<u>Tutorial 4 : PCA and cluster</u>



I hope this tutorial will help you better understand PCA and clustering with R! Don't hesitate to ask me questions if you need further clarification. That's what I'm here for.

<u>Learning objectives:</u>

Familiarize yourself with R's functions for PCA and clustering

# ACP

## Part 0 : Packages downloading and loading
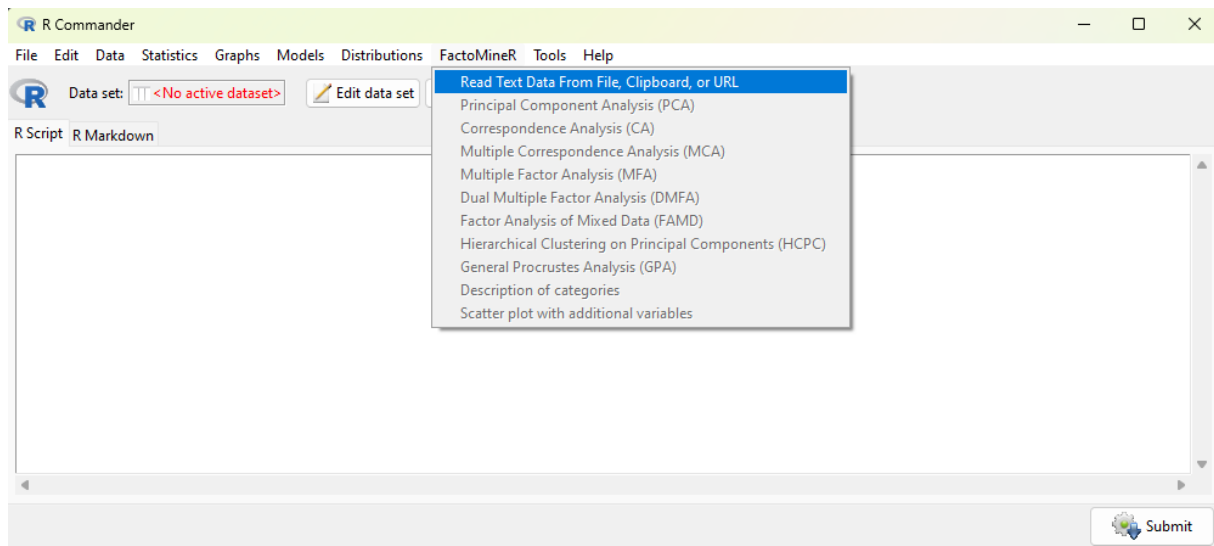I let you install and load the following packages :

- Rcmdr
- FactoMineR
- RcmdrPlugin.FactoMineR

## Part1 : Data loading
Download the dataset « data_PCA_got.csv » with a small change that is the name of rows need to be the characters' name. It will be better for graphs.
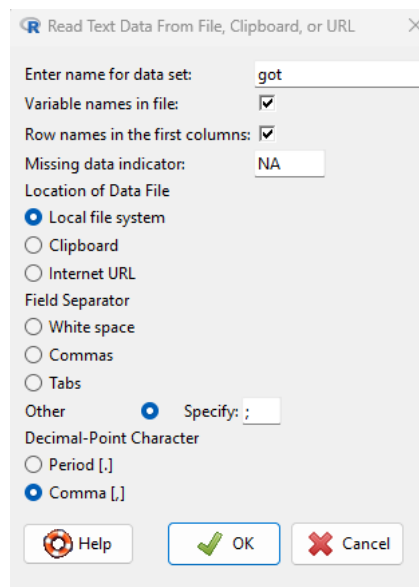
CLICK BUTTON TEAM
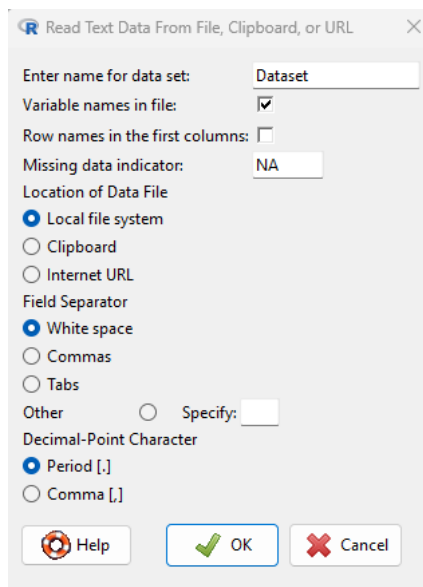
In the FactoMineR tab of Rcmdr.

Change dataset name, indicate 'got'.

Change field separator: check 'other' and specify semicolon '; ' as it is a csv file.

Change filed decimal-point character if needed.

Also check that the1st line corresponds to the characters' names (not compulsory, but visually more interesting).



Go and find your file in the directory structure of your computer.

