



UNIVERSIDAD
NACIONAL DEL OESTE

Explotación de Datos

ACTIVIDAD N^o 6

Clustering de Criminalidad entre Provincias

PROFESORES:

*Dejean, Gustavo
Españadero, Juan
Mendoza, Dante*

INTEGRANTES GRUPO B:

*Benitez, Nicolas
Garcia Ravlic, Ignacio Agustin
Rechimon, Pablo Hernan
Rodríguez, Miguel Ángel*

FECHA DE ENTREGA:

21 de Noviembre de 2020

Resumen

Analizamos la distribución de las Tasas Criminales por provincia en Argentina y aplicamos distintas técnicas de agrupamiento, con el fin de formar grupos de provincias con características similares, darles una clase y que sea una base para futuros análisis en el que se incluyan otro tipo de variables.

Para esto contamos con las herramientas que nos provee el lenguaje R y los datos obtenidos en el Ministerio de Seguridad acerca de Delitos por Provincia, también utilizamos un pequeño dataset donde se detalla para cada provincia la cantidad de habitantes.

Palabras Clave:

*agrupamiento - cluster - dendograma - distancia - Euclídea - Manhattan - Maxima
- Minkowski*

Índice

1	Introducción	1
1.1	Problemática	1
1.2	Datos a utilizar	1
1.3	Objetivo	1
2	Desarrollo	2
2.1	Preparación de los datos	2
2.2	Análisis de los datos	2
2.3	Cálculo de Distancias	3
2.4	Número Óptimo de Clusters	7
2.5	Agrupamiento por Provincia	9
2.6	Agrupamiento por Crimen	14
3	Conclusión	15
4	Anexo	16
4.1	Código de limpieza en R	16
4.2	Código de clustering en R	17

Graficos

Fig. 1	Boxplot de Tasas Criminales	2
Fig. 2	Correlaciones entre los Crimenes	3
Fig. 3	Distancia euclidia	4
Fig. 4	Distancia Manhattan	5
Fig. 5	Distancia máxima	6
Fig. 6	Distancia Minkowski	7
Fig. 7	Numero optimo de clusters - Metodo wss	8
Fig. 8	Numero optimo de clusters - Metodo Silouette	8
Fig. 9	Dendograma	9
Fig. 10	K-means clustering	10
Fig. 11	Pam clustering	11
Fig. 12	Mapa de Calor	12
Fig. 13	PCA dinamico	13
Fig. 14	Diferencias por crimen entre grupos	13
Fig. 15	Dendograma por crimenes	14

1 Introducción

1.1 Problemática

Las tasas criminales a lo largo de las provincias de Argentina presentan notables diferencias, un análisis de las similitudes entre estas tasas facilitaría el desarrollo de otros análisis donde se incluyan variables culturales y socioeconómicas en busca de correlaciones.

1.2 Datos a utilizar

Utilizamos el lenguaje R y el dataset de Estadísticas Criminales de la República Argentina, por provincia, obtenido en la página del ministerio de seguridad de la nación.

1.3 Objetivo

Decidimos realizar un análisis y aplicar diferentes técnicas de agrupamiento con el fin de obtener grupos de provincias con menos diferencias en las tasas de los diferentes crímenes.

2 Desarrollo

2.1 Preparación de los datos

Los datos han sido procesados con anterioridad, contábamos con mediciones por año, por provincia, de diferentes tipos de crímenes, donde aplicamos un primer filtro para quedarnos con las muestras tomadas en el 2019, luego otro filtro para quedarnos con determinados crímenes (ya que en un principio eran 29). Realizamos unos cambios de nombres, eliminamos unas columnas y mergeamos con los datos de los habitantes por provincia para calcular la tasa del crimen por cada cien mil habitantes de dicha provincia.

2.2 Análisis de los datos

Realizamos un Boxplot de los datos, podemos visualizar como se distribuyen los valores de las tasas en los diferentes delitos, podemos ver que los robos se dan una mayor medida, seguido por las amenazas, y que los homicidios son los que se dan en menor medida, junto con las tentativas de homicidio y las violaciones.

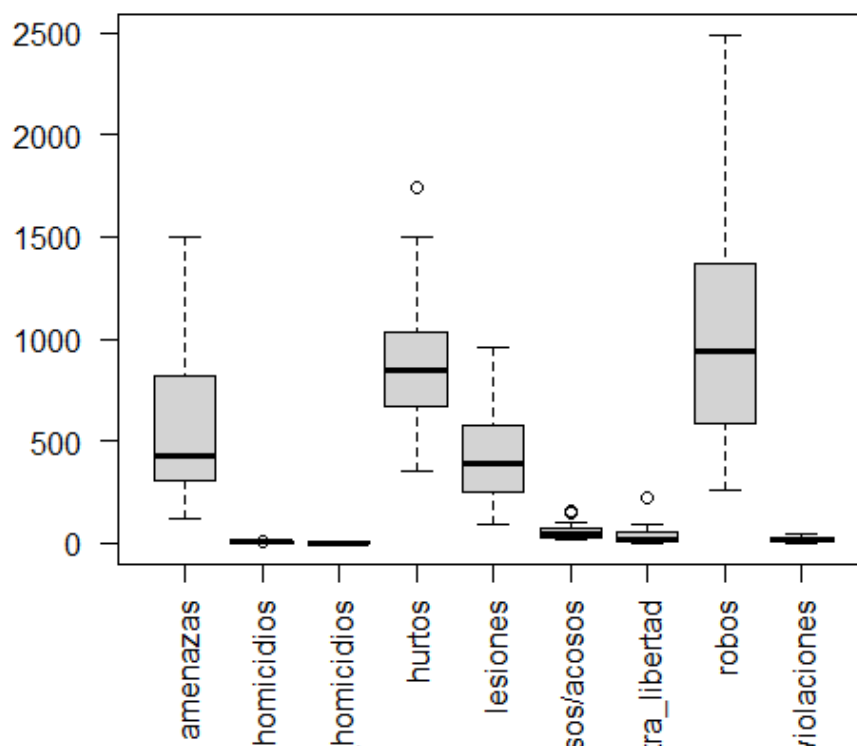


Fig. 1: Boxplot de Tasas Criminales

Graficamos las correlaciones, aquí podemos observar las distintas relaciones entre variables, siendo amenazas es la variable con mayor grado de correlación positiva

con casi todas las demás, y que en general el grado de correlación es mayor a .3 exceptuando algunos casos.

