

Cluster course: Hierarchical and non hierarchical clustering R code

By: S. Civit

Contents

| | |
|--|----------|
| Hierarchical clustering | 1 |
| Dataset | 1 |
| Distància | 1 |
| Dendograma | 2 |
| Possible exercici: | 6 |
| | |
| Non hierarchical clustering: Kmeans | 6 |
| Dataset | 6 |
| Distància | 7 |
| Graph | 7 |
| Conexió amb altres tècniques multivariants PCA | 12 |

Hierarchical clustering

Dataset

Fangataufa és un atolò en “format rectangular” amb una superfície total de 45 km². Originalment no tenia cap pas a la llacuna interior, però les forces armades franceses van volar 400 m d’esculls per obrir un pas que facilités el programa de proves nuclears. És una zona militar amb accés prohibit sense autorització.

Entre el 1966 i 1996 s’hi van fer 5 explosions nuclears atmosfèriques i 10 subterrànies a una profunditat entre 500 i 700 m sota la llacuna.

Les dades procedeixen d’un estudi que pretenia detectar canvis temporals en l’estructura de la comunitat gasteròpoda en els esculls del atolò. Els mol·luscs dels esculls van ser parcialment o totalment esborrats per la calor de les proves nuclears i van ser recolonitzats per larves oceàniques. En tots els esculls, la composició de la comunitat abans de les proves era molt diferent de la que va evolucionar fins a després.

Cal doncs entendre les dades i per tant la **matriu multivariant** de dades

- Files zones del mostreig, que correspon a zones del atolò
- Columnes: espècies de gasteròpods

```
data<-read.table("fangataufa.txt",sep="\t",header=TRUE)
```

Distància

- Elecció de la distància com element clau. Comencem de forma incorrecta (en aquest cas amb la distància euclidia)

```
data2<-scale(data[,2:14])# Standardise the variables  
is.matrix(data2)
```

```
## [1] TRUE
```

```
data.D1 = dist(data2, method="eucl") # Compute Euclidean distance (for illustrative purposes)
```

Dendograma

- Representem mitjançant el mètode UPGM i mètode de Ward
- Determimen la correlació coefenètica (distorsió de les dades originals a les observades en el dendrograma)

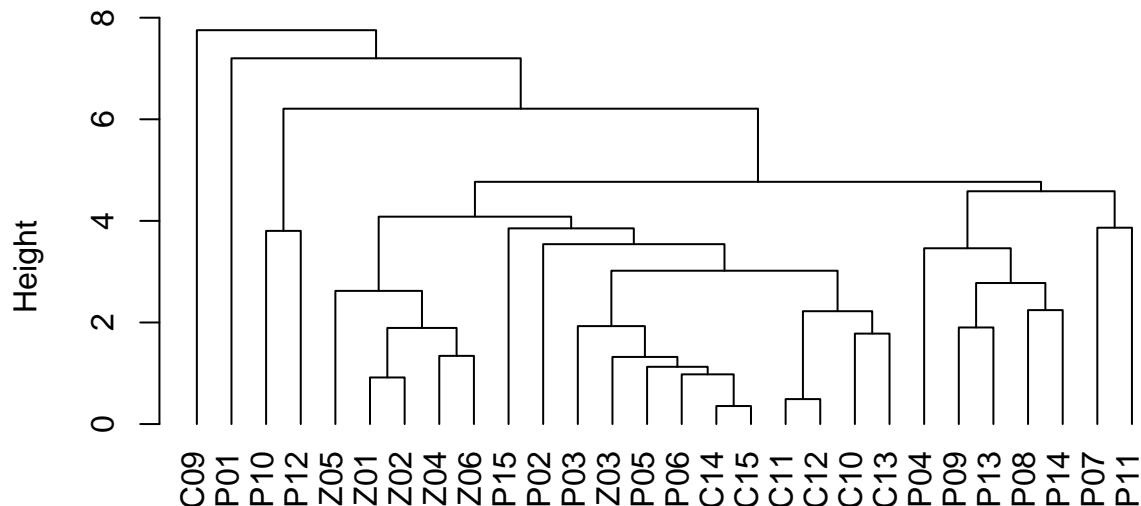
```
# Agglomerative clustering, UPGMA method:
```

```
clusterAV = hclust(data.D1, method="average")#average=UPGMA  
clusterW<-hclust(data.D1^2, method="ward.D2")#ward
```

```
# Plot the dendrogram
```

```
plot(clusterAV, hang=-1, labels=data[,1])
```

