

Éléments de classification

Christophe Chesneau

<http://www.math.unicaen.fr/~chesneau/>

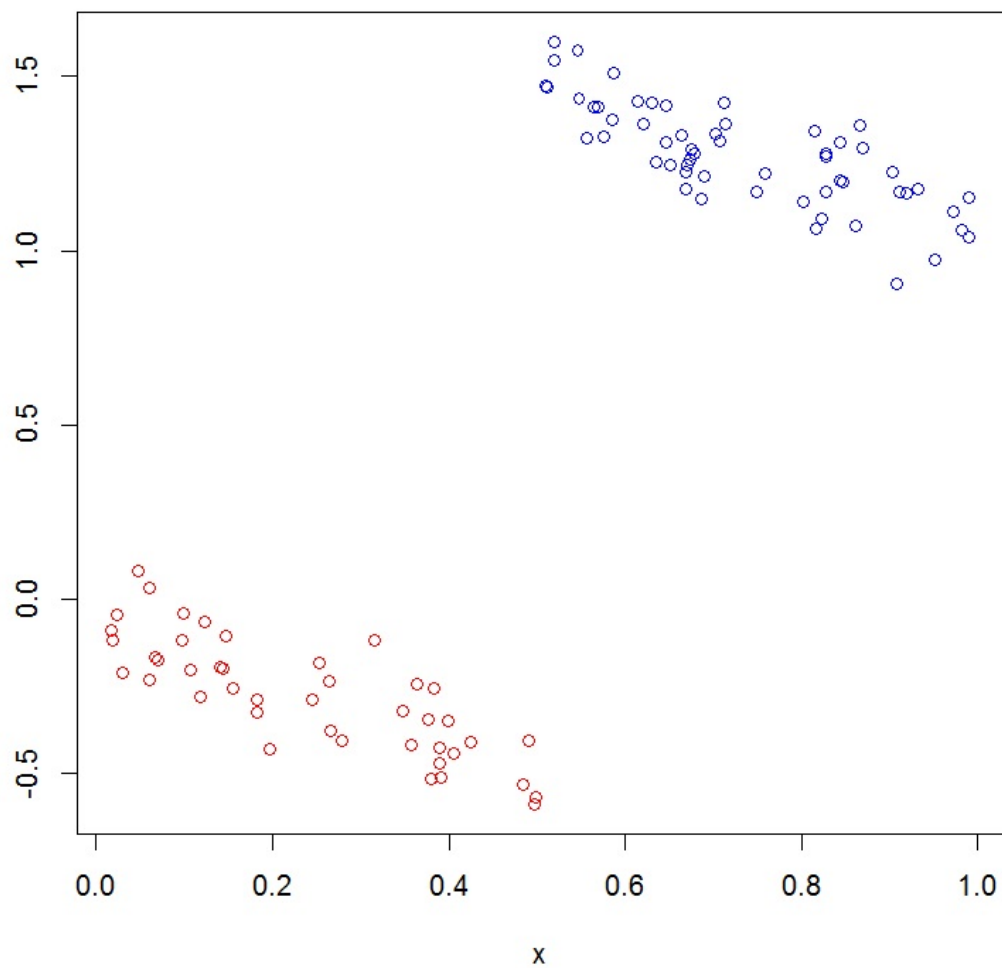


Table des matières

1	Présentation	4
2	Introduction	5
2.1	Classification non-supervisée	5
2.2	Classification supervisée	6
2.3	Les métiers	8
3	Enjeux de la classification non-supervisée	9
4	Étude de la ressemblance	10
4.1	Nuage de points	10
4.2	Distances	15
4.3	Écartés	17
5	Algorithme de Classification Ascendante Hiérarchique (CAH)	22
5.1	Introduction	22
5.2	Description de l'algorithme	22
5.3	Dendrogramme	24
5.4	Quelques commandes R	28
6	CAH et méthode de Ward ; compléments	32
7	Qualité d'une partition	40
8	ACP et CAH	44
9	Caractérisation des groupes	46
10	Algorithme des centres mobiles (k means)	48
11	Consolidation de l'algorithme de CAH	60
12	CAH avec des caractères qualitatifs	61

13 Enjeux de la classification supervisée	67
14 Méthode des k plus proches voisins	69
15 Modèle de mélange de densités	75
16 Régression logistique	79
Index	82

1 Présentation

Ce document résume les notions abordées dans le cours *Éléments de classification* du Master 1 MIASHS de l'université de Caen.

Un des objectifs est de donner des pistes de réflexion au regroupement/classification des individus à partir de données.

Les méthodes statistiques y sont décrites de manière concise, avec les commandes R associées.

N'hésitez pas à me contacter pour tout commentaire :

`christophe.chesneau@gmail.com`

Bonne lecture.

2 Introduction

On présente ici les enjeux de la classification non-supervisée et de la classification supervisée.

2.1 Classification non-supervisée

Contexte : On considère n individus extraits au hasard d'une population. Pour chacun d'entre eux, on dispose de p valeurs de p caractères X_1, \dots, X_p .

Objectif : Partant des données, l'objectif est de regrouper/classer les individus qui se ressemblent le plus/qui ont des caractéristiques semblables.

Ce regroupement peut avoir des buts divers : tenter de séparer des individus appartenant à des sous-populations distinctes, décrire les données en procédant à une réduction du nombre d'individus pour communiquer, simplifier, exposer les résultats. . .

Exemple : Dans une classe, un professeur souhaite faire des binômes constitués d'élèves ayant des compétences semblables. Parmi ceux-ci, 6 élèves ont obtenu les notes suivantes :

	Maths	Physique	Ed Mus	Art Plas
Boris	20	20	0	0
Mohammad	8	8	12	12
Stéphanie	20	20	0	0
Jean	0	0	20	20
Lilly	10	10	10	10
Annabelle	2	2	18	18

Tous les élèves ont une moyenne de 10/20 mais, vu les notes,

- Boris et Stéphanie ont un profil similaire,
- Mohammad et Lilly ont un profil similaire,
- Jean et Annabelle ont un profil similaire.

Finalement, le professeur décide de faire 2 groupes cohérents de 3 élèves avec ces 6 élèves. Lesquels proposez-vous ?

En comparant les notes par matière, on propose :

- *Groupe 1* : Boris, Stéphanie et Lilly,
- *Groupe 2* : Mohammad, Jean et Annabelle.

De plus, par exemple, le profil de Jean est plus proche de celui de Lilly, que celui de Stéphanie. Bien entendu, cette analyse intuitive n'est pas possible si, par exemple, on a 30 élèves à classer par groupes de 3 et on considère 12 matières. C'est pourquoi des méthodes mathématiques ont été mises en place.

Applications : Quelques exemples d'applications de la classification non-supervisée sont présentés ci-dessous.

- *Application 1* : En biologie, on veut regrouper les espèces suivant leurs caractéristiques et donc leurs origines communes.
- *Application 2* : En psychologie, on veut classer les individus selon leur type de personnalités.
- *Application 3* : En chimie, on veut classer des composés selon leurs propriétés.
- *Application 4* : Dans l'industrie, on veut
 - analyser des résultats d'enquêtes,
 - identifier les clients potentiels d'une entreprise,
 - identifier les clients susceptibles de partir à la concurrence,
 - déterminer des lieux de ventes (pose de distributeurs de billets...),
 - analyser, identifier les risques (dégâts des eaux...),
 - analyser des données textuelles.

2.2 Classification supervisée

Contexte : On considère une population divisée en q groupes d'individus différents. Ces groupes sont distinguables suivant les valeurs de p caractères X_1, \dots, X_p , sans que l'on ait connaissance des valeurs de X_1, \dots, X_p les caractérisant. On dispose

- de n individus avec, pour chacun d'entre eux, les valeurs de X_1, \dots, X_p et son groupe d'appartenance,

- d'un individu ω_* de la population avec ses valeurs de X_1, \dots, X_p , sans connaissance de son groupe d'appartenance.

Objectif : Partant des données, l'objectif est de déterminer à quel groupe l'individu ω_* a le plus chance d'appartenir.

Exemple : Dans une classe, un professeur considère deux groupes d'élèves, $G1$ et $G2$, en fonction de leur compétence. On dispose uniquement des notes et de l'affectation de 6 élèves :

	Maths	Physique	Ed Mus	Art Plas	Groupe
Boris	20	20	0	0	G1
Mohammad	8	8	12	12	G2
Stéphanie	20	20	0	0	G1
Jean	0	0	20	20	G2
Lilly	10	10	10	10	G1
Annabelle	2	2	18	18	G2

D'autre part, un étudiant de la classe, Bob, a les résultats suivants :

	Maths	Physique	Ed Mus	Art Plas	Groupe
Bob	9	15	13	11	inconnu

À partir de ses notes, à quel groupe Bob a le plus de chances d'appartenir ?

Autrement écrit, quel est la probabilité que Bob appartienne au groupe $G1$ sachant qu'il a obtenu les notes (Maths, Physique, Ed Mus, Art Plas) = (9, 15, 13, 11) ?

On peut écrire cette probabilité comme :

$$\mathbb{P}(\text{Bob} \in G1 / (9, 15, 13, 11)).$$

La réponse n'est pas immédiate ; c'est pourquoi des méthodes mathématiques ont été mises en place.

Applications : Quelques exemples d'applications de la classification supervisée sont présentés ci-dessous.

- *Application 1 :* Un archéologue cherche à déterminer si des restes humains sont ceux d'un homme ou d'une femme.

- *Application 2* : Dans une banque, une commission de crédit doit décider, à partir de paramètres financiers, si on accorde ou non un prêt à un particulier.
- *Application 3* : Étant donné un ensemble de symptômes, un médecin doit poser un diagnostic.
- *Application 4* : Dans l'industrie, on veut
 - identifier des visages, des empreintes digitales,
 - identifier des objets dans des séquences vidéos,
 - rechercher des clients potentiels dans des bases de données,
 - rapprocher un ou plusieurs mots de manière pertinente au texte le plus pertinent.

2.3 Les métiers

Il y a de nombreux métiers où la classification est utilisée dont

- responsable logistique du traitement et de l'analyse des études,
- chargé d'études junior : prise en charge de la documentation, codage des questionnaires, traitement statistiques simples,
- chargé d'études senior, assistant du chargé d'étude : prise en main d'une étude de marché,
- analyste statisticien, études quantitatives, aide à la décision, expert en statistiques,
- chef de projet.

3 Enjeux de la classification non-supervisée

Contexte : On considère n individus $\Gamma = \{\omega_1, \dots, \omega_n\}$ extraits au hasard d'une population. Pour chacun d'entre eux, on dispose de p valeurs de p caractères X_1, \dots, X_p . Dans un premier temps, on suppose que les caractères étudiés sont quantitatifs.

Les données sont donc de la forme :

	X_1	\dots	X_p
ω_1	$x_{1,1}$	\dots	$x_{p,1}$
\vdots	\vdots	\dots	\vdots
ω_n	$x_{1,n}$	\dots	$x_{p,n}$

où, pour tout $(i, j) \in \{1, \dots, n\} \times \{1, \dots, p\}$, $x_{j,i} = X_j(\omega_i)$ est l'observation du caractère X_j sur l'individu ω_i .

Objectif : Partant des données, l'objectif est de regrouper/classer les individus qui se ressemblent le plus/qui ont des caractéristiques semblables.

Méthodes : Pour atteindre l'objectif, plusieurs méthodes sont possibles. Parmi elles, il y a

- l'algorithme de Classification Ascendante Hiérarchique (CAH),
- l'algorithme des centres mobiles,
- l'algorithme de Classification Descendante Hiérarchique (CDH),
- la méthode des nuées dynamiques (partitionnement autour d'un noyau),
- la méthode de classification floue,
- la méthode de classification par voisinage dense.

Ce document aborde quelques aspects des deux premiers points.

4 Étude de la ressemblance

4.1 Nuage de points

Matrice de données : On appelle matrice de données associées à Γ la matrice \mathbf{X} définie par

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & \dots & x_{p,1} \\ \vdots & \dots & \vdots \\ x_{1,n} & \dots & x_{p,n} \end{pmatrix}$$

Nuage de points : Pour tout $i \in \{1, \dots, n\}$, l'individu ω_i peut être représenté dans \mathbb{R}^p par un point m_i de coordonnées $(x_{1,i}, \dots, x_{p,i})$. On appelle nuage de points la représentation graphique de l'ensemble de ces points. Il est noté $\mathcal{N} = \{m_1, \dots, m_n\}$.

Ressemblance : On dira que des individus se ressemblent si les points associés sont proches les uns des autres/si les distance qui les séparent sont petites.

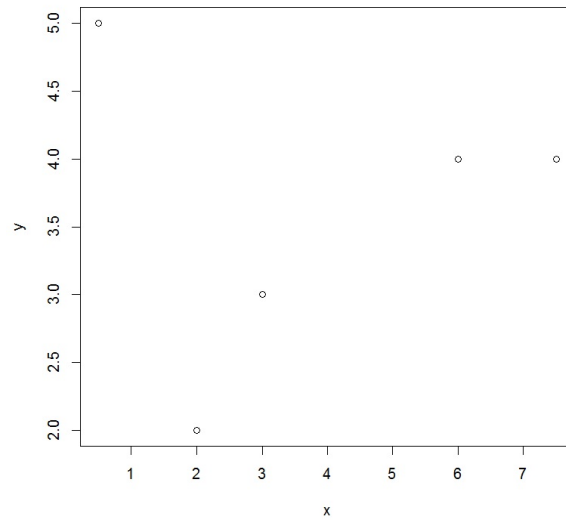
Ainsi, on souhaite rechercher dans \mathcal{N} les zones denses pouvant correspondre à des groupes d'individus qu'il s'agira d'interpréter par la suite.

Exemple 1 : On considère la matrice de données \mathbf{X} associée à 5 individus, $\Gamma = \{\omega_1, \dots, \omega_5\}$, définie par

$$\mathbf{X} = \begin{pmatrix} 2 & 2 \\ 7.5 & 4 \\ 3 & 3 \\ 0.5 & 5 \\ 6 & 4 \end{pmatrix}.$$

Implicitement, on considère donc 2 caractères X_1 et X_2 . Par exemple, l'individu ω_2 a pour caractéristiques $X_1 = 7.5$ et $X_2 = 4$.

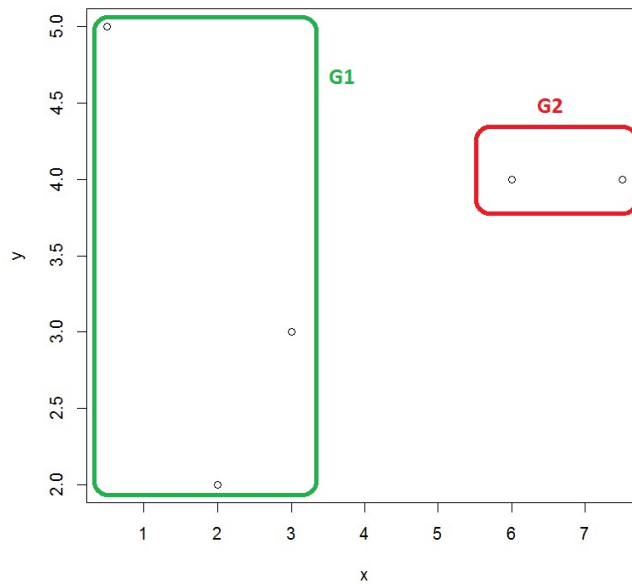
Le nuage de point associé est :



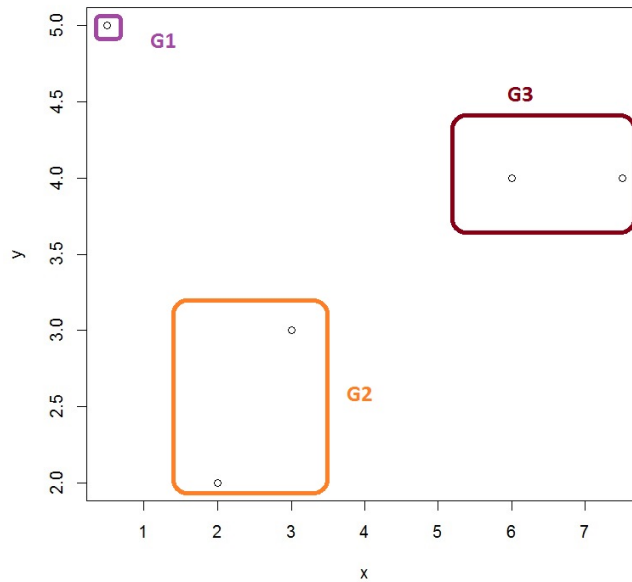
La problématique est la suivante : comment regrouper ces individus en 2 ou 3 groupes, par exemple, en fonction de leur position dans \mathbb{R}^2 ?

Visuellement, en fonction des zones denses, on peut envisager

◦ les 2 groupes suivants :



◦ les 3 groupes suivants :



Exemple 2 : On considère un tableau de notes de 6 élèves :

	Maths	Physique	Ed Mus	Art Plas
Boris	19	17	2	8
Mohammad	7	8	12	12
Stéphanie	20	19	9	9
Jean	1	6	18	17
Lilly	10	11	12	12
Annabelle	2	12	18	18

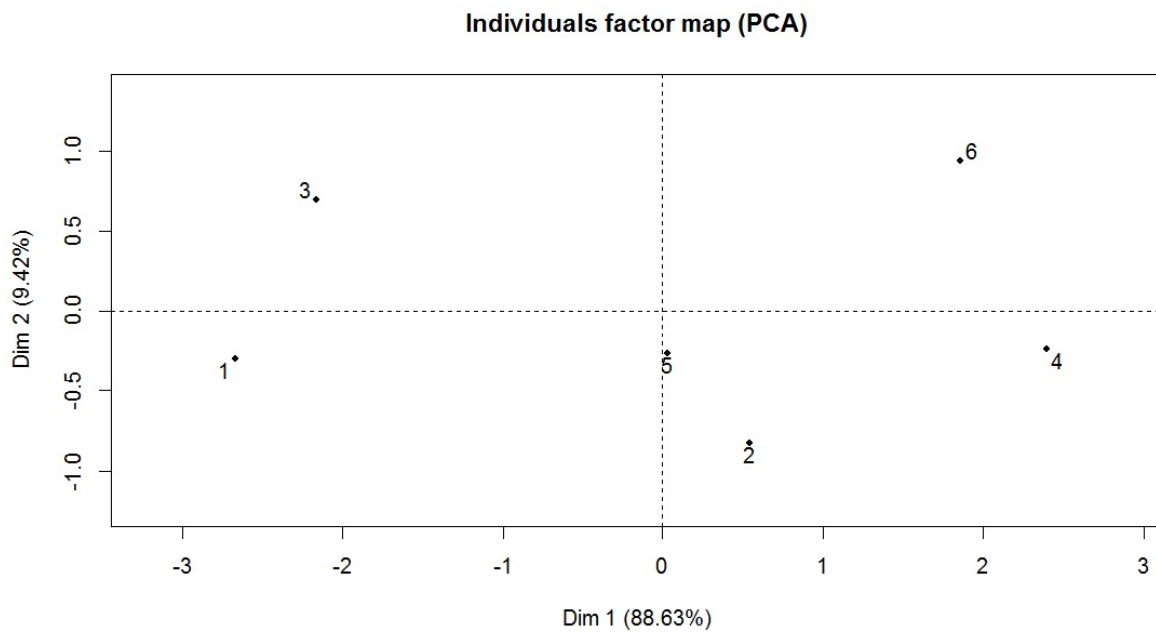
La matrice de données \mathbf{X} associée est

$$\mathbf{X} = \begin{pmatrix} 19 & 17 & 2 & 8 \\ 7 & 8 & 12 & 12 \\ 20 & 19 & 9 & 9 \\ 1 & 6 & 18 & 17 \\ 10 & 11 & 12 & 12 \\ 2 & 12 & 18 & 18 \end{pmatrix}$$

La problématique est la suivante : comment regrouper ces individus en 2 ou 3 groupes, par exemple, en fonction de leur position dans \mathbb{R}^4 ?

