), \forall k \in \{1, \dots, g\}$  que nous stockons dans un vecteur et retourne le numéro de la classe auquel  $x_i$  appartient avec la fonction *which.max()*.

Quant à *graph\_coloree*, nous donnons le numéro du groupe le plus probable pour tout point de  $\mathbb{R}^2$  dans le pavé défini par les deux paramètres d'entrées *xlim* et *ylim*. Grâce à notre *appartenance*, nous stockons ensuite ces valeurs dans  $A$  afin de pouvoir afficher le graphe avec la fonction *filled.contour* de RStudio.

### 4.3 Frontière entre les groupes

En observant le graphique ci-dessous, on constate que les frontières entre les groupes ont des formes particulières. Nous distinguons une hyperbole (en Jaune et Bleu) et une ellipse (en vert) qui s'ajoute au graphique. Les caractéristiques mathématiques des frontières correspondent à l'égalité des probabilités conditionnelles de classement d'un individu dans un des groupes.



Nous pouvons dans le cas d'un mélange à deux lois, retrouver les équations théoriques régissant les frontières des classes en distinguant deux cas : un premier cas homoscédastique et le second hétéroscédastique.

L'égalité des probabilités conditionnelles s'écrit comme suit :

$$\mathbb{P}(Z_{i1} = 1|X_i = x) = \mathbb{P}(Z_{i2} = 1|X_i = x)$$

La formule de Bayes nous donne ensuite :  $\frac{\mathbb{P}(Z_{i1}=1 \cap X_i=x)}{\mathbb{P}(X_i=x)} = \frac{\mathbb{P}(Z_{i2}=1 \cap X_i=x)}{\mathbb{P}(X_i=x)}$

Il en résulte les équivalences suivantes :

$$\begin{aligned} & \mathbb{P}(Z_{i1} = 1 \cap X_i = x) = \mathbb{P}(Z_{i2} = 1 \cap X_i = x) \\ \Leftrightarrow & \mathbb{P}(X_i = x|Z_{i1}).\mathbb{P}(Z_{i1}) = \mathbb{P}(X_i = x|Z_{i2}).\mathbb{P}(Z_{i2}) \\ \Leftrightarrow & \pi_1.f_1(x) = \pi_2.f_2(x) \\ \Leftrightarrow & \pi_1 \cdot |\Sigma_1|^{-1} \cdot e^{-\frac{1}{2}(x-\mu_1)^t \cdot \Sigma_1^{-1} \cdot (x-\mu_1)} = \pi_2 \cdot |\Sigma_2|^{-1} \cdot e^{-\frac{1}{2}(x-\mu_2)^t \cdot \Sigma_2^{-1} \cdot (x-\mu_2)} \\ \Leftrightarrow & \log\left(\frac{\pi_1}{\pi_2}\right) + \log[\det(\Sigma_1^{-1}) \cdot \det(\Sigma_2)] = \frac{1}{2}(x - \mu_1)^t \cdot \Sigma_1^{-1} \cdot (x - \mu_1) - \frac{1}{2}(x - \mu_2)^t \cdot \Sigma_2^{-1} \cdot (x - \mu_2) \end{aligned}$$

— Si  $\Sigma_1 = \Sigma_2 = \Sigma$  l'équation devient :

$$\begin{aligned} 2 \ln\left(\frac{\pi_1}{\pi_2}\right) &= (x - \mu_1)^t \cdot \Sigma^{-1} \cdot (x - \mu_1) - (x - \mu_2)^t \cdot \Sigma^{-1} \cdot (x - \mu_2) \\ &= (x - \mu_1 - \mu_2 + \mu_2)^t \cdot \Sigma^{-1} \cdot (x - \mu_1 - \mu_2 + \mu_2) - (x - \mu_2)^t \cdot \Sigma^{-1} \cdot (x - \mu_2) \\ &= (x - \mu_2)^t \cdot \Sigma^{-1} \cdot (\mu_2 - \mu_1) + (\mu_2 - \mu_1)^t \cdot \Sigma^{-1} \cdot (\mu_2 - \mu_1) \\ &= 2 \cdot x^t \cdot \Sigma^{-1} \cdot (\mu_2 - \mu_1) - (\mu_1 + \mu_2) \cdot \Sigma^{-1} \cdot (\mu_2 - \mu_1) \end{aligned}$$

L'équation décrivant la frontière entre les deux lois est alors une droite :

$$\boxed{\ln\left(\frac{\pi_1}{\pi_2}\right) = (\mu_2 - \mu_1) \cdot \Sigma^{-1} \cdot x - \frac{1}{2} \cdot (\mu_1 - \mu_2) \cdot \Sigma^{-1} \cdot (\mu_2 + \mu_1)}$$

— Si  $\Sigma_1 \neq \Sigma_2$  (cas hétéroscédastique). L'équation devient une conique :

$$\log\left(\frac{\pi_1}{\pi_2}\right) + \log[\det(\Sigma_1^{-1}) \cdot \det(\Sigma_2)] = \frac{1}{2}(x - \mu_1)^t \cdot \Sigma_1^{-1} \cdot (x - \mu_1) - \frac{1}{2}(x - \mu_2)^t \cdot \Sigma_2^{-1} \cdot (x - \mu_2)$$

