# Lab10.R

rstudio-user

2021-04-12

```r
#Heirarchical Clustering

#1.Load the necessary packages for clustering
#install.packages("tidyverse")
#install.packages("cluster")
#install.packages("factoextra")
#install.packages("dendextend")

library(tidyverse)  # data manipulation
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.0 --
```

```
## v ggplot2 3.3.3     v purrr   0.3.4
## v tibble  3.0.5     v dplyr   1.0.4
## v tidyr   1.1.3     v stringr 1.4.0
## v readr   1.4.0     v forcats 0.5.1
```

```
## -- Conflicts ------------------------------------------ tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```r
library(cluster)      # clustering algorithms
library(factoextra) # clustering visualization
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```r
library(dendextend) # for comparing two dendrograms
```

```
##
## ---------------------
## Welcome to dendextend version 1.14.0
## Type citation('dendextend') for how to cite the package.
##
## Type browseVignettes(package = 'dendextend') for the package vignette.
## The github page is: https://github.com/talgalili/dendextend/
##
## Suggestions and bug-reports can be submitted at: https://github.com/talgalili/dendextend/issues
## Or contact: <tal.galili@gmail.com>
##
##  To suppress this message use:  suppressPackageStartupMessages(library(dendextend))
## ---------------------
```

```
##
## Attaching package: 'dendextend'
```

```
## The following object is masked from 'package:stats':
```

```
##
##       cutree
```

```r
#Hierarchical Clustering Algorithms
## Agglomerative Clustering
## Divisive hierarchical clustering

#Reading file
df <- USArrests
df <- na.omit(df)
head(df)
```

```
##            Murder Assault UrbanPop Rape
## Alabama      13.2     236       58 21.2
## Alaska       10.0     263       48 44.5
## Arizona       8.1     294       80 31.0
## Arkansas      8.8     190       50 19.5
## California    9.0     276       91 40.6
## Colorado      7.9     204       78 38.7
```
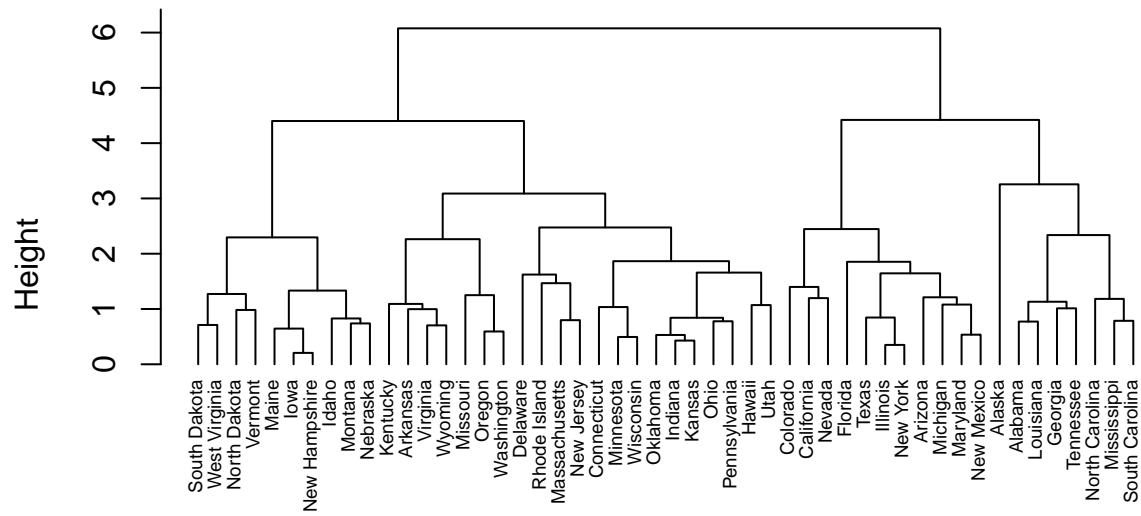
```r
#3. Scaling/Standardizing
df <- scale(df)
head(df)
```

```
##                 Murder    Assault   UrbanPop         Rape
## Alabama     1.24256408 0.7828393 -0.5209066 -0.003416473
## Alaska      0.50786248 1.1068225 -1.2117642  2.484202941
## Arizona     0.07163341 1.4788032  0.9989801  1.042878388
## Arkansas    0.23234938 0.2308680 -1.0735927 -0.184916602
## California  0.27826823 1.2628144  1.7589234  2.067820292
## Colorado    0.02571456 0.3988593  0.8608085  1.864967207
```

```r
#4. Perform Agglomerative Hierarchical Clustering by computing dissimilarity
#values and perform any hierarchical clustering method like complete linkage and
#then plot the dendogram.

##Agglomerative
# Dissimilarity matrix
d <- dist(df, method = "euclidean")
# Hierarchical clustering using Complete Linkage
hc1 <- hclust(d, method = "complete" )
# Plot the obtained dendrogram
plot(hc1, cex = 0.6, hang = -1, main="Dendogram")
```

## Dendogram



d
hclust (*, "complete")

```
#5.Determine optimal number of clusters
# methods to assess
m <- c( "average", "single", "complete", "ward")
names(m) <- c( "average", "single", "complete", "ward")

# function to compute coefficient
ac <- function(x) {
  agnes(df, method = x)$ac
}

map_dbl(m, ac)

##   average    single  complete      ward
## 0.7379371 0.6276128 0.8531583 0.9346210
```