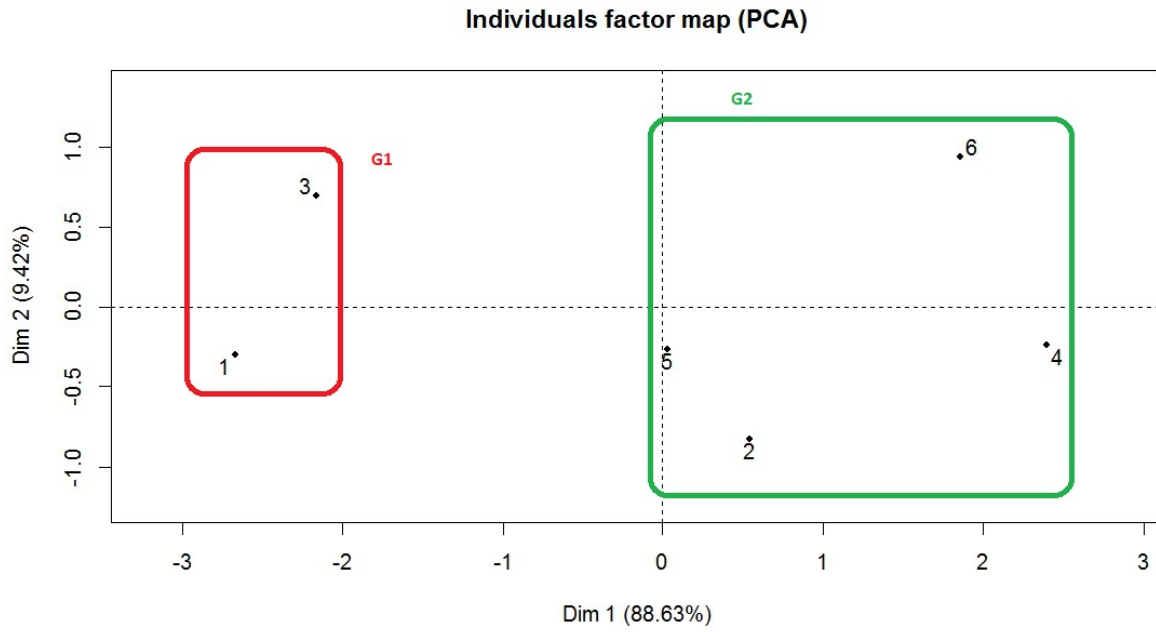
La représentation graphique dans \mathbb{R}^2 n'est donc pas possible.

Une solution est de considérer le plan principal d'une analyse en composante principale (ACP), méthode statistique qui ne sera pas développée ici. Cela donne la représentation graphique suivante :

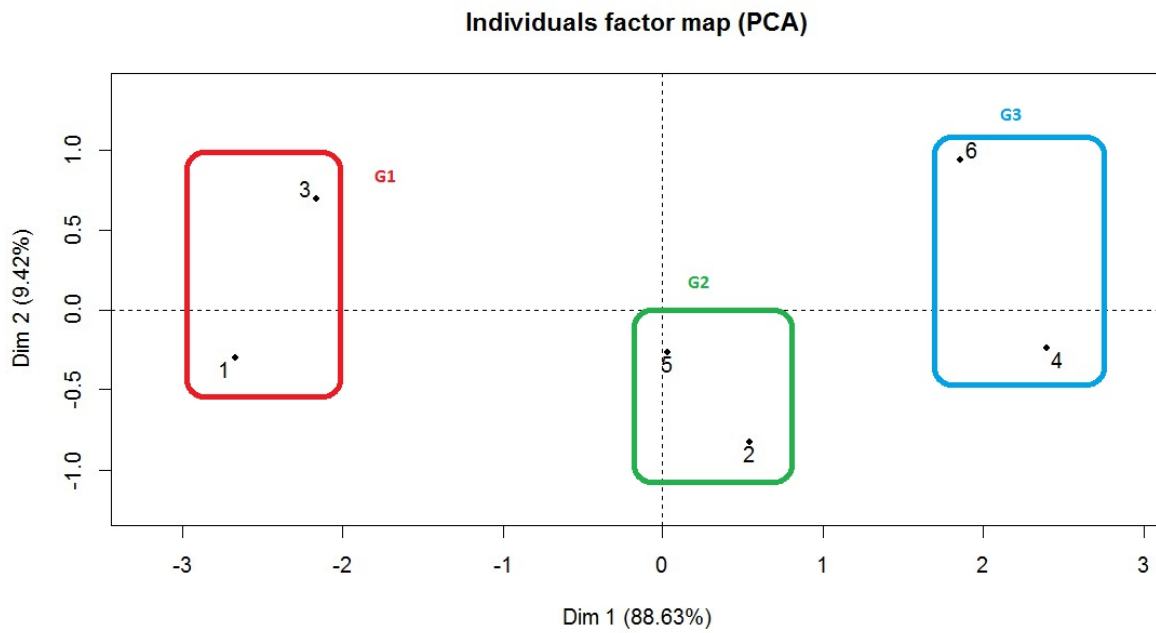


Visuellement, en fonction des zones denses, on peut envisager

- les 2 groupes suivants :



- les 3 groupes suivants :



4.2 Distances

Distances : On peut donc aborder le problème de la ressemblance entre individus par le biais de la notion de distance. On appelle distance sur un ensemble M toute application $d : M^2 \rightarrow [0, \infty[$ telle que

- pour tout $(x, y) \in M^2$, on a $d(x, y) = 0$ si, et seulement si, $x = y$,
- pour tout $(x, y) \in M^2$, on a $d(x, y) = d(y, x)$,
- pour tout $(x, y, z) \in M^3$, on a

$$d(x, y) \leq d(x, z) + d(z, y).$$

Exemple 1 : distance euclidienne : Soient $m \in \mathbb{N}^*$, $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ et $y = (y_1, \dots, y_m) \in \mathbb{R}^m$. On appelle distance euclidienne entre x et y la distance :

$$d(x, y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}.$$

Exemple 2 : distance de Manhattan : Soient $m \in \mathbb{N}^*$, $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ et $y = (y_1, \dots, y_m) \in \mathbb{R}^m$. On appelle distance de Manhattan entre x et y la distance :

$$d(x, y) = \sum_{i=1}^m |x_i - y_i|.$$

Exemple 3 : distance de Minkowski : Soient $m \in \mathbb{N}^*$, $q \geq 1$, $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ et $y = (y_1, \dots, y_m) \in \mathbb{R}^m$. On appelle distance de Minkowski entre x et y la distance :

$$d(x, y) = \left(\sum_{i=1}^m |x_i - y_i|^q \right)^{\frac{1}{q}}.$$

Exemple 4 : distance de Canberra : Soient $m \in \mathbb{N}^*$, $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ et $y = (y_1, \dots, y_m) \in \mathbb{R}^m$. On appelle distance de Canberra entre x et y la distance :

$$d(x, y) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|}.$$

Exemple 5 : distance maximum : Soient $m \in \mathbb{N}^*$, $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ et $y = (y_1, \dots, y_m) \in \mathbb{R}^m$. On appelle distance maximum entre x et y la distance :

$$d(x, y) = \sup_{i \in \{1, \dots, m\}} |x_i - y_i|.$$

Quelques commandes R : Quelques commandes R associées à ces distances sont :

```
x = c(1, 16, 2, 9, 10, 16, 1)
y = c(14, 9, 9, 12, 4, 3, 13)
z = rbind(x, y)
dist(z, method = "euclidean")
dist(z, method = "manhattan")
dist(z, method = "minkowski", p = 6)
dist(z, method = "maximum")
```

Dorénavant, pour raison de simplicité et de popularité, seule la distance euclidienne sera considérée.

Distance entre 2 individus : Pour tout $(u, v) \in \{1, \dots, n\}^2$ avec $u \neq v$, la distance euclidienne entre les individus ω_u et ω_v est

$$d(\omega_u, \omega_v) = \sqrt{\sum_{j=1}^p (x_{j,u} - x_{j,v})^2}.$$

Tableau des distances : Soit d une distance. On appelle tableau des distances associées aux individus $(\omega_1, \dots, \omega_n)$ le tableau :

$$D = \begin{array}{c|ccccc} & \omega_1 & \omega_2 & \dots & \omega_{n-1} & \omega_n \\ \hline \omega_1 & 0 & d_{1,2} & \dots & d_{1,n-1} & d_{1,n} \\ \omega_2 & d_{2,1} & 0 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \omega_{n-1} & d_{n-1,1} & \dots & \dots & 0 & d_{n-1,n} \\ \omega_n & d_{n,1} & \dots & \dots & d_{n,n-1} & 0 \end{array}$$

où, pour tout $(u, v) \in \{1, \dots, n\}^2$ avec $u \neq v$,

$$d_{u,v} = d(\omega_u, \omega_v) = \sqrt{\sum_{j=1}^p (x_{j,u} - x_{j,v})^2}.$$

Exemple : On considère la matrice de données \mathbf{X} définie par

$$\mathbf{X} = \begin{pmatrix} 2 & 2 \\ 7.5 & 4 \\ 3 & 3 \\ 0.5 & 5 \\ 6 & 4 \end{pmatrix}$$

En prenant 2 chiffres après la virgule, on a, par exemple,

$$d(\omega_1, \omega_2) = \sqrt{(2 - 7.5)^2 + (2 - 4)^2} = 5.85.$$

En procédant de même, on obtient le tableau des distances :

$$\mathbf{D} = \begin{array}{c|ccccc} & \omega_1 & \omega_2 & \omega_3 & \omega_4 & \omega_5 \\ \hline \omega_1 & 0 & 5.85 & 1.41 & 3.35 & 4.47 \\ \omega_2 & 5.85 & 0 & 4.60 & 7.07 & 1.50 \\ \omega_3 & 1.41 & 4.60 & 0 & 3.20 & 3.16 \\ \omega_4 & 3.35 & 7.07 & 3.20 & 0 & 5.59 \\ \omega_5 & 4.47 & 1.50 & 3.16 & 5.59 & 0 \end{array}$$

4.3 Écarts

Écarts : En notant $\mathcal{P}(\Gamma)$ l'ensemble des parties de Γ , on appelle écart toute application $e : \mathcal{P}(\Gamma)^2 \rightarrow [0, \infty[$ définie à partir d'une distance et évaluant la ressemblance entre deux groupes d'individus.

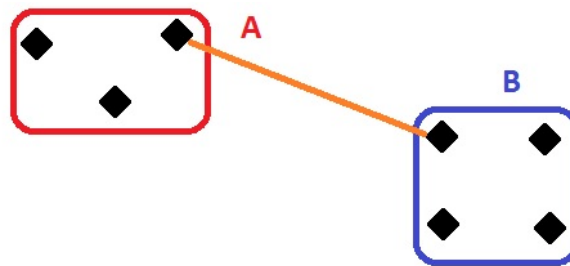
Règle centrale : Plus l'écart entre deux éléments est petit, plus ils se ressemblent.

Écarts usuels : Parmi les écarts usuels entre deux groupes A et B /méthodes usuelles mesurant la ressemblance entre deux groupes A et B , il y a :

- **Écart simple (single linkage)/Méthode du plus proche voisin** :

$$e(A, B) = \min_{(\omega, \omega_*) \in A \times B} d(\omega, \omega_*).$$

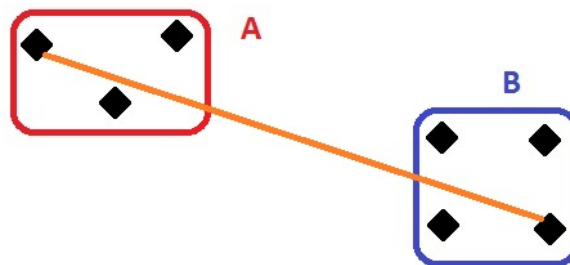
L'écart entre deux groupes A et B est caractérisé par la distance la plus faible entre un point de A et un point de B :



- **Écart complet (complete linkage)/Méthode du voisin le plus éloigné** :

$$e(A, B) = \max_{(\omega, \omega_*) \in A \times B} d(\omega, \omega_*).$$

L'écart entre deux groupes A et B est caractérisé par la distance la plus forte entre un point de A et un point de B :

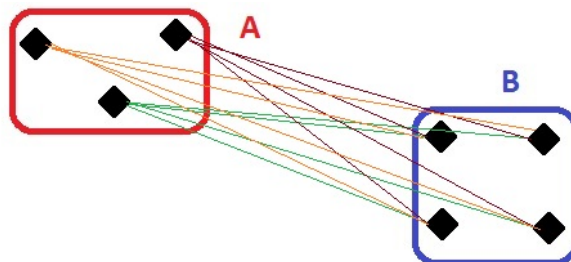


- **Écart moyen (average linkage)/Méthode de la distance moyenne** :

$$e(A, B) = \frac{1}{n_A n_B} \sum_{\omega \in A} \sum_{\omega_* \in B} d(\omega, \omega_*),$$

où n_A est le nombre d'individus dans le groupe A , et n_B le nombre d'individus dans le groupe B .

L'écart entre deux groupes A et B est caractérisé par la distance moyenne entre les points de A et B :



- o **Écart de Ward** : Soit d la distance euclidienne. La méthode de Ward considère l'écart :

$$e(A, B) = \frac{n_A n_B}{n_A + n_B} d^2(g_A, g_B),$$

où g_A est le centre de gravité de A , et g_B celui de B . On rappelle que g_A est le point de coordonnées $(\bar{x}_{1,A}, \dots, \bar{x}_{p,A})$, où, pour tout $j \in \{1, \dots, p\}$, $\bar{x}_{j,A}$ désigne la moyenne des valeurs observées du caractère X_j sur les n_A individus du groupe A . De même pour g_B .

Cette méthode prend en compte à la fois la dispersion à l'intérieur d'un groupe et la dispersion entre les groupes. Elle est utilisée par défaut dans la plupart des programmes informatiques. Elle fera l'objet d'un chapitre à venir.

Tableau des écarts : Soit e un écart défini par une des méthodes précédentes. On appelle tableau des écarts associé aux groupes d'individus (A_1, \dots, A_n) le tableau :

$$\mathbf{E} = \begin{array}{c|cccccc} & A_1 & A_2 & \dots & A_{n-1} & A_n \\ \hline A_1 & 0 & e_{1,2} & \dots & e_{1,n-1} & e_{1,n} \\ A_2 & e_{2,1} & 0 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ A_{n-1} & e_{n-1,1} & \dots & \dots & 0 & e_{n-1,n} \\ A_n & e_{n,1} & \dots & \dots & e_{n,n-1} & 0 \end{array}$$

où, pour tout $(u, v) \in \{1, \dots, n\}^2$ avec $u \neq v$,

$$e_{u,v} = e(A_u, A_v).$$

Exemple : On considère la matrice de données \mathbf{X} dans \mathbb{R}^2 définie par

$$\mathbf{X} = \begin{pmatrix} 2 & 2 \\ 7.5 & 4 \\ 3 & 3 \\ 0.5 & 5 \\ 6 & 4 \end{pmatrix}$$

On considère la méthode du voisin le plus éloigné munie de la distance euclidienne.

Le tableau des écarts associé à $(\{\omega_1\}, \dots, \{\omega_5\})$ est en fait le tableau des distances :

$$\mathbf{E} = \begin{array}{c|ccccc} & \omega_1 & \omega_2 & \omega_3 & \omega_4 & \omega_5 \\ \hline \omega_1 & 0 & 5.85 & 1.41 & 3.35 & 4.47 \\ \omega_2 & 5.85 & 0 & 4.60 & 7.07 & 1.50 \\ \omega_3 & 1.41 & 4.60 & 0 & 3.20 & 3.16 \\ \omega_4 & 3.35 & 7.07 & 3.20 & 0 & 5.59 \\ \omega_5 & 4.47 & 1.50 & 3.16 & 5.59 & 0 \end{array}$$

Soit A le couple d'individus : $A = \{\omega_1, \omega_3\}$. Par la même méthode, on obtient

$$e(\omega_2, A) = \max(e(\omega_2, \omega_1), e(\omega_2, \omega_3)) = \max(5.85, 4.60) = 5.85,$$

$$e(\omega_4, A) = \max(e(\omega_4, \omega_1), e(\omega_4, \omega_3)) = \max(3.35, 3.20) = 3.35$$

et

$$e(\omega_5, A) = \max(e(\omega_5, \omega_1), e(\omega_5, \omega_3)) = \max(4.47, 3.16) = 4.47.$$

Le tableau des écarts associé à $(\{\omega_2\}, \{\omega_4\}, \{\omega_5\}, A)$ est

$$\mathbf{E} =$$

	ω_2	ω_4	ω_5	A
ω_2	0	7.07	1.50	5.85
ω_4	7.07	0	5.59	3.35
ω_5	1.50	5.59	0	4.47
A	5.85	3.35	4.47	0

5 Algorithme de Classification Ascendante Hiérarchique (CAH)

5.1 Introduction

CAH : L'idée de l'algorithme de Classification Ascendante Hiérarchique (CAH) est de créer, à chaque étape, une partition de $\Gamma = \{\omega_1, \dots, \omega_n\}$ en regroupant les deux éléments les plus proches. Le terme "élément" désigne aussi bien un individu qu'un groupe d'individus.

Objectif : On veut

- mettre en relief les liens hiérarchiques entre les individus ou groupe d'individus,
- détecter les groupes d'individus qui se démarquent le plus.

