

TP Analyse de données - Apprentissage non supervisé

INTRODUCTION

Nous allons utiliser le logiciel R ([documentation](#) ; possibilité d'utiliser Rstudio). Dans un dossier *AnalyseDeDonnees* créer deux sous dossiers :

- Code - dans lequel vous placerez vos fichiers de code
- Data - dans lequel vous placerez les fichiers de données (à télécharger via Moodle)

Les 3 TP *Analyse de données* présentent différentes méthodes d'analyse de données : le but n'est pas de finir les TP le plus vite possible mais d'analyser les résultats ! Un rapport de 1 page **maximum** vous est demandé pour chacun des 3 TP : ne choisir que les résultats les plus intéressants et les **commenter**.

Ne pas hésiter à utiliser l'aide de R grâce à la commande :

```
1 help(...)
```

ANALYSE EN COMPOSANTES PRINCIPALES

0. Télécharger le cours sur l'Analyse en Composantes Principales.

1. Créer un fichier *ACP.R* dans le dossier *Code*. Charger les packages d'intérêt en ajoutant dans le fichier :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../AnalyseDonnees/TP/TP/Code")
3 # Packages utilisés dans la suite
4 library(GGally)
5 library("FactoMineR")
```

2. Charger les données et les manipuler :

```
1 # Chargement des données
2 load("../Data/eaux.rda")
3 # Création de trois objets
4 X <- data
5 n <- nrow(X) # nombre d'observations
6 p <- ncol(X) # nombre de variables
7 # Manipulation des données
8 X[, 1:3] # les 3 premières colonnes
9 X[, c(1,5)] # la première et la 5ème colonne
10 X[, -c(1,5)] # toutes les colonnes sauf la première et la 5ème
11 X[1:3, ] # les 3 premières lignes
12 class(X) # Un data.frame est une liste dont les éléments sont les colonnes
13 X$saveur.amère # On récupère les éléments d'une liste avec le $
14 X[,1]
```

3. Calculer la moyenne et l'écart type de chacune des variables :

```
1 # Moyenne de la variable saveur amère
2 sum(X$saveur.amère)/n
3 mean(X$saveur.amère)
4 # Ecart-type de la variable saveur amère
5 x <- X$saveur.amère
6 sqrt(1/(n-1)*sum((x-mean(x))^2))
7 sd(x)
8 # Moyenne et écart-type de toutes les variables
9 moy <- apply(X, 2, mean) # 1=lignes, 2=colonnes
10 sigma <- apply(X,2, sd) #standard deviation
```

4. Centrer-réduire (standardiser) les variables :

```
1 # Création des données centrées ...
2 Y <- sweep(X, 2, moy, "-")
3 apply(Y, 2, mean) # les colonnes sont bien centrées
4 # ... et réduites
5 Z <- sweep(Y, 2, sigma, "/")
6 apply(Z, 2, sd) # les colonnes sont bien de variance 1
7 # ou de manière équivalente
8 Z <- scale(X)
9 # ou avec l'écart-type non corrigé (comme en ACP)
10 Z <- scale(X)*sqrt(n/(n-1))
```

5. Description bivariée des 5 premières variables :

```
1 # Avec la fonction ggpairs du package GGally
2 ggpairs(X[,1:5])
3 # Calculer et visualiser la matrice de corrélation
4 z1 <- Z[, 1] # variable saveur amère standardisée
5 z2 <- Z[, 2] # variable saveur sucrée standardisée
6 sum(z1*z2)/n # corrélation entre les deux variables
7 cor(X$saveur.amère, X$saveur.sucrée)
8 cor(X[,1:5]) # matrice des corrélations entre les 5 premières variables
9 ggcorr(X[,1:5])
```

6. Calculer la matrice des distances entre les individus et les corrélations entre les variables :

```
1 # Matrice des distances entre les individus
2 dist(X) # données brutes
3 dist(Y) # données centrées
4 dist(Z) # données centrées-réduites
5 # Corrélation entre les variables
6 cor(X)
7 # ou encore
8 t(Z) %*% Z/n # %*% est le produit matriciel
```

7. Faire l'Analyse en Composantes Principales non normalisée (sur matrice de covariances) :