Après simplification, on trouve l'équation d'une conique :

$$\boxed{\log\left(\frac{\pi_1}{\pi_2}\right) + \log[\det(\Sigma_1^{-1}) \cdot \det(\Sigma_2)] = \frac{1}{2} \cdot x^t (\Sigma_1^{-1} - \Sigma_2^{-1}) \cdot x + x [\Sigma_2^{-1} \cdot \mu_2 - \Sigma_1^{-1} \cdot \mu_1] \frac{1}{2} (\mu_1 \cdot \Sigma_1^{-1} \cdot \mu_1 - \mu_2 \cdot \Sigma_2^{-1} \cdot \mu_2)}$$

## 4.4 Les deux problèmes pour classer un individu

Deux problèmes quant à l'utilisation de cette méthode apparaissent pour classer une nouvelle observation. En pratique, le paramètre  $\theta$  nous est inconnu et on l'estime avec une des méthodes en classification non supervisée. Outre les limitations de EM/CEM pour des lois trop imbriquées, le fait de vouloir classer un individu trop proche de la frontière ou trop éloigné des moyennes des groupes peut être problématique.

En effet, pour comprendre la difficulté de classement aux frontières, reprenons l'exemple du mélange à deux lois. La figure ci-dessus nous montre la frontière théorique et celle de droite représente la frontière obtenue avec l'algorithme EM et *grahpe\_coloree*. Nous voyons nettement la différence entre les deux qui est due à l'erreur commise lors de l'approximation de  $\theta$ . Nous prendrons alors la liberté de ne pas classer ce type d'individus. On parle dans ce cas de *rejet d'ambiguïté*.

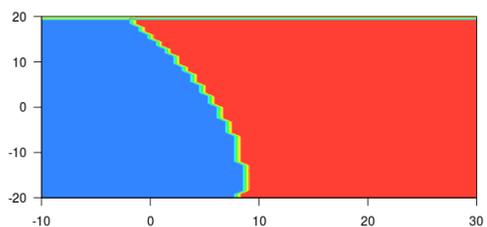
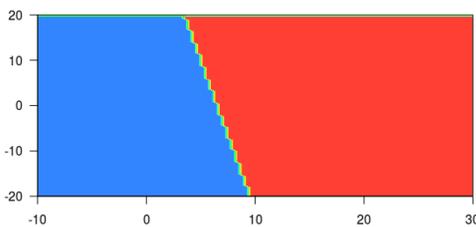


FIGURE 9 – Graphique des frontières théoriques    FIGURE 10 – Graphique des frontières avec EM

Une deuxième forme de rejet consiste à ne pas classer une nouvelle observation dite atypique, c'est-à-dire très en dehors des zones de confiance associées aux groupes. La figure ci-dessus affiche en rouge les ellipses de confiances au seuil de 95 %. Si nous prenons par exemple un nouvel individu  $x_{n+1}$  de coordonnées  $(-5, -5)$ . D'après la figure 10, la règle du MAP classerait sans hésiter  $x_{n+1}$  dans le groupe de gauche. Cependant ce nouveau point étant très éloigné des régions de confiances (fig. ??), il n'est pas très pertinent de l'affecter à l'une des classes car il pourrait provenir d'une nouvelle classe pas encore observée jusqu'à présent. Ce type de rejet est appelé *rejet de distance*.

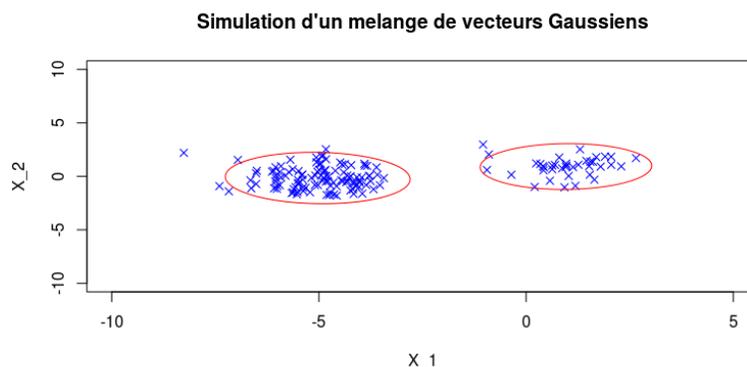


FIGURE 11 – Mélange deux lois avec ellipses confiances à 95%

## 5 Exemple avec le jeu de données *Iris*

Nous récupérons le jeu de données *Iris* qui contient notamment 5 variables : *Sepal.length* / *Sepal.width* / *Petal.length* / *Petal.width* / *Species*. Nous pouvons tracer les deux graphes des largeurs en fonction des longueurs pour les pétales et les sépales. De plus, avec un code couleur nous pouvons voir la réelle appartenance de chaque fleur a son espèce.