5.2 Description de l'algorithme

Algorithme CAH : L'algorithme de CAH est décrit ci-dessous :

- On choisit un écart. On construit le tableau des écarts pour la partition initiale des n individus de Γ :

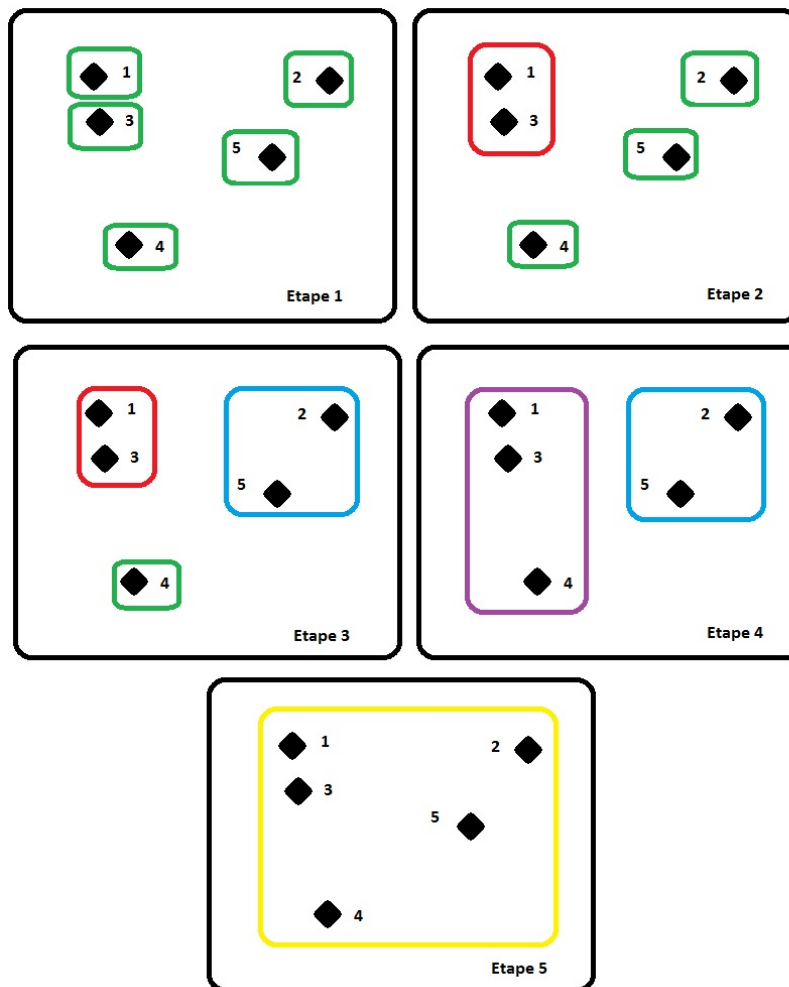
$$\mathcal{P}_0 = (\{\omega_1\}, \dots, \{\omega_n\}).$$

Chaque individu constitue un élément.

- On parcourt le tableau des écarts pour identifier le couple d'individus ayant l'écart le plus petit. Le regroupement de ces deux individus forme un groupe A . On a donc une partition de Γ de $n - 1$ éléments : A et les $n - 2$ individus restants.
- On calcule le tableau des écarts entre les $n - 1$ éléments obtenus à l'étape précédente et on regroupe les deux éléments ayant l'écart le plus petit (cela peut être deux des $n - 2$ individus, ou un individu des $n - 2$ individus restants avec A). On a donc une partition de Γ de $n - 2$ éléments.
- On itère la procédure précédente jusqu'à ce qu'il ne reste que deux éléments.

- On regroupe les deux éléments restants. Il ne reste alors qu'un seul élément contenant tous les individus de Γ .

Exemple graphique : Ci-dessous, un exemple graphique des étapes de l'algorithme CAH :



Quelques commandes R : `hclust` et `agnes` : On peut aussi utiliser la commande `hclust` (pour Hierarchical CLUSTERing) :

```
x = c(2.4, 7.1, 3.8, 1.2, 6.5, 2.1, 4.3, 3, 5.1, 4.3)
m = matrix(x, ncol = 2, nrow = 5)
d = dist(m, method = "euclidean")
cah = hclust(d, "complete")
cah$merge
```

Si on a directement affaire au tableau des distances, la commande est `as.dist` :

```
M = matrix(c(0, 23, 15, 22, 30, 26, 20, 23, 0, 26, 25, 16, 25, 33, 15, 26, 0, 28, 37, 28, 20, 22,
25, 28, 0, 22, 7, 28, 30, 16, 37, 22, 0, 20, 22, 26, 25, 28, 7, 20, 0, 18, 20, 33, 20, 28, 22, 18, 0),
byrow = T, ncol = 7)
rownames(M) = c("A","B","C","D","E","F","G")
colnames(M) = c("A","B","C","D","E","F","G")
d = as.dist(M)
cah = hclust(d, "single")
cah$merge
```

Alternativement à `hclust`, on peut utiliser la commande `agnes` (pour AGglomerative NESTing) qui offre plus de possibilités :

```
x = matrix(c(1, 16, 2, 9, 10, 16, 1, 17, 15, 2, 1, 37, 0, 14, 9, 9, 12, 4, 3, 13), ncol = 5, nrow =
4)
library(cluster)
ag = agnes(x, method = "average")
ag$merge
```

Avec `agnes`, on peut soit travailler directement avec la matrice de données (en précisant `diss = F` si la matrice est carrée, sinon la commande comprend), soit avec le tableau des distances (en précisant `diss = T`).

5.3 Dendrogramme

Dendrogramme : Les partitions de Γ faites à chaque étape de l'algorithme de la CAH peuvent se visualiser via un arbre appelé dendrogramme. Sur un axe apparaît les individus à regrouper et sur l'autre axe sont indiqués les écarts correspondants aux différents niveaux de regroupement. Cela se fait graphiquement par le biais de branches et de nœuds.

Une partition naturelle se fait en coupant l'arbre au niveau du plus grand saut de nœuds.

Exemple : On considère la matrice de données \mathbf{X} dans \mathbb{R}^2 définie par

$$\mathbf{X} = \begin{pmatrix} 2 & 2 \\ 7.5 & 4 \\ 3 & 3 \\ 0.5 & 5 \\ 6 & 4 \end{pmatrix}$$

On va regrouper les individus avec l’algorithme CAH et la méthode du voisin le plus éloigné munie de la distance euclidienne.

◦ Le tableau des écarts associé à $\mathcal{P}_0 = (\{\omega_1\}, \dots, \{\omega_5\})$ est

	ω_1	ω_2	ω_3	ω_4	ω_5
ω_1	0	5.85	1.41	3.35	4.47
ω_2	5.85	0	4.60	7.07	1.50
ω_3	1.41	4.60	0	3.20	3,16
ω_4	3.35	7.07	3.20	0	5.59
ω_5	4.47	1.50	3.16	5.59	0

Les éléments (individus) ω_1 et ω_3 ont l’écart le plus petit : ce sont les éléments les plus proches.

On les rassemble pour former le groupe : $A = \{\omega_1, \omega_3\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_1 = (\{\omega_2\}, \{\omega_4\}, \{\omega_5\}, A).$$

◦ Le tableau des écarts associé à \mathcal{P}_1 est

	ω_2	ω_4	ω_5	A
ω_2	0	7.07	1.50	5.85
ω_4	7.07	0	5.59	3.35
ω_5	1.50	5.59	0	4.47
A	5.85	3.35	4,47	0

On a

$$e(\omega_2, A) = \max(e(\omega_2, \omega_1), e(\omega_2, \omega_3)) = \max(5.85, 4.60) = 5.85,$$

$$e(\omega_4, A) = \max(e(\omega_4, \omega_1), e(\omega_4, \omega_3)) = \max(3.35, 3.20) = 3.35$$

et

$$e(\omega_5, A) = \max(e(\omega_5, \omega_1), e(\omega_5, \omega_3)) = \max(4.47, 3.16) = 4.47.$$

Les éléments (individus) ω_2 et ω_5 sont les plus proches. On les rassemble pour former le groupe : $B = \{\omega_2, \omega_5\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_2 = (\{\omega_4\}, A, B).$$

◦ Le tableau des écarts associé à \mathcal{P}_2 est

	ω_4	A	B
ω_4	0	3.35	7.07
A	3.35	0	5.85
B	7.07	5.85	0

On a

$$e(B, \omega_4) = \max(e(\omega_2, \omega_4), e(\omega_5, \omega_4)) = \max(7.07, 5.59) = 7.07$$

et

$$e(B, A) = \max(e(\omega_2, A), e(\omega_5, A)) = \max(5.85, 4.47) = 5.85.$$

Les éléments ω_4 et A sont les plus proches. On les rassemble pour former le groupe : $C = \{\omega_4, A\} = \{\omega_1, \omega_3, \omega_4\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_3 = (B, C).$$

- Le tableau des écarts associé à \mathcal{P}_3 est

	<i>B</i>	<i>C</i>
<i>B</i>	0	7.07
<i>C</i>	7.07	0

On a

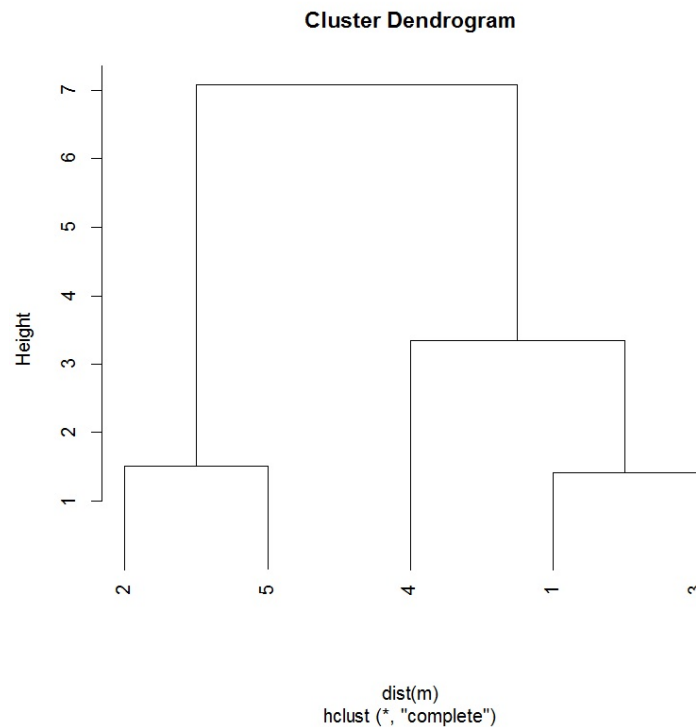
$$e(C, B) = \max(e(\omega_4, B), e(A, B)) = \max(7.07, 5.85) = 7.07.$$

Il ne reste plus que 2 éléments, *B* et *C* ; on les regroupe. On obtient la partition $\mathcal{P}_5 = \{\omega_1, \dots, \omega_5\} = \Gamma$. Cela termine l'algorithme de CAH.

Au final,

- les éléments $\{\omega_1\}$ et $\{\omega_3\}$ ont été regroupés avec un écart de 1.41,
- les éléments $\{\omega_2\}$ et $\{\omega_5\}$ ont été regroupés avec un écart de 1.50,
- les éléments $A = \{\omega_1, \omega_3\}$ et $\{\omega_4\}$ ont été regroupés avec un écart de 3.35,
- les éléments $C = \{\omega_4, A\}$ et $B = \{\omega_2, \omega_5\}$ ont été regroupés avec un écart de 7.07.

On peut donc construire le dendrogramme associé :



Comme le plus grand saut se situe entre les éléments B et C (on a $7.07 - 3.35 = 3.72$), on propose les deux groupes : B et C .

5.4 Quelques commandes R

Avec la commande `hclust` :

- D'abord, on met les données dans une matrice et on trace le nuage de points :

```
x = c(2, 7.5, 3, 0.5, 6, 2, 4, 3, 5, 4)
m = matrix(x, ncol = 2, nrow = 5)
plot(m)
```

- On calcule les distances euclidiennes :

```
dist(m)
```

- On met en œuvre l'algorithme CAH avec la méthode du voisin le plus éloigné (complete linkage) :

```
hc = hclust(dist(m), "complete")
```

On affiche les regroupements :

```
hc$merge
```

Cela renvoie :

```
      [,1] [,2]
[1,] -1  -3
[2,] -2  -5
[3,] -4   1
[4,]  2   3
```

Ainsi, à la première étape, les individus ω_1 et ω_3 ont été regroupés, formant ainsi le groupe 1, à la deuxième étape, ω_2 et ω_5 ont été regroupés, formant ainsi le groupe 2, à la troisième étape ω_4 et le groupe 1, ont été regroupés, formant ainsi le groupe 3, et pour finir, les groupes 2 et 3 ont été regroupés.

- On affiche les écarts de regroupements :

```
hc$height
```

Cela renvoie : 1.414214 1.500000 3.354102 7.071068, rejoignant ainsi la conclusion de l'exercice, à savoir :

- les éléments $\{\omega_1\}$ et $\{\omega_3\}$ ont été regroupés avec un écart de 1.41,
- les éléments $\{\omega_2\}$ et $\{\omega_5\}$ ont été regroupés avec un écart de 1.50,
- les éléments $A = \{\omega_1, \omega_3\}$ et $\{\omega_4\}$ ont été regroupés avec un écart de 3.35,
- les éléments $C = \{\omega_4, A\}$ et $B = \{\omega_2, \omega_5\}$ ont été regroupés avec un écart de 7.07.

- On trace le dendrogramme :

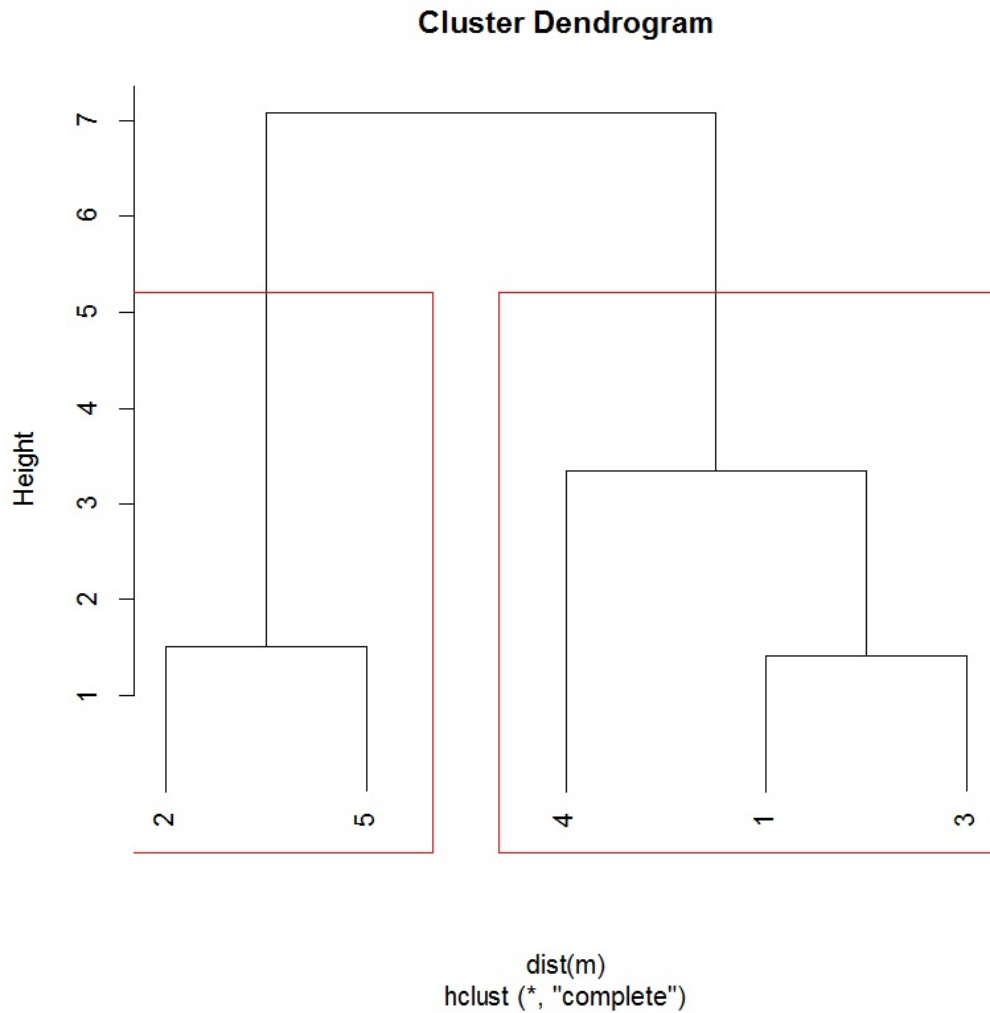
```
plot(hc, hang = -1)
```

- On peut demander à quel groupe chaque individu appartient suivant la hauteur des sauts avec la commande `cutree`. Avec 2 groupes, on a :

```
cutree(hc, k = 2)
```

- On peut alors afficher clairement les groupes :

```
rect.hclust(hc, 2)
```

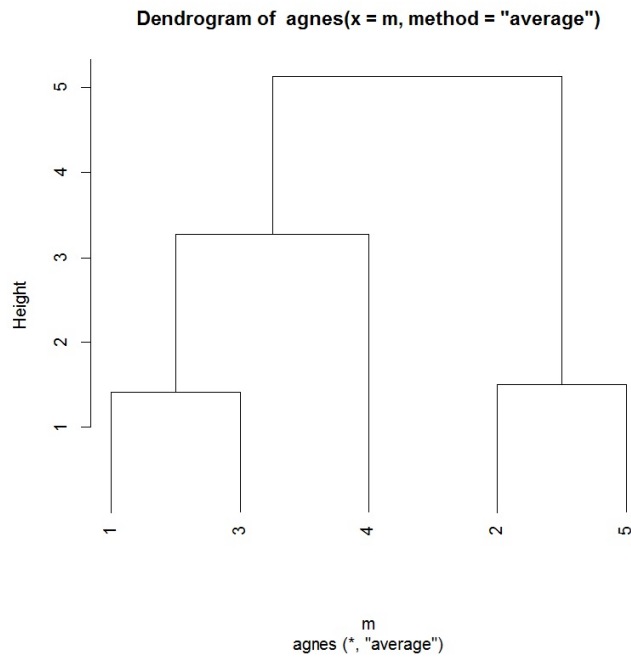


Avec la commande agnes : Avec la commande `agnes`, on propose :

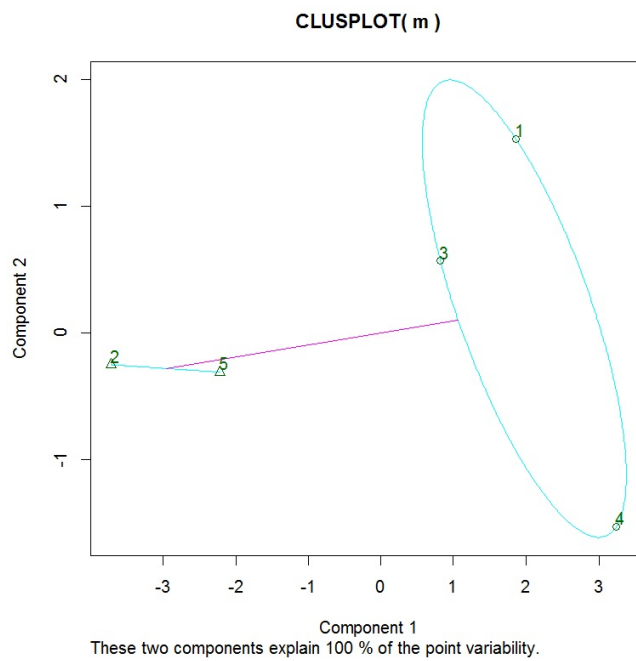
```
x = c(2, 7.5, 3, 0.5, 6, 2, 4, 3, 5, 4)
m = matrix(x, ncol = 2, nrow = 5)
library(cluster)
ag = agnes(m, method = "average")
ag$merge
ag$height
cutree(ag, k = 2)
pltree(ag, hang = -1)
clusplot(m, cutree(ag, k = 2), labels = 3)
```

On obtient

o le dendrogramme :



o une représentation graphique des regroupements obtenus (utilisant l'ACP) :



6 CAH et méthode de Ward ; compléments

Centre de gravité : On appelle centre de gravité du nuage de points $\mathcal{N} = \{m_1, \dots, m_n\}$ le point g de coordonnées $(\bar{x}_1, \dots, \bar{x}_p)$, où, pour tout $j \in \{1, \dots, p\}$,

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{j,i}.$$

Pour raison de simplicité, on dira que g est le centre de gravité associé à $\Gamma_n = \{\omega_1, \dots, \omega_n\}$; on ne se ramènera pas toujours au nuage de point associé.