Cluster Dendrogram



```
data.D1  
hclust (*, "average")
```

```
# Cophenetic distances of the dendrogram
```

```
Ours.coph = cophenetic(clusterAV)  
cor(data.D1, Ours.coph)
```

```
## [1] 0.851541
```

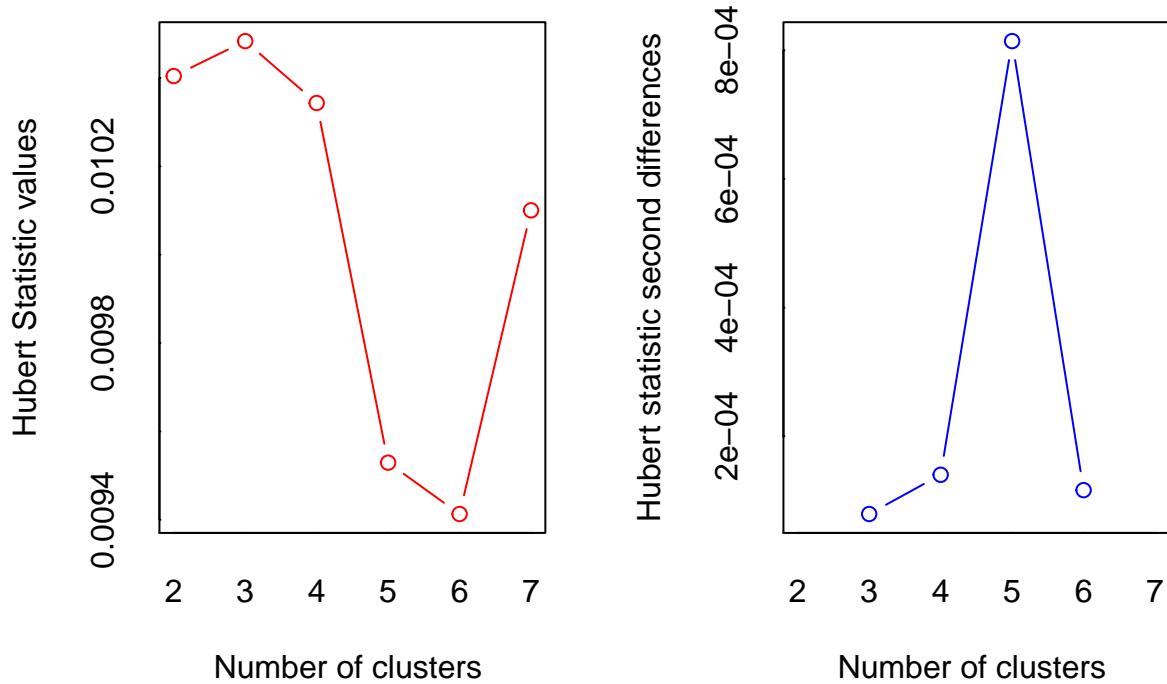
Nombre de grups?

Funció **NbClust ()** [al paquet NbClust R] (Charrad et al. 2014): proporciona 30 índexs per determinar el nombre rellevant de clúster a partir de les combinacions de nombre de clústers, mesures de distància i mètodes d'agrupació.