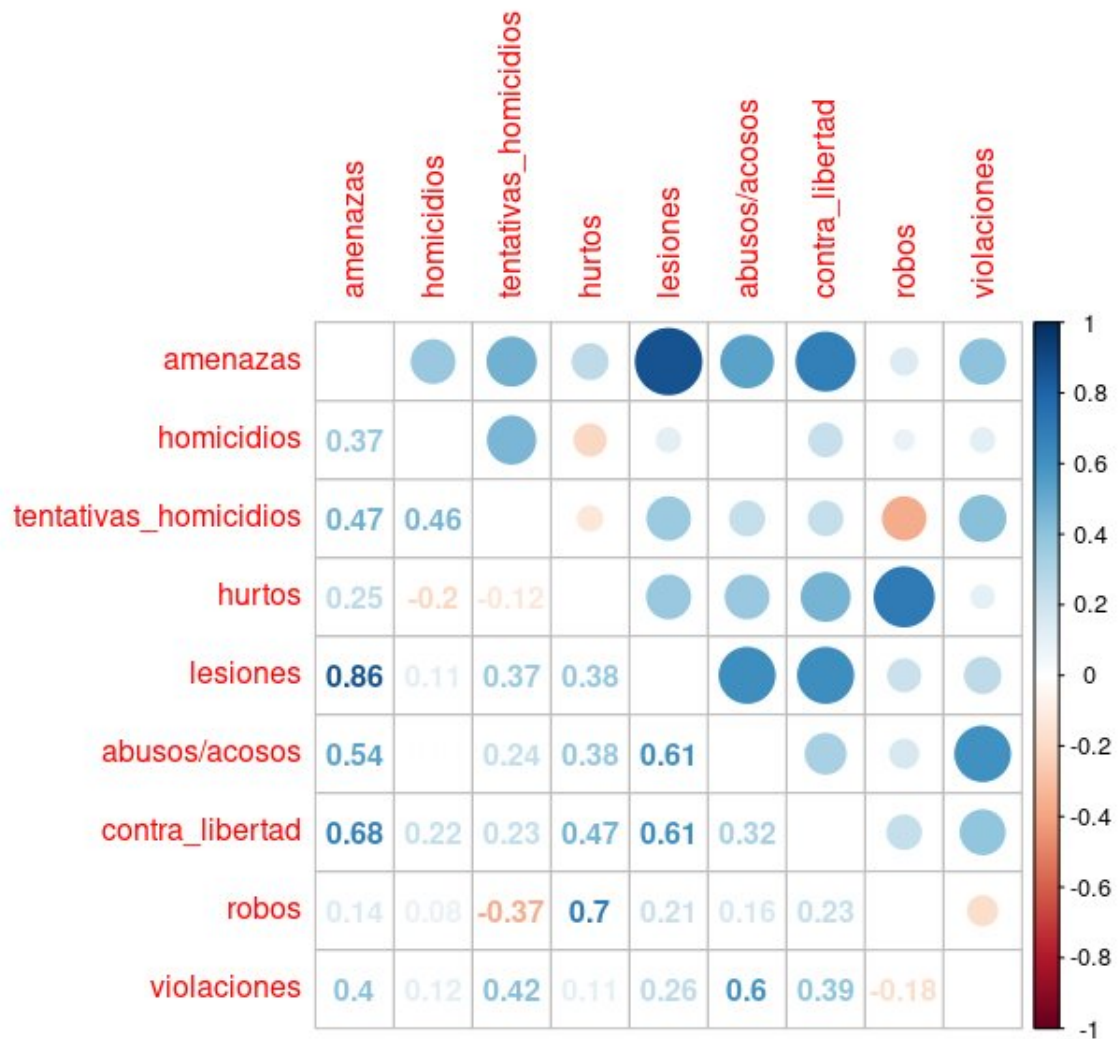


Fig. 2: Correlaciones entre los Crímenes

2.3 Cálculo de Distancias

La distancia euclídea es la raíz cuadrada de la sumatoria de la diferencia de cuadrados de cada coordenada; en el siguiente gráfico podemos ver la matriz de distancia generada por el algoritmo, mientras mas rojo sea, menor distancia habra, y mientras mas azul, mayor distancia.

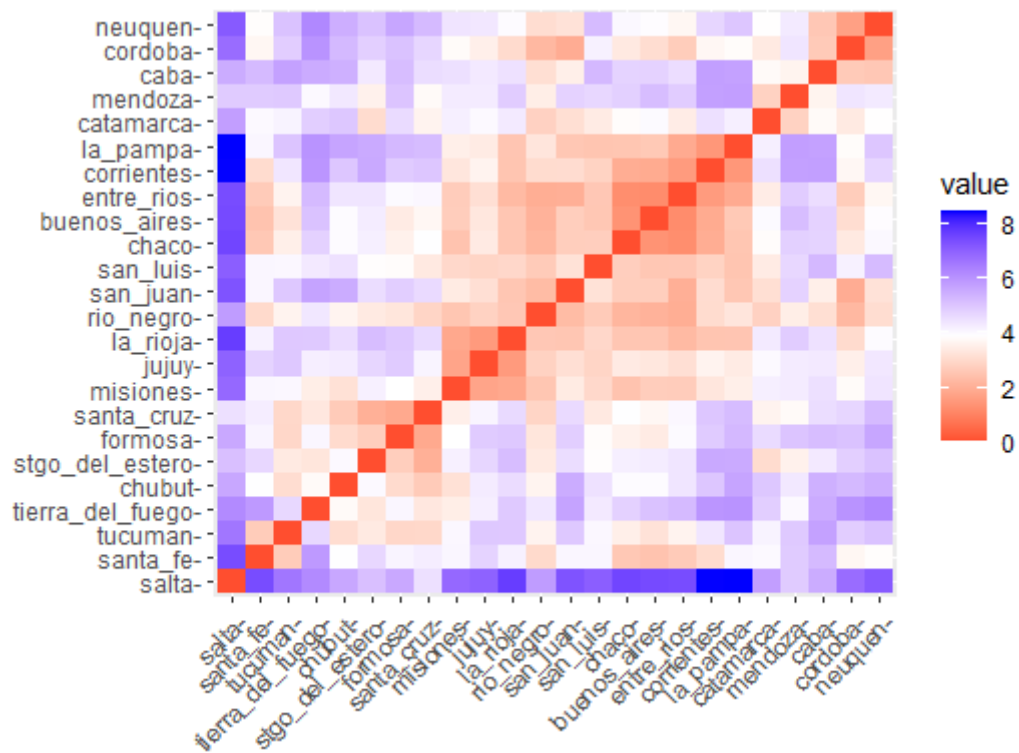


Fig. 3: Distancia euclidia

En la distancia Manhattan es igual a la sumatoria de las diferencias absolutas de cada coordenada. A continuación visualizamos el gráfico de distancias:

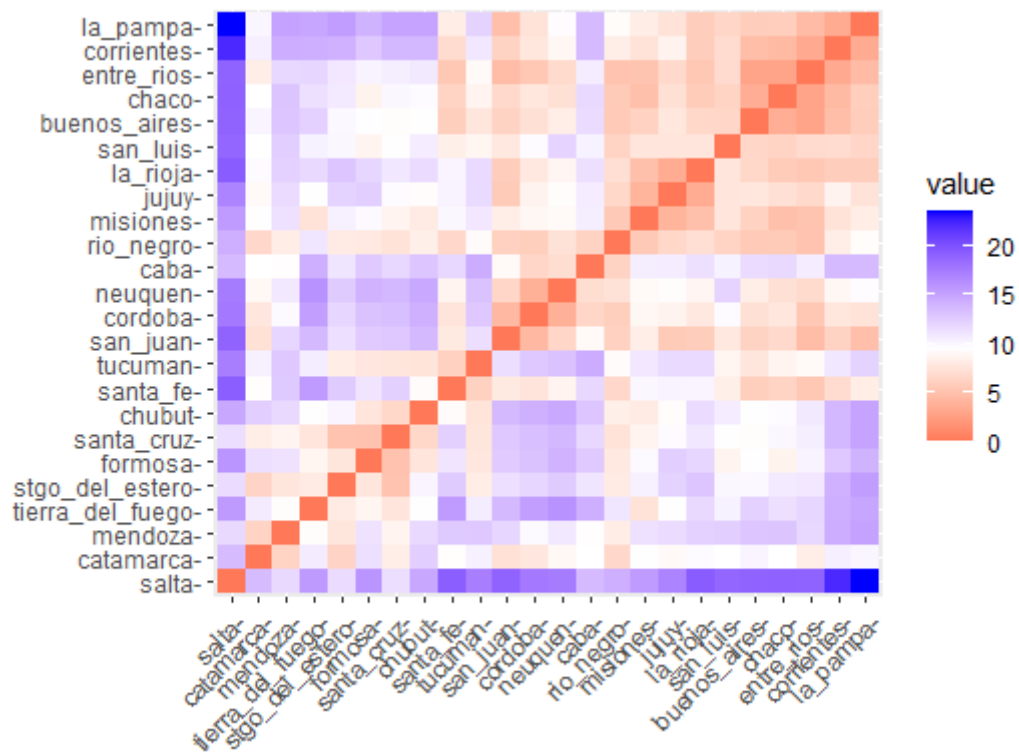


Fig. 4: Distancia Manhattan

La distancia máxima es el valor máximo hallado al hacer la diferencia absoluta, muy similar a la distancia de Manhattan con la diferencia de que no es la sumatoria, sino que es el máximo.

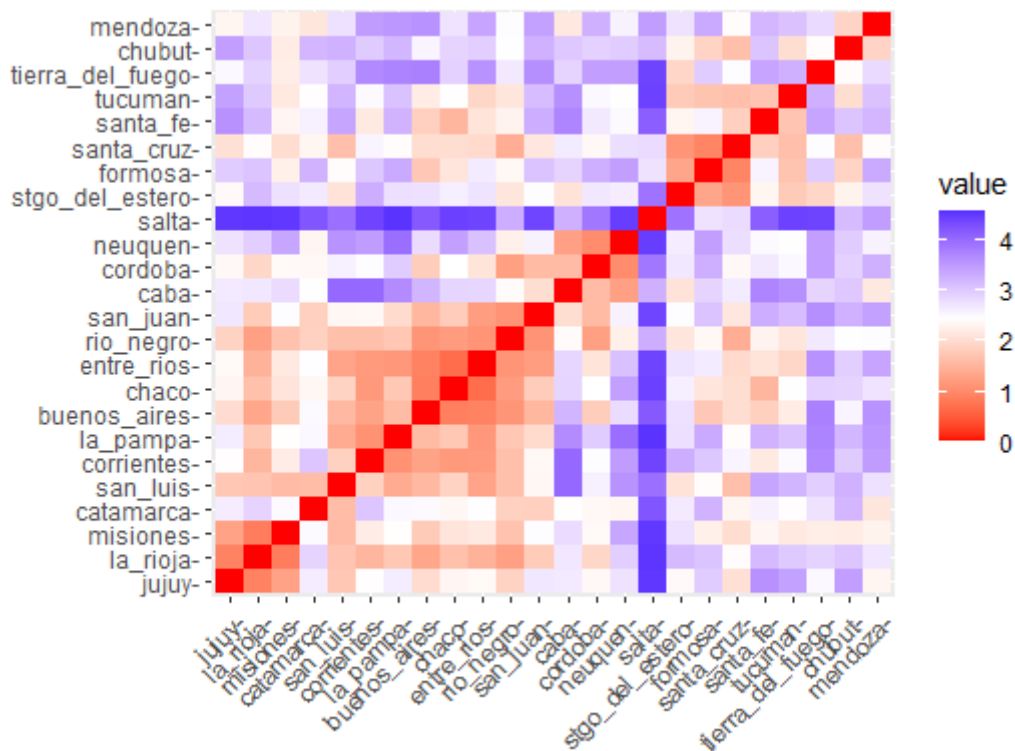


Fig. 5: Distancia máxima

La distancia de Minkowski es una generalización de las formulas para calcular la distancia Euclidea y la de Manhattan, donde le agrega un valor p , el cual si vale 1 se estara haciendo la distancia Manhattan, y si es 2 sera la de Euclidea.