Inertie totale : On appelle inertie totale de \mathcal{N} autour de son centre de gravité g le réel

$$\mathcal{I}_{tot} = \frac{1}{n} \sum_{i=1}^n d^2(\omega_i, g).$$

On peut remarquer que

$$\mathcal{I}_{tot} = \sum_{j=1}^p s_j^2, \quad s_j^2 = \frac{1}{n} \sum_{i=1}^n (x_{j,i} - \bar{x}_j)^2.$$

L'inertie de \mathcal{N} est une mesure de l'homogénéité de \mathcal{N} .

Inertie d'un sous-nuage de points : Soient $h \in \{1, \dots, n\}$ et $\mathcal{P} = (\mathcal{N}_\ell)_{\ell \in \{1, \dots, h\}}$ une partition de \mathcal{N} . Ainsi, pour tout $\ell \in \{1, \dots, h\}$, \mathcal{N}_ℓ est un sous-nuage de points de \mathcal{N} . On note

- n_ℓ le nombre d'individus représentés par \mathcal{N}_ℓ ,
- g_ℓ le centre de gravité de \mathcal{N}_ℓ , donc le point de coordonnées $(\bar{x}_{1,\ell}, \dots, \bar{x}_{p,\ell})$, où, pour tout $j \in \{1, \dots, p\}$, $\bar{x}_{j,\ell}$ désigne la moyenne des valeurs observées du caractère X_j sur les n_ℓ individus du sous-nuage \mathcal{N}_ℓ .
- **Inertie totale :** On appelle inertie totale de \mathcal{N}_ℓ autour de son centre de gravité g_ℓ le réel

$$\mathcal{I}(\mathcal{N}_\ell) = \frac{1}{n_\ell} \sum_{i \in \mathcal{N}_\ell} d^2(\omega_i, g_\ell).$$

- Inertie intra-classes : On appelle inertie intra-classes le réel

$$\mathcal{I}_{intra}(\mathcal{P}) = \sum_{\ell=1}^h \frac{n_{\ell}}{n} \mathcal{I}(\mathcal{N}_{\ell}) \quad \left(= \frac{1}{n} \sum_{j=1}^p \sum_{\ell=1}^h \sum_{i \in \mathcal{N}_{\ell}} (x_{j,i} - \bar{x}_{j,\ell})^2 \right).$$

L'inertie intra-classes mesure l'homogénéité de l'ensemble des sous-nuages de la partition.

- Inertie inter-classes : On appelle inertie inter-classes le réel

$$\mathcal{I}_{inter}(\mathcal{P}) = \sum_{\ell=1}^h \frac{n_{\ell}}{n} d^2(g_{\ell}, g) \quad \left(= \frac{1}{n} \sum_{j=1}^p \sum_{\ell=1}^h \sum_{i \in \mathcal{N}_{\ell}} (\bar{x}_{j,\ell} - \bar{x}_j)^2 \right).$$

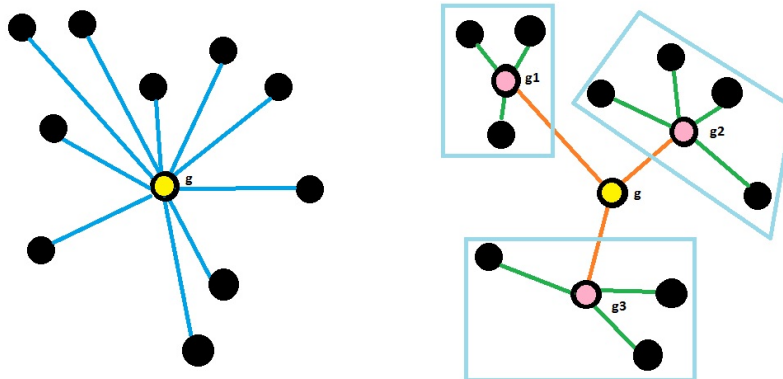
L'inertie inter-classes mesure la séparation entre les sous-nuages de la partition.

Décomposition de Huygens : Pour toute partition \mathcal{P} de \mathcal{N} , on a

$$\mathcal{I}_{tot} = \mathcal{I}_{intra}(\mathcal{P}) + \mathcal{I}_{inter}(\mathcal{P}).$$

On constate que minimiser l'inertie intra-classes est équivalent à maximiser l'inertie inter-classes.

Cette décomposition est illustrée par les schémas ci-dessous :



Le point g est le centre de gravité du nuage de points, g_1 est celui du sous-nuage de points à gauche, g_2 est celui du sous-nuage de points à droite et g_3 est celui du sous-nuage de points en bas. Les traits de couleurs représentent les distances entre les points et les centres de gravité.

Alors la somme des distances des traits bleus au carré est égale à la somme des distances des traits verts au carré plus la somme des traits orange au carré. Ce résultat est une conséquence du théorème de Pythagore.

Sur l'écart de Ward : L'utilisation de l'algorithme de CAH avec la méthode de Ward est justifiée par le résultat suivant :

Soient $\Gamma_n = \{\omega_1, \dots, \omega_n\}$ n individus et g le centre de gravité associé. Soient A et B deux groupes d'individus

- d'effectifs respectifs n_A et n_B ,
- de centres de gravité associés respectifs g_A et g_B .

Le regroupement de A et B , noté $A \cup B$, a pour centre de gravité

$$g_{A \cup B} = \frac{n_A g_A + n_B g_B}{n_A + n_B}.$$

La perte d'inertie inter-classes lors du regroupement de A et B est égale à $1/n$ multiplié par

$$n_A d^2(g_A, g) + n_B d^2(g_B, g) - (n_A + n_B) d^2(g_{A \cup B}, g) = \frac{n_A n_B}{n_A + n_B} d^2(g_A, g_B).$$

On reconnaît alors l'écart de Ward qui est donc une mesure de la perte d'inertie inter-classes lors du regroupement de A et B . Ainsi, à chaque étape de l'algorithme de CAH, on veut regrouper des éléments dont le regroupement provoque une perte minimale de l'inertie inter-classes.

Dendrogramme associé à l'écart de Ward : On peut procéder comme à l'accoutumée. Toutefois, une variante est souvent implémentée dans divers programmes dont ceux inhérents à la commande `agnes` : pour 2 éléments A et B qui se regroupent, la hauteur de la branche correspondante est donnée par la formule :

$$\sqrt{2e(A, B)}.$$

Cela ne change rien quant à la hiérarchie de la classification.

Exemple : On considère la matrice de données \mathbf{X} dans \mathbb{R}^2 définie par

$$\mathbf{X} = \begin{pmatrix} 2 & 2 \\ 7.5 & 4 \\ 3 & 3 \\ 0.5 & 5 \\ 6 & 4 \end{pmatrix}$$

On fait l'algorithme de CAH avec la méthode de Ward.

◦ Le tableau des écarts associé à $\mathcal{P}_0 = (\{\omega_1\}, \dots, \{\omega_5\})$ est

	ω_1	ω_2	ω_3	ω_4	ω_5
ω_1	0	17.12	1	5.62	10
ω_2	17.12	0	10.62	25	1,12
ω_3	1	10.62	0	5.12	5
ω_4	5.62	25	5.12	0	15.62
ω_5	10	1.12	5	15.62	0

Par exemple, on a

$$e(\omega_1, \omega_2) = \frac{1 \times 1}{1 + 1} ((2 - 7.5)^2 + (2 - 4)^2) = 17.12.$$

Les éléments (individus) ω_1 et ω_3 ont l'écart le plus petit : ce sont les éléments les plus proches.

On les rassemble pour former le groupe : $A = \{\omega_1, \omega_3\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_1 = (\{\omega_2\}, \{\omega_4\}, \{\omega_5\}, A).$$

L'inertie intra-classes de \mathcal{P}_1 est

$$\mathcal{I}_{intra}(\mathcal{P}_1) = \frac{1}{5} \times 1 = 0.2.$$

- Le centre de gravité associé à A est le point g_A de coordonnées :

$$\left(\frac{2+3}{2}, \frac{2+3}{2}\right) = (2.5, 2.5).$$

Le tableau des écarts associé à \mathcal{P}_1 est

	ω_2	ω_4	ω_5	A
ω_2	0	25	1.12	18.16
ω_4	25	0	15.62	6.83
ω_5	1.12	15.62	0	9.66
A	18.16	6.83	9.66	0

Par exemple, on a

$$e(\omega_2, A) = \frac{1 \times 2}{1 + 2} ((7.5 - 2.5)^2 + (4 - 2.5)^2) = 18.16.$$

Les éléments (individus) ω_2 et ω_5 sont les plus proches. On les rassemble pour former le groupe : $B = \{\omega_2, \omega_5\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_2 = (\{\omega_4\}, A, B).$$

L'inertie intra-classes de \mathcal{P}_2 est

$$\mathcal{I}_{intra}(\mathcal{P}_2) = 0.2 + \frac{1}{5} \times 1.12 = 0.424.$$

- Le centre de gravité associé à B est le point g_B de coordonnées $(6.75, 4)$.

Le tableau des écarts associé à \mathcal{P}_2 est

	ω_4	A	B
ω_4	0	6.83	26.7
A	6,83	0	20.31
B	26,7	20.31	0

On a, par exemple,

$$e(B, A) = \frac{2 \times 2}{2 + 2} ((6.75 - 2.5)^2 + (4 - 2.5)^2) = 20.31.$$

Les éléments ω_4 et A sont les plus proches. On les rassemble pour former le groupe : $C = \{\omega_4, A\}$. On a une nouvelle partition de Γ :

$$\mathcal{P}_3 = (B, C).$$

L'inertie intra-classes de \mathcal{P}_3 est

$$\mathcal{I}_{intra}(\mathcal{P}_3) = 0.424 + \frac{1}{5} \times 6.83 = 1.79.$$

- Le centre de gravité associé à C est le point g_C de coordonnées :

$$\left(\frac{2 + 3 + 0.5}{3}, \frac{2 + 3 + 5}{3} \right) = (1.83, 3.33).$$

Le tableau des écarts associé à \mathcal{P}_3 est

	B	C
B	0	29.58
C	29.58	0

On a

$$e(B, C) = \frac{2 \times 3}{2 + 3} ((6.75 - 1.83)^2 + (4 - 3.33)^2) = 29.58.$$

Il ne reste plus que 2 éléments, B et C ; on les regroupe. Cela donne la partition $\mathcal{P}_4 = \{\omega_1, \dots, \omega_5\} = \Gamma$.

L'inertie intra-classes de \mathcal{P}_4 est

$$\mathcal{I}_{intra}(\mathcal{P}_4) = 1.79 + \frac{1}{5} \times 29.58 = 7.706.$$

Cela termine l'algorithme de CAH.

Un indicateur qu'aucune erreur n'a été commise est l'égalité : $\mathcal{I}_{intra}(\mathcal{P}_4) = \mathcal{I}_{tot}$. On vérifie alors cela : on a

$$\mathcal{I}(\mathcal{N}) = \sigma_1^2 + \sigma_2^2,$$

avec

$$\sigma_1 = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{1,i} - \bar{x}_1)^2} = 2.5807, \quad \sigma_2 = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{2,i} - \bar{x}_2)^2} = 1.0198.$$

Donc

$$\mathcal{I}(\mathcal{N}) = 2.5807^2 + 1.0198^2 = 7.701.$$

On admet alors l'égalité (en prenant en compte les arrondis de décimales).

Au final,

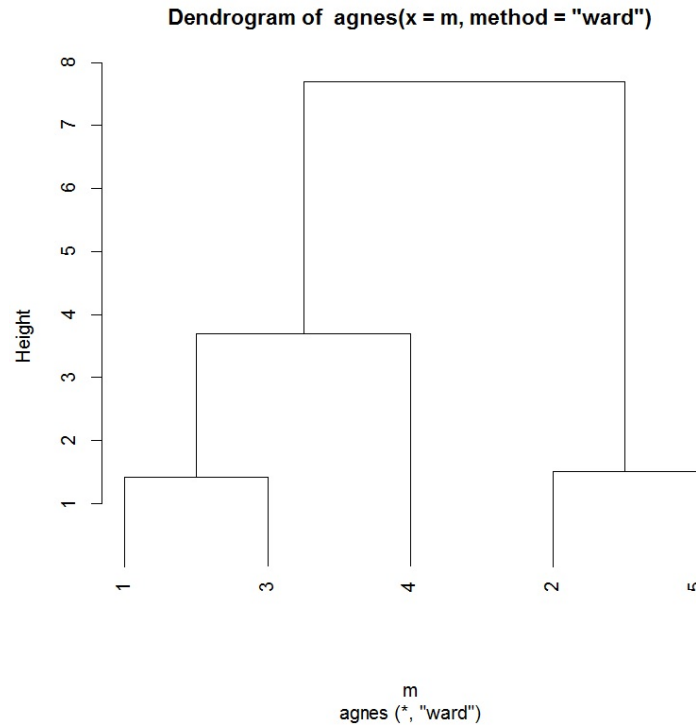
- les éléments $\{\omega_1\}$ et $\{\omega_3\}$ ont été regroupés avec un écart de 1,
- les éléments $\{\omega_2\}$ et $\{\omega_5\}$ ont été regroupés avec un écart de 1.12,
- les éléments $A = \{\omega_1, \omega_3\}$ et $\{\omega_4\}$ ont été regroupés avec un écart de 6.83,
- les éléments $B = \{\omega_2, \omega_5\}$ et $C = \{\omega_4, A\}$ ont été regroupés avec un écart de 29.58.

On peut donc construire le dendrogramme associé.

Quelques commandes R : Avec la commande `agnes`, on propose :

```
x = c(2, 7.5, 3, 0.5, 6, 2, 4, 3, 5, 4)
m = matrix(x, ncol = 2, nrow = 5)
library(cluster)
ag = agnes(m, method = "ward")
pltree(ag, hang = -1)
```

On obtient :



Comme le plus grand écart se situe entre les éléments B et C , on envisage de considérer ces deux groupes.

Notons que la formule " $\sqrt{2e(A, B)}$ " a été utilisée pour les hauteurs des branches du dendrogramme : on a $\sqrt{2 \times 1} = 1.41$, $\sqrt{2 \times 1.12} = 1.49$, $\sqrt{2 \times 6.83} = 3.69$ et $\sqrt{2 \times 29.58} = 7.69$.

7 Qualité d'une partition

Coefficient d'agglomération : On appelle coefficient d'agglomération le réel :

$$AC = \frac{1}{n} \sum_{i=1}^n \left(1 - \frac{e(\omega_i, A_i)}{e(Q, R)} \right),$$

où

- pour tout $i \in \{1, \dots, n\}$, A_i désigne le premier élément avec lequel ω_i a été regroupé,
- Q et R désignent les deux derniers groupes rassemblés à l'étape finale de l'algorithme.

On a $AC \in]0, 1[$.

Plus AC est proche de 1, plus les individus sont fortement structurés en plusieurs groupes. Une valeur proche de 0 signifie que les individus appartiennent tous à un même groupe.

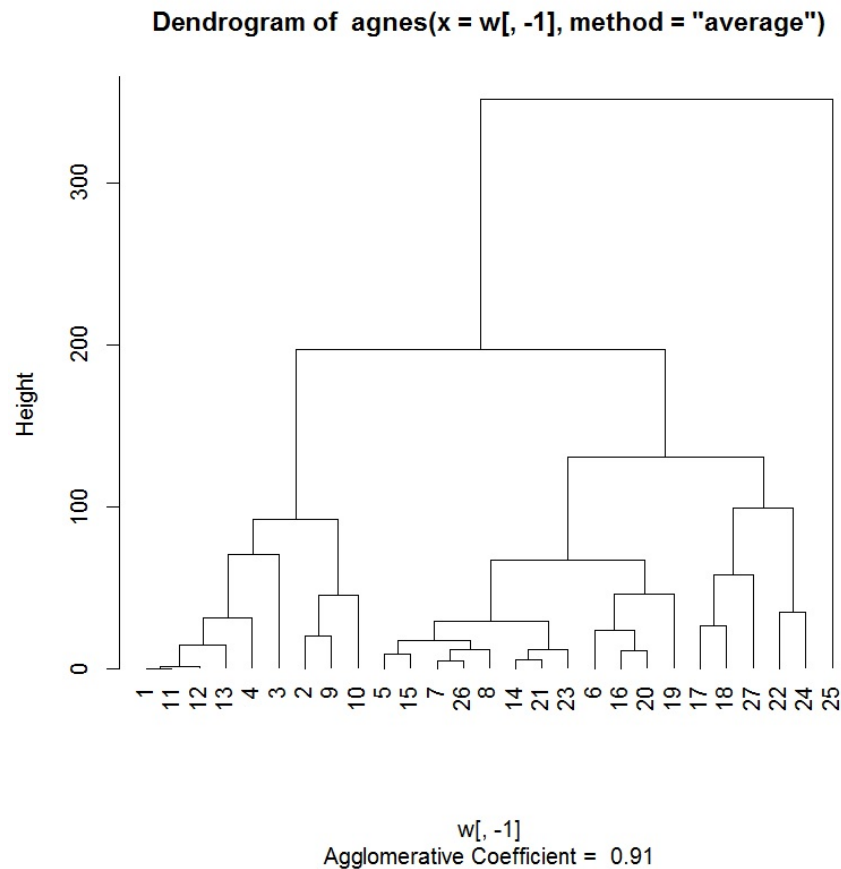
Quelques commandes R : On considère le jeu de données "aliments" dont voici l'entête :

	Individus	X1	X2	X3	X4	X5
1	BB	340	20	28	9	2.60
2	HR	245	21	17	9	2.70
3	BR	420	15	39	7	2.00
4	BS	375	19	32	9	2.50
5	BC	180	22	10	17	3.70
6	CB	115	20	3	8	1.40

Un exemple de CAH et AC avec la commande `agnes` est :

```
w = read.table("http://www.math.unicaen.fr/~chesneau/aliments.txt", header = T)
w
attach(w)
library(cluster)
ag = agnes(w[, -1], method = "average")
ag$ac
plot(ag, which = 2, hang = -1)
```


Cela renvoie le coefficient d'agglomération $AC = 0.9054413$ et le dendrogramme :



On constate alors une bonne structure de groupes, confirmée par le coefficient d'agglomération proche de 1.

