

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 du répertoire "Data" (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("FactoMineR")
5 library(PCAmixdata)
```

2. Charger les données et les afficher :

```
1 # Chargement des données
2 load("../Data/eaux.RData")
3 # Affichage des données
4 print(data,digits=4)
```

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 mean <- apply(data,2,mean)
3 std <- apply(data,2,sd) #standard deviation
4 stat <- rbind(mean,std)
5 # Affichage
6 print(stat,digits=4)
```

4. Afficher les données centrées-réduites :

```
1 # Création des données centrées ...
2 datanorm <- sweep(data,2,mean,"-")
3 # ... et réduites
4 datanorm <- sweep(datanorm,2,std,"/")
5 # Affichage des données centrées - réduites
6 print(datanorm,digits=4)
```

5. Visualiser la description bivariée des 5 premières variables :

```
1 # Visualisation des données en description bivariée
2 pairs(data[,1:5])
3 # Afficher la matrice de corrélation
4 ggcorr(data[,1:5])
5 # Aller encore plus loin avec ggpairs
6 ggpairs(data[,1:5])
```

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

```
1 # Matrice des distances entre les individus
2 dist(data)
3 # Corrélation entre les variables
4 cor(data[,1:5])
```

7. Faire l'Analyse en Composantes Principales sur les données d'origine :

```
1 # Analyse en composantes principales sur les données d'origine
2 # (scale.unit=FALSE)
3 res <- PCA(data,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. Faire l'Analyse en Composantes Principales sur les données centrées-réduites :

```
1 # Analyse en composantes principales sur les données centrées-réduites
2 # (par défaut: scale.unit=TRUE)
3 resnorm <- PCA(data,graph=FALSE)
4 # Figure individus
5 plot(resnorm,choix="ind",cex=1.5,title="")
6 # Figure variables
7 plot(resnorm,choix="var",cex=1.5,title="")
```

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 resnorm$eig
3 barplot(resnorm$eig[,1])
```

10. Interprétation des résultats : qualité de la projection des individus (angle entre l'individu et les composantes principales)

```
1 # Projection des individus
2 resnorm$ind$cos2
3 # Somme avec les 2 premières
4 resnorm$ind$cos2[,1]+resnorm$ind$cos2[,2]
5 # Et les 3 premières ?
6 resnorm$ind$cos2[,1]+resnorm$ind$cos2[,2]+resnorm$ind$cos2[,3]
```

11. Interprétation des résultats : contribution des individus (à regarder en même temps que le graphe de projection)

```
1 # Contribution des individus
2 resnorm$ind$contrib
```

12. Interprétation des résultats : qualité de la projection et contribution des variables

```
1 # Projection des variables
2 resnorm$var$cos2
3 # Somme avec les 2 premières
4 resnorm$var$cos2[,1]+resnorm$var$cos2[,2]
5 # Et les 3 ?
6 resnorm$var$cos2[,1]+resnorm$var$cos2[,2]+resnorm$var$cos2[,3]
7 # Contribution des variables
8 resnorm$var$contrib
```

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.

PARTITIONNEMENT

0. Télécharger le cours sur le partitionnement.

Partitionnement kmeans

1. Créer un fichier *partitionnement.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")
```

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 print(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 mean <- apply(X,2,mean)
3 std <- apply(X,2,sd) #standard deviation
4 stat <- rbind(mean,std)
5 # Affichage
6 print(stat,digits=4)
```

4. Créer et afficher les données centrées-réduites :

```
1 # Création des données centrées ...
2 Xnorm <- sweep(X,2,mean,"-")
3 # ... et réduites
4 Xnorm <- sweep(Xnorm,2,std,"/")
5 # Affichage des données centrées - réduites
6 print(Xnorm,digits=4)
```

5. Fixer le nombre de clusters souhaité (faire varier ce paramètre à la fin de l'exercice!)

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

6. Appliquer l'algorithme des kmeans sur les données brutes ... puis sur les données centrées réduites. étudier les résultats obtenus :

```
1 ## KMEANS
2 # Algorithme des kmeans (avec affichage)
3 km <- kmeans(X,numcluster,nstart=50)
4 print(km)
5 # Algorithme des kmeans sur données centrées-réduites (avec affichage)
6 kmnorm <- kmeans(Xnorm,numcluster,nstart=50)
7 print(kmnorm)
```

7. Comme il y a plus de 2 variables dans cet exemple nous ne pouvons pas visualiser facilement les résultats ... L'ACP peut nous y aider : on va afficher les clusters sur les premières directions principales. Tout d'abord concaténer aux données le résultat obtenu par l'algorithme des kmeans pour les 2 cas considérées :

```
1 # Concatenation des données avec leur résultat de cluster
2 cluster <- as.factor(km$cluster)
3 clusternorm <- as.factor(kmnorm$cluster)
4 XplusCluster <- data.frame(X,cluster=cluster)
5 XnormplusCluster <- data.frame(Xnorm,cluster=clusternorm)
6 colclust <- length(X)+1
7 print(XplusCluster)
8 print(XnormplusCluster)
```

