

Travaux Pratiques n° 1 : Approches non-supervisées

*Objectifs : mise en pratique sous R des concepts de l'apprentissage supervisé vus en cours - Réduction de dimensionalité, Clustering*

## 1 Avant goût

### 1.1 ACP

Considerer la matrice  $X$  de données brutes :

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 2 \\ 2 & 2 & 1 \\ 1 & 0 & 0 \\ 2 & 3 & 2 \end{pmatrix}$$

- (i) Créer un fichier `.txt` avec la matrice  $X$  et importer les données, e.g. :

```
X = read.table("Nom.txt", header = FALSE)
```

- (ii) Calculer le centre de gravité (individu moyen), la matrice  $Y$  des données centrées.  
(iii) Calculer la matrice  $V$  des variances-covariances.  
(iv) Calculer la variance de la première variable avec la commande `var`. Commenter.  
(v) Utiliser la commande `eigen` pour obtenir les vecteurs propres de  $V$ .  
(vi) Projeter les individus sur le plan factoriel principal. Utiliser la commande `plot` pour les visualiser.  
(vii) Charger la librairie "FactoMineR" et utiliser la commande `PCA` pour faire une ACP sur  $X$ .

### 1.2 Analyse sémantique latente

Le but dans cette section est d'introduire le paquet `lsa` et de son utilisation. Commencez par télécharger celui et exécutez pas à pas le code suivant :

```
library(lsa)
# create some files
td = tempfile()
dir.create(td)
```

```

write( c("dog", "cat", "mouse"), file=paste(td, "D1", sep="/"))
write( c("hamster", "mouse", "sushi"), file=paste(td, "D2", sep="/"))
write( c("dog", "monster", "monster"), file=paste(td, "D3", sep="/"))
write( c("dog", "mouse", "dog"), file=paste(td, "D4", sep="/"))

# read files into a document-term matrix
myMatrix = textmatrix(td, minWordLength=1)

# create the latent semantic space
myLSAspace = lsa(myMatrix, dims=dimcalc_raw())

# display it as a textmatrix again
round(as.textmatrix(myLSAspace),2) # should give the original

# create the latent semantic space
myLSAspace = lsa(myMatrix, dims=dimcalc_share())

# display it as a textmatrix again
myNewMatrix = as.textmatrix(myLSAspace)
myNewMatrix # should look be different!
# compare two terms with the cosine measure
cosine(myNewMatrix["dog",], myNewMatrix["cat",])
# compare two documents with pearson
cor(myNewMatrix[,1], myNewMatrix[,2], method="pearson")

##Find closest terms

# calc associations for mouse
associate(myNewMatrix, "mouse")

# clean up
unlink(td, recursive=TRUE)

```

### 1.3 K-means

Le but ici est de vous donner un exemple simple d'utilisation de l'algorithme k-moyennes :

```

x = c(-2,-2,0,2,-2,3)
y= c(2, -1,-1,2,3,0)
don = matrix(data=c(x,y), nr=6, nc=2)
ctre = c(-1,2,-1,3)
ctre1 =matrix(data=ctre, nr=2, nc=2)
cl1 = kmeans(don,ctre1,algorithm="Lloyd")
plot(don, col = cl1$cluster)
points(cl1$centers, col = 1 : 2, pch = 8, cex=2)

```

### 1.4 CAH

Même objectif que pour la k-means pour en utilisant l'algorithme CAH. La suggestion dans le code suivant est d'utiliser l'indice du lien moyen. Vous pouvez tester différents

types d'indices tels que le lien minimal ou maximal.

```
> data(iris)
> don=iris[,1 :4]
#
# Classification par la C.A.H.
#
> hclust(dist(don), "ave")
> plot(hc)
> plot(hc, hang = -1)}
```

## 2 La vraie vie : Text Mining

Nous allons considérer dans ce TP l'analyse d'un corpus de pages wiki, extrait à partir de <http://edutechwiki.unige.ch>. Commencez par récupérer l'archive `TP_TextMining.zip` envoyée par mail ou disponible sur `MyCourse`

1. Commencer par décompresser l'archive envoyée par email dans votre espace de travail local (e.g. `~/DataMining/TP1/`)
2. Changer d'espace de travail : `setwd("~/DataMining/TP1/")`
3. Commencer par télécharger et charger les librairies :

```
library(tm)
library(tm.plugin.webmining)
library(SnowballC)
library(corrplot)
```

