

Rédigé par : Astrid Jourdan

Ref :

A l'intention de : Elèves d'ING1-GI

Créé le : 16/09/2016

Exercice 1

Dans \mathbb{R} , on considère 5 points $x_1=1$, $x_2=2$, $x_3=9$, $x_4=12$ et $x_5=20$.

1) Appliquer l'algorithme des k-means avec les valeurs de k et les points de départ suivants. Calculer le pourcentage d'inertie expliquée par la partition obtenue.

- a) $k=2$, $g_1=1$ et $g_2=20$
- b) $k=2$, $g_1=2$ et $g_2=9$
- c) $k=3$, $g_1=1$, $g_2=9$ et $g_3=12$
- d) Peut-on utiliser l'inertie inter-classes pour déterminer le meilleur regroupement des deux?

2) Appliquer une méthode de classification hiérarchique ascendante en utilisant la distance minimale comme critère de dissimilarité entre classes. Tracer le dendrogramme (avec en ordonnée la distance minimale). Quel regroupement vous paraît correct ?

Exercice 2

L'objectif de cet exercice est de tester les algorithmes k-means et CAH sur des jeux de données simulés,

Test_Clusters_Distincts.txt

Test_Clusters_Random.txt

Test_Clusters_Melanges.txt

Test_Clusters_Atypiques.txt

Pour cela, on utilisera le langage R avec ses fonctions `kmeans` et `hclust`.

- 1) Algorithme des Kmeans
 - a) Tester l'algorithme des kmeans sur les données *Test_Clusters_Distincts.txt*. Essayer plusieurs nombres de classes et choisir le meilleur.
 - b) Tester l'algorithme des kmeans sur les données *Test_Clusters_Distincts.txt*, *Test_Clusters_Melanges.txt* et *Test_Clusters_Random.txt*. Constater l'évolution de l'inertie expliquée.
 - c) Tester l'algorithme des kmeans sur les données *Test_Clusters_Corr.txt*. Que pourrait-on faire pour améliorer le résultat.
 - d) Tester l'algorithme des kmeans sur les données *Test_Clusters_Atypique.txt* avec les individus n°1 et n°1499 pour initialisation
- 2) Classification ascendante hiérarchique
 - a) Tester l'algorithme CAH sur les données *Test_Clusters_Distincts.txt*, *Test_Clusters_Melanges.txt* et *Test_Clusters_Random.txt*. Constater l'évolution du dendrogramme.
 - b) Tester l'algorithme CAH sur les données *Test_Clusters_Atypique.txt* avec la méthode « ward.D2 » et la méthode « average ».
 - c) Comparer les résultats obtenus entre CAH et Kmeans sur les données *Test_Clusters_Distincts.txt*, *Test_Clusters_Melanges.txt* et *Test_Clusters_Random.txt*.

Exercice 3

Déterminer des clusters dans le jeu de données « iris ». Est-ce que les clusters correspondent aux trois types de fleurs ?

Fonction kmeans

`kmeans(x,centers,nstart,...)`

Entrées :

- `x` = données de type matrice
- `centers` = soit le nombre de classes, soit une matrice contenant les coordonnées des points initiaux
- `nstart`=nombre d'initialisations

Sorties :

- `$cluster` = vecteur d'entiers de indiquant le numéro de la classe de chaque individu
- `$centers` = matrice des distances entre les individus et les centres de chaque classe
- `$size` = vecteur indiquant la taille de chaque classe
- `$iter` = nombre d'itérations

`# a 2-dimensional example`

```
x=rbind(matrix(rnorm(100, sd = 0.3), ncol = 2),      # rnorm génère une matrice de
matrix(rnorm(100, mean = 1, sd = 0.3), ncol = 2))    # réalisations d'une loi normale
# rbind concatène des matrices
```

```
x=scale(x) #centre et réduit
```

```
cl= kmeans(x,center=2,nstart=5) #clustering à 2 classes
```

```
print(cl) # affiche les résultats
```

```
plot(x, col = cl$cluster) # affiche les points avec une couleur différente par classe
# indexée par le numéro de la classe
```

```
points(cl$centers, col = 1:2, pch = 8, cex = 2) # ajoute les centres des classes
```

Fonction hclust pour la construction de l'arbre

```
hclust(d, method = "ward.D2",...)
```

Entrées :

- d=structure de dissimilarité entre les individus générée par la fonction dist
- method = mesure de dissimilarité entre les classes

Sorties : plusieurs attributs décrivant l'arbre. On retient

- \$height = vecteur indiquant la valeur du critère à chaque branche

Fonction cutree pour la construction des classes

```
cutree(tree, k = 2,...)
```

Entrées :

- tree=arbre résultant de hclust
- k = entier indiquant le nombre de classes

Sorties :

- un vecteur indiquant le numéro des classes

Fonctions pour représenter le dendrogramme

```
plot(tree)      affiche le dendrogramme
```

```
rect.hclust(tree,k=nclusters)    ajoute les classes sur le dendrogramme
```

```
# a 2-dimensional example
```

```
x=rbind(matrix(rnorm(100, sd = 0.3), ncol = 2),      # rnorm génère une matrice de  
         matrix(rnorm(100, mean = 1, sd = 0.3), ncol = 2))    # réalisations d'une loi normale  
                                                                 # rbind concatène des matrices
```

```
x=scale(x) #centre et réduit
```

```
distance=dist(x,"euclidean") #crée une structure de distance entre les individus
```

```
h=hclust(distance, "ward.D2") # crée l'arbre
```

```
plot(h$height) # affiche l'évolution du critère de dissimilarité entre classes
```

```
plot(h) # affiche le dendrogramme
```

```
rect.hclust(h,k=2) # ajoute les classes
```

```
c=cutree(h,k=2) # crée les classes
```

```
plot(x, col = c) # affiche les points avec une couleur différente par classe indexée  
                 # par le numéro de la classe
```

Quelques rappels sur R

R est un langage et un environnement logiciel open source pour les calculs statistiques et graphiques. R devient incontournable dans le traitement exploratoire et statistique des données.

Site Internet : Statistiques avec R

http://zoonek2.free.fr/UNIX/48_R_2004/all.html

- # pour écrire un commentaire
- `help(fonction)` # pour obtenir l'aide sur une fonction
- `ls()` # permet d'afficher tous les objets de la session de travail
- `rm()` # efface les objets de la session de travail
- `q()` # permet de quitter R

La session de travail avec tous les objets qui auront été définis dedans (matrices, dataframes, fonctions, ...) peut être sauvegardée puis rechargée à chaque utilisation :

- `getwd()` # affiche le répertoire courant
- `setwd("K:/ING1/GI/Statistiques Descriptives")` # permet de fixer le chemin d'accès au répertoire. attention d'utiliser / et non \
- `save.image(file="nom.Rdata")` # sauvegarde la session de travail dans le fichier nom.Rdata du répertoire courant
- `load("nom.Rdata")` # recharge la session de travail

Lecture d'un fichier

- `tab=read.table("Nom Fichier",header=TRUE ou FALSE si noms des colonnes indiqués dans le fichier,sep=" séparateur de champs",dec="symbole pour la décimale")`

Quelques instructions

- `as.matrix(tab)` # transforme tab en matrice
- `plot(matrice de points)` # graphique des points
- `points(matrice de points)` # pour ajouter des points à un graphique existant
- `X11()` pour ouvrir une nouvelle fenêtre graphique
- `tab[i,j]` # élément (i,j) du tableau
- `table(V1,V2)` # crée le tableau de contingence des variables V1 et V2