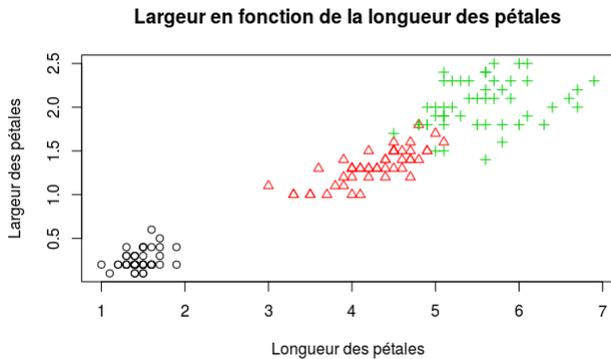


FIGURE 12 – Graphique des individus classés théoriques pour les pétales

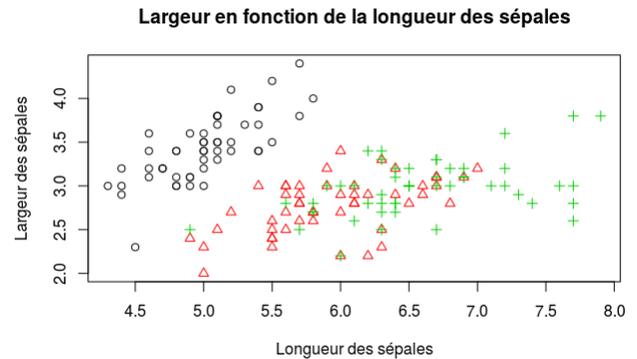


FIGURE 13 – Graphique des individus classés théoriques pour les sépales

Notre objectif est maintenant de retrouver cette classification. Tout d’abord, n’ayant que les valeurs des longueurs et des largeurs, nous cherchons le nombre de classes de l’échantillon. En utilisant *critere\_ICL* sur les pétales et les sépales, nous remarquons que celui-ci nous retourne un nombre de classes qui est égale à deux ou trois. De plus, lorsque nous recherchons le nombre de groupe avec les variables des sépales, nous obtenons régulièrement deux. Cela peut s’expliquer par le fait que les classes représentées en rouge et en vert sur la figure 13 sont imbriquées.

En gardant la valeur de trois pour le nombre de groupe, nous appliquons dans un premier temps l’algorithme *MULTI\_EM* puis *graph\_EM* pour nos deux jeux de données.

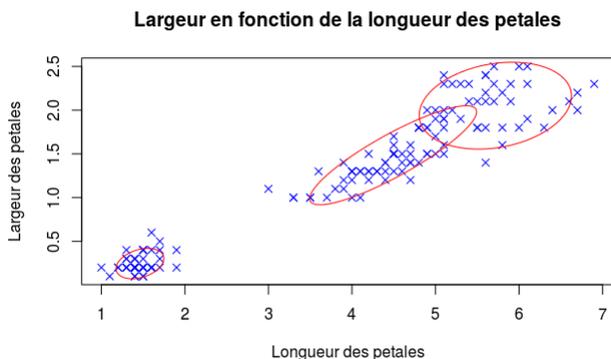


FIGURE 14 – Graphique avec les ellipses de confiance pour les pétales

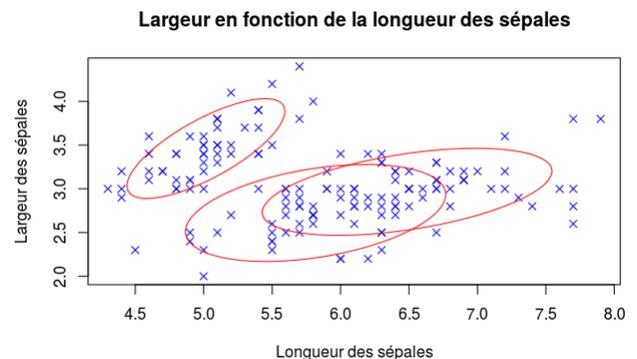


FIGURE 15 – Graphique avec les ellipses de confiance pour les sépales

Nous pouvons dire que l’algorithme *EM* a bien retrouvé les différents groupes de nos ensembles. De plus, nous remarquons d’autant plus le fait que deux classes sont imbriquées au vu du résultat présent sur la figure 15.

Nous pouvons maintenant nous intéresser à l'appartenance de chaque individu à une classe avec l'algorithme *MULTI\_CEM* et *graph\_CEM*.

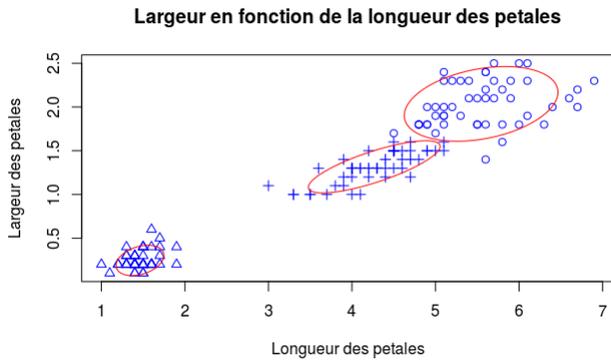


FIGURE 16 – Graphique avec les ellipses de confiance et classification de chaque individu pour les pétales