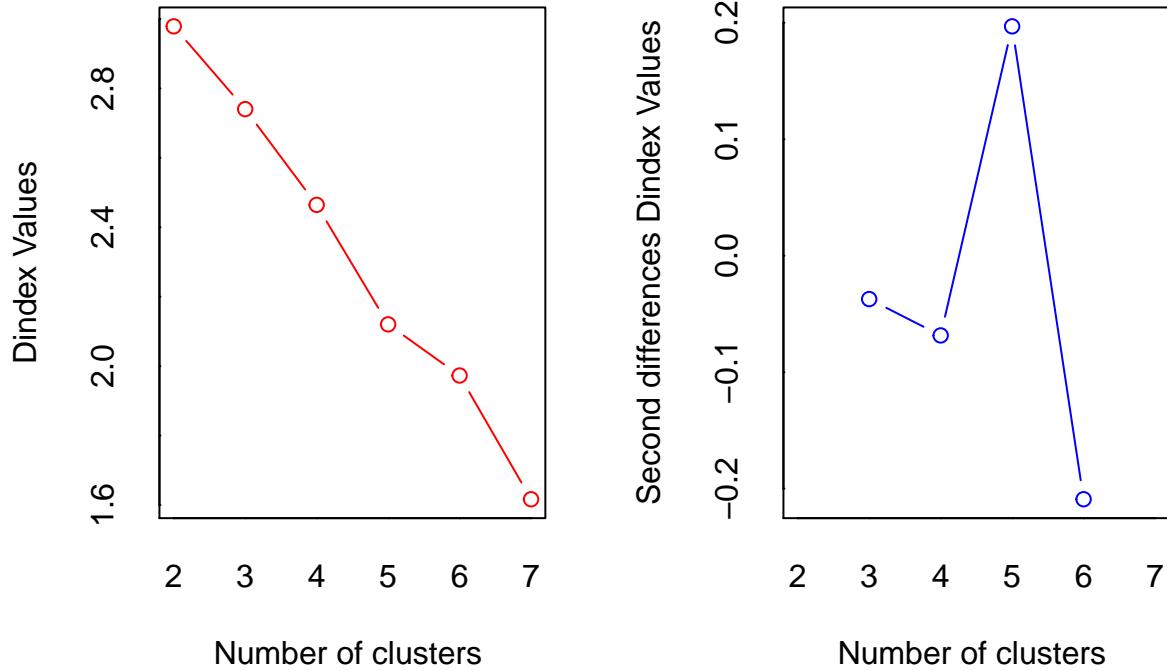
```
#install.packages("NbClust")
library(NbClust)

NbClust(data = data2, diss = NULL, distance = "euclidean",
        min.nc = 2, max.nc = 7, method = "average")

## Warning in pf(beale, pp, df2): NaNs produced
```



```
## *** : The Hubert index is a graphical method of determining the number of clusters.
## In the plot of Hubert index, we seek a significant knee that corresponds to a
## significant increase of the value of the measure i.e the significant peak in Hubert
## index second differences plot.
```

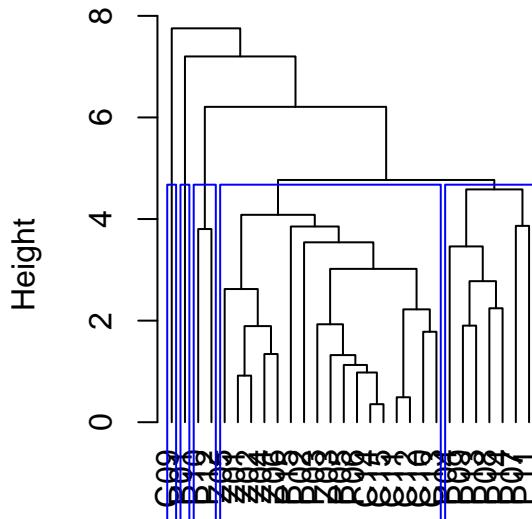


```

## *** : The D index is a graphical method of determining the number of clusters.
## In the plot of D index, we seek a significant knee (the significant peak in Dindex
## second differences plot) that corresponds to a significant increase of the value of
## the measure.
##
## *****
## * Among all indices:
## * 6 proposed 2 as the best number of clusters
## * 5 proposed 3 as the best number of clusters
## * 3 proposed 4 as the best number of clusters
## * 3 proposed 5 as the best number of clusters
## * 2 proposed 6 as the best number of clusters
## * 5 proposed 7 as the best number of clusters
##
## ***** Conclusion *****
##
## * According to the majority rule, the best number of clusters is 2
##
## *****
## $All.index
##      KL      CH Hartigan      CCC      Scott      Marriot      TrCovW      TraceW
## 2 0.5406  4.1846  4.0962 -2.8921 127.1457 3.437719e+13 1223.8504 302.3397
## 3 0.5686  4.2981  5.8870 -3.6442 195.8445 6.651199e+12 977.0556 261.1903
## 4 0.7255  5.2824  9.1556 -3.2498 292.6396 3.727648e+11 699.4205 211.4077
## 5 3.1174  7.4386  3.5009 -0.8025 366.3150 4.192873e+10 337.1872 153.0295

