

## **Tarea 4**

Estudiante

**John Daniel hoyos Arias  
Ivan Santiago Rojas Martinez  
Genaro Alfonso Aristizabal Echeverri**

Docente

**Cesar Augusto Gomez Velez**

Asignatura

**Analitica de datos**



UNIVERSIDAD  
**NACIONAL**  
DE COLOMBIA

Sede Medellín  
29 de Noviembre del 2022

# Índice

<b>1. Ejercicio 2</b>	<b>4</b>
1.1. a) . . . . .	6
1.2. b) . . . . .	6
1.3. c) . . . . .	8
1.4. d) . . . . .	8

## 1. Ejercicio 2

Se considera el conjunto de datos **USArrests**. En este ejercicio se agruparán los estados en **USArrests** con agrupamiento jerárquico. Este conjunto de datos contiene estadísticas, en arrestos por cada 100,000 residentes por agresión, asesinato y violación en cada uno de los 50 estados de EE. UU. en 1973. También se proporciona el porcentaje de la población que vive en áreas urbanas.

Primeramente, se procede a cargar la base de datos **USArrests** y examinar sus características:

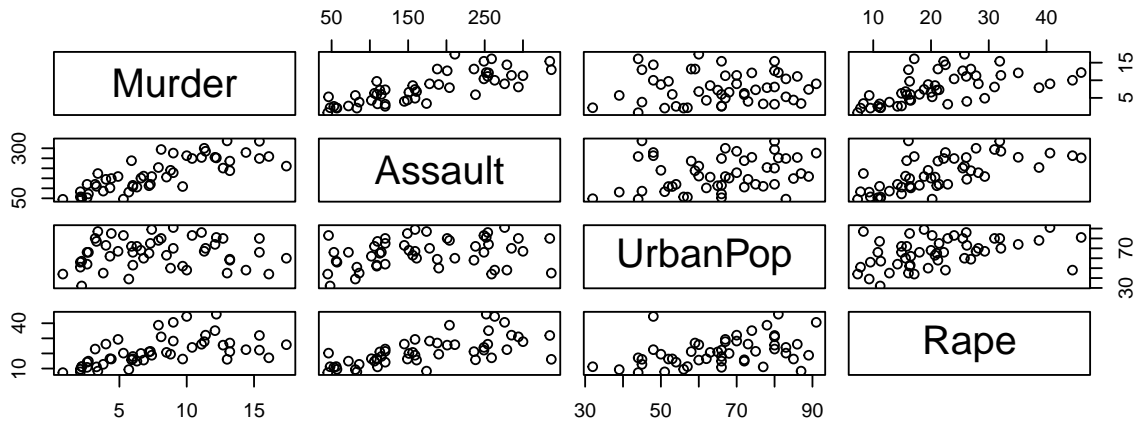
```
##           Murder Assault UrbanPop Rape
## Alabama      13.2      236        58 21.2
## Alaska       10.0      263        48 44.5
## Arizona       8.1      294        80 31.0
## Arkansas      8.8      190        50 19.5
## California    9.0      276        91 40.6
## Colorado      7.9      204        78 38.7

## 'data.frame':   50 obs. of  4 variables:
## $ Murder   : num  13.2 10 8.1 8.8 9 7.9 3.3 5.9 15.4 17.4 ...
## $ Assault  : int  236 263 294 190 276 204 110 238 335 211 ...
## $ UrbanPop: int  58 48 80 50 91 78 77 72 80 60 ...
## $ Rape     : num  21.2 44.5 31 19.5 40.6 38.7 11.1 15.8 31.9 25.8 ...
```

Se observa como la base cuenta con 50 observaciones y 4 variables las cuales todas son numericas y su descripción se presenta seguidamente:

- **Murder:** Arrestos por asesinato (por 100.000).
- **Assault:** Arrestos por asalto (por 100.000).
- **UrbanPop:** Porcentaje de población urbana.
- **Rape:** Arrestos por violaciones (por 100.000).