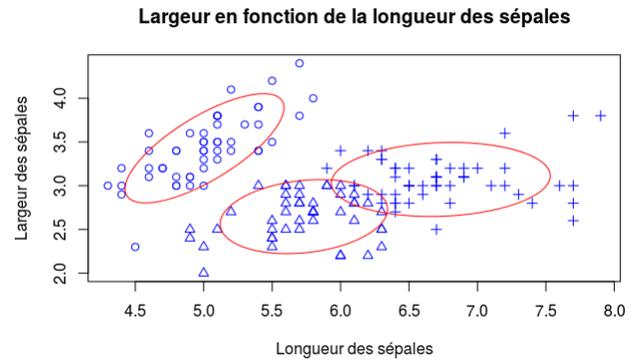


FIGURE 17 – Graphique avec les ellipses de confiance et classification de chaque individu pour les sépales

Cette fois ci, nous constatons que les groupes obtenus par l'algorithme *CEM* sont beaucoup moins imbriqués que pour *EM*. Si nous comparons ces résultats avec les figure 12 et 13, nous voyons que la partition estimée pour les pétales correspond bien au cas théorique. Il doit y avoir environ trois ou quatre erreurs de classement. En revanche pour les sépales, les classes étant plus fortement imbriquées, nous remarquons une erreur de classement plus importante : environ dix individus mal classés.

Finalement, nous nous intéressons à réaliser la classification supervisée de nos deux ensembles grâce à *graph\_colore*.

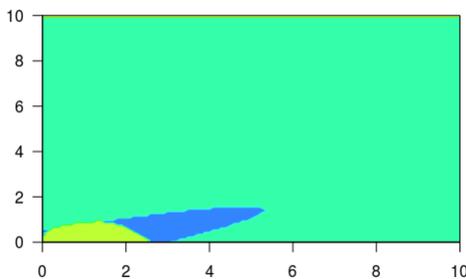


FIGURE 18 – Graphique des frontières calculées empiriques

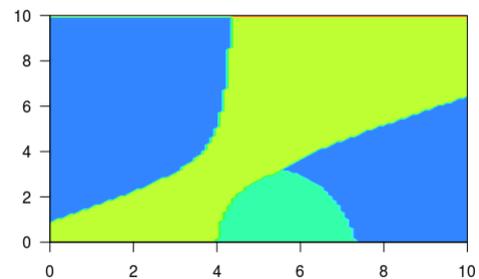


FIGURE 19 – Graphique des frontières calculées empiriques

Concernant les pétales, nous constatons que l'un des groupes, en cyan, est extrêmement dominant. On peut néanmoins être critique sur ce résultat en disant de ne pas classer un individu trop éloigné de notre échantillon original (rejet de distance). Si nous devons classer plusieurs individus supplémentaires loin de notre ensemble, ces individus pourraient former une autre classe. Pour les sépales, il est plus facile de savoir à quel groupe un individu supplémentaire peut appartenir.

## 6 Conclusion

Pour conclure, nous avons introduit divers méthodes et algorithmes pouvant réaliser une classification d'éléments d'un jeu de données. Cette présentation reste néanmoins qu'une introduction. Les algorithmes réalisés ne sont pas optimisés mais restent efficaces pour des jeux de données de taille raisonnable. De plus, lors de cette présentation nous avons fait des choix de méthodes aussi bien que pour l'initialisation de *EM* ou *CEM* que pour l'importance accordée à l'inertie des groupes.

Il existe d'autres méthodes faisant d'autres choix qui peuvent s'avérer plus efficaces dans certains cas. Nous avons décidé de nous intéresser à la *classification supervisée* et *non supervisée* mais il existe justement un juste milieu en faisant une *classification semi-supervisée* pour des données de très haute dimension (plusieurs milliers de variables) ou encore pour des données partiellement étiquetées.

## Références

- [1] Wikipédia : Algorithme espérance-maximisation. [https://fr.wikipedia.org/wiki/Algorithme\\_espérance-maximisation](https://fr.wikipedia.org/wiki/Algorithme_espérance-maximisation).
- [2] Christophe Biernacki. Pourquoi les modèles de mélange pour la classification? *Université Lille 1*, oct 2009.
- [3] Jean Christophe Breton. Estimation d'un mélange avec l'algorithme em. *Université Rennes 1*, 2008.
- [4] J. de Leeuw. Akaike(1973) information theory and an extension of the maximum likelihood principle. *University of California at Los Angeles*, 1992.
- [5] Tristan Mary-Huard Emilie Lebarbier. Le critère bic : fondements théoriques et interprétation. *Inria*, sept 2004. <https://hal.inria.fr/inria-00070685>.
- [6] Yoann Gelineau. Vecteurs gaussiens. *Université Lyon 1*, dec 2007.
- [7] Daniel Li. Vecteurs aléatoires gaussiens. *Université d'Artois*.
- [8] Ruocong Zhang Marc Weber. Formulaire de dérivation matricielle. *Inria*, Oct 2009.

## A Annexe

### Intervalle de confiance en dimension 1

Dans cette partie, nous allons prouver que pour une variable aléatoire  $X \sim N(m, \sigma^2)$  son intervalle de confiance est un intervalle centré autour de la moyenne. Sans perte de généralités, on peut prendre une variable aléatoire  $X \sim N(0, 1)$ .