```


Hcluster Dendrogram (UPGMA method)



```
data.D1  
hclust (*, "average")
```

Caracterització dels grups

```
#number of elements and elements within groups  
table(groups.3)
```

```
## groups.3  
## 1 2 3 4 5  
## 1 17 1 7 2  
g1=data[groups.3==1,]  
  
# Mean profile G1  
G1mean<-round(apply(g1[-1], 2, mean), 3)
```

Possible exercici:

- 1. Repeat the analysis plotting Hierarchical clustering dendrogram and cofenetic correlation according to “real distance” for this dataset.
- 2. Decide how many groups and characterize them and show a 2D plot

Non hierarchical clustering: Kmeans

Dataset

In this data set we observe the composition of different wines. Given a set of observations (x_1, x_2, \dots, x_n), where each observation is a d -dimensional real vector,

```
#install.packages('rattle')
data(wine, package='rattle')
head(wine)

##   Type Alcohol Malic  Ash Alkalinity Magnesium Phenols Flavanoids Nonflavanoids
## 1    1   14.23  1.71 2.43        15.6       127     2.80      3.06      0.28
## 2    1   13.20  1.78 2.14        11.2       100     2.65      2.76      0.26
## 3    1   13.16  2.36 2.67        18.6       101     2.80      3.24      0.30
## 4    1   14.37  1.95 2.50        16.8       113     3.85      3.49      0.24
## 5    1   13.24  2.59 2.87        21.0       118     2.80      2.69      0.39
## 6    1   14.20  1.76 2.45        15.2       112     3.27      3.39      0.34
##   Proanthocyanins Color  Hue Dilution Proline
## 1              2.29  5.64 1.04      3.92     1065
## 2              1.28  4.38 1.05      3.40     1050
## 3              2.81  5.68 1.03      3.17     1185
## 4              2.18  7.80 0.86      3.45     1480
## 5              1.82  4.32 1.04      2.93      735
## 6              1.97  6.75 1.05      2.85     1450
```

Distància

```
wine.stand <- scale(wine[-1]) # To standarize the variables
```

Graph

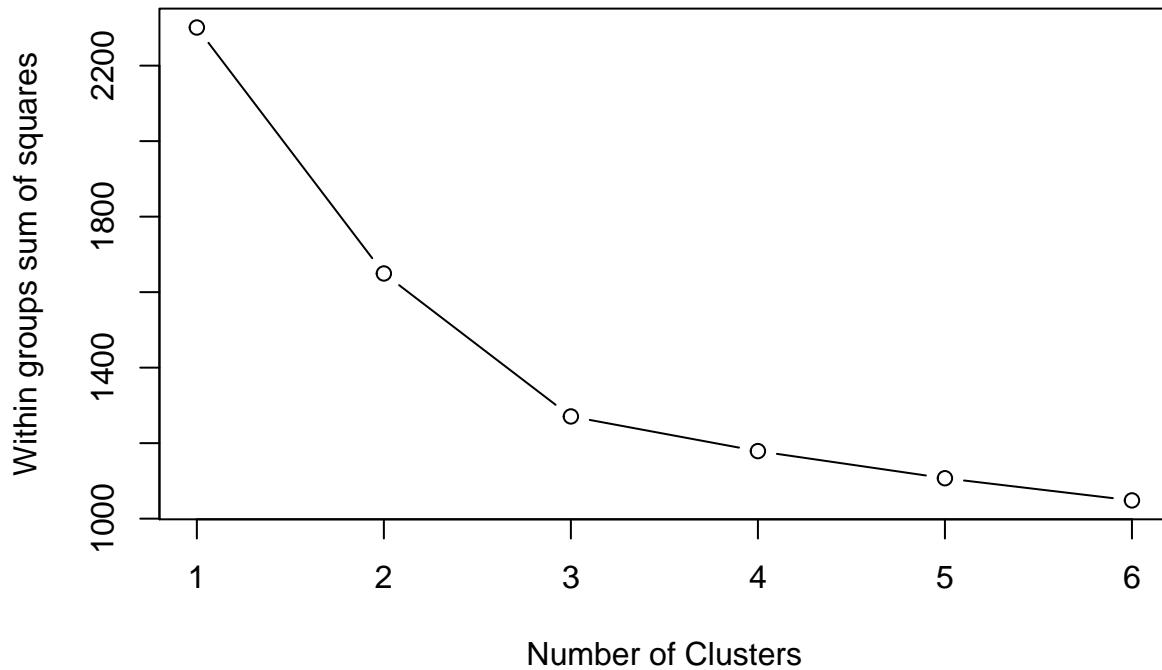
Nombre de grups?

