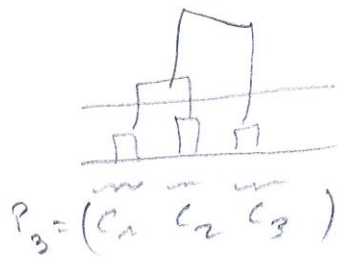


TP classification: méthodologie

Exercice 1

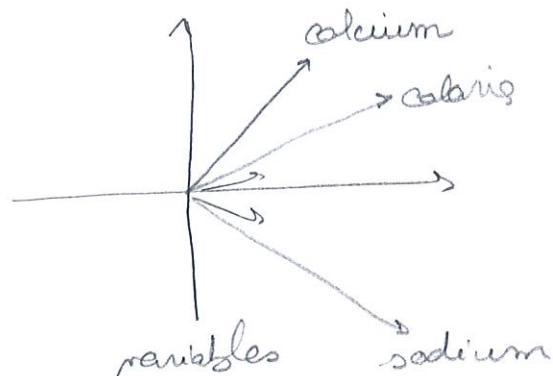
Matrice de données quantitatives $X: (29, 9)$

- 1) Calories, sodium et calcium ont de \oplus grandes variances $\Rightarrow \oplus$ grande variabilité
- 2) On effectue une classificat^o hiérarchique de Ward sur X et on se réfère la partition en trois classes



$$\Rightarrow X_n = \begin{pmatrix} 1 & \dots & 9 \\ \vdots & & \vdots \\ n & & \vdots \end{pmatrix} \begin{pmatrix} P_3 \\ \vdots \\ 1 \\ 2 \\ 3 \\ \vdots \end{pmatrix}$$

- 3) Interprétation de la partition en 3 classes via l'ACP non normalisée (i.e. ACP sur données centrées sur matrice de covariance)



Interprétation:

- seuls les variables de fortes variances contribuent aux axes (car ACP non normalisée).

- Les fromages projetés à droite (classes 1 et 2) ont des valeurs fortes en calories-calcium-sodium (pâtes durs et fromage gras), à gauche \oplus faible (fromage frais)

(2)
⇒ On ne peut interpréter les classes qu'en j° des variables de forte variance

Explication: Lorsqu'on calcule la distance entre deux fromages sur les données brutes, on donne implicitement plus d'importance aux variables de forte variance.

	Calories	Sodium	Calcium	lipides	retards	
Ex: Emmentaler	378	60	308,2	29,4	56,3
Fr. gras	80	41	146,3	3,5	50	
Ecart au carré	$(378-80)^2$	$(60-41)^2$	$(308,2-146)^2$	$(29,4-3,5)^2$	$(56,3-50)^2$
carre	88804	361	26211,4	670	39,69
distance ² =	88804	+	361	+	26211,4	+
	<div style="display: flex; justify-content: space-around; align-items: center;"> <div>↙</div> <div>↘</div> <div>⊕ "lourd"</div> </div>					

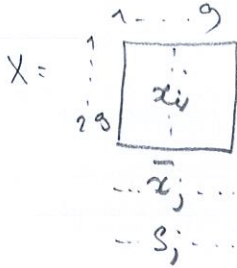
⇒ distance ne reflète que la ressemblance sur ces variables de forte variance.

⇒ On va regrouper les fromages dans des classes en fonction de ces variables uniquement

⇒ Contrer-réduire les données pour donner le même poids à toutes les variables dans le calcul des distances.