Indice de silhouette : Pour tout $i \in \{1, \dots, n\}$, on appelle indice de silhouette associé à l'individu ω_i le réel :

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

où

- a_i est la moyenne des distances entre ω_i et les individus de son groupe,
- b_i est la moyenne des distances entre ω_i et les individus du groupe le plus proche de celui auquel il appartient.

On a $S(i) \in] - 1, 1[$.

Plus $S(i)$ est proche de 1, plus l'appartenance de ω_i a son groupe est justifiée.

Ainsi, les individus ayant des grands indices de silhouette sont bien regroupés.

Si l'indice de silhouette d'un individu est négatif, l'individu n'est pas dans le bon groupe et pourrait être déplacé dans le groupe le plus proche.

Largeur de silhouette : On appelle largeur de silhouette de la partition le réel :

$$S = \frac{1}{n} \sum_{i=1}^n S(i).$$

On a alors l'interprétation suivante :

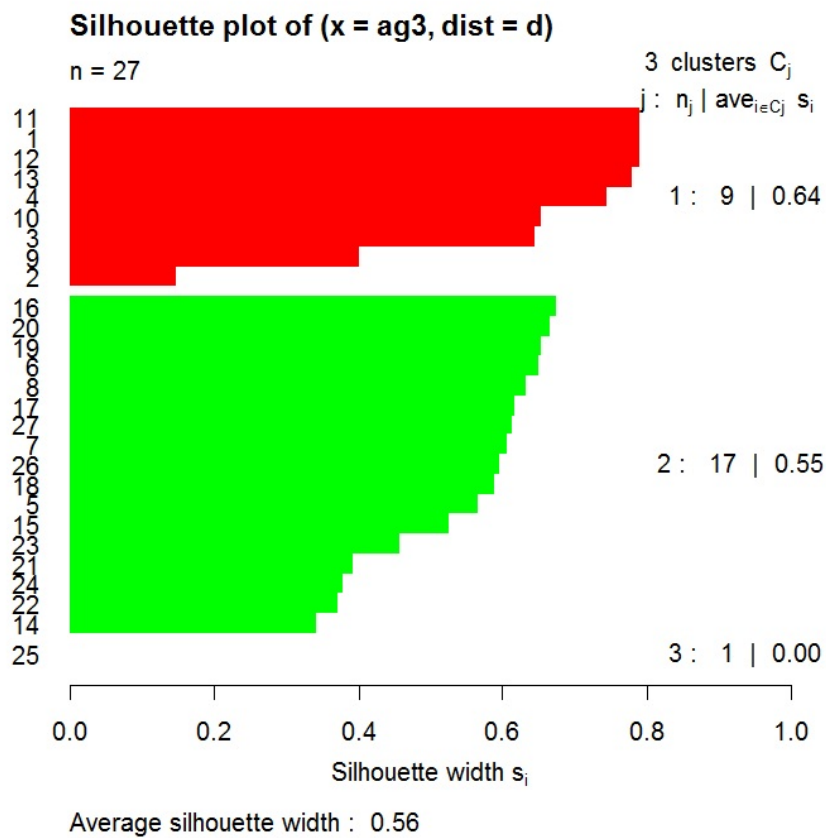
Valeur de S	Nature de la structure
$\in]0.51, 1]$	Forte
$\in]0.31, 0.50]$	Raisnable
$\in [0, 0.30[$	Faible
$\in [-1, 0[$	Inexistante

On peut également calculer S pour les individus d'un groupe.

Quelques commandes R : Ci-dessous un exemple utilisant les indices de silhouette :

```
w = read.table("http://www.math.unicaen.fr/~chesneau/aliments.txt", header = T)
w
attach(w)
d = dist(w[, -1], method = "euclidean")
library(cluster)
ag = agnes(d, method = "average")
ag3 = cutree(ag, 3)
si = silhouette(ag3, d)
plot(si, col = c("red", "green", "blue"))
```

Cela renvoie le graphique :



Cela renvoie une largeur de silhouette de 0.56, soit une structure forte de la partition, un individu isolé dans le troisième groupe (ω_5) et pas d'individu mal regroupé ; aucun indice de silhouette n'est négatif.

Remarques : D'autres indices de qualité existent. Il y a notamment :

- l'indice d'inertie,
- l'indice de connectivité,
- l'indice de Dunn,
- le Cubic Clustering Criterion (CCC).

8 ACP et CAH

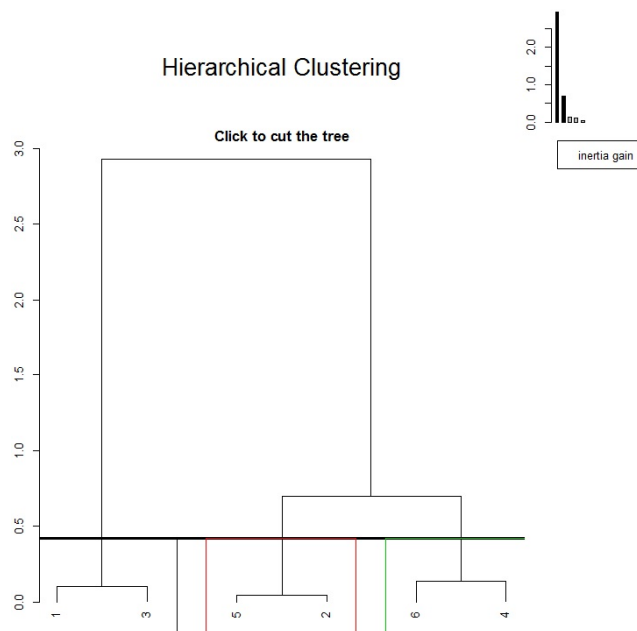
Idée : Lorsque l'on travaille avec plus de 3 variables quantitatives, donc $p \geq 3$, on peut faire une analyse en composantes principales (ACP) et considérer les coordonnées des individus sur le plan principal.

Quelques commandes R : Un exemple de commandes utilisant le package FactoMineR et l'écart de Ward est présenté ci-dessous :

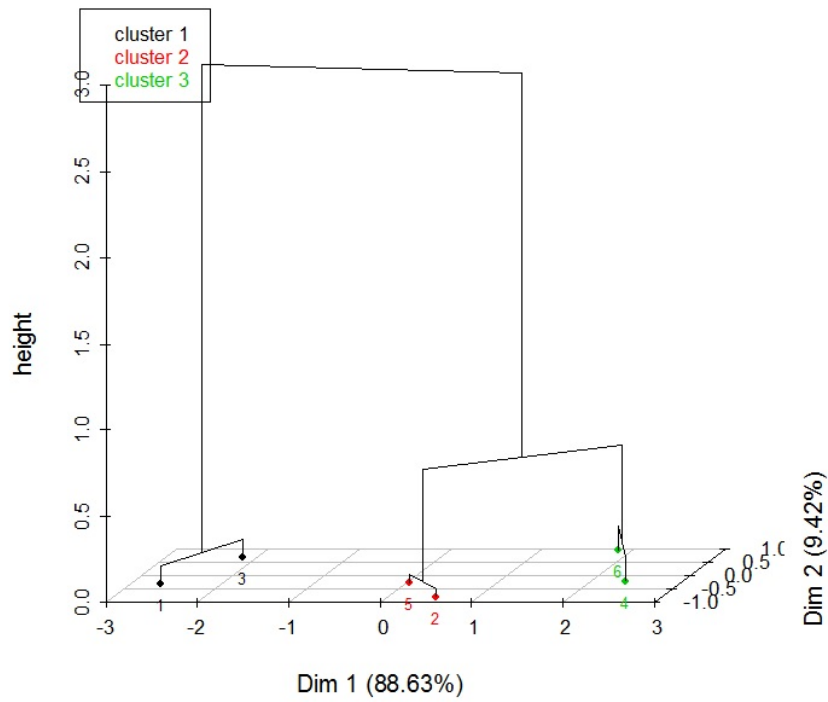
```
library(FactoMineR)
x = c(19, 7, 20, 1, 10, 2, 17, 8, 19, 6, 11, 12, 2, 12, 9, 18, 12, 18, 8, 12, 9, 17, 12, 18)
m = matrix(x, ncol = 4, nrow = 6)
m
acp = PCA(m, ncp = 2, graph = F)
res = HCPC(acp)
```

On décide de faire 3 groupes (la coupure se fait interactivement sur le dendrogramme affiché).

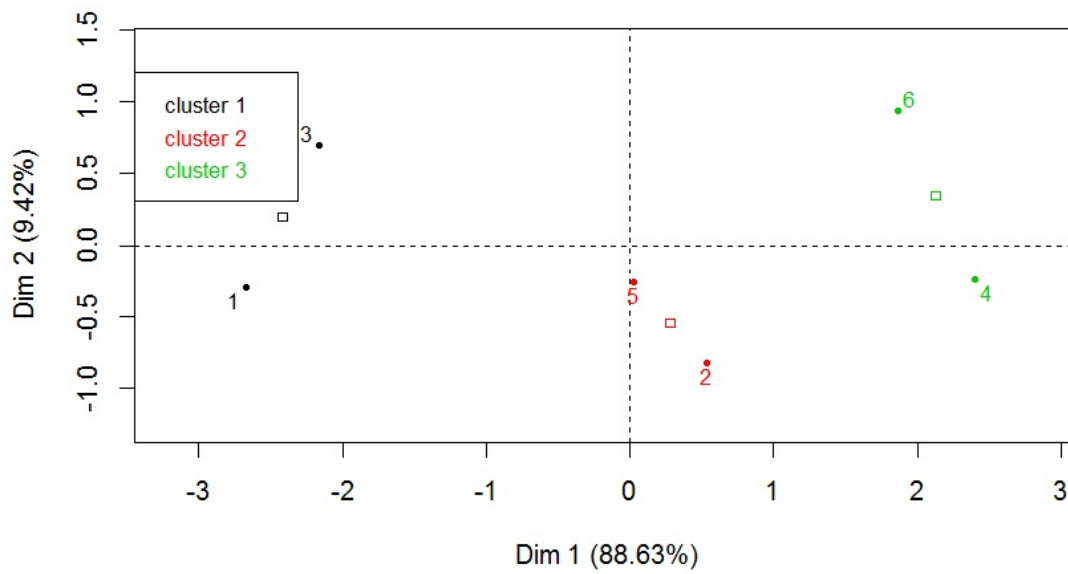
On obtient les graphiques :



Hierarchical clustering on the factor map



Factor map



9 Caractérisation des groupes

Parangons : Pour chaque groupe formé, on appelle parangon l'individu dont les coordonnées sont les plus proches du centre de gravité du groupe. Le profil de cet individu caractérise alors le groupe auquel il appartient.

Caractères dominants dans la classification : Pour connaître les caractères qui jouent un rôle important dans la classification, on peut faire p ANOVA à 1 facteur associées aux p caractères considérés. Plus précisément, pour tout $j \in \{1, \dots, p\}$, on fait une ANOVA à 1 facteur avec :

- le facteur G ayant pour modalités les q groupes formés : G_1, \dots, G_q ,
- le caractère X_j (variable quantitative).

Pour chacun des p tests d'hypothèses, le test de Fisher renvoie alors une p-valeur évaluant l'influence du facteur sur le caractère considéré. Ainsi, les caractères associés aux p-valeurs les plus petites sont ceux qui importent le plus dans la classification obtenue.

Caractères dominants d'un groupe : On peut déterminer les caractères dominants pour chacun des groupes formés. Pour se faire, pour chacun des caractères, on peut faire un test d'hypothèses reposant sur loi normale. Soient G_1, \dots, G_q les q groupes formés. Pour tout $g \in \{1, \dots, q\}$ et tout $j \in \{1, \dots, p\}$, on calcule :

- $\bar{x}_{j,g}$: la moyenne des valeurs du caractère X_j pour les individus du groupe g ,
- \bar{x}_j : la moyenne des valeurs du caractère X_j ,
- n_g : le nombre d'individus dans le groupe g ,
- s_j : l'écart-type corrigé des valeurs du caractère X_j ,
- le z_{obs} :

$$z_{obs} = \frac{\bar{x}_{j,g} - \bar{x}_j}{\sqrt{\frac{s_j^2}{n_g} \left(\frac{n - n_g}{n - 1} \right)}}.$$

On considère alors la p-valeur :

$$\text{p-valeur} = \mathbb{P}(|Z| \geq |z_{obs}|), \quad Z \sim \mathcal{N}(0, 1).$$

Ainsi, pour tout $g \in \{1, \dots, q\}$, on obtient p p-valeurs qu'il convient de classer par ordre croissant. Pour chaque groupe, les plus petites correspondent aux caractères qui importent le plus dans la constitution de ce groupe.

Quelques commandes R : On considère le jeu de données "zebu" dont voici l'entête :

	vif	carc	qsup	tota	gras	os
1	395	224	35.10	79.10	6.00	14.90
2	410	232	31.90	73.40	9.70	16.40
3	405	233	30.70	76.50	7.50	16.50
4	405	240	30.40	75.30	8.70	16.00
5	390	217	31.90	76.50	7.80	15.70
6	405	243	32.10	77.40	7.10	15.50

On décrit ci-dessous des exemples de commandes R avec FactoMineR :

```
library(FactoMineR)
w = read.table("http://www.math.unicaen.fr/~chesneau/zebu.txt", header = T)
w
attach(w)
acp = PCA(w, ncp = 5, graph = F)
res = HCPC(acp, consol = F)
```

On décide de faire 2 groupes.

◦ Classification :

```
res$clust
```

◦ Parangons (donnés par les premiers noms de chaque liste) :

```
res$desc.ind
```

L'individu ω_{13} est un parangon pour le premier groupe et ω_3 est un parangon pour le deuxième.

◦ Étude des caractères dominants dans la classification et des caractères dominants d'un groupe :

```
res$desc.var
```

On remarque que *gras* (caractère X_3), *tota* (caractère X_4) et *qsup* (caractère X_5) caractérisent le mieux la partition et ce sont les caractères dominants des groupes.

10 Algorithme des centres mobiles (k means)

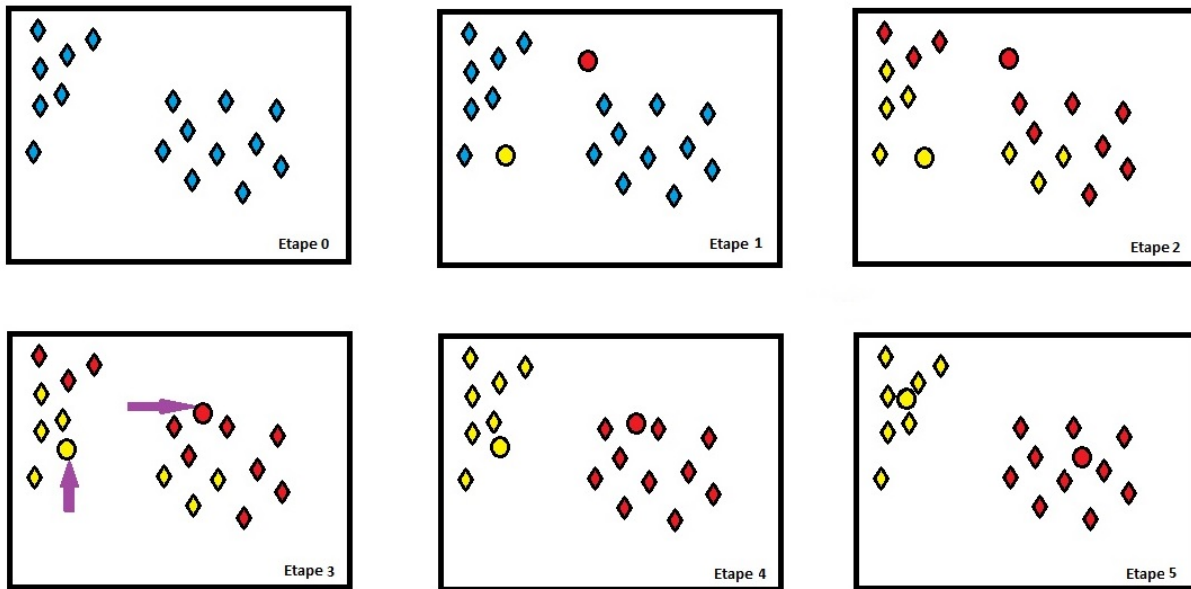
Algorithme des centres mobiles (k means) : L'algorithme des centres mobiles vise à classer une population Γ en q classes. Cela se fait de manière automatique ; il n'y a pas de lien hiérarchique dans les regroupements contrairement à l'algorithme CAH. Il est le mieux adapté aux très grands tableaux de données.

L'algorithme des centres mobiles avec la méthode de Lloyd (la plus standard) est décrit ci-dessous :

- On choisit q points au hasard dans \mathbb{R}^p . Ces points sont appelés centres.
- On calcule le tableau de distances entre tous les individus et les q centres.
- On forme alors q groupes de la manière suivante : chaque groupe est constitué d'un centre et des individus les plus proches de ce centre que d'un autre. On obtient une partition \mathcal{P}_1 de Γ .
- On calcule le centre de gravité de chacun des q sous-nuages de points formés par les q groupes. Ces q centres de gravité sont nos nouveaux q centres.
- On calcule le tableau de distances entre tous les individus et les nouveaux q centres.
- On forme alors q groupes, chaque groupe étant constitué d'un centre et des individus les plus proches de ce centre que d'un autre. On a une nouvelle partition \mathcal{P}_2 de Γ .
- On itère la procédure précédente jusqu'à ce que deux itérations conduisent à la même partition.

Remarque importante : La classification des individus dépend du choix des centres initiaux. Plusieurs méthodes existent pour choisir judicieusement ces centres.

Illustration : Une illustration de l'algorithme des centres mobiles est présentée ci-dessous :



Exemple : Dans une étude industrielle, on a étudié 2 caractères : X_1 et X_2 , sur 6 individus $\omega_1, \dots, \omega_6$.

Les données recueillies sont :

	X_1	X_2
ω_1	-2	2
ω_2	-2	-1
ω_3	0	-1
ω_4	2	2
ω_5	-2	3
ω_6	3	0

1. Dans un premier temps, on fait une classification par l'algorithme des centres mobiles avec, pour centres initiaux, c_1^0 de coordonnées $(-1, -1)$ et c_2^0 de coordonnées $(2, 3)$.
2. Dans un deuxième temps, on fait de même avec, pour centres initiaux, c_1^0 de coordonnées $(-1, 2)$ et c_2^0 de coordonnées $(1, 1)$.
1. \circ On considère les centres initiaux c_1^0 de coordonnées $(-1, -1)$ et c_2^0 de coordonnées $(2, 3)$.