Adicionalmente, se presentan un análisis descriptivos de estas variables:

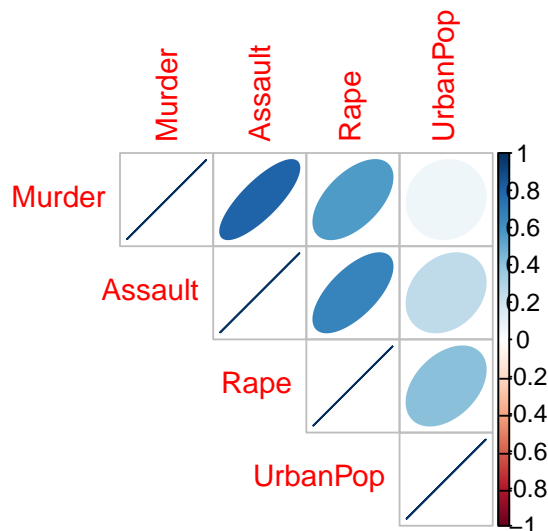


Del anterior gráfico de dispersión entre las variables, se observa como entre cada par de combinación de variables, existe una relación creciente. Lo cual en primera instancia podría ser un indicativo de que posiblemente en los estados donde se presente mayor porcentaje de población urbana también se puede presentar mayores casos de arrestos por asalto, asesinato o violación.

Por otro lado, se presenta una matriz de correlación entre las cuatro variables:

```
##           Murder  Assault  UrbanPop   Rape
## Murder    1.00000000 0.8018733 0.06957262 0.5635788
## Assault    0.80187331 1.0000000 0.25887170 0.6652412
## UrbanPop   0.06957262 0.2588717 1.00000000 0.4113412
## Rape       0.56357883 0.6652412 0.41134124 1.0000000
```

También, se presenta un gráfico de correlaciones de estas variables:



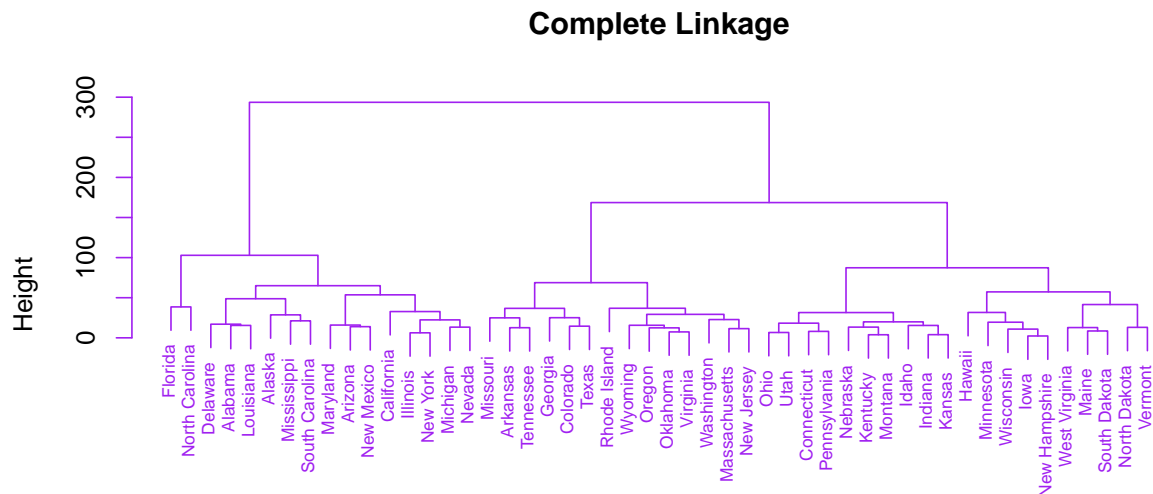
De los resultados obtenidos anteriormente, se observa como:

- Existe una alta correlación positiva entre los arrestos por asesinato y los arrestos por asaltos, la cual es de un 0.8018733, Esto puede indicar que, así como pueden aumentar los arrestos por asalto en un estado de USA, también puede aumentar los arrestos por asesinato en ese mismo estado.
- Existe una alta correlación positiva entre los arrestos por asalto y los arrestos por violaciones, la cual es de un 0.6652412, Esto puede indicar que, así como pueden aumentar los arrestos por asalto en un estado de USA, también puede aumentar los arrestos por violaciones en ese mismo estado.
- Se observa en la matriz de correlaciones como, no existe una aparente correlación significativa entre los arrestos por asesinatos y el porcentaje de población urbana.