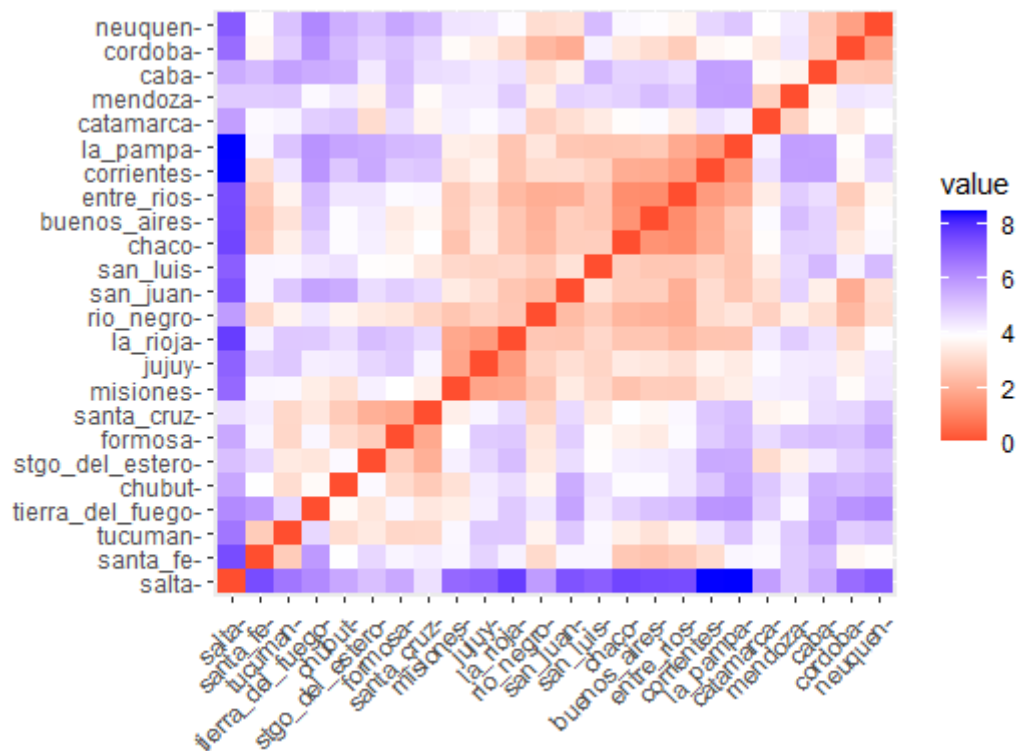


Fig. 6: Distancia Minkowski

2.4 Número Óptimo de Clusters

Antes de realizar los agrupamientos debemos hacernos una idea aproximada de cuál es la cantidad de grupos, siendo esta la que mejor representaría a los datos, para esto realizamos distintas pruebas que nos brinda la biblioteca *factoextra*.

En este primer gráfico, siguiendo la regla del pliegue del codo o rodilla podemos ver que el número óptimo se encuentra entre 4 y 6 grupos.

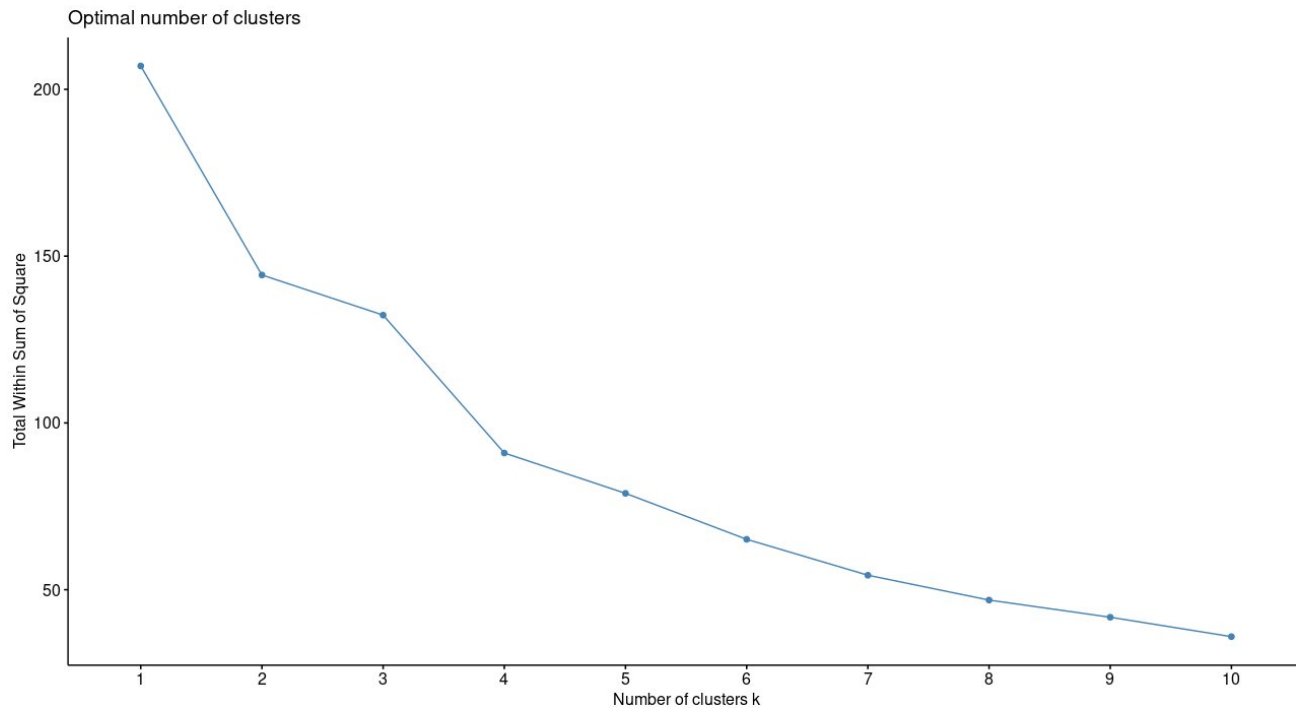


Fig. 7: Numero optimo de clusters - Metodo wss

En este otro gráfico, si bien la linea punteada se encuentra en el valor 2, podemos ver que los puntos antes de los primeros valles son el 3 y el 5, siendo estos candidatos a ser k