Le tableau des distances entre les individus et ces centres est

	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6
c_1^0	3.16	1	1	4.24	4.12	4.12
c_2^0	4.12	5.66	4.47	1	4	3.16

Exemple de calcul : $d(\omega_1, c_1^0) = \sqrt{(-2 - (-1))^2 + (2 - (-1))^2} = 3.16$.

D'où les deux groupes :

$$A = \{\omega_1, \omega_2, \omega_3\}, \quad B = \{\omega_4, \omega_5, \omega_6\}.$$

- On considère deux nouveaux centres, c_1^1 et c_2^1 , lesquels sont les centres de gravité des deux groupes A et B . Donc c_1^1 a pour coordonnées $(\frac{-2-2+0}{3}, \frac{2-1-1}{3}) = (-1.33, 0)$ et c_2^1 a pour coordonnées $(\frac{2-2+3}{3}, \frac{2+3+0}{3}) = (1, 1.67)$.

Le tableau des distances entre les individus et ces centres est

	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6
c_1^1	2.11	1.20	1.66	3.88	3.07	4.33
c_2^1	3.02	4.02	2.85	1.05	3.28	2.61

D'où les deux groupes :

$$A = \{\omega_1, \omega_2, \omega_3, \omega_5\}, \quad B = \{\omega_4, \omega_6\}.$$

- On considère deux nouveaux centres, c_1^2 et c_2^2 , lesquels sont les centres de gravité des deux groupes A et B . Donc c_1^2 a pour coordonnées $(\frac{-2-2+0-2}{4}, \frac{2-1-1+3}{4}) = (-1.5, 0.75)$ et c_2^2 a pour coordonnées $(\frac{2+3}{2}, \frac{2+0}{2}) = (2.5, 1)$.

Le tableau des distances entre les individus et ces centres est

	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6
c_1^2	1.35	1.82	2.30	3.72	2.30	4.56
c_1^2	4.61	4.92	3.20	1.12	4.92	1.12

D'où les deux groupes :

$$A = \{\omega_1, \omega_2, \omega_3, \omega_5\}, \quad B = \{\omega_4, \omega_6\}.$$

On retrouve la même classification que l'étape précédente, on arrête l'algorithme.

2. Considérons maintenant les centres initiaux c_1^0 de coordonnées $(-1, 2)$ et c_2^0 de coordonnées $(1, 1)$.

- On considère les centres initiaux c_1^0 de coordonnées $(-1, 2)$ et c_2^0 de coordonnées $(1, 1)$.

Le tableau des distances entre les individus et ces centres est

	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6
c_1^0	1	3.16	3.16	3	1.41	4.47
c_2^0	3.16	3.60	2.24	1.41	3.60	2.24

D'où les deux groupes :

$$A = \{\omega_1, \omega_2, \omega_5\}, \quad B = \{\omega_3, \omega_4, \omega_6\}.$$

- On considère deux nouveaux centres, c_1^1 et c_2^1 , lesquels sont les centres de gravité des deux groupes A et B. Donc c_1^1 a pour coordonnées $(\frac{-2-2-2}{3}, \frac{2-1+3}{3}) = (-2, 1.33)$ et c_2^1 a pour coordonnées $(\frac{0+2+3}{3}, \frac{-1+2+0}{3}) = (1.67, 0.33)$.

Le tableau des distances entre les individus et ces centres est

	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6
c_1^1	0.67	2.33	3.07	4.06	1.67	5.17
c_2^1	4.03	3.90	2.13	1.70	4.54	1.37

D'où les deux groupes :

$$A = \{\omega_1, \omega_2, \omega_5\}, \quad B = \{\omega_3, \omega_4, \omega_6\}.$$

On retrouve la même classification que l'étape précédente, on arrête l'algorithme.

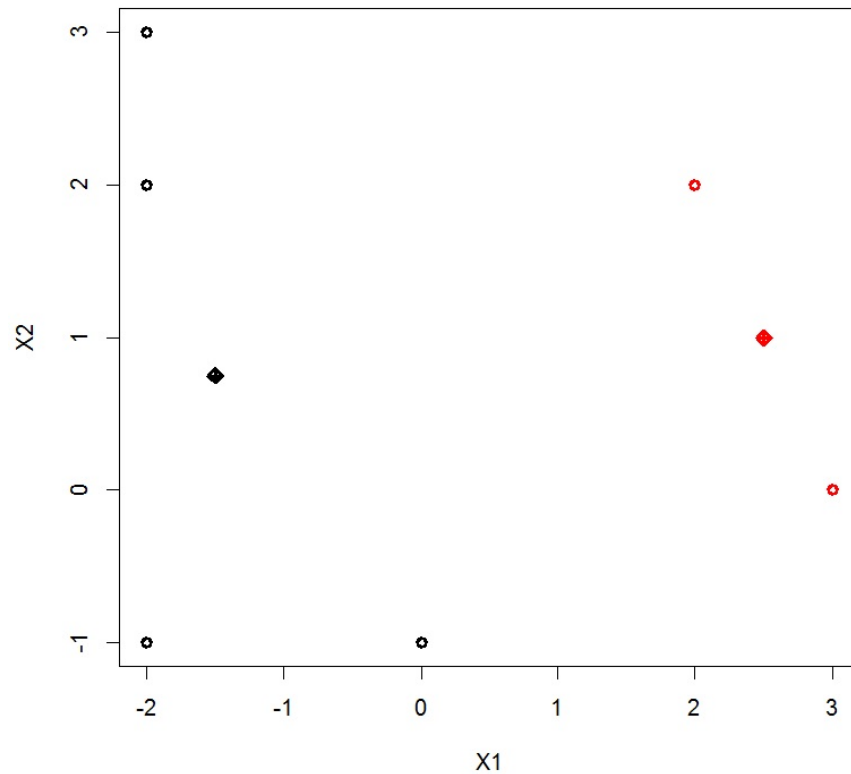
Remarque : On obtient deux classifications différentes suivant les choix des centres initiaux.

Commandes R de l'exemple :

Pour le 1., les commandes R associées sont :

```
x = c(-2, -2, 0, 2, -2, 3, 2, -1, -1, 2, 3, 0)
m = matrix(x, ncol = 2, nrow = 6)
clus = kmeans(m, centers=rbind(c(-1, -1), c(2, 3)), algorithm = "Lloyd")
clus$cluster
clus$centers
plot(m, col = clus$cluster, pch = 1, lwd = 3, xlab = "X1", ylab = "X2")
points(clus$centers, col = 1 :2, pch = 9, lwd = 3)
```

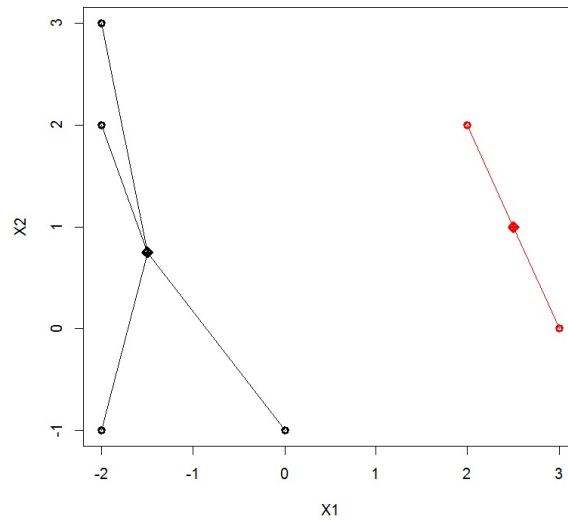
Cela renvoie les groupes d'affectation de chaque individu (`clus$cluster`), les coordonnées des centres de gravité de chaque groupe (`clus$centers`) et le graphique :



On peut rejoindre les individus au centre de gravité dans chaque groupe avec la commande `segments` :

```
segments(m[clus$cluster == 1, ][ ,1], m[clus$cluster == 1, ][ ,2], clus$centers[1, 1],
clus$centers[1, 2])
segments(m[clus$cluster == 2, ][ ,1], m[clus$cluster == 2, ][ ,2], clus$centers[2, 1],
clus$centers[2, 2], col = 2)
```

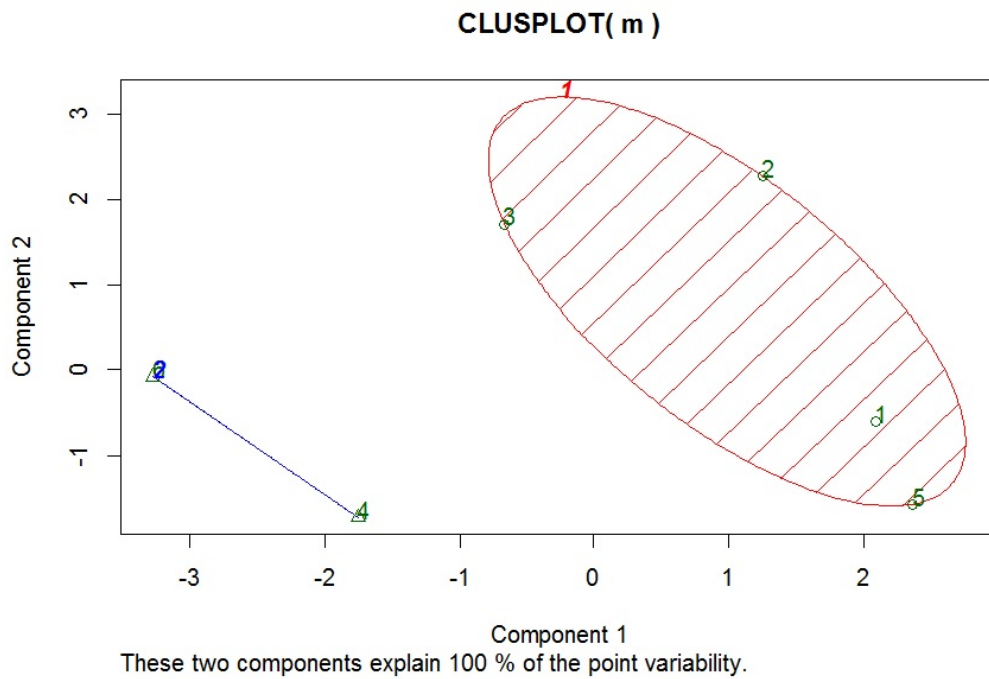
Cela renvoie :



On peut aussi utiliser `clusplot` pour la visualisation des groupes :

```
library(cluster)
clusplot(m, cluster, color = T, shade = T, labels = 2, lines = 0)
```

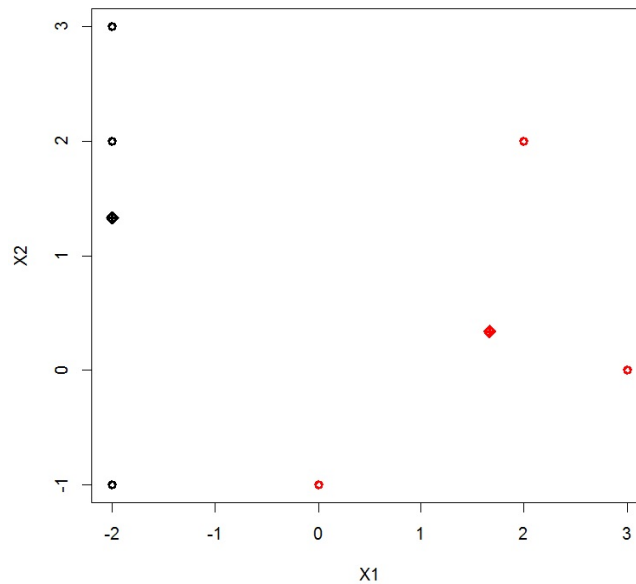
Cela renvoie :



Pour le 2., les commandes R associées sont :

```
x = c(-2, -2, 0, 2, -2, 3, 2, -1, -1, 2, 3, 0)
m = matrix(x, ncol = 2, nrow = 6)
clus = kmeans(m, centers=rbind(c(-1, 2), c(1, 1)), algorithm = "Lloyd")
clus$cluster
clus$centers
plot(m, col = clus$cluster, pch = 1, lwd = 3, xlab = "X1", ylab = "X2")
points(clus$centers, col = 1 :2, pch = 9, lwd = 3)
```

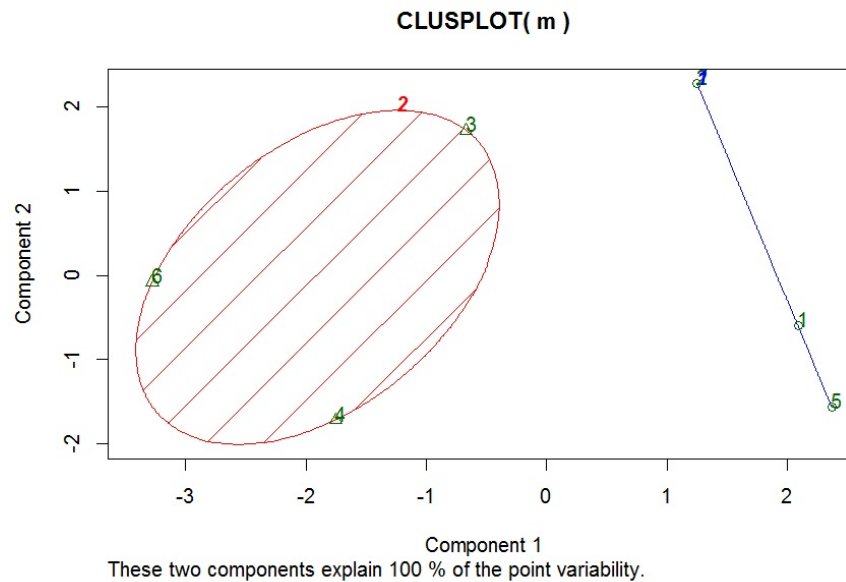
Cela renvoie le graphique :



On peut aussi utiliser `clusplot` pour la visualisation des groupes :

```
library(cluster)
clusplot(m, clus$cluster, color = T, shade = T, labels = 2, lines = 0)
```

Cela renvoie :



Présentation du jeu de données iris : Une célèbre jeu de données étudié par le statisticien Fisher en 1936 est "les iris de Fisher". Pour 3 variétés d'iris : Setosa, Versicolor, Virginica, et pour 150 iris par variété, on considère 4 caractères quantitatifs :

- X_1 la longueur en cm d'un pétale,
- X_2 la largeur en cm d'un pétale,
- X_3 la longueur en cm d'un sépale,
- X_4 la largeur en cm d'un sépale.

Ce sont les variables explicatives X_1, X_2, X_3 et X_4 . La variable à expliquer Y est une variable qualitative dont les modalités sont les espèces d'iris $\{setosa, versicolor, virginica\}$.

Voici l'entête du jeu de données "iris" :

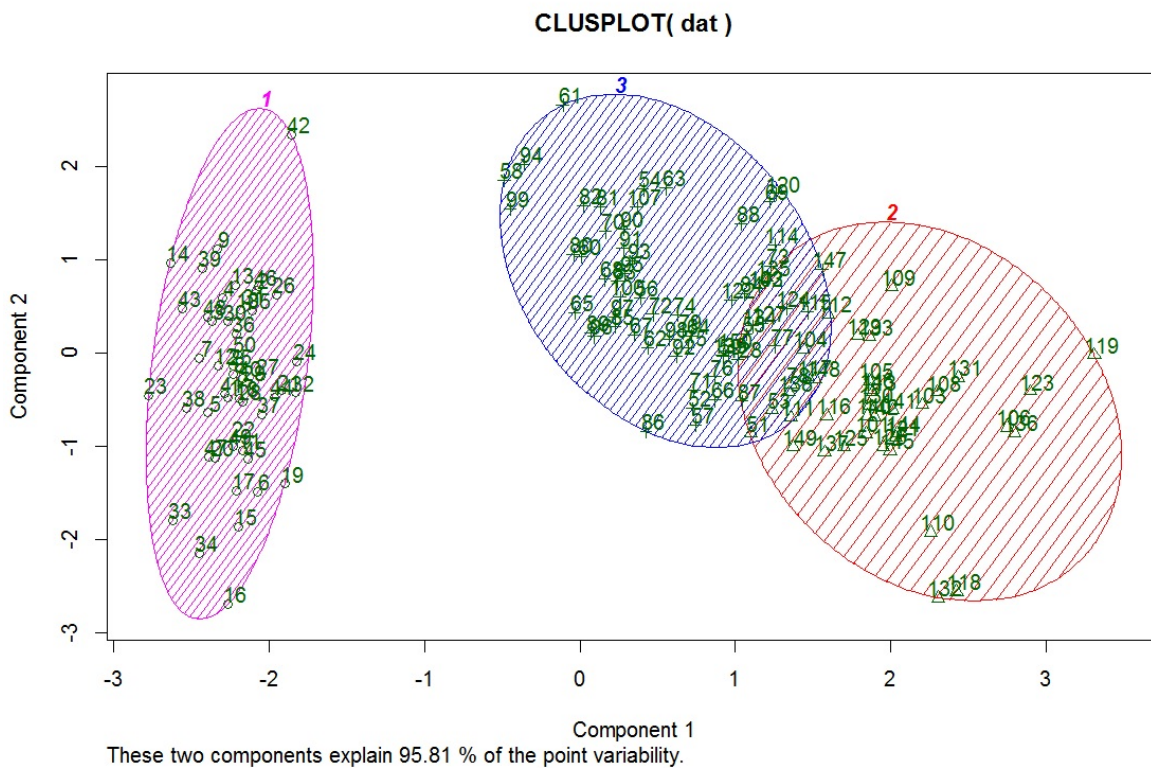
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.10	3.50	1.40	0.20	setosa
2	4.90	3.00	1.40	0.20	setosa
3	4.70	3.20	1.30	0.20	setosa
4	4.60	3.10	1.50	0.20	setosa
5	5.00	3.60	1.40	0.20	setosa
6	5.40	3.90	1.70	0.40	setosa

Quelques commandes R : Un exemple de commandes R utilisant le jeu de données iris et l’algorithme des centres mobiles est présenté ci-dessous :

```
dat = iris [,1 :4]
library(stats)
clus = kmeans(dat, centers = dat[c(15, 135, 65), ], algorithm = "Lloyd")
library(cluster)
clusplot(dat, clus$cluster, color = T, shade = T, labels = 2, lines = 0)
```

Dans cet exemple, on a donc considéré l’algorithme des centres mobiles avec 3 centres initiaux qui sont les individus correspondants aux lignes 15, 135 et 65 du jeu de données iris.