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$$

$$\text{ai) } S_j = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$



$$Z = \frac{x_i - \bar{x}}{s_i}$$

mean - 0
variance - 1

$1, \dots, j, \dots, P$
 x_i x_i'
 x_i x_i''

- $$d_M^2(x_i, x_{i'}) = \sum_{j=1}^p \frac{1}{s_j^2} (x_{ij} - x_{i'j})^2 \quad \text{avec } M = \begin{pmatrix} 1/s_1^2 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & 1/s_p^2 \end{pmatrix}$$
- \nwarrow pondère les écarts par la variance de j

$$= \sum_{j=1}^n \left(\frac{x_{ij} - \bar{x}_j - (x'_{ij} - \bar{x}_j)}{s_j} \right)^2$$

$$= \sum_{j=1}^n (z_j - z_{i_j})^2 = d^2(z_i, z_{i'})$$

distance euclidean simple.

sur les données contraires réduites * Exemple 14

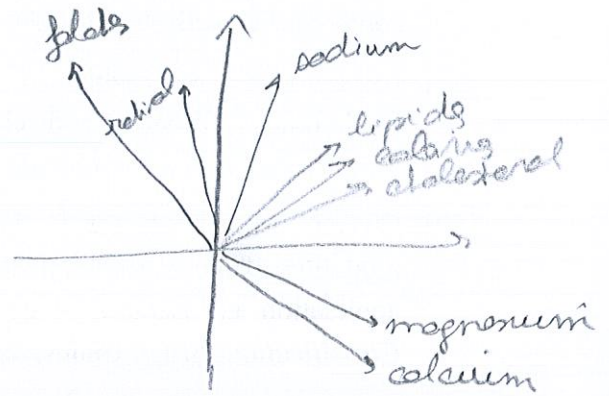
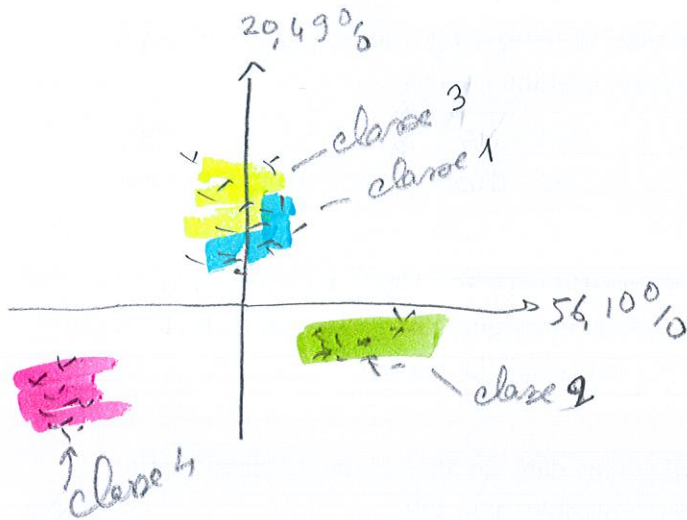
-

$$P_L = (\overset{\sim}{C_1} \quad \overset{\sim}{C_2} \quad \overset{\sim}{C_3} \quad \overset{\sim}{C_4})$$

$$\Rightarrow X_2 = \begin{bmatrix} x \end{bmatrix} \begin{pmatrix} p_1 \\ 1 \\ 4 \\ 4 \end{pmatrix}$$

3) Interpretation de la partition en 4 classes via

l'ACP normée (i.e. ACP aux données centrées-réduites sur matrice de corrélations)



* 2 dim^o → 76,59% de l'inertie expliquée + chute de la valeur propre.

* Toutes les variables jouent un rôle

→ classe 4: peu gras et calcaire (fromages frais)

classe 2: riches en magnésium et calcium (pâtes dures)

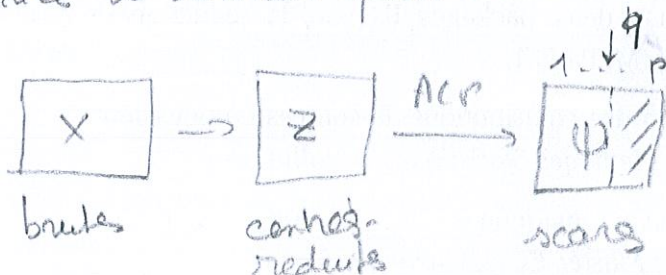
classe 3: riches en plats

→ Aller voir sur classe 3 ?

classe 1: ④ salés

		calcaire	sodium	calcium	lipides	retrait...
④	Emmental	378	60	308	29	56
	Fromage frais	80	41	146	3,5	50
Exercice 3	carte au canard	$(378-80)^2$	$(60-41)^2$	$(308-146)^2$	$(29-3,5)^2$	
	reavance	8448	1181	3260	66	

Matrice de données quantitatives : Ψ ← Scores de l'ACP de X.