### 1.1. a)

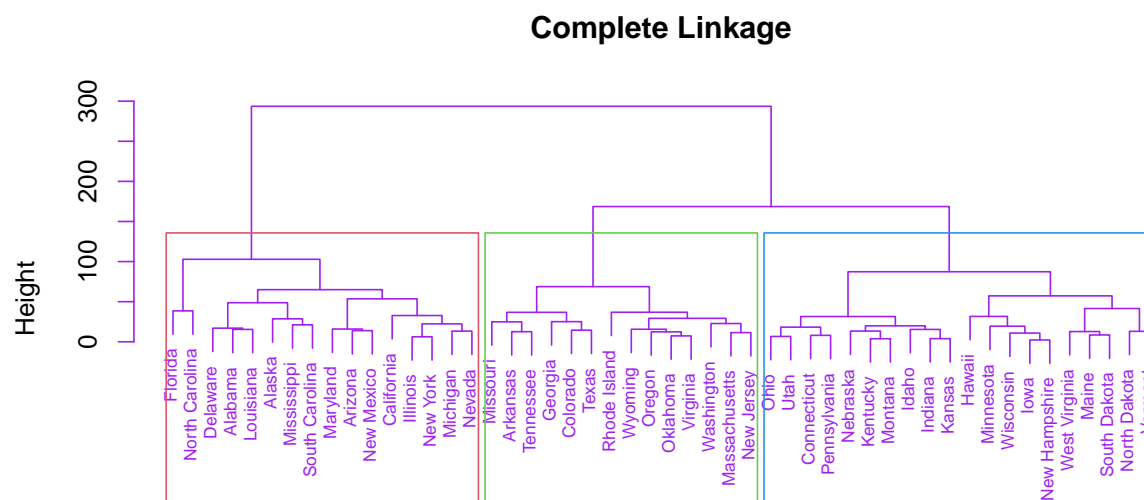
Se utiliza agrupación jerárquica con enlace completo y distancia euclidiana, para agrupar los estados, de la siguiente forma:

Luego se presenta el **dendrograma** de dicho enlace completo:



### 1.2. b)

Se procede a separar en el **dendrograma** a una altura que dé como resultado 3 *clusters*.



Luego, usando la función **cutree()** se puede observar las etiquetas a las que pertenece cada estado según el cluster al que se le asigno.

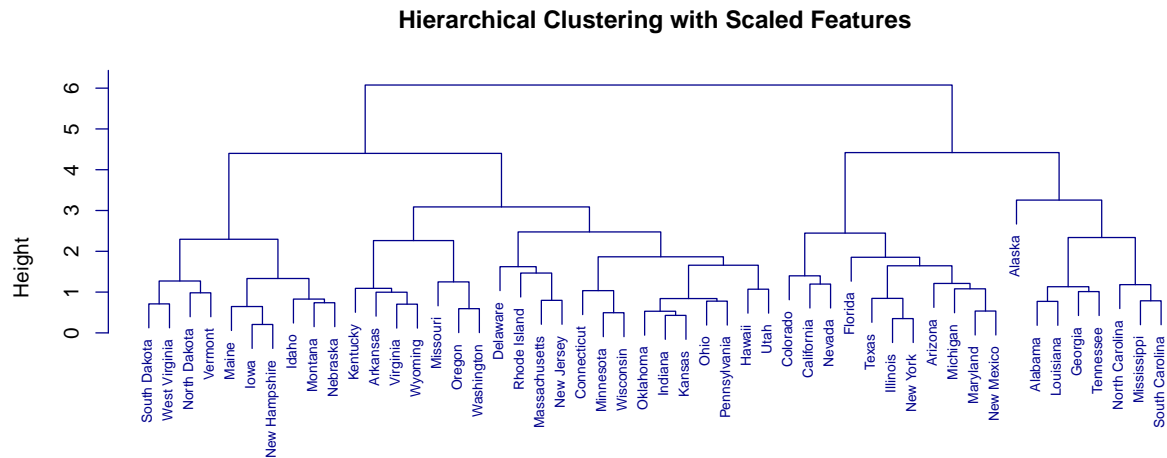
##	Alabama	Alaska	Arizona	Arkansas	California
##	1	1	1	2	1
##	Colorado	Connecticut	Delaware	Florida	Georgia
##	2	3	1	1	2
##	Hawaii	Idaho	Illinois	Indiana	Iowa
##	3	3	1	3	3
##	Kansas	Kentucky	Louisiana	Maine	Maryland
##	3	3	1	3	1
##	Massachusetts	Michigan	Minnesota	Mississippi	Missouri
##	2	1	3	1	2
##	Montana	Nebraska	Nevada	New Hampshire	New Jersey
##	3	3	1	3	2
##	New Mexico	New York	North Carolina	North Dakota	Ohio
##	1	1	1	3	3
##	Oklahoma	Oregon	Pennsylvania	Rhode Island	South Carolina
##	2	2	3	2	1
##	South Dakota	Tennessee	Texas	Utah	Vermont
##	3	2	2	3	3
##	Virginia	Washington	West Virginia	Wisconsin	Wyoming
##	2	2	3	3	2