Una qüestió fonamental és com determinar el valor del paràmetre k . Si observem el percentatge de variància explicat en funció del nombre de clústers: s'ha de triar un nombre de clústers de manera que afegir un altre clúster no doni una modelització molt millor de les dades. Més exactament, si es traça el percentatge de variància explicat pels clústers contra el nombre de clústers, els primers clústers afegiran molta informació (expliquen molta variància), però en algun moment el guany marginal caurà, donant un angle al gràfic.

- opció 1: En aquest moment es tria el nombre de clústers, d'aquí el “criteri del colze”.

```
wssplot <- function(data, nc=15, seed=1234){
  wss <- (nrow(data)-1)*sum(apply(data,2,var))
  for (i in 2:nc){
    set.seed(seed)
    wss[i] <- sum(kmeans(data, centers=i)$withinss)}
  plot(1:nc, wss, type="b", xlab="Number of Clusters",
       ylab="Within groups sum of squares")}

wssplot(wine.stand, nc=6)
```



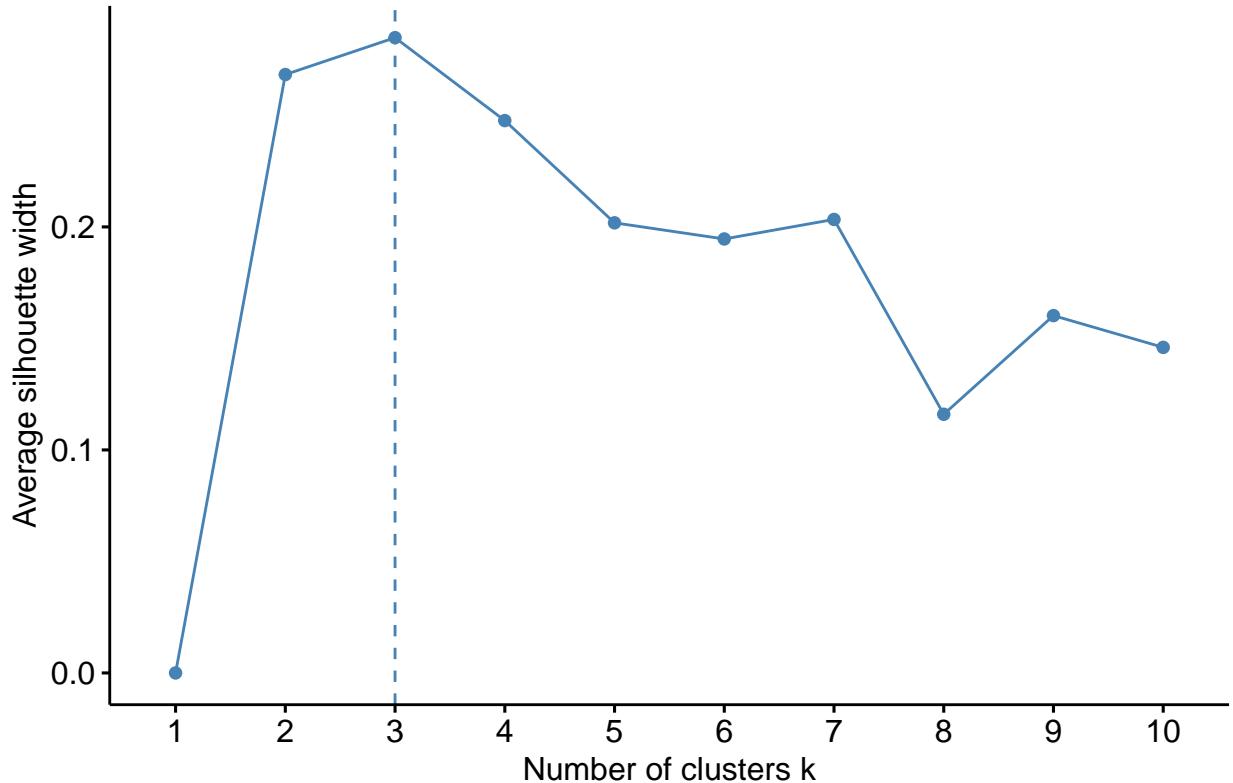
- Opció 2: Ús de la library **factoextra**

```
#install.packages(factoextra)
library(factoextra)

## Warning: package 'factoextra' was built under R version 3.6.3
## Loading required package: ggplot2
## Warning: package 'ggplot2' was built under R version 3.6.3
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
library(ggplot2)

fviz_nbclust(wine.stand, kmeans, method = c("silhouette", "wss", "gap_stat"))
```