On cherche les coefficients  $a < 0$  et  $b > 0$  tel que :

$$\mathbb{P}(X \in [a; b]) = \alpha > 0.5$$

On note  $f$  la densité de la variable  $X$ .

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{\sigma^2}}$$

Supposons que  $[a, b]$  est le plus petit intervalle vérifiant  $\mathbb{P}(X \in [a; b]) = \alpha$  tel que  $|a| > |b|$ .

$$\begin{aligned} \mathbb{P}(X \in [a; b]) &= \int_a^b f(x) dx \\ &= \int_0^b f(x) dx + \int_a^0 f(x) dx \\ &= \int_{\frac{|a|+|b|}{2}}^b f(x) dx + \int_0^{\frac{|a|+|b|}{2}} f(x) dx + \int_{-\frac{|a|+|b|}{2}}^0 f(x) dx + \int_a^{-\frac{|a|+|b|}{2}} f(x) dx \\ &= \int_{-\frac{|a|+|b|}{2}}^{\frac{|a|+|b|}{2}} f(x) dx - \int_b^{\frac{|a|+|b|}{2}} f(x) dx + \int_a^{-\frac{|a|+|b|}{2}} f(x) dx \end{aligned}$$

On cherche maintenant le signe de  $-\int_b^{\frac{|a|+|b|}{2}} f(x) dx + \int_a^{-\frac{|a|+|b|}{2}} f(x) dx$ .

$$-\int_b^{\frac{|a|+|b|}{2}} f(x) dx + \int_a^{-\frac{|a|+|b|}{2}} f(x) dx = -\int_b^{\frac{|a|+|b|}{2}} f(x) dx + \int_{\frac{|a|+|b|}{2}}^{-a} f(x) dx \quad \text{car } f \text{ est pair}$$

On sait que  $f$  est décroissante sur  $[0, +\infty]$  ( $f$  est même unimodale) et  $|a| > |b|$ .

$$\begin{aligned} \text{Donc } &\int_b^{\frac{|a|+|b|}{2}} f(x) dx > \int_{\frac{|a|+|b|}{2}}^{-a} f(x) dx \\ &-\int_b^{\frac{|a|+|b|}{2}} f(x) dx + \int_a^{-\frac{|a|+|b|}{2}} f(x) dx < 0 \\ \mathbb{P}(X \in [a; b]) &< \int_{-\frac{|a|+|b|}{2}}^{\frac{|a|+|b|}{2}} f(x) dx = \mathbb{P}\left(X \in \left[-\frac{|a|+|b|}{2}; \frac{|a|+|b|}{2}\right]\right) \end{aligned}$$

## Courbes de niveau

Dans cette partie, nous allons expliciter les courbes de niveau de la fonction de densité d'un vecteur gaussien de dimension  $d$ .

Pour un vecteur gaussien de dimension  $d$ , la fonction de densité est :

$$f(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

On cherche donc une expression en fonction de  $x$  tel que :

$$\begin{aligned} f(x) &= r \\ \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} &= r \\ e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} &= r(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}} \\ (x-\mu)^T \Sigma^{-1} (x-\mu) &= -2 \log(r(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}) \\ (x-\mu)^T \Sigma^{-1} (x-\mu) &= R \end{aligned}$$

Comme vu précédemment, avec cette dernière ligne de calcul on a une équation d'une ellipse si  $R$  est positif ou nul. En effet :

$$\begin{aligned} r &\in \left] 0; \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \right] \\ r(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}} &\in ]0; 1] \\ \log(r(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}) &\leq 0 \\ -2 \log(r(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}) &= R \geq 0 \end{aligned}$$

# Simulation

```
r_VG <- function(MU,SIG){
  # --- nos parametres ---
  d <- length(MU)
  x <- matrix(0,nrow = d,ncol = 1)
  # -----
  # --- Calcul de l'echantillon ---
  x <- rnorm(d,0,sqrt(eigen(SIG)$values))
  x <- eigen(SIG)$vectors %*% x + MU
  # -----
  return(x)
}
d_VG <- function(x,MU,SIG){
  # --- nos parametres ---
  d <- length(x)
  P <- eigen(SIG)$vectors
  y <- t(eigen(SIG)$vectors) %*% (x-MU)
  # -----
  # --- Calcul de la densite ---
  densit <- dnorm(y,0,sqrt(eigen(SIG)$values))
  # -----
  return(prod(densit))
}
r_MelangeGauss <- function(n,PI,MU,SIG){
  # --- nos parametres ---
  d <- length(MU[1,])
  g <- length(PI)
  x <- matrix(0,nrow = n,ncol = d)
  # -----
  # --- Calcul de l'echantillon ---
  for (i in 1:n) {
    # tirage de la loi multinomiale
    m = rmultinom(1,1,PI)
    for (j in 1:g) {
      if (m[j]) {
        x[i,] <- t( r_VG( MU[j,], SIG[((j-1)*d+1):(j*d),] ) )
      }
    }
  }
  # -----
  return(x)
}
```