### 1.3. c)

Ahora se procede a escalar las variables a fin de tener una desviación estándar de uno y luego se realiza la respectiva agrupación jerárquica usando un enlace completo y la distancia euclidiana.

La función **scale()** permite escalar las variables, así como el proceso de agrupación jerárquica se muestra a continuación:

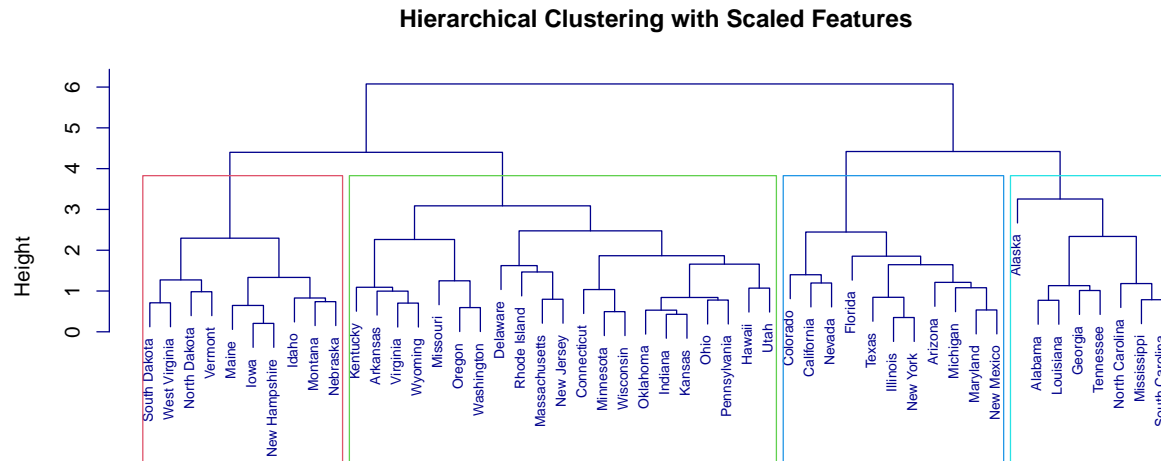
Luego se presenta el **dendrograma** de dicho enlace completo:



### 1.4. d)

Se observa como, en el **dendrograma** correspondiente a las agrupaciones jerárquicas con las variables escaladas, es notablemente distinto a la agrupación jerárquica generada sin las variables escaladas. Dado que, si bien estamos tratando con los mismos datos, la escalación de variables hace que en cada sub rama existan agrupaciones más uniformes. Inicialmente con las variables sin escalar se observaba claramente una distinción entre tres grupos o *clusters* distintos. En cambio, realizando el escalado de las variables se puede observar una posible distinción entre 4 *clusters*.

A continuación se presenta una posible agrupación entre 4 *clusters*:



Aunque también podría ser una agrupación de tres *clusters*.

En definitiva, se considera que las variables deben ser escaladas previamente, ya que proporciona una mejor estabilidad a la hora de hacer agrupaciones jerárquicas. Porque es bien sabido que la distancia euclidiana no tiene en cuenta el tipo de escala en la cual se encuentran las variables. Lo cual hace que, en ocasiones, usar esta medida no sea del todo preciso y como se pudo observar en los literales anteriores, se obtuvieron *clusters* y **dendrogramas** distintos.