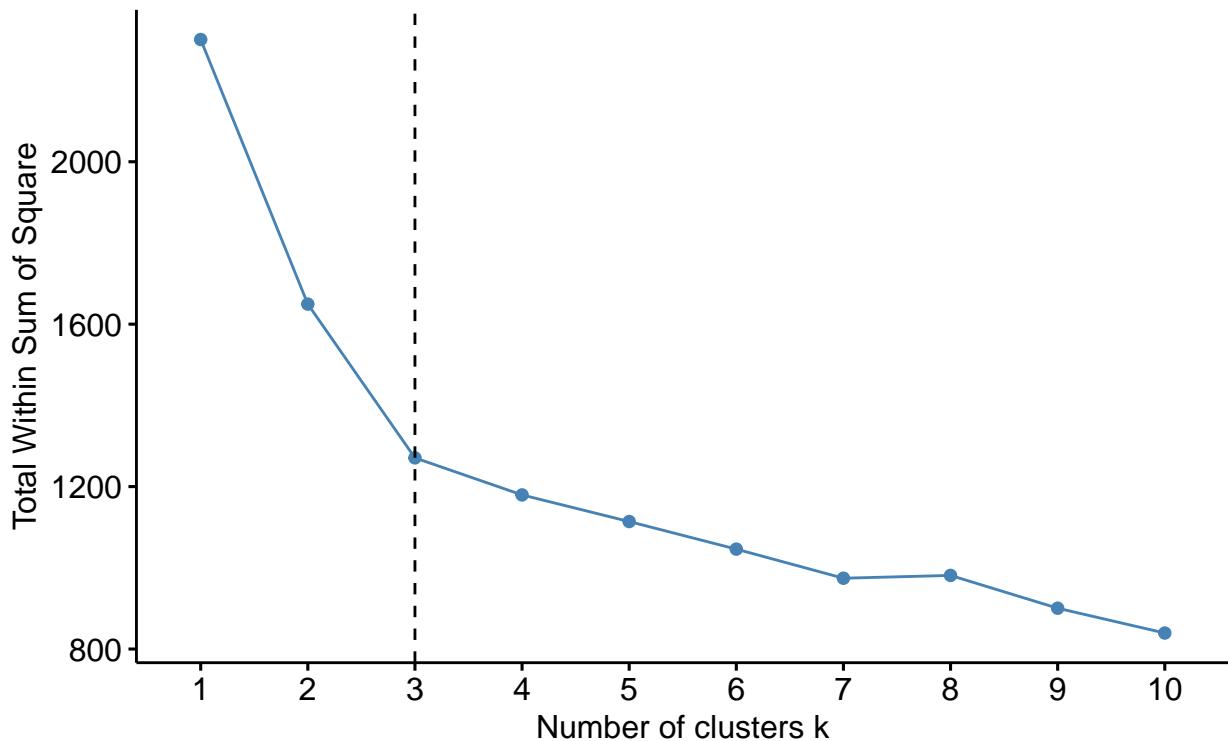
Optimal number of clusters



```
# Elbow method
fviz_nbclust(wine.stand, kmeans, method = "wss") +
  geom_vline(xintercept = 3, linetype = 2) +
  labs(subtitle = "Elbow method")
```