8. Mettre en place l'ACP sur les données brutes et afficher les résultats. Les couleurs correspondent aux différents clusters. Interpréter les résultats.

```
1 # ACP sur les données brutes
2 rPCA <- PCA(XplusCluster,scale.unit=FALSE,graph=FALSE,quali.sup=colclust)
3
4 # Nuage des individus et des variables dans le premier plan factoriel
5 par(mfrow=c(1,2))
6 plot.PCA(rPCA,axes=c(1,2),choix="ind",habillage=colclust,invisible="quali")
7 plot.PCA(rPCA,axes=c(1,2),choix="var")
8 # Nuage des individus et des variables dans le deuxième plan factoriel
9 par(mfrow=c(1,2))
10 plot.PCA(rPCA,axes=c(1,3),choix="ind",habillage=colclust,invisible="quali")
11 plot.PCA(rPCA,axes=c(1,3),choix="var")
```

9. Mettre en place l'ACP sur les données centrées-réduites et afficher les résultats. Les couleurs correspondent aux différents clusters. Interpréter les résultats.

```
1 # ACP sur les données centrées-réduites
2 rPCAnorm <- PCA(XnormplusCluster,graph=FALSE,quali.sup=colclust)
3
4 # Nuage des individus et des variables dans le premier plan factoriel
5 par(mfrow=c(1,2))
6 plot.PCA(rPCAnorm,,axes=c(1,2),choix="ind",habillage=colclust,invisible="quali")
7 plot.PCA(rPCAnorm,axes=c(1,2),choix="var")
8 # Nuage des individus et des variables dans le deuxième plan factoriel
9 par(mfrow=c(1,2))
10 plot.PCA(rPCAnorm,axes=c(1,3),choix="ind",habillage=colclust,invisible="quali")
11 plot.PCA(rPCAnorm,axes=c(1,3),choix="var")
```

Partitionnement hiérarchique (distance de Ward)

10. Mettre en place la classification hiérarchique de Ward sur les données brutes ... :

```
1 #Classification hiérarchique de Ward sur données brutes
2 d <- dist(X)
3 tree <- hclust(d^2,method="ward.D2")
4 par(mfrow=c(1,1))
5 plot(tree)
```

11. ... ainsi que sur les données centrées-réduites :

```
1 #Classification hiérarchique de Ward sur données centrées-réduites
2 dnorm <- dist(Xnorm)
3 treenorm <- hclust(dnorm^2,method="ward.D2")
4 plot(treenorm)
```

12. Concaténer les données avec leurs résultats de cluster en vue de l'ACP :

```
1 # Concatenation des données avec leur résultat de cluster
2 clusterW <- as.factor(cutree(tree,numcluster))
3 XplusClusterW <- data.frame(X,cluster=clusterW)
4 print(XplusClusterW)
5 clusternormW <- as.factor(cutree(treenorm,numcluster))
6 XnormplusClustW <- data.frame(Xnorm,cluster=clusternormW)
7 print(XnormplusClustW)
```

13. Mise en place de l'ACP sur les données brutes :

```
1 # ACP sur les données brutes
2 rPCAW <- PCA(XplusClusterW,scale.unit=FALSE,graph=FALSE,quali.sup=colclust)
3 # Nuage des individus et des variables dans le premier plan factoriel
4 par(mfrow=c(1,2))
5 plot.PCA(rPCAW,axes=c(1,2),choix="ind",habillage=colclust,invisible="quali")
6 plot.PCA(rPCAW,axes=c(1,2),choix="var")
7 # Nuage des individus et des variables dans le deuxième plan factoriel
8 par(mfrow=c(1,2))
9 plot.PCA(rPCAW,axes=c(1,3),choix="ind",habillage=colclust,invisible="quali")
10 plot.PCA(rPCAW,axes=c(1,3),choix="var")
```

14. Mise en place de l'ACP sur les données centrées-réduites :

```
1 # ACP sur les données centrées-réduites
2 rPCAnormW <- PCA(XnormplusClustW,scale.unit=FALSE,graph=FALSE,quali.sup=colclust)
3 # Nuage des individus et des variables dans le premier plan factoriel
4 par(mfrow=c(1,2))
5 plot.PCA(rPCAnormW,axes=c(1,2),choix="ind",habillage=colclust,invisible="quali")
6 plot.PCA(rPCAnormW,axes=c(1,2),choix="var")
7 # Nuage des individus et des variables dans le deuxième plan factoriel
8 par(mfrow=c(1,2))
9 plot.PCA(rPCAnormW,axes=c(1,3),choix="ind",habillage=colclust,invisible="quali")
10 plot.PCA(rPCAnormW,axes=c(1,3),choix="var")
```