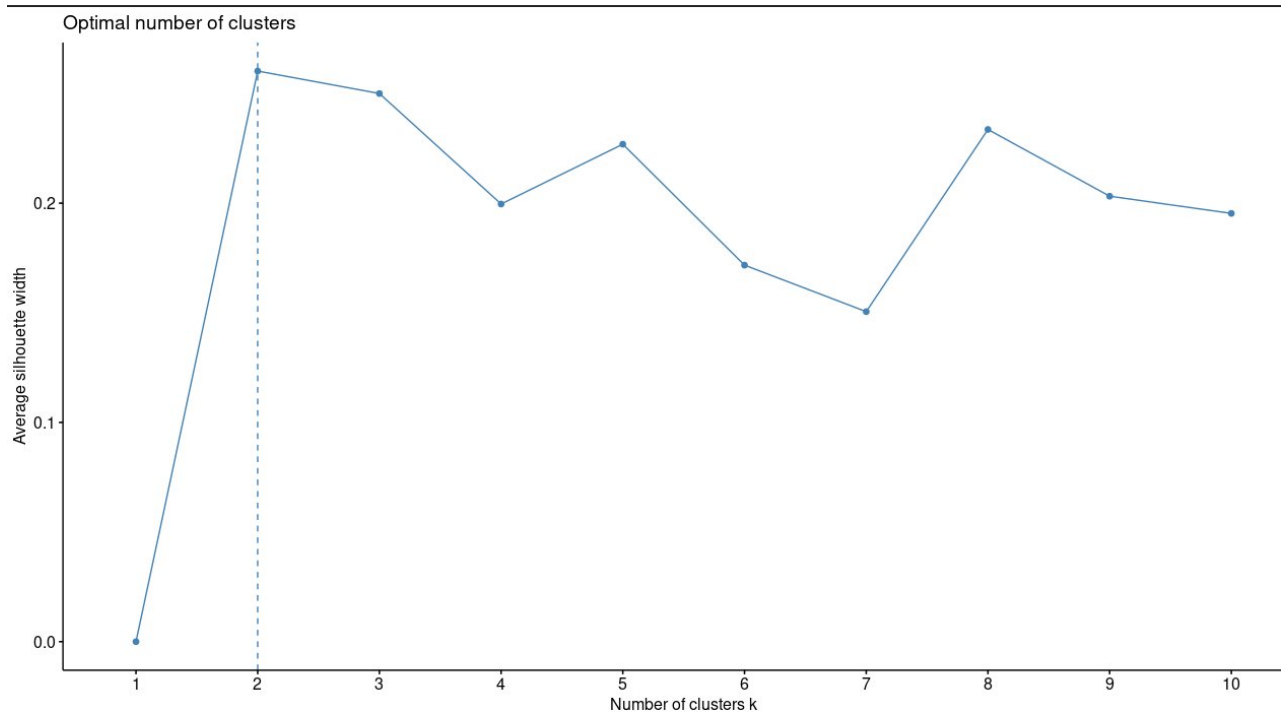


Fig. 8: Numero optimo de clusters - Metodo Silouette

A partir de los resultados obtenidos y diferentes experimentos podemos concluir

que un buen número de clusters es 5.

2.5 Agrupamiento por Provincia

Realizamos el primer agrupamiento a partir del método Jerárquico, en un principio lo visualizamos sin la división por grupos, donde pudimos ver que los grupos a realizar podrían ser 5 o 6. En esta visualización podemos ver 5 grupos.

Cluster Dendrogram

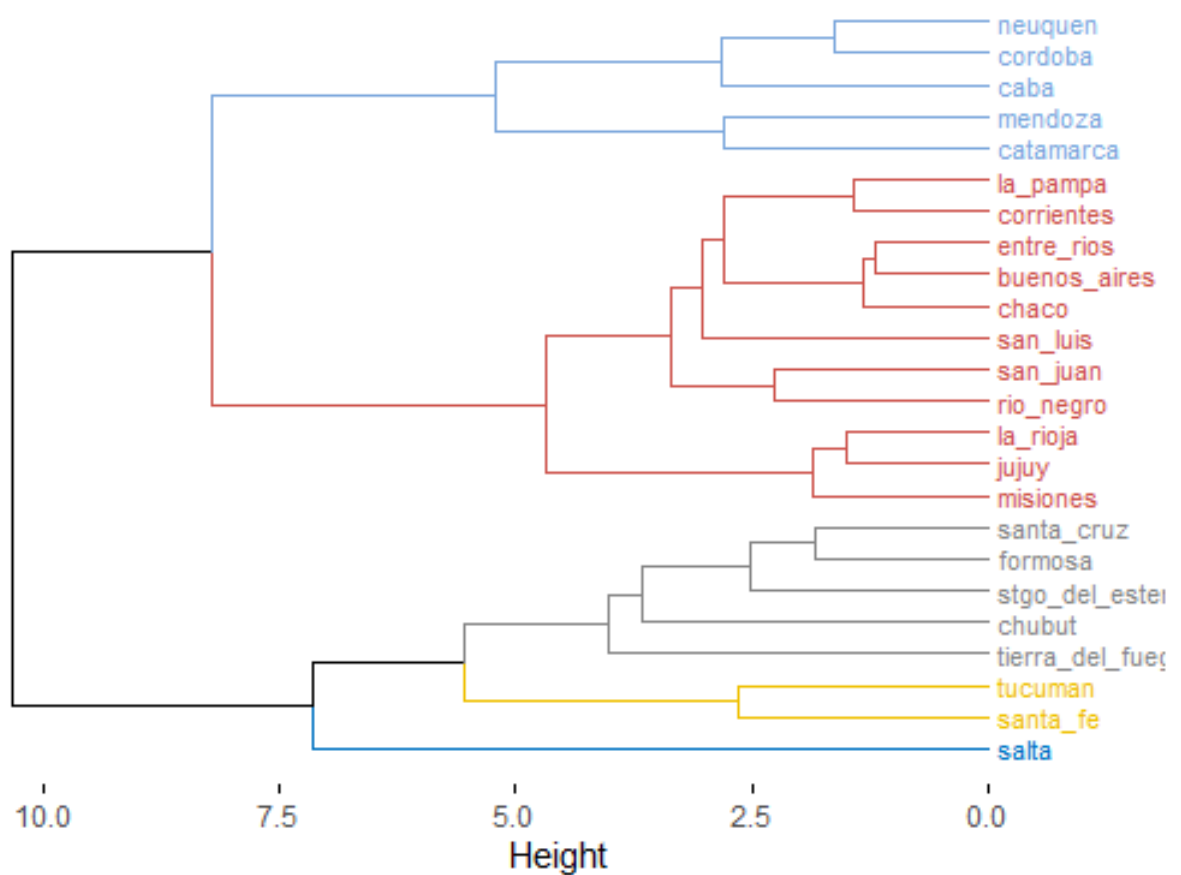


Fig. 9: Dendrograma

Siguiendo con los experimentos, realizamos un agrupamiento con el método del *k-means*, con 5 grupos.