Optimal number of clusters

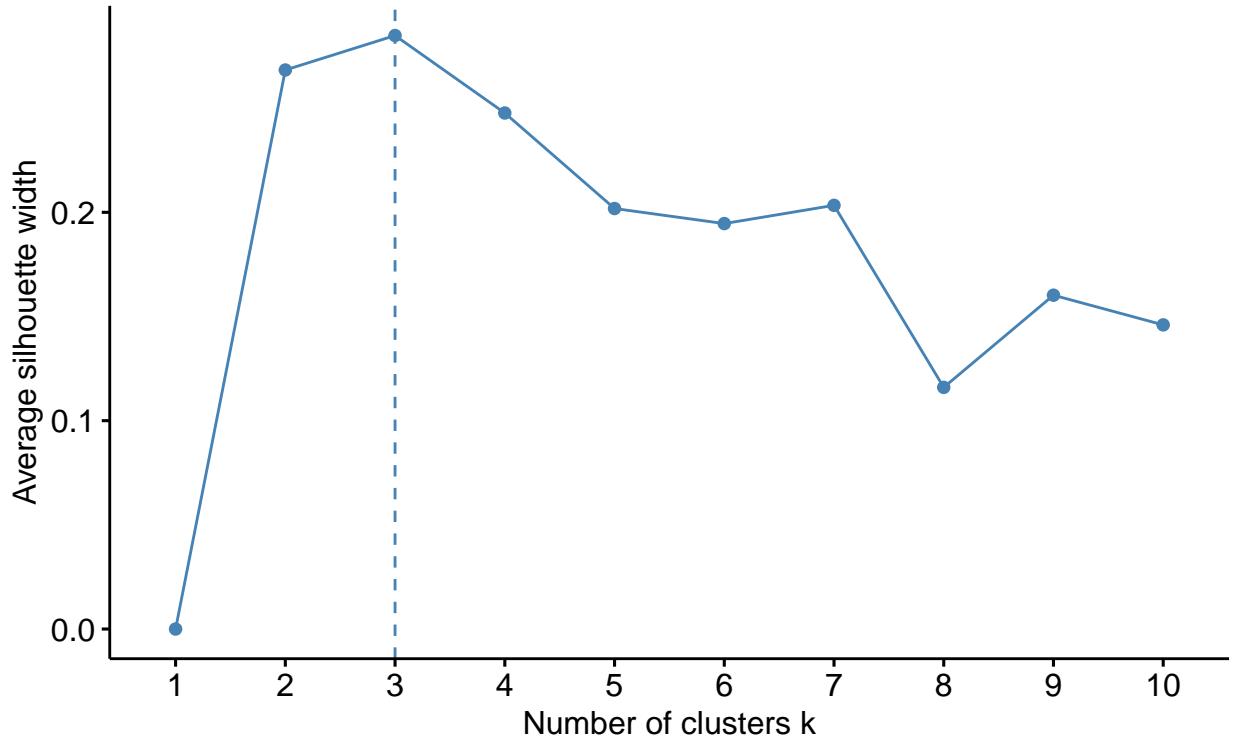
Elbow method



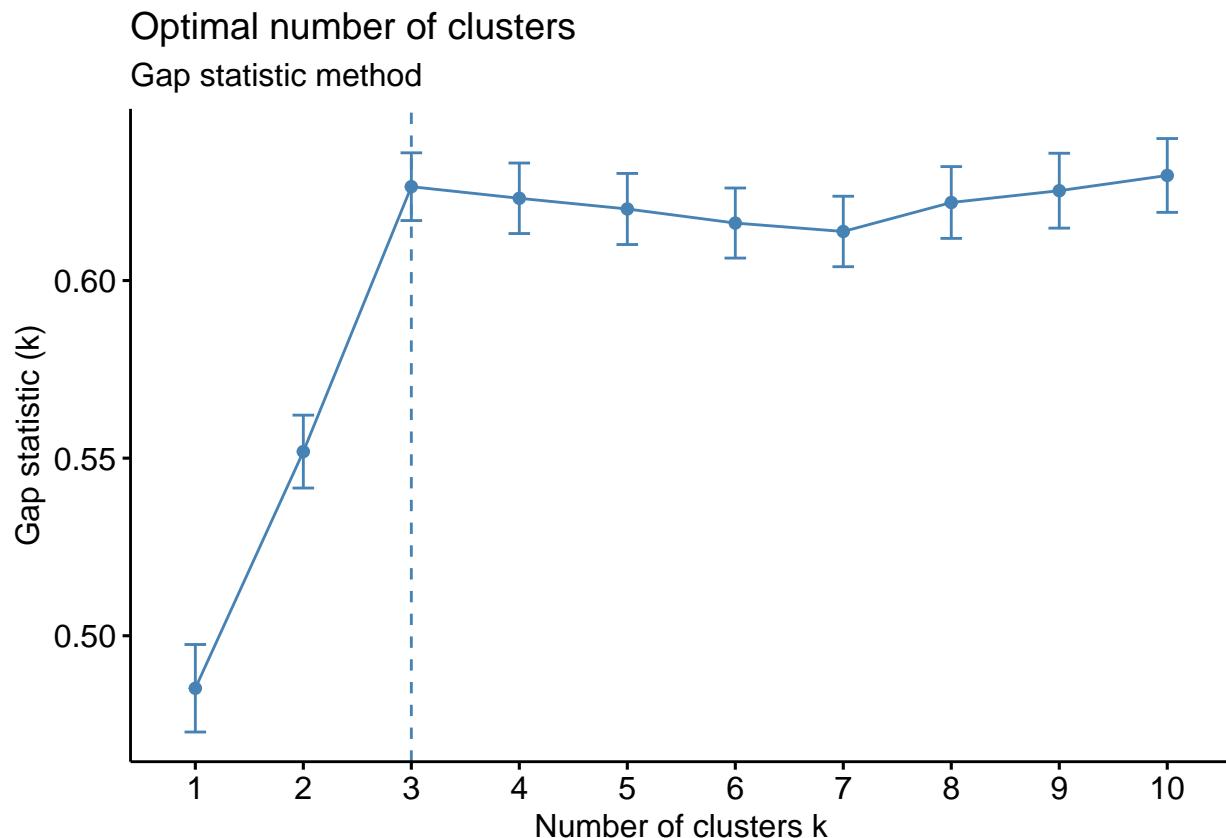
```
# Silhouette method  
fviz_nbclust(wine.stand, kmeans, method = "silhouette") +  
  labs(subtitle = "Silhouette method")
```

Optimal number of clusters

Silhouette method



```
# Gap statistic
# nboot = 50 to keep the function speedy.
# recommended value: nboot= 500 for your analysis.
# Use verbose = FALSE to hide computing progression.
set.seed(123)
fviz_nbclust(wine.stand, kmeans, nstart = 25, method = "gap_stat", nboot = 50) +
  labs(subtitle = "Gap statistic method")
```



Caracterització dels grups

```
# Centroids:
k.means.fit$centers
```

```
##          Alcohol      Malic       Ash Alkalinity   Magnesium   Phenols
## 1  0.1644436  0.8690954  0.1863726  0.5228924 -0.07526047 -0.97657548
## 2  0.8328826 -0.3029551  0.3636801 -0.6084749  0.57596208  0.88274724
## 3 -0.9234669 -0.3929331 -0.4931257  0.1701220 -0.49032869 -0.07576891
##          Flavanoids Nonflavanoids Proanthocyanins      Color        Hue Dilution
## 1 -1.21182921     0.72402116     -0.77751312  0.9388902 -1.1615122 -1.2887761
## 2  0.97506900    -0.56050853      0.57865427  0.1705823  0.4726504  0.7770551
## 3  0.02075402    -0.03343924      0.05810161 -0.8993770  0.4605046  0.2700025
##          Proline
## 1 -0.4059428
## 2  1.1220202
## 3 -0.7517257
```

```
# Cluster size:
k.means.fit$size
```

```
## [1] 51 62 65
```

Conexió amb altres tècniques multivariants PCA

library cluster ens permeten representar (amb l'ajut de PCA) la solució de clúster en 2 dimensions

```

library(cluster)
clusplot(wine.stand, k.means.fit$cluster, main='2D representation of the Cluster solution',
          color=TRUE, shade=TRUE,
          labels=2, lines=0)

```