# Outils pour les algorithmes EM et CEM

```
H <- function(x,PI,MU,SIG){
  # --- nos parametres ---
  n <- length(x[,1])
  g <- length(PI)
  h <- matrix(0,nrow = n,ncol = g)
  # -----
  # --- Calcul de la probabilite conditionnelle ---
  # quand g = 1
  if (g==1){
    h <- matrix(1,nrow = n,ncol = g)
  } else {
    d <- dim(x)[2]
    df <- 0
    for (i in 1:n){
      df <- 0
      for (j in 1:g){
        df <- df + ( PI[j] * d_VG(x[i,],MU[j,],SIG[((j-1)*d+1):(j*d),]) )
      }
      for (j in 1:g){
        h[i,j] <- (PI[j]*d_VG(x[i,],MU[j,],SIG[((j-1)*d+1):(j*d),])) / df
      }
    }
  }
  # -----
  return(h)
}

norme <- function(x){
  return(sqrt(sum(x**2)))
}

Inertie_intra <- function(x,z,MU,g,n){
  res <- 0
  for(i in 1:n){
    for (j in 1:g){
      res <- res + z[i,j]*norme(x[i,]-MU[j,])**2
    }
  }
  return(res)
}
```

# L'algorithme EM

```
EM <- function(x,g,eps,init_p){
  # --- nos parametres ---
  n <- length(x[,1])
  d <- length(x[1,])
  p <- vector()
  MU <- matrix(0,nrow = g,ncol = d)
  SIG <- matrix(0,nrow = g*d,ncol = d)
  # -----
  # --- Initilisation ---
  p <- init_p
  u <- matrix(0,nrow = d,ncol = d)
  for (j in 1:g){
    MU[j,] <- x[sample.int(n,size = 1),]
  }
  for (i in 1:g){
    for (k1 in 1:d){
      u[k1,k1] <- sd(x[,k1])
    }
    SIG[((i-1)*d+1):(i*d),] <- u
  }
  p_prec <- rep(0,g)
  MU_prec <- matrix(0,nrow = g,ncol = d)
  SIG_prec <- matrix(0,nrow = g*d,ncol = d)
  # -----
  # --- Boucle de l'agorithme EM ---
  k <- 0
  while((norme(p_prec-p)>eps)&&(norme(MU_prec-MU)>eps)&&(norme(SIG_prec-SIG)>eps)){
    p_prec <- p
    MU_prec <- MU
    SIG_prec <- SIG
    h <- H(x,p,MU,SIG)
    SIG = matrix(0,nrow = g*d,ncol = d)
    for (j in 1:g){
      p[j] = mean(h[,j])
      MU[j,] = colSums(h[,j] * x) / sum(h[,j])
      SIG[((j-1)*d+1):(j*d),] = t((h[,j]*t((t(x)-MU[j,])))%*%t(t(x)-MU[j,]))
      SIG[((j-1)*d+1):(j*d),] = SIG[((j-1)*d+1):(j*d),] / sum(h[,j])
    }
    k <- k+1
  }
  # -----
  # --- Calcul de l'inertie ---
  z <- matrix(0,nrow=n,ncol=g)
  h <- H(x,p,MU,SIG)
  for (i in 1:n){
    z[i,which.max(h[i,])] = 1
  }
  I <- Inertie_intra(x,z,MU,g,n)
  return(list(Poids=p,Moyenne=MU,Variance=SIG,Nb_iteration=k,Inertie=I))
}
```

# L'algorithme CEM

```
CEM <- function(x,g,init_p){
  # --- nos parametres ---
  n <- length(x[,1])
  d <- length(x[1,])
  p <- rep(0,g)
  MU <- matrix(0,nrow = g,ncol = d)
  SIG <- matrix(0,nrow = g*d,ncol = d)
  z <- matrix(0,nrow=n,ncol=g)
  # -----
  # --- Initilisation ---
  p <- init_p
  u <- matrix(0,nrow = d,ncol = d)
  MU <- x[sample.int(n,size = g),]
  for (j in 1:g){
    for (k1 in 1:d){
      u[k1,k1] <- sd(x[,k1])
    }
    SIG[((j-1)*d+1):(j*d),] <- u
  }
  p_prec <- rep(0,g)
  MU_prec <- matrix(0,nrow = g,ncol = d)
  SIG_prec <- matrix(0,nrow = g*d,ncol = d)
  # -----
  # --- Boucle de L'agorithme CEM ---
  k <- 0
  while((norme(p_prec-p)>0)&&(norme(MU_prec-MU)>0)&&(norme(SIG_prec-SIG)>0)){
    p_prec <- p
    MU_prec <- MU
    SIG_prec <- SIG
    h <- H(x,p,MU,SIG)
    z <- matrix(0,nrow=n,ncol=g)
    for (i in 1:n){
      z[i,which.max(h[i,])] = 1
    }
    # Quand une classe ne contient aucun element
    if(min(colSums(z))==0){
      return(list(Inertie=Inf))
    }
    for (j in 1:g){
      p[j] = mean(z[,j])
      MU[j,] = colSums(x[which(z[,j]==1),,drop=FALSE]) / sum(z[,j])
      SIG[((j-1)*d+1):(j*d),] = (t(x[which(z[,j]==1),])%*%t(t(x[which(z[,j]==1),])-
MU[j,]))
      SIG[((j-1)*d+1):(j*d),] = SIG[((j-1)*d+1):(j*d),] / sum(z[,j])
    }
    k <- k+1
  }
  # -----
  # --- Calcul de L'inertie ---
  I <- Inertie_intra(x,z,MU,g,n)
  # -----
  return(list(Poids=p,Moyenne=MU,Variance=SIG,Classification=z,Nb_iteration=k,Inertie=I))
}
```