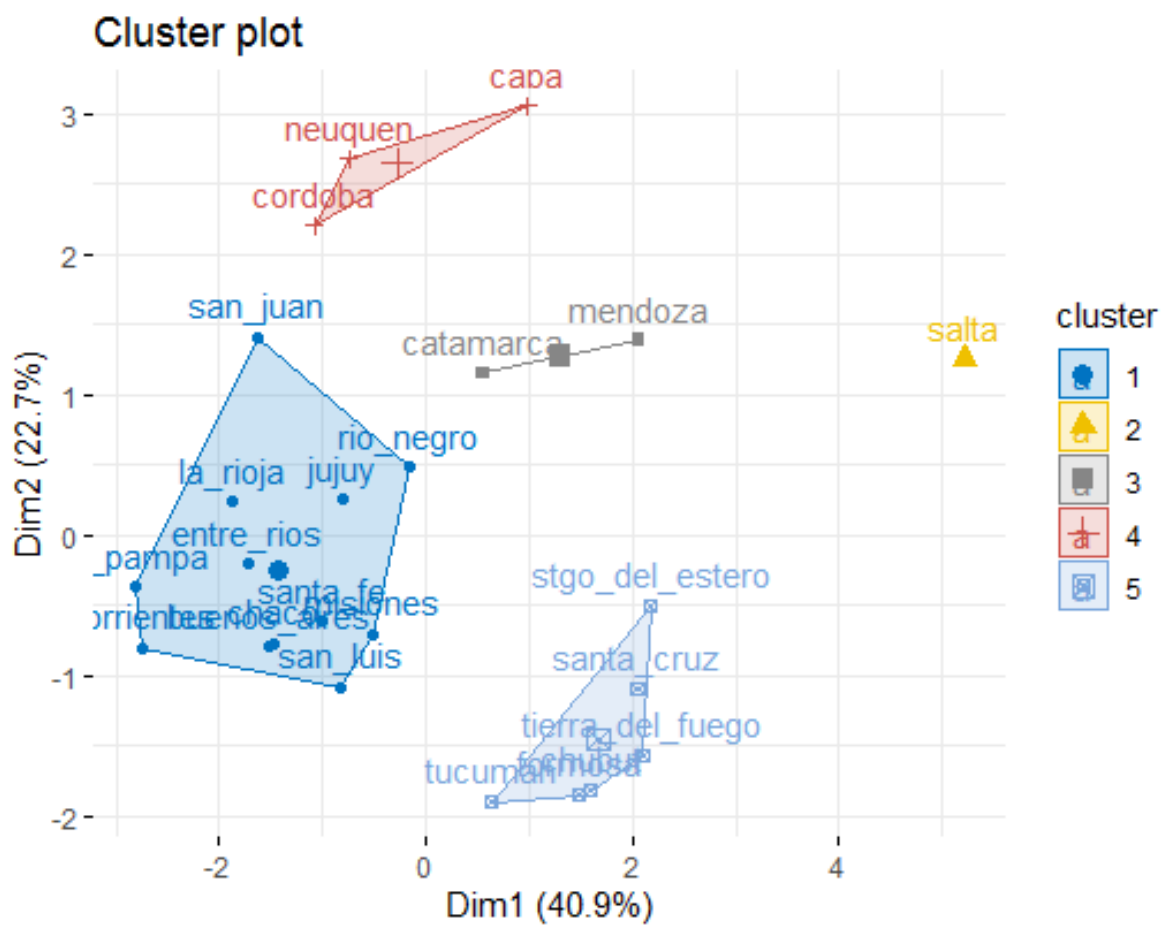


Fig. 10: K-means clustering

Y por último un experimento a partir del agrupamiento del método del *k-medoids* o PAM, el cual en vez de calcular una media selecciona un punto existente como centroide.