4. Importer les données dans une variable "corpus", e.g. :

```
corpus <- Corpus(DirSource("./", encoding="UTF8"),
                 readerControl = list(language="lat"))
```

5. Vous pouvez visualiser le contenu d'un élément du corpus (ex. le premier) :

```
corpus[[1]]$content
```

6. Nettoyer le corpus pour pouvoir le traiter :

```
— Mettre en minuscule
wiki.cl1 <- tm_map(corpus, content_transformer(tolower))
— Tuer les balises
```

```
wiki.cl2 <- tm_map(wiki.cl1, content_transformer(extractHTMLStrip),
                  encoding="UTF-8")
```

```
— Enlever les chiffres
wiki.cl2 <- tm_map(wiki.cl2, removeNumbers)
— Autres nettoyage (facultatif)
```

```
(kill_chars
  <- content_transformer (function(x, pattern) gsub(pattern, " ", x)))

tm_map (wiki.cl2, kill_chars, "\u2019")
tm_map (wiki.cl2, kill_chars, "'")
tm_map (wiki.cl2, kill_chars, "[«»'“\"]")
tm_map (wiki.cl2, kill_chars, "\\[modifier\\]")
— Enlever les ponctuations qui restent :
```

```
wiki.cl3
  <- tm_map (wiki.cl2, removePunctuation, preserve_intra_word_dashes = TRUE)
— Enlever les mots fréquents :
wiki.essence <- tm_map (wiki.cl3, removeWords, stopwords("french"))
— Extraire les racines
wiki.racines <- tm_map (wiki.essence, stemDocument, language="french")
— Enlever les blancs s'il en reste
wiki.racines <- tm_map (wiki.racines, stripWhitespace)
— Test
wiki.racines[[2]]
class(wiki.racines)
```

#### 7. Créer la matrice documents-termes

```
wiki.mots <- Corpus(VectorSource(wiki.racines))
matrice_termes_docs <- DocumentTermMatrix(wiki.mots)
```

#### 8. Réduction de la matrice

```
inspect(removeSparseTerms(matrice_termes_docs, 0.4))
inspect(removeSparseTerms(matrice_termes_docs, 0.6))
```

#### 9. Visualisation de la matrice termes-documents

```
# Créer une DTM avec des poids normalisés
mtd.norm <- as.matrix(removeSparseTerms(
  TermDocumentMatrix(wiki.mots, control=list(weighting=weightTf)),
  0.2))
corrplot (mtd.norm, is.corr=FALSE)
```

#### 10. Création de la matrice TFidf

```
mtd.TfIdf2 <- as.matrix(removeSparseTerms(
  TermDocumentMatrix(wiki.mots, control=list(weighting=weightTfIdf)),
  0.2))
# Plot simple
corrplot (mtd.TfIdf2, is.corr=FALSE)
```

#### 11. Application d'une classification hiérarchique

```
m2 <- as.matrix(mtd.TfIdf2)
fit <- hclust(distMatrix, method = "ward")
plot(fit) # On dessine
rect.hclust(fit, k = 6) # couper le dendrogramme en 6 clusters
```

#### 12. Application du k-means sur les documents

```
m3 <- t(m2) # transposer la matrice pour segmenter les documents (pages wiki)
set.seed(122) # fixer un germe aléatoire
k <- 6 # nombre de clusters
kmeansResult <- kmeans(m3, k)
round(kmeansResult$centers, digits = 3) # centres de clusters
```

13. Visualisation des groupes de documents

```
for (i in 1:k) {
  cat(paste("cluster ", i, ": ", sep = ""))
  s <- sort(kmeansResult$centers[i, ], decreasing = T)
  cat(names(s)[1:5], "nn")
}
```

14. Appliquer une Décomposition en Valeurs Singulières à la matrice `mtd.TfIdf2`. Expliquez ? (commande `svd`)
15. Appliquer une ACP à cette matrice de données. Interprétez !
16. Appliquer la méthode k-moyennes sur la matrice `mtd.TfIdf2` après application de la Décomposition en Valeurs Singulières. Commentez ! Rque : Vous pouvez aussi passer par le paquet LSA.
17. Appliquer la méthode k-moyennes sur la matrice `mtd.TfIdf2` t après application de l'ACP. Commentez !