# MULTI\_EM et MULTI\_CEM

```
MULTI_EM <- function(x,g,eps){
  # --- nos parametres ---
  res_EM_final <- list()
  I <- 0
  I_min <- 0
  init_p <- rep(0,g)
  s <- seq(1/g,1,1/g)
  s <- s/sum(s)
  # --- Premiere iteration avec poids identique pour chaque classe ---
  res_EM_final <- EM(x,g,eps,rep(1/g,g))
  I_min <- res_EM_final[[5]]
  # --- Boucle avec differentes initialisation pour p ---
  for(k in 1:g){
    for(j in 1:g){
      init_p[j] <- s[(j+k)%g+1]
    }
    res_EM <- EM(x,g,eps,init_p)
    I <- res_EM[[5]]
    if(I<I_min){
      I_min <- I
      res_EM_final <- res_EM
    }
  }
  return(res_EM_final)
}
MULTI_CEM <- function(x,g){
  # --- nos parametres ---
  res_CEM_final <- list()
  I <- 0
  I_min <- 0
  init_p <- rep(0,g)
  s <- seq(1/g,1,1/g)
  s <- s/sum(s)
  # --- Premiere iteration avec poids identique pour chaque classe ---
  res_CEM_final <- CEM(x,g,rep(1/g,g))
  I_min <- res_CEM_final$Inertie
  # --- Boucle avec differentes initialisation pour p ---
  for(k in 1:g){
    for(j in 1:g){
      init_p[j] <- s[(j+k)%g+1]
    }
    res_CEM <- CEM(x,g,init_p)
    I <- res_CEM$Inertie
    if(I<I_min){
      I_min <- I
      res_CEM_final <- res_CEM
    }
  }
  return(res_CEM_final)
}
```

# Graphique EM et CEM

```
ellipse_dim2 <- function(x, g, THETA, prec){
  # --- nos parametres ---
  SIG <- THETA$Variance
  MU <- THETA$Moyenne
  n <- dim(x)[1]
  # -----
  # --- Boucle pour tracer une ellipse pour chaque groupe ---
  for (j in 1:g){
    v <- eigen(THETA$Variance[((j-1)*2+1):(j*2),])
    a <- sqrt(-2*log(1-prec)*v$values[1])
    b <- sqrt(-2*log(1-prec)*v$values[2])
    t <- seq(0,2*3.2,0.1)
    cor <- matrix(0,nrow = 2,ncol = 65)
    cor[1,] <- a*cos(t)
    cor[2,] <- b*sin(t)
    cor <- v$vectors %%% cor
    cor <- cor + MU[j,]
    lines(t(cor), col="red")
  }
  # -----
}

graph_EM <- function(x,g,THETA,prec){ # prec =0.95
  # --- Representation des points de notre echantillon ---
  plot(x, type = "p",pch = 4,main = "Simulation d'un melange de vecteurs Gaussiens",xlab = "X_1",ylab = "X_2", col="blue")
  # -----
  # --- Representation de nos ellipses ---
  ellipse_dim2(x, g, THETA, prec)
  # -----
}

graph_CEM <- function(x,g,THETA,prec){
  # --- Representation des points de notre echantillon avec un code pour la representation de la partition ---
  plot(x[which(THETA$Classification[,1]==1),], type = "p",pch = 1,main = "Simulation d'un melange de vecteurs Gaussiens",xlab = "X_1",ylab = "X_2", col="blue",xlim = c(min(x[,1]),max(x[,1])),ylim = c(min(x[,2]),max(x[,2])))
  for (j in 2:g){
    points(x[which(THETA$Classification[,j]==1),],pch = j,col="blue")
  }
  # -----
  # --- Representation de nos ellipses ---
  ellipse_dim2(x, g, THETA, prec)
  # -----
}
```

# Outils pour critere de selection du modele

```
log_vraisemblance <- function(donnees,THETA,g){
  # --- nos parametres ---
  ret <- 0
  n <- length(x[,1])
  # -----
  # --- Calcul de notre Log vraisemblance ---
  for(i in 1:n){
    res1 <- 0
    for(j in 1:g){
      res1 <- res1 + THETA$Poids[j] * d_VG(donnees[i,],THETA$Moyenne[j,],THETA$Variance[((j-
1)*2+1):(j*2),])
    }
    ret <- ret + log(res1)
  }
  # -----
  return(ret)
}
Entropie <- function(n,g,z,h){
  ret = 0
  # --- Calcul de notre valeur d'entropie pour Le critere ICL ---
  for(i in 1:n){
    ret <- ret - log(h[i,which.max(z[i,])])
  }
  # -----
  return(ret)
}
```