```
1 # Fonction PCA du package FactoMineR
2 # (scale.unit=FALSE)
3 res <- PCA(X, graph = FALSE, scale.unit = FALSE)
4 # Figure individus
5 plot(res, choix = "ind", cex = 1.5, title = "")
6 # Figure variables
7 plot(res, choix = "var", cex = 1.5, title = "")
8 # Equivalent à la décomposition en valeurs propres de la matrice des covariances
9 Y <- as.matrix(Y)
10 C <- (t(Y) %*% Y)/n # matrice des covariances
11 eigen(C)$values
12 res$eig[, 1]
```

8. Faire l'Analyse en Composantes Principales normalisée (sur matrice des corrélations) :

```
1 # Analyse en composantes principales normalisée (sur matrice des corrélations)
2 # (par défaut: scale.unit=TRUE)
3 res <- PCA(X, graph=FALSE)
4 # Figure individus
5 plot(res, choix = "ind", cex = 1.5, title = "") # plan 1-2
6 plot(res, choix = "ind", axes=c(2,3), cex = 1.5, title = "") # plan 2-3
7 # Figure variables
8 plot(res, choix = "var", cex = 1.5, title = "") # plan 1-2
9 plot(res, choix = "var", axes=c(2,3), cex = 1.5, title = "") # plan 2-3
10 # Equivalent à la décomposition en valeurs propres de la matrice des corrélations
11 R <- (t(Z) %*% Z)/n # matrice des corrélations
12 eigen(R)$values
13 res$eig[, 1]
14 # Récupérer les 2 premières composantes principales
15 F <- res$ind$coord[, 1:2]
16 plot(F, pch = 16)
17 text(F, rownames(X), pos = 3) # on retrouve la figure des individus
18 # Récupérer les loadings (corrélations aux deux premières CP)
19 A <- res$var$coord[, 1:2]
20 plot(A, pch=16)
21 text(A, colnames(X), pos = 3) # on retrouve la figure des variables
22 A[1, , drop=FALSE] # corrélations entre saveur amère et les 2 premières CP
23 cor(F, X$saveur.amère)
24 # Interprétation du premier plan des individus en fonction des variables ?
```

9. Interprétation des résultats : combien de composantes peut-on retenir ? Utiliser les règles de Kaiser et du coude.

```
1 # Inertie (variance) des composantes principales
2 apply(F, 2, var)*(n-1)/n # variances des 2 premières CP
3 res$eig[, 1]
4 sum(res$eig[, 1])
5 res$eig
```

10. Interprétation des résultats : qualité de la projection des individus

```
1 # Qualité de la projection des individus sur les axes
2 res$ind$cos2
3 # Qualité de la projection des individus sur le premier plan
4 apply(res$ind$cos2, 1, sum)
5 # Interprétation du premier plan factoriel des individus ?
```

11. Interprétation des résultats : qualité de la projection des variables

```
1 # Qualité de la projection des variables sur les axes
2 res$var$cos2
3 # Qualité de la projection des variables sur le premier plan
4 apply(res$var$cos2, 1, sum) # ou regarder la longueur des flèches !
5 # Interprétation du premier plan factoriel des variables ?
```

En plus - Quand vous faites de l'ACP il y a deux erreurs à éviter :

- Attention aux données très asymétriques : par exemple beaucoup de très petites valeurs et quelques très grandes (dans ce cas là une transformation des données peut être utile ...).
- Attention à l'effet *taille* (toutes les variables ont des contributions positives sur un axe) quand les données sont corrélées entre elles.

CLUSTERING

0. Télécharger le cours sur le clustering.

1. Créer un fichier *clustering.R* dans le dossier *Code*. Charger les packages d'intérêt en ajoutant dans le fichier :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../AnalyseDonnees/TP/TP/Code")
3 # Packages utilisés dans la suite
4 library("FactoMineR")
```

Les kmeans

2. Charger les données et les afficher :

```
1 # Données sur les fromages
2 X <- read.table("../Data/fromage.txt", sep="", header=TRUE, row.names=1)
3 n <- nrow(X)
4 p <- ncol(X)
```

3. Calculer la moyenne et l'écart type de chacune des variables :

```
1 # Calcul de la moyenne et de l'écart type des variables
2 moy <- apply(X, 2, mean)
3 sigma <- apply(X, 2, sd)*sqrt((n-1)/n) #écart-type non corrigé
```

4. Standardiser les données :

```
1 # Création des données centrées ...
2 Y <- sweep(X, 2, moy, "-")
3 # ... et réduites
4 Z <- sweep(Y, 2, sigma, "/")
```