Le graphique obtenu est :



Méthodes alternatives : Il existe de nombreuses méthodes autres que celle de Lloyd. Il y a notamment :

- la méthode de Forgy : les centres initiaux sont tirés au hasard parmi ceux associés aux individus de Γ .

```
clus = kmeans(dat, 3, algorithm = "Forgy")
```

- la méthode de MacQueen : les centres sont recalculés à chaque réaffectation d'un seul individu.

```
clus = kmeans(dat, 3, algorithm = "MacQueen")
```

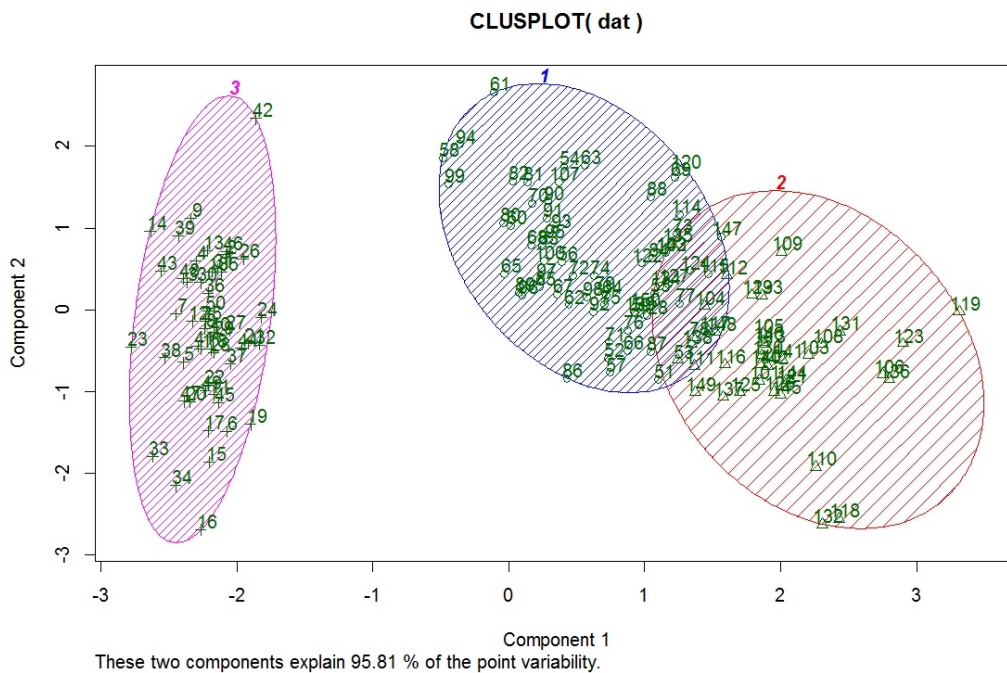
- la méthode de Hartigan-Wong : c'est la méthode par défaut de la commande kmeans. Elle est considérée comme la plus robuste de toutes.

```
clus = kmeans(dat, 3)
```

Quelques commandes R : Un exemple de commandes R pour utiliser l'algorithme des centres mobiles avec la méthode de Hartigan-Wong :

```
dat = iris [,1 :4]
clus = kmeans(dat, 3)
library(cluster)
clusplot(dat, clus$cluster, color = T, shade = T, labels = 2, lines = 0)
```

Le graphique obtenu est :



Le résultat est identique à celui obtenu avec la méthode de Lloyd ; cela est un hasard.

Alternatives à l'algorithme des centres mobiles : Il existe de nombreuses alternatives à l'algorithme des centres mobiles. Il y a notamment :

- la méthode PAM (Partition Around Medoids) : cette méthode a la particularité de marcher aussi avec un tableau des distances et d'être moins sensible que l'algorithme des centres mobiles aux individus atypiques.

```
library(cluster)
clus = pam(dat, 3)
plot(clus)
```

- la méthode CLARA (Clustering LARge Application) :

```
library(cluster)
clus = clara(dat, 3)
plot(clus)
```

- la méthode FANNY :

```
library(cluster)
clus = fanny(dat, 3)
plot(clus)
```

11 Consolidation de l'algorithme de CAH

Idée : On peut consolider/améliorer les regroupements obtenus via l'algorithme de CAH en utilisant l'algorithme des centres mobiles. On prend alors pour centres initiaux les parangons obtenus lors de la CAH. Il est donc possible que des individus changent de groupes.

Quelques commandes R : Un exemple de code R utilisant la librairie FactoMineR et la commande `consol = T` est :

```
library(FactoMineR)
w = read.table("http://www.math.unicaen.fr/~chesneau/zebu.txt", header = T)
w
attach(w)
acp = PCA(w, ncp = 5, graph = F)
res = HCPC(acp, consol = T)
res$data.clust
```

En fait, à la base, on dispose du jeu de données `zebu` avec une classification des individus en deux groupes. Celle-ci est visible par :

```
w2 = read.table("http://www.math.unicaen.fr/~chesneau/zebu-g.txt", header = T)
w2
```

On constate alors que la classification obtenue avec la CAH consolidée est exacte. Ce n'est pas le cas sans consolidation :

```
res2 = HCPC(acp, consol = F)
res2$data.clust
```

12 CAH avec des caractères qualitatifs

Indice de similarité : Soit $\Gamma = \{\omega_1, \dots, \omega_n\}$. On appelle indice de similarité toute application $s :$

$\Gamma^2 \rightarrow [0, \infty[$ telle que, pour tous individus ω et ω_* dans Γ , on a

- $s(\omega, \omega_*) = s(\omega_*, \omega)$,
- $s(\omega, \omega_*) \leq s(\omega, \omega)$,
- on a $s(\omega, \omega_*) = s(\omega, \omega)$ si, et seulement si, $\omega = \omega_*$.

Règle centrale : Plus l'indice de similarité entre deux individus est élevé, plus ils se ressemblent.

Cas des caractères quantitatifs : Quand les caractères X_1, \dots, X_p sont quantitatifs, on peut choisir comme indice de similarité la fonction s telle que

$$s(\omega, \omega_*) = d_{max} - d(\omega, \omega_*),$$

où d désigne une distance et $d_{max} = \max_{\omega \times \omega_* \in \Gamma^2} d(\omega, \omega_*)$.

Tableau disjonctif complet : Dans le cas où les caractères X_1, \dots, X_p sont qualitatifs, on peut présenter les données sous la forme d'un tableau disjonctif complet (TDC) de dimension $n \times r$, où r est le nombre total de modalités des p caractères considérés. Pour tout $i \in \{1, \dots, n\}$, la i -ème ligne du tableau est constituée du vecteur $(n_{1,i}, \dots, n_{k,i}, \dots, n_{r,i})$, avec

$$n_{k,i} = \begin{cases} 1 & \text{si } i \text{ possède la modalité } k, \\ 0 & \text{sinon.} \end{cases}$$

Il n'y a donc que des 0 et 1 dans le tableau.

Valeurs intermédiaires : Pour tout $(u, v) \in \{1, \dots, n\}^2$, on pose

- le nombre de (1, 1) aux (u, v) -ème lignes du TDC,
- le nombre de (1, 0) aux (u, v) -ème lignes du TDC,
- le nombre de (0, 1) aux (u, v) -ème lignes du TDC,
- le nombre de (0, 0) aux (u, v) -ème lignes du TDC.

Notons que

$$a_{u,v} + b_{u,v} + c_{u,v} + d_{u,v} = r.$$

Indices de similarité usuels : Les indices de similarité les plus utilisés sont les suivants :

- Indice de Russel et Rao :

$$s(\omega_u, \omega_v) = \frac{a_{u,v}}{r}.$$

- Indice de Jaccard :

$$s(\omega_u, \omega_v) = \frac{a_{u,v}}{a_{u,v} + b_{u,v} + c_{u,v}} = \frac{a_{u,v}}{r - d_{u,v}}.$$

- Indice de Dice :

$$s(\omega_u, \omega_v) = \frac{2a_{u,v}}{2a_{u,v} + b_{u,v} + c_{u,v}}.$$

- Indice de d'Anderberg :

$$s(\omega_u, \omega_v) = \frac{a_{u,v}}{a_{u,v} + 2(b_{u,v} + c_{u,v})}.$$

- Indice de Rogers et Tanimoto :

$$s(\omega_u, \omega_v) = \frac{a_{u,v} + d_{u,v}}{a_{u,v} + d_{u,v} + 2(b_{u,v} + c_{u,v})}.$$

- Indice de Pearson :

$$s(\omega_u, \omega_v) = \frac{a_{u,v}d_{u,v} - b_{u,v}c_{u,v}}{\sqrt{(a_{u,v} + b_{u,v})(a_{u,v} + c_{u,v})(d_{u,v} + b_{u,v})(d_{u,v} + c_{u,v})}}.$$

- Indice de Yule :

$$s(\omega_u, \omega_v) = \frac{a_{u,v}d_{u,v} - b_{u,v}c_{u,v}}{a_{u,v}d_{u,v} + b_{u,v}c_{u,v}}.$$

Distances à partir d'un indice de similarité et CAH : À partir d'un indice de similarité s , on définit une application $d_* : \Gamma^2 \rightarrow [0, \infty[$ par

$$d_*(\omega_u, \omega_v) = s_{max} - s(\omega_u, \omega_v),$$

où $s_{max} = s(\omega_u, \omega_u) (= s(\omega_v, \omega_v))$.

Cette application est appelée dissimilarité.

On peut alors faire de la CAH avec cette dissimilarité d_* au lieu de d et l'écart de son choix.

Sur l'indice de Jaccard : Si s est l'indice de Jaccard, alors $s_{max} = 1$ et on peut prendre la dissimilarité :

$$d_*(\omega_u, \omega_v) = 1 - s(\omega_u, \omega_v) = 1 - \frac{a_{u,v}}{r - d_{u,v}}.$$

Quelques commandes R : Ci-dessous, des exemples de commandes R illustrant plusieurs dissimilarités :

```
m = matrix(sample(c(0, 1), 100, replace = T), ncol = 10)
m
library(arules)
d = dissimilarity(m, method = "jaccard")
d
d = dissimilarity(m, method = "pearson")
d
```

Exemple : On interroge 6 individus en leur demandant leur sexe X_1 (F : femme, H : homme), leur type de logement X_2 (R : rural, U : urbain) et leur état civil X_3 (C : célibataire, M : marié, A : autre). On obtient :

	X_1	X_2	X_3
ω_1	H	U	C
ω_2	F	U	C
ω_3	F	R	M
ω_4	F	U	A
ω_5	H	R	M
ω_6	H	R	A

1. En considérant l'indice de Jaccard, calculer $s(\omega_1, \omega_2)$ et $s(\omega_3, \omega_6)$.
2. Est-ce que ω_1 est plus proche de ω_2 , que ω_3 de ω_6 ?

Solution : 1. Le TDC associé est

	<i>F</i>	<i>H</i>	<i>R</i>	<i>U</i>	<i>C</i>	<i>M</i>	<i>A</i>
ω_1	0	1	0	1	1	0	0
ω_2	1	0	0	1	1	0	0
ω_3	1	0	1	0	0	1	0
ω_4	1	0	0	1	0	0	1
ω_5	0	1	1	0	0	1	0
ω_6	0	1	1	0	0	0	1

On a $a_{1,2} = 2$, $b_{1,2} = 1$ et $c_{1,2} = 1$ (et $d_{1,2} = 3$). Donc

$$s(\omega_1, \omega_2) = \frac{a_{1,2}}{a_{1,2} + b_{1,2} + c_{1,2}} = \frac{2}{2 + 1 + 1} = 0,5.$$

On a $a_{3,6} = 1$, $b_{3,6} = 2$ et $c_{3,6} = 2$ (et $d_{3,6} = 2$). Donc

$$s(\omega_3, \omega_6) = \frac{a_{3,6}}{a_{3,6} + b_{3,6} + c_{3,6}} = \frac{1}{1 + 2 + 2} = 0,2.$$

2. Comme

$$s(\omega_1, \omega_2) > s(\omega_3, \omega_6),$$

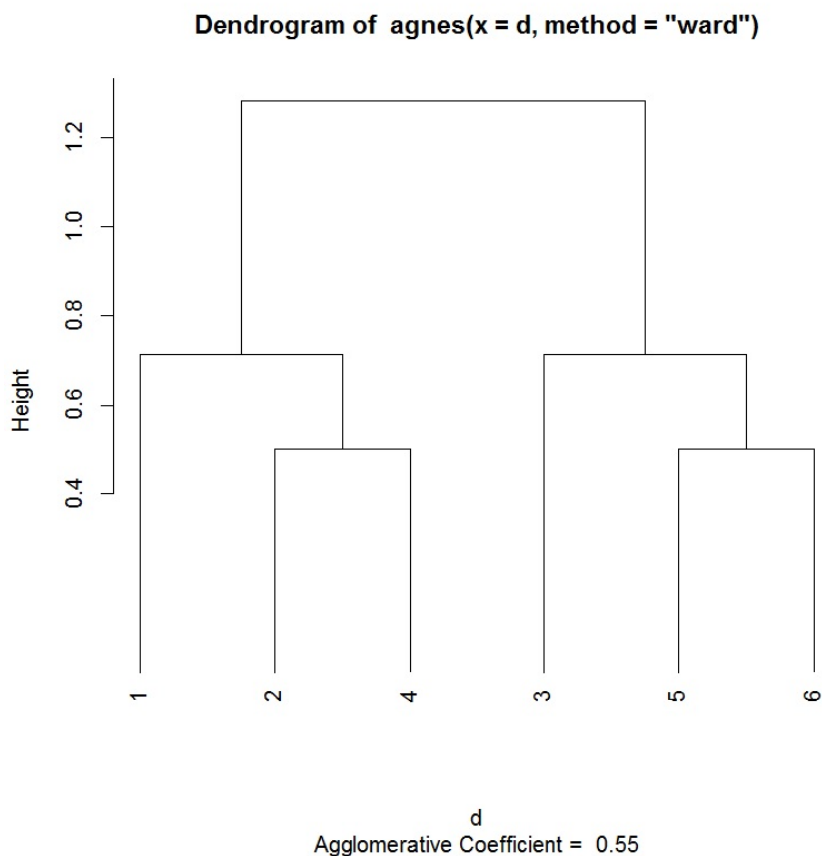
ω_1 est plus proche de ω_2 que ω_3 de ω_6 .

On peut aller plus loin en calculant les distances entre tous les individus et faire une CAH.

Commandes R de l'exemple : Les commandes R ci-dessous renvoient les dissimilarités entre tous les individus avec l'indice de Jaccard et l'algorithme de CAH est mis en œuvre.

```
x = c(0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0,
1, 0, 0, 0, 0, 1, 0, 1)
m = matrix(x, nrow = 6)
library(arules)
d = dissimilarity(m, method = "jaccard")
d
library(cluster)
ag = agnes(d, method = "ward")
cutree(ag, k = 2)
plot(ag, which = 2, hang = -1)
```

Cela renvoie le graphique :



"Distance" du Chi-deux : À partir du tableau disjonctif complet, on appelle "distance" du Chi-deux entre ω_u et ω_v la distance :

$$d(\omega_u, \omega_v) = \sqrt{\sum_{k=1}^r \frac{1}{\rho_k} (f_{u,k} - f_{v,k})^2},$$

où

$$f_{u,k} = \frac{n_{u,k}}{n_{u,\bullet}}, \quad n_{u,\bullet} = \sum_{k=1}^r n_{u,k}, \quad \rho_k = \frac{n_{\bullet,k}}{n}, \quad n_{\bullet,k} = \sum_{i=1}^n n_{i,k}.$$

On peut aussi utiliser cette "distance" pour mettre en œuvre l'algorithme de CAH.

Caractères de natures différentes : Si certains caractères sont qualitatifs et d'autres quantitatifs, on peut toujours transformer les caractères qualitatifs en quantitatifs en introduisant des classes de valeurs et en les considérant comme des modalités.

13 Enjeux de la classification supervisée

Contexte : On considère une population divisée en q groupes d'individus différents $\{G_1, \dots, G_q\}$.

Ces groupes sont distinguables suivant les valeurs de p caractères X_1, \dots, X_p , sans que l'on ait connaissance des valeurs de X_1, \dots, X_p les distinguant. Soit Y le caractère égal au groupe dans lequel appartient un individu extrait au hasard dans la population. On dispose de n individus avec, pour chacun d'entre eux, les valeurs de Y, X_1, \dots, X_p .

Les données sont donc de la forme :

	Y	X_1	\dots	X_p
ω_1	y_1	$x_{1,1}$	\dots	$x_{p,1}$
\vdots	\vdots	\vdots	\dots	\vdots
ω_n	y_n	$x_{1,n}$	\dots	$x_{p,n}$

où, pour tout $(i, j) \in \{1, \dots, n\} \times \{1, \dots, p\}$, $x_{j,i} = X_j(\omega_i)$ est l'observation du caractère X_j sur l'individu ω_i et $y_i = Y(\omega_i)$ est le groupe dans lequel appartient ω_i .

Objectif : On s'intéresse à un individu ω_* de la population avec ses valeurs $x = (x_1, \dots, x_p)$ de X_1, \dots, X_p , sans connaissance de son groupe d'appartenance. Partant des données, l'objectif est de déterminer à quel groupe l'individu ω_* a le plus chance d'appartenir. En terme mathématique, ce groupe inconnu est

$$G = \underset{g \in \{G_1, \dots, G_q\}}{\text{Argmax}} \mathbb{P}(\{Y = g\} / \{(X_1, \dots, X_p) = x\}).$$

Méthodes : Pour estimer G , plusieurs méthodes sont possibles. Parmi elles, il y a

- la méthode des k plus proches voisins (kNN pour K Nearest Neighbors),
- le modèle de mélange de densités,
- le modèle de régression logistique (pour $q = 2$),
- les arbres de décision,
- les réseaux de neurone,

- le Support Vector Machine,
- les forêts aléatoires.

Ce document aborde quelques aspects des trois premiers points.