# Critere AIC ou critere BIC

```
critere_AIC_BIC <- function(x,mod,nb_test,multi){
  # --- nos parametres ---
  n <- length(x[,1])
  d <- length(x[1,])
  valeur <- rep(0,nb_test)
  # -----
  # --- Selection de notre critere ---
  if(mod=="BIC"){
    penalite = 0.5*log(n)
  } else if (mod=="AIC"){
    penalite = 1
  } else {
    stop("la fonction ne reconnait pas ce critere")
  }
  # -----
  # --- Avec un seul groupe ---
  Poids <- 1
  Moyenne <- matrix(colMeans(x),nrow = 1,ncol = d)
  Variance <- var(x)
  THETA <- list(Poids=Poids,Moyenne=Moyenne,Variance=Variance)
  lg_prec <- log_vraisemblance(x,THETA,1) - (2+3)*1* penalite
  valeur[1] <- lg_prec
  nb_groupe_modele <- 1
  # -----
  # --- Avec plusieurs groupes ---
  for(g in 2:nb_test){
    if (multi==TRUE){
      THETA <- MULTI_EM(x,g,0.01)
    }else{
      init_p <- runif(g)
      init_p <- init_p/sum(init_p)
      THETA <- EM(x,g,0.01,init_p)
    }
    lg <- 0
    lg <- log_vraisemblance(x,THETA,g) - (2+3)*g * penalite
    valeur[g] <- lg
    if (lg>lg_prec){
      nb_groupe_modele <- g
      lg_prec <- lg
    }
  }
  # -----
  return(list(nb_groupe_modele=nb_groupe_modele,valeur_critere=valeur))
}
```

# Critere ICL

```
critere_ICL <- function(x,nb_test){
  # --- nos parametres ---
  n <- length(x[,1])
  d <- length(x[1,])
  valeur <- vector(length = nb_test)
  # -----
  # --- Avec un seul groupe ---
  Poids <- 1
  Moyenne <- matrix(colMeans(x),nrow = 1,ncol = d)
  Variance <- var(x)
  Classification <- matrix(1,nrow = n,ncol =1)
  THETA <- list(Poids=Poids,Moyenne=Moyenne,Variance=Variance,Classification=Classification)
  lg_prec <- log_vraisemblance(x,THETA,1) - (2+3)*1*0.5*log(n)
  # Comme h = 1 pour tous element => Entropie est nulle
  valeur[1] <- lg_prec
  lg <- 0
  nb_groupe_modele <- 1
  # -----
  # --- Avec plusieurs groupes ---
  for(g in 2:nb_test){
    THETA <- MULTI_CEM(x,g)
    lg <- log_vraisemblance(x,THETA,g) - (2+3)*g*0.5*log(n)
    h <- H(x,THETA$Poids,THETA$Moyenne,THETA$Variance)
    lg <- lg - Entropie(n,g,THETA$Classification,h)
    if (length(lg)){
      valeur[g] <- lg
      if(lg>lg_prec){
        nb_groupe_modele <- g
        lg_prec = lg
      }
    }else{
      valeur[g]<- -Inf
    }
  }
  # -----
  return(list(nb_groupe_modele=nb_groupe_modele,valeur_critere=valeur))
}
```

# Graphique colore

```
appartenance <- fonction(x,g,THETA){
  # --- nos parametres ---
  MU <- THETA$Moyenne
  SIG <- THETA$Variance
  PI <- THETA$Poids
  h <- rep(0,g)
  df <- 0
  # -----
  # --- Calcul des probabilites d'appartenance d'un point a chaque groupe ---
  for (j in 1:g){
    df <- df + ( PI[j] * d_VG(x,MU[j,],SIG[((j-1)*d+1):(j*d),]) )
  }
  for (j in 1:g){
    h[j] <- (PI[j]*d_VG(x,MU[j,],SIG[((j-1)*d+1):(j*d),])) / df
  }
  # -----
  # --- Selection de La valeur maximal ---
  return(which.max(h))
}
graph_colore <- fonction(g,THETA,xlim,ylim,N){
  # --- Calcul de L'appartenance de chaque poids du plan discretise ---
  A <- matrix(0,nrow = N+1,ncol = N+1)
  for(i in 0:N+1){
    for (j in 0:N+1){
      A[i,(N+1)-j] = appartenance(c(ylim[1]+(i-1)*(ylim[2]-ylim[1])/N,xlim[2]+(j-1)*(xlim[1]-
xlim[2])/N),g,THETA)
    }
  }
  # -----
  # --- Limite de notre plan et pas de discretisation ---
  dx <- (xlim[2]-xlim[1])/N
  dy <- (ylim[2]-ylim[1])/N
  X <- seq(xlim[1],xlim[2],dx)
  Y <- seq(ylim[1],ylim[2],dy)
  # -----
  # --- Representation du graphique colore ---
  filled.contour(X,Y,A)
  # -----
}
```