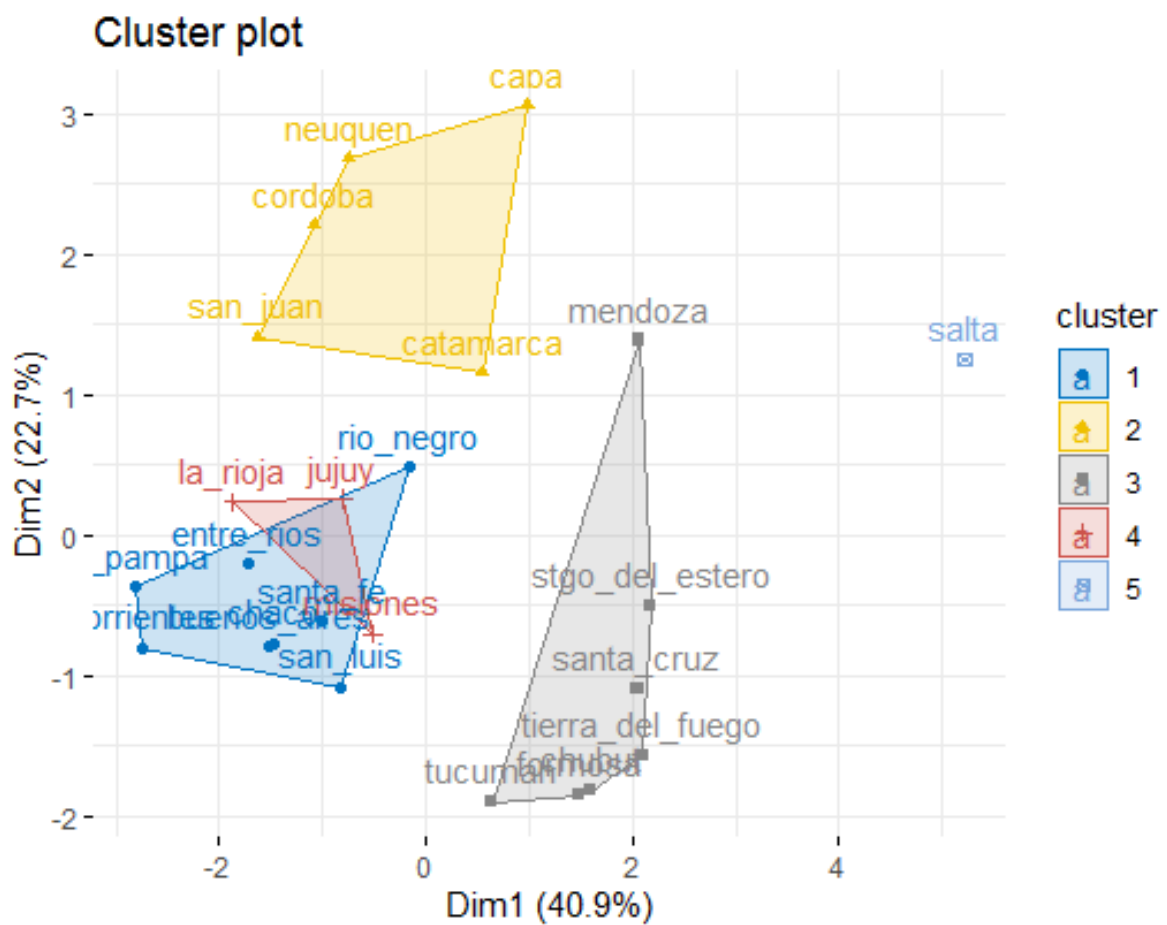


Fig. 11: Pam clustering

Como agrupamiento final realizamos un doble dendrograma con mapa de calor, donde podemos observar la interseccion entre cada provincia y cada delito, detallando su valor a partir de la paleta de colores que podemos observar a la derecha.

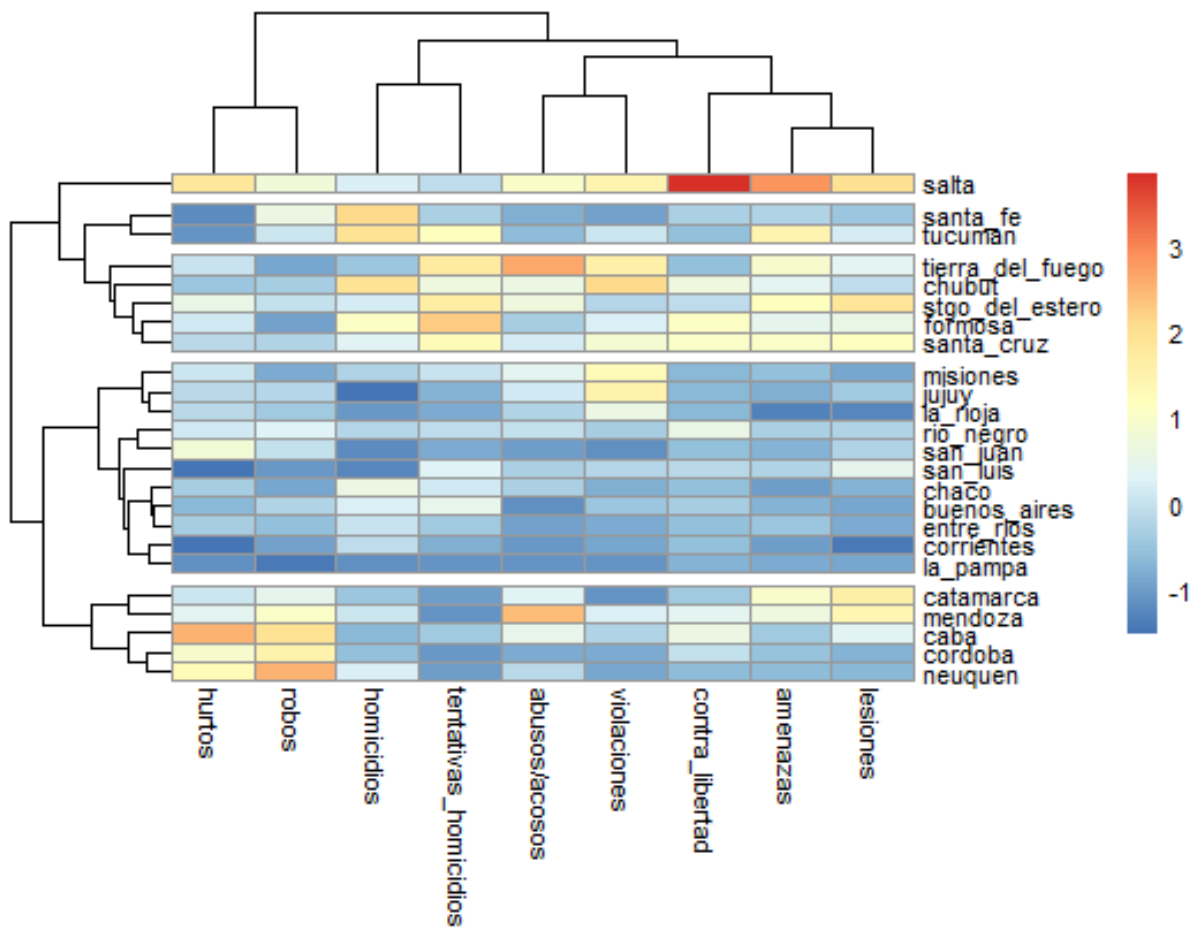


Fig. 12: Mapa de Calor

A partir de esta visualización podemos decir que hay 5 grupos bien marcados, donde las similitudes de las tasas criminales dentro de un grupo son coherentes, de manera que vemos un grupo donde abundan los robos y los hurtos, otro en el que en general hay muchos crímenes, otro en el que hay una enorme cantidad de delitos en contra de la libertad, otro en el que hay mayor número de homicidios y tentativas, y por último un grupo que posee en casi todas las tasas unos valores muy bajos.

Con el resultado del agrupamiento jerárquico hecho se ha asignado a cada provincia una clase. Permittiéndonos hacer un Análisis de Componentes Principales y visualizar como se distribuyen los puntos con el énfasis de que cada punto tendrá el color de su clase.