5. Fixer le nombre de clusters souhaité :

```
1 # Nombre de clusters souhaité
2 K <- 5
```

6. Appliquer l'algorithme des kmeans sur les données standardisées :

```
1 # kmeans en 5 classes sur données standardisées
2 km <- kmeans(Z, centers=K, nstart=50)
3 km$cluster # la partition
4 table(km$cluster) # effectifs des classes
5 km$withinss # inertie de chaque classe (within sum of squares)
6 km$betweenss # inertie inter-classe (between sum of squares)
7 km$totss # inertie totale (total sum of squares)
8 sum(km$withinss) + km$betweenss # inertie total= intra + inter
9 km$betweenss/km$totss # pourcentage d'inertie expliquée par la partition
10 # Vérifier que ce pourcentage augmente quand K augmente
```

7. Cette partition en 5 classes est une nouvelle variable qualitative :

```
1 # La partition est une nouvelle variable qualitative
2 part <- as.factor(km$cluster) # object de classe factor
3 levels(part) <- paste("cluster", 1:K, sep="") # modalités de la variable
4 # Un data.frame mélange les colonnes numériques et qualitatives !
5 Xplus <- data.frame(X, part)
6 View(Xplus)
```

8. Interpréter cette partition en visualisant les classes sur une ACP normée et avec la fonction `catdes` du package `FactoMineR` :

```
1 # ACP normée en mettant la variable "part" en illustrative
2 res <- PCA(Xplus, quali.sup="part", graph = FALSE)
3 res$eig
4 plot(res, choix = "ind", habillage = "part", invisible = "quali")
5 plot.PCA(res, choix = "var")
6 # Interprétation des classes avec la fonction catdes
7 des <- catdes(Xplus, num.var = 10)
8 des$quanti$cluster1
9 des$quanti$cluster2
10 # etc...
```

La CAH de Ward

9. Appliquer la classification ascendante hiérarchique de Ward sur les données standardisées :

```
1 #Classification hiérarchique de Ward sur données centrées-réduites
2 D <- dist(Z)
3 tree <- hclust(D^2/(2*n), method = "ward.D")
4 plot(tree, hang=-1, main="Dendrogramme de Ward", xlab="", sub = "")
5 tree$height
6 sum(tree$height) # Inertie totale
```

10. Partition en 5 classes de Ward et comparaison avec la partition en 5 classes des kmeans :

```
1 # Partition en 5 classes avec la CAH de Ward
2 K <- 5
3 plot(tree, hang=-1, main="", sub="", xlab="")
4 rect.hclust(tree, k=K)
5 part_ward <- cutree(tree, k=K) # partition en 5 classes
6 # Comparaison avec la partition en 5 classes des kmeans
7 table(part_ward) # effectif des classes de la partition de Ward
8 table(km$cluster) # effectifs des classes de la partition des kmeans
9 table(part_ward, km$cluster) # tableau de contingence entre les deux partitions
10 W <- sum(tree$height[1:(n-K)]) # inertie intra-classe de la partition de Ward
11 (1-W/9)*100 # Pourcentage d'inertie expliquée à comparer avec les kmeans
```

11. Combiner Ward et kmeans :

```
1 # 1. Calculer les centres des 5 classes de Ward
2 sp <- split(Z, part_ward)
3 centres_ward <- t(sapply(sp, colMeans))
4 # 2. Appliquer les kmeans en partant des centres des classes de Ward
5 km <- kmeans(Z, centers = centres_ward) # pas besoin de nstart !
6 table(part_ward, km$cluster) # deux fromages ont changés de classe
7 km$betweenss/km$totss # Le pourcentage d'inertie expliquée est un peu meilleur
```

Clustering de données qualitatives

12. Charger les données sur les races de chiens. Les recoder via l'ACM.

```
1 load("../Data/chiens.rda")
2 res <- MCA(chiens, graph = FALSE) # Multiple Correspondance Analysis
3 res$eig # 10 valeurs propres non nulles
4 res <- MCA(chiens, ncp=10, graph = FALSE) # 10 composantes principales
5 F <- res$ind$coord # données recodées
```

13. Appliquer la CAH de Ward aux données recodées.

```
1 tree <- hclust(dist(F)^2/(2*nrow(F)), method="ward.D")
2 plot(tree, hang=-1)
```

14. Couper le dendrogramme en conséquence et interpréter la partition avec la fonction `catdes`.

```
1 part <- cutree(tree, k=4) # partition en 4 classes
2 part <- as.factor(part)
3 levels(part) <- paste("cluster", 1:4, sep="")
4 des <- catdes(data.frame(part, chiens), num.var=1)
5 des$category
```