## CODING TEAM

row.names = 'Name' means that the column 'Name' will be taken as row names.

stringsAsFactors = TRUE convert all the columns of text into factors, which allow using some functions.

```
got<-read.csv2("data_PCA_got.csv",row.names = "Name", stringsAsFactors = TRUE)
got
```

```
summary(got)
```
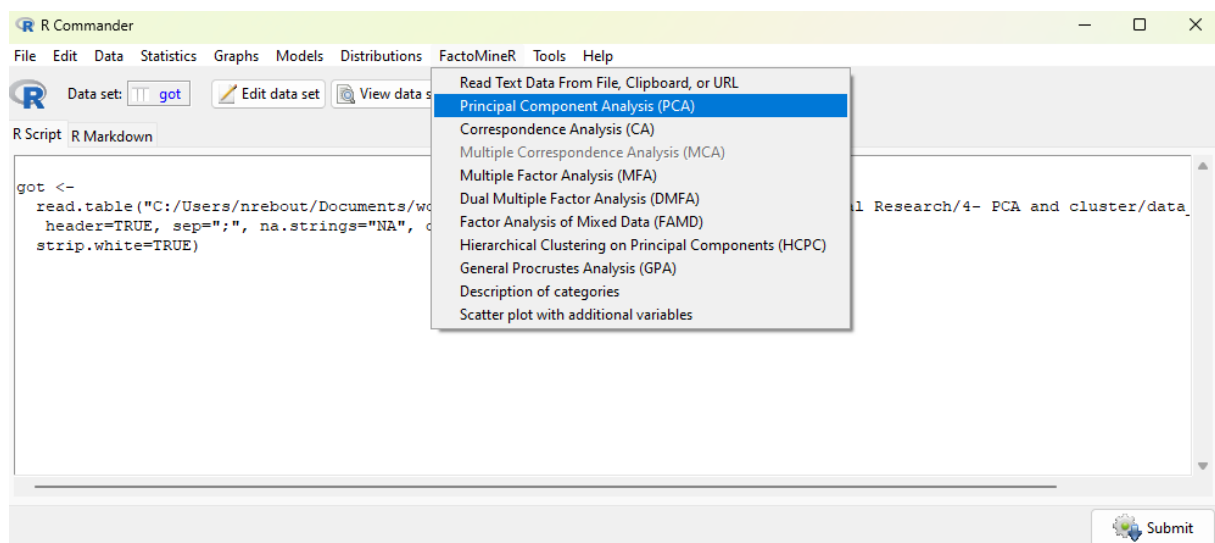
# Part2 : Principal Component Analysis (PCA) :
## 1) Creation of PCA and first conclusions

You will run PCA on the got dataset and store the result in a variable you'll call 'res' (for result).

We'll consider the following quantitative variables: number of chapters prepared ('Nb_chp_narr'), Gender, Married, Age, number of deaths in relationships ('Nb_death_relationships'), Popularity.

==We don't consider here survival ('IsAlive').==

CLICK BUTTON TEAM



Then select the variables you wish to include in the PCA. Note: they must be quantitative. Then click on OK.

The disadvantage of Rcmdr is that it is gonna write all the lines of code at once and will superimpose all the graphs.

```
got <-
  read.table("C:/Users/nrebout/Documents/worshop Easy Statistics to Shine in Biological Research/4- PCA and cluster/data",
    header=TRUE, sep=";", na.strings="NA", dec=",", row.names=1,
  strip.white=TRUE)
got.PCA<-got[, c("Nb_chp_narr", "Gender", "Married", "Age",
  "Nb_death_relationships", "Popularity", "isAlive")]
res<-PCA(got.PCA , scale.unit=TRUE, ncp=5, graph = FALSE)
print(plot.PCA(res, axes=c(1, 2), choix="ind", habillage="none",
  col.ind="black", col.ind.sup="blue", col.quali="magenta", label=c("ind",
  "ind.sup", "quali"),new.plot=TRUE))
print(plot.PCA(res, axes=c(1, 2), choix="var", new.plot=TRUE,
  col.var="black", col.quanti.sup="blue", label=c("var", "quanti.sup"),
  lim.cos2.var=0))
summary(res, nb.dec = 3, nbelements=10, nbind = 10, ncp = 3, file="")
remove(got.PCA)
```

I encourage you to isolate blocks of code as following :

```
got <-
  read.table("C:/Users/nrebout/Documents/worshop Easy Statistics to Shine in Biological Research/4- PCA and cluster/data
   header=TRUE, sep=";", na.strings="NA", dec=",", row.names=1,
  strip.white=TRUE)

got.PCA<-got[, c("Nb_chp_narr", "Gender", "Married", "Age",
 "Nb_death_relationships", "Popularity", "isAlive")]

res<-PCA(got.PCA , scale.unit=TRUE, ncp=5, graph = FALSE)

print(plot.PCA(res, axes=c(1, 2), choix="ind", habillage="none",
  col.ind="black", col.ind.sup="blue", col.quali="magenta", label=c("ind",
 "ind.sup", "quali"),new.plot=TRUE))

print(plot.PCA(res, axes=c(1, 2), choix="var", new.plot=TRUE,
  col.var="black", col.quanti.sup="blue", label=c("var", "quanti.sup"),
  lim.cos2.var=0))

summary(res, nb.dec = 3, nbelements=10, nbind = 10, ncp = 3, file="")

remove(got.PCA)
```

And then run the blocks of code again to get the different graphs.