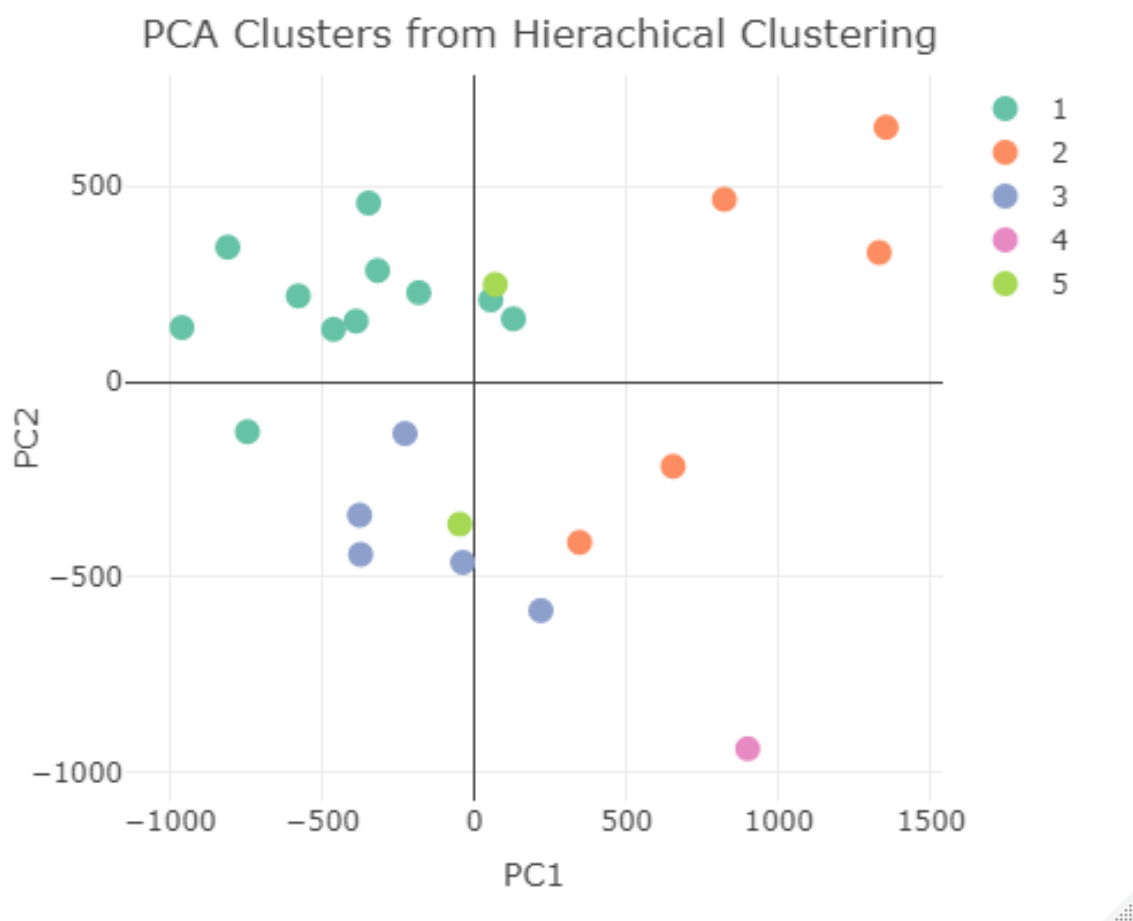


Fig. 13: PCA dinamico

Visualizamos las diferencias entre cada tasa de cada grupo.

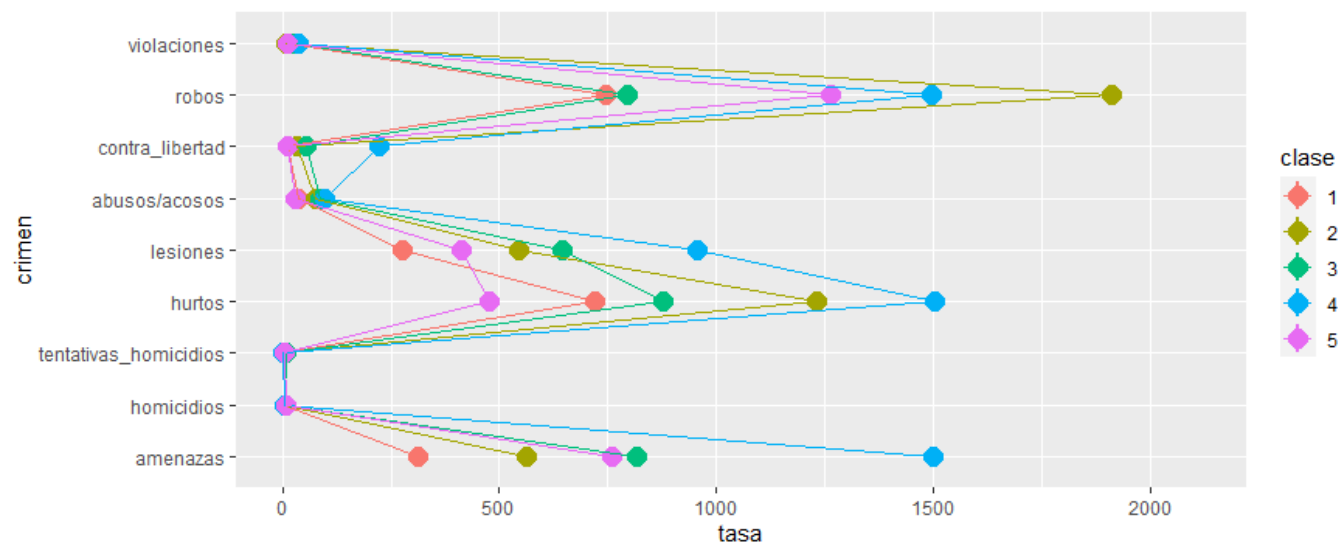


Fig. 14: Diferencias por crimen entre grupos

2.6 Agrupamiento por Crimen

Por ultimo realizamos un agrupamiento de los distintos tipos de crímenes y realizamos su dendrograma, vemos como el agrupamiento que hace es sumamente lógico y confiable.

Cluster Dendrogram