14 Méthode des k plus proches voisins

Méthode des k plus proches voisins (kNN pour K Nearest Neighbors) : Soient d une distance et $k \in \{1, \dots, n\}$. En utilisant d , on considère l'ensemble \mathcal{U}_k des k individus de $\Gamma = \{\omega_1, \dots, \omega_n\}$ les plus proches de ω_* . Ainsi, pour tout $i \in \mathcal{U}_k$ et tout $j \in \Gamma - \mathcal{U}_k$, on a

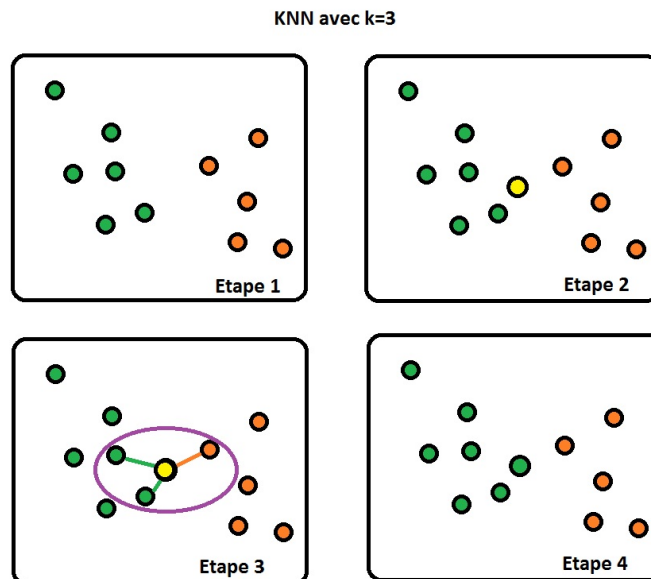
$$d(\omega_*, \omega_i) < d(\omega_*, \omega_j).$$

Par la méthode des k plus proches voisins, une estimation ponctuelle du groupe dans lequel ω_* a le plus de chances d'appartenir est

$$\hat{G} = \underset{g \in \{G_1, \dots, G_q\}}{\text{Argmax}} \sum_{i \in \mathcal{U}_k} \mathbb{I}_{\{y_i = g\}},$$

$$\text{où } \mathbb{I}_{\{y_i = g\}} = \begin{cases} 1 & \text{si } y_i = g \text{ (c'est-à-dire } \omega_i \text{ appartient au groupe } g), \\ 0 & \text{sinon.} \end{cases}$$

Illustration : Une illustration de la méthode des k plus proches voisins avec $k = 3$ est présentée ci-dessous :



Quelques commandes R : Un exemple simple de commandes R est décrit ci-dessous.

On introduits 3 individus $A1$, $A2$ et $A3$ qui vont former un groupe A :

```
A1 = c(0.1, 0.4)
A2 = c(0.8, 0.9)
A3 = c(3, 3.5)
```

On introduits 3 individus $B1$, $B2$ et $B3$ qui vont former un groupe B :

```
B1 = c(5.7, 6.1)
B2 = c(5.5, 6.8)
B3 = c(6.5, 4.9)
```

On considère la matrice de données correspondante et on spécifie l'appartenance des individus aux groupes A et B :

```
train = rbind(A1, A2, A3, B1, B2, B3)
cl = factor(c(rep("A", 3), rep("B", 3)))
```

On s'intéresse à un nouvel individu ω_* de caractéristiques $X_1 = 4.1$ et $X_2 = 3.8$ et on trace le nuage de points :

```
point = c(4.1, 3.8)
plot(rbind(train, point))
```

On évalue le groupe dans lequel ω_* a le plus de chance d'appartenir avec la méthode des k plus proche voisins avec $k = 1$:

```
library(class)
knn(train, point, cl, k = 1)
```

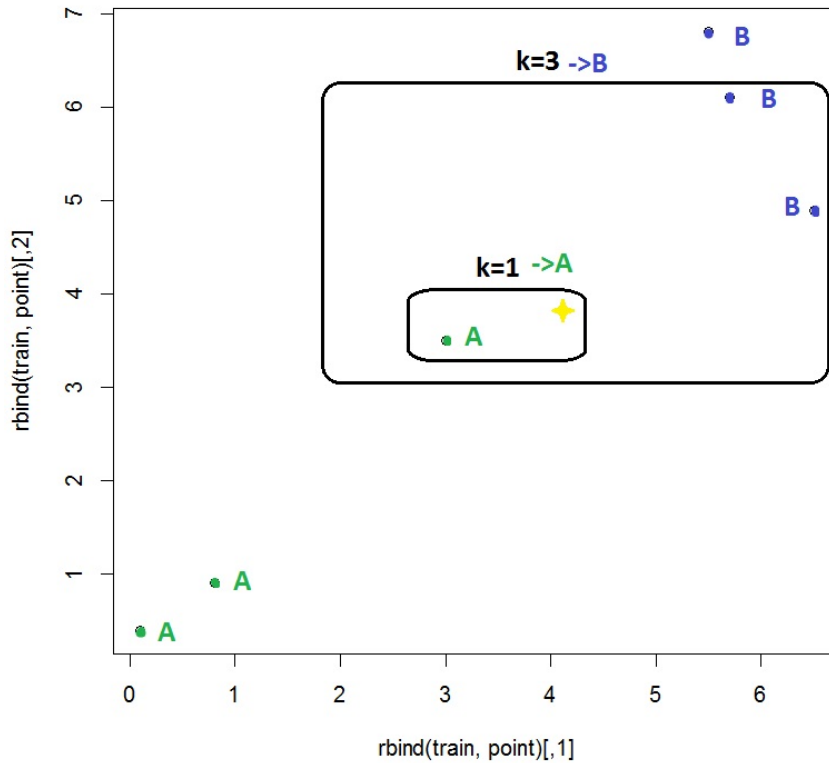
Cela renvoie A .

On fait la même chose avec $k = 3$:

```
knn(train, point, cl, k = 3)
```

Cela renvoie B .

On se rend compte de ce qui se passe par un graphique :



On considère un autre exemple portant sur des mesures du crane associées aux chiens et aux loups. L'entête du jeu de données "loups-g.txt" associé est :

	LCB	LMS	LPM	LP	LM	LAM	RACE
1	129	64	95	17.50	11.20	13.80	CHIEN
2	154	74	76	20.00	14.20	16.50	CHIEN
3	170	87	71	17.90	12.30	15.90	CHIEN
4	188	94	73	19.50	13.30	14.80	CHIEN
5	161	81	55	17.10	12.10	13.00	CHIEN
6	164	90	58	17.50	12.70	14.70	CHIEN

On propose les commandes R suivantes :

```
w = read.table("http://www.math.unicaen.fr/~chesneau/loups-g.txt", header = T)
attach(w)
w
cl = factor(w[,7])
point = c(210, 200, 76, 22, 12, 15)
library(class)
knn(w[,1:6], point, cl, k = 3)
```

Cela renvoie CHIEN. Ainsi, l'individu ω_* de caractéristiques $X_1 = 210$, $X_2 = 200$, $X_3 = 76$, $X_4 = 22$, $X_5 = 12$ et $X_6 = 15$ appartient au groupe des chiens.

kNN avec validation croisée : On peut aussi évaluer la qualité de l'algorithme des k plus proches voisins à l'aide d'une validation croisée. Cela consiste à extraire un petit groupe d'individus du jeu de données dont on connaît parfaitement leur groupe d'affectation et de faire l'algorithme des k plus proches voisins sur ceux-ci. On peut ainsi voir le nombre de fois où l'algorithme se trompe.

Un exemple avec le jeu de données iris est présenté ci-dessous :

```
data(iris)
u = iris[,-5]
library(class)
class = as.factor(iris[,5])
results = knn.cv(u, class, 1:length(class))
levels(results) = levels(class)
table(results, class)
```

Cela renvoie :

	setosa	versicolor	virginica
setosa	50	0	0
versicolor	0	47	3
virginica	0	3	47

On voit alors qu'il y a eut 6 iris qui ont été mal affectés.

Taux d'erreur de classification : Partant de l'algorithme des k plus proches voisins avec validation croisée, le taux moyen d'erreur de classification, noté t , est donné par le nombre d'individus mal affectés sur le nombre total d'individus.

Plus t est proche de 0, meilleur est la qualité prédictive du modèle.

On convient que la qualité de la classification est mauvaise lorsque $t > 0.5$.

Exemple : Sur l'exemple du jeu de données iris, ce taux est bon :

$$t = \frac{3 + 3}{50 + 47 + 47 + 3 + 3} = 0.04.$$

15 Modèle de mélange de densités

Hypothèses de gaussianité : On adopte le contexte de la classification supervisée. Le modèle de mélange de densités (gaussiennes) repose sur l'hypothèse que (X_1, \dots, X_p) est un vecteur aléatoire réel de densité :

$$f(x, \mu, \Sigma, r) = \sum_{k=1}^q r_k \phi(x, \mu_k, \Sigma_k), \quad x = (x_1, \dots, x_p) \in \mathbb{R}^p,$$

$r_k = \mathbb{P}(Y = G_k)$, $r = (r_1, \dots, r_q)$, $\mu = (\mu_1, \dots, \mu_p) \in \mathbb{R}^p$, $\Sigma = (\Sigma_1, \dots, \Sigma_p)$ est un vecteur de matrices de covariances de dimension $p \times p$ et $\phi(x, \mu_k, \Sigma_k)$ est la densité associée à la loi $\mathcal{N}_p(\mu_k, \Sigma_k)$:

$$\phi(x, \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{p/2} \sqrt{\det(\Sigma_k)}} \exp\left(-\frac{1}{2}(x - \mu_k)^t \Sigma_k^{-1} (x - \mu_k)\right), \quad x \in \mathbb{R}^p.$$

Les paramètres μ , Σ et p sont inconnus.

Enjeu : Partant des données, on souhaite estimer la probabilité inconnue qu'un individu ω_* vérifiant $(X_1, \dots, X_p) = x$ appartienne au groupe G_k :

$$p_{G_k}(x) = \mathbb{P}(\{Y = G_k\} / \{(X_1, \dots, X_p) = x\}).$$

Par la règle de Bayes et les hypothèses de départ, on peut exprimer cette probabilité comme

$$p_{G_k}(x) = \frac{r_k \phi(x, \mu_k, \Sigma_k)}{f(x, \mu, \Sigma, r)}.$$

Par conséquent, une estimation des paramètres μ , Σ et p à l'aide des données donne une estimation de $p_{G_k}(x)$ et, a fortiori, une estimation du groupe d'appartenance de ω_* .

Estimateurs du maximum de vraisemblance : Pour estimer les paramètres μ , Σ et r à partir des données, on utilise la méthode du maximum de vraisemblance qui donne les estimateurs $\hat{\mu}$, $\hat{\Sigma}$ et \hat{r} .

Dans le cas général, il n'y a pas d'expression analytique pour $\hat{\mu}$, $\hat{\Sigma}$ et \hat{r} ; ils peuvent être approchés avec l'algorithme de Newton-Raphson ou l'algorithme EM (Expectation-Maximization) (ou mieux : l'algorithme Classification EM ou l'algorithme Stochastique EM).

Estimation : Une estimation de $p_{G_k}(x)$ avec $x = (x_1, \dots, x_p)$ est

$$\hat{p}_{G_k}(x) = \frac{\hat{r}_k \phi(x, \hat{\mu}_k, \hat{\Sigma}_k)}{f(x, \hat{\mu}, \hat{\Sigma}, \hat{r})}.$$

Prédiction du groupe : On appelle prédiction du groupe d'un individu ω_* vérifiant $(X_1, \dots, X_p) = x$ la réalisation de

$$\hat{G} = \underset{g \in \{G_1, \dots, G_q\}}{\text{Argmax}} \hat{p}_g(x).$$

Quelques commandes R : Rappel du jeu de données iris : pour 3 variétés d'iris : Setosa, Versicolor, Virginica, et pour 150 iris par variété, on considère 4 caractères quantitatifs :

- X_1 la longueur en cm d'un pétale,
- X_2 la largeur en cm d'un pétale,
- X_3 la longueur en cm d'un sépale,
- X_4 la largeur en cm d'un sépale.

Ce sont les variables explicatives X_1 , X_2 , X_3 et X_4 . La variable à expliquer Y est une variable qualitative dont les modalités sont les espèces d'iris $\{\textit{setosa}, \textit{versicolor}, \textit{virginica}\}$.

Voici un problème de classification supervisée possible : on dispose d'un iris vérifiant :

$$X_1 = 2.1, \quad X_2 = 3, \quad X_3 = 2.3, \quad X_4 = 4.3.$$

À l'aide des mesures effectuées, à quelle variété a-t'il le plus de chances d'appartenir ?

Les commandes ci-dessous, dont `lda`, apportent une réponse en utilisant le modèle de mélange des densités :

```
data(iris)
library(MASS)
results = lda(Species ~ ., iris, prior = c(1, 1, 1) / 3)
library(MASS)
newiris = data.frame(Sepal.Length = 2.3, Sepal.Width = 4.3, Petal.Length = 2.1, Petal.Width = 3)
plda = predict(results, newiris)
plda
```

Cela renvoie `Virginica` avec une probabilité de 0.9980522. On est donc presque sûr que l'iris observé est de l'espèce `Virginica`.

Validation croisée : On peut aussi évaluer la qualité du modèle de mélange de densités à l'aide d'une validation croisée. Cela consiste à extraire un petit groupe d'individus du jeu de données dont on connaît parfaitement leur groupe d'affectation et les tester avec le modèle de mélange des densités. On peut ainsi voir le nombre de fois où l'algorithme se trompe.

Un exemple avec le jeu de données `iris` et la commande `lda` est présenté ci-dessous :

```
data(iris)
library(MASS)
results = lda(Species ~ ., iris, prior = c(1, 1, 1) / 3, CV = T)
table(iris$Species, results$class)
```

Cela renvoie :

	setosa	versicolor	virginica
setosa	50	0	0
versicolor	0	48	2
virginica	0	1	49

On voit alors qu'il y a eut 3 iris qui ont été mal affectés.

Taux d'erreur de classification : Partant du modèle de mélange de densités, le taux moyen d'erreur de classification, noté t , est donné par le nombre d'individus mal affectés sur le nombre total d'individus.

Plus t est proche de 0, meilleur est la qualité prédictive du modèle.

On convient que la qualité de la classification est mauvaise lorsque $t > 0.5$.

Exemple : Sur l'exemple du jeu de données iris, ce taux est :

$$t = \frac{1 + 2}{50 + 48 + 49 + 1 + 2} = 0.02.$$

Cela est très correct.

16 Régression logistique

Enjeu : On suppose que la population est divisée en $q = 2$ groupes : G_1 et G_2 et on adopte le contexte de la classification supervisée. Partant des données, on souhaite estimer la probabilité inconnue qu'un individu ω_* vérifiant $(X_1, \dots, X_p) = x$ appartienne au groupe G_1 :

$$p(x) = \mathbb{P}(\{Y = G_1\} | \{(X_1, \dots, X_p) = x\}).$$

Si $p(x) \geq 0.5$, alors ω_* a plus de chances d'appartenir à G_1 qu'à G_2 .

Transformation logit : On appelle transformation logit la fonction :

$$\text{logit}(y) = \log\left(\frac{y}{1-y}\right) \in \mathbb{R}, \quad y \in]0, 1[.$$

Son inverse est la fonction :

$$\text{logit}^{-1}(y) = \frac{\exp(y)}{1 + \exp(y)} \in]0, 1[, \quad y \in \mathbb{R}.$$

Régression logistique : On appelle régression logistique la modélisation :

$$\text{logit}(p(x)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p,$$

où β_0, \dots, β_p désigne $p + 1$ réels inconnus.

Ainsi, $p(x)$ et $\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ sont liés par la transformation logit ; on parle de lien logit.

On en déduit l'expression de $p(x)$:

$$p(x) = \text{logit}^{-1}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p) = \frac{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)}{1 + \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)},$$

Estimateurs du maximum de vraisemblance : Notre objectif est d'estimer les coefficients inconnus β_0, \dots, β_p à partir des données. Pour ce faire, on utilise la méthode du maximum de vraisemblance qui donne les estimateurs $\hat{\beta}_0, \dots, \hat{\beta}_p$.

Dans le cas général, il n'y a pas d'expression analytique pour $\hat{\beta}_0, \dots, \hat{\beta}_p$; ils peuvent être approchés avec l'algorithme de Newton-Raphson.

Estimation : Une estimation de $p(x)$ avec $x = (x_1, \dots, x_p)$ est

$$\hat{p}(x) = \text{logit}^{-1}(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p)}.$$

Prédiction du groupe : On appelle prédiction du groupe d'un individu ω_* vérifiant $(X_1, \dots, X_p) = x$ la réalisation de

$$\hat{G} = \begin{cases} G_1 & \text{si } \hat{p}(x) \geq 0.5, \\ G_2 & \text{sinon.} \end{cases}$$

Quelques commandes R : On considère le jeu de données "puits" dont voici l'entête :

Y	X1	X2
1	2.36	16.83
1	0.71	47.32
0	2.07	20.97
1	1.15	21.49
1	1.10	40.87
1	3.90	69.52

Un exemple de commandes R associées est donné ci-dessous :

```
w = read.table("http://www.math.unicaen.fr/~chesneau/puits.txt", header = T)
attach(w)
w
library(stats)
reg = glm(Y ~ X1 + X2, family = binomial)
pred.prob = predict.glm(reg, data.frame(X1 = 1, X2 = 60), type = "response")
pred.mod = factor(ifelse(pred.prob > 0.5, "G1", "G2"))
pred.mod
```

Cela renvoie $G2$.

Taux d'erreur de classification : Partant du modèle de régression logistique, le taux moyen d'erreur de classification noté t est donné par le nombre d'individus mal affectés sur le nombre total d'individus.

Plus t est proche de 0, meilleur est la qualité prédictive du modèle.

On convient que la qualité de la classification est mauvaise lorsque $t > 0.5$.

Quelques commandes R : Un exemple de commandes R associées est donné ci-dessous :

```
pred.prob = predict(reg, type = "response")
pred.mod = factor(ifelse(pred.prob > 0.5, "G1", "G2"))
mc = table(Y, pred.mod)
t = (mc[1, 2] + mc[2, 1]) / sum(mc)
t
```

Cela renvoie $t = 0.1364238$, ce qui est très faible.

Plus d'éléments seront donnés en Master 2. Voir, par exemple, les documents :

<http://www.math.unicaen.fr/~chesneau/Reg-M2.pdf>

<http://www.math.unicaen.fr/~chesneau/etudes-reg.pdf>

Index

- hclust, 23
- FactoMineR, 44, 47
- HCPC, 44
- agnes, 24, 65
- kmeans, 57
- lda, 76
- silhouette, 42

- ACP et CAH, 44

- CAH, 22
- CAH caractères qualitatifs, 61
- Classification non-supervisée, 5
- Classification supervisée, 6, 67
- Coefficient d'agglomération, 40
- Consolidation CAH, 60

- Dendrogramme, 24
- Distance du Chi-deux, 66
- Distance entre 2 individus, 16
- Distance euclidienne, 15
- Distances, 15
- Décomposition de Huygens, 33

- Ecart de Ward, 19
- Écarts, 17

- Indice de Dice, 62
- Indice de Jaccard, 62, 63

- Indice de silhouette, 41
- Inertie inter-classes, 33
- Inertie intra-classes, 33
- Inertie totale, 32

- k means, 48
- kNN, 69

- Largeur de silhouette, 42

- Matrice de données, 10
- Modèle de mélange de densités, 74
- Méthode de Forgy, 57
- Méthode de Hartigan-Wong, 58
- Méthode de la distance moyenne, 18
- Méthode de MacQueen, 58
- Méthode de Ward, 19, 32
- Méthode du plus proche voisin, 18
- Méthode du voisin le plus éloigné, 18

- PAM, 59
- Parangons, 46

- Ressemblance, 10
- Régression logistique, 78

- Tableau des écarts, 19
- Tableau disjonctif complet, 61
- Taux d'erreur de classification, 73, 77, 80
- Transformation Logit, 78