1) On effectue une classification hiérarchique de Ward sur Ψ (en conservant toute les composantes principales)

On effectue une classification hiérarchique de Ward sur z

⇒ Même dendrogramme

⇒ Même indices

⇒ Même hiérarchie

2) On effectue une classification hiérarchique de Ward sur les q premières colonnes de φ (en conservant les q premières composantes principales)

→ Choix de q : 2 ou 3

→ On retrouve la même partition à 1 individu près

→ Stratégie de SPAD: type faire une ACP et "monter" la hiérarchie sur les q premières composantes.

(+) de SPAD: } Outil d'interprétation des classes

```

#-----TP Classification : méthodologie

#Attention : retirer les accents dans les noms de colonne !!
X<-read.table(file.choose(),sep=" ",header=TRUE,row.names=1)
dim(X)

#-----Exercice 1-----

#Variance des variables
apply(X,2,var)

#Classification hiérarchique de Ward sur donnees brutes
d <- dist(X)
tree <- hclust(d^2,method="ward")
plot(tree)

P3 <- as.factor(cutree(tree,3))
levels(P3) <- paste("Classe",1:3)
X1 <- data.frame(X,P3=P3)

#Interprétation de la partition
require(FactoMineR)

#ACP non normee
pca1 <- PCA(X1,scale=FALSE,quali.sup=10,graph=FALSE)
pca1$eig

# Nuage des individus et des variables dans le premier plan factoriel
par(mfrow=c(1,2))
plot.PCA(pca1,axes=c(1,2),choix="ind",habillage=10,invisible="quali")
plot.PCA(pca1,axes=c(1,2),choix="var")

#distance entre Emmental et Fromage frais
X[c(13,16),]
(X[13,]-X[16,])^2

#-----Exercice 2-----

#donnees centrees-reduire
n <- nrow(X)
Z<-scale(X,center=TRUE,scale=TRUE)*sqrt((n)/(n-1))

apply(Z,2,mean)
apply(Z,2,sd)*sqrt((n-1)/n)

# Classification sur donnees centrees-reduites:

d <- dist(Z)
tree <- hclust(d^2,method="ward")
plot(tree)

P4 <- as.factor(cutree(tree,4))
levels(P4)<-paste("Classe",1:4)
X2<-data.frame(X,P4=P4)

#ACP normee
pca2 <- PCA(X2,scale=TRUE,quali.sup=10,graph=FALSE)
pca2$eig

#distance normalisee par l'inverse des variances entre Emmental et Fromage frais
X[c(13,16),]
apply(X,2,var)
(X[13,]-X[16,])^2

par(mfrow=c(1,2))
plot.PCA(pca2,axes=c(1,2),choix="ind",habillage=10,invisible="quali")
plot.PCA(pca2,axes=c(1,2),choix="var")

plot.PCA(pca2,axes=c(1,3),choix="ind",habillage=10,invisible="quali")
plot.PCA(pca2,axes=c(1,3),choix="var")

#-----Exercice 3-----

#ACP normee en conservant TOUTES les composantes principales
pca2 <- PCA(X,ncp=9,graph=FALSE)
Psi <- pca2$ind$scoord #matrice des scores de l'ACP normee

#Ward sur donnees centrees-reduites
d <- dist(Z)
tree <- hclust(d,method="ward")

```