TEAM CODAGE :

First, select the quantitative variables in your dataset. To do this, we create a sub-dataset from got, called got.PCA.

```
got.PCA<-got[, c("Nb_chp_narr", "Gender", "Married", "Age",
   "Nb_death_relationships", "Popularity")]
```

Then run the PCA function and store the result in 'res' (the calculations of PCA).

```
res<-PCA(got.PCA, scale.unit=TRUE, ncp=5, graph =  FALSE)
```

Then build the graphs based on the calculations :

```
    # plot of individuals
    plot(res,choix="ind")

    # plot of variables
    plot(res,choix="var")
```

What do the % values on the PCA axes correspond to ? What can we conclude?

Observe the cloud of variables. What can we conclude?

Observe the cloud of individuals. What can we conclude?

Considering the two types of information above, what can we conclude?


2) Because customers/researchers want « significant results », want ρ-values.

Use the dim.desc() function directly in the console.

```
dimdesc(res)
```

In the function, you specify the result of your PCA (here called 'res').

> The dimdesc() function helps to interpret the dimensions of a PCA, MCA, CFA, MFA or HFA.
>
> It allows dimensions to be described by both categorical and continuous variables.
>
> It shows which variables the axes are most related to: which continuous variables are most correlated with each axis, and which categorical variables and modalities best describe each axis.
>
> In the first part of the output, the correlation coefficient between each continuous variable and each PCA dimension is calculated. Correlation coefficients significantly different from zero are kept and sorted.

Conclude.

## 3) BONUS : use of FactoShiny

What is Factoshiny?

It's a graphical interface that lets you parameterize methods and modify graphs interactively.

No programming skills required.

Modifications to graphical options are instantly visible on graphs

Results (graphs and listings) are modified instantly.

Graphs can be downloaded, along with the lines of code needed to redo the analysis.

Factoshiny's result object can be reused to modify graphs. The interface is reopened with the latest settings and graphical options that can be modified.

First, install the Factoshiny package.

Then run this line of code:

```
RES<-PCAshiny(as.data.frame(got))
RES
```

An 'Internet' application window opens. All you have to do is explore and click.

Warnings : Don't forget to quit the application when you're finished, as R waits for you to do so before continuing.

# Partie 3 : Hierarchical ascending classification

CLICK BUTTON TEAM

Let's go back to the creation of the ACP.



Then click on 'perform clustering after PCA'.



Then click on OK.

Once again, you may be confronted with an overlay of all the graphics.

```
got.PCA<-got[, c("Nb_chp_narr", "Gender", "Married", "Age",
  "Nb_death_relationships", "Popularity")]
res<-PCA(got.PCA , scale.unit=TRUE, ncp=5, graph = FALSE)
res.hcpc<-HCPC(res ,nb.clust=-1,consol=FALSE,min=3,max=10,graph=TRUE)
print(plot.PCA(res, axes=c(1, 2), choix="ind", habillage="none",
  col.ind="black", col.ind.sup="blue", col.quali="magenta", label=c("ind",
  "ind.sup", "quali"),new.plot=TRUE))
print(plot.PCA(res, axes=c(1, 2), choix="var", new.plot=TRUE,
  col.var="black", col.quanti.sup="blue", label=c("var", "quanti.sup"),
  lim.cos2.var=0))
summary(res, nb.dec = 3, nbelements=10, nbind = 10, ncp = 3, file="")
remove(got.PCA)
```

I advise you to isolate blocks of code.

```
got.PCA<-got[, c("Nb_chp_narr", "Gender", "Married", "Age",
  "Nb_death_relationships", "Popularity")]

res<-PCA(got.PCA , scale.unit=TRUE, ncp=5, graph = FALSE)

res.hcpc<-HCPC(res ,nb.clust=-1,consol=FALSE,min=3,max=10,graph=TRUE)

print(plot.PCA(res, axes=c(1, 2), choix="ind", habillage="none",
  col.ind="black", col.ind.sup="blue", col.quali="magenta", label=c("ind",
  "ind.sup", "quali"),new.plot=TRUE))

print(plot.PCA(res, axes=c(1, 2), choix="var", new.plot=TRUE,
  col.var="black", col.quanti.sup="blue", label=c("var", "quanti.sup"),
  lim.cos2.var=0))

summary(res, nb.dec = 3, nbelements=10, nbind = 10, ncp = 3, file="")

remove(got.PCA)
```

## CODING TEAM

```
res.hcpc<-HCPC(res ,nb.clust=-1,consol=FALSE,min=3,max=10,graph=TRUE)

plot(res.hcpc)
plot(res.hcpc,choice="tree")
plot(res.hcpc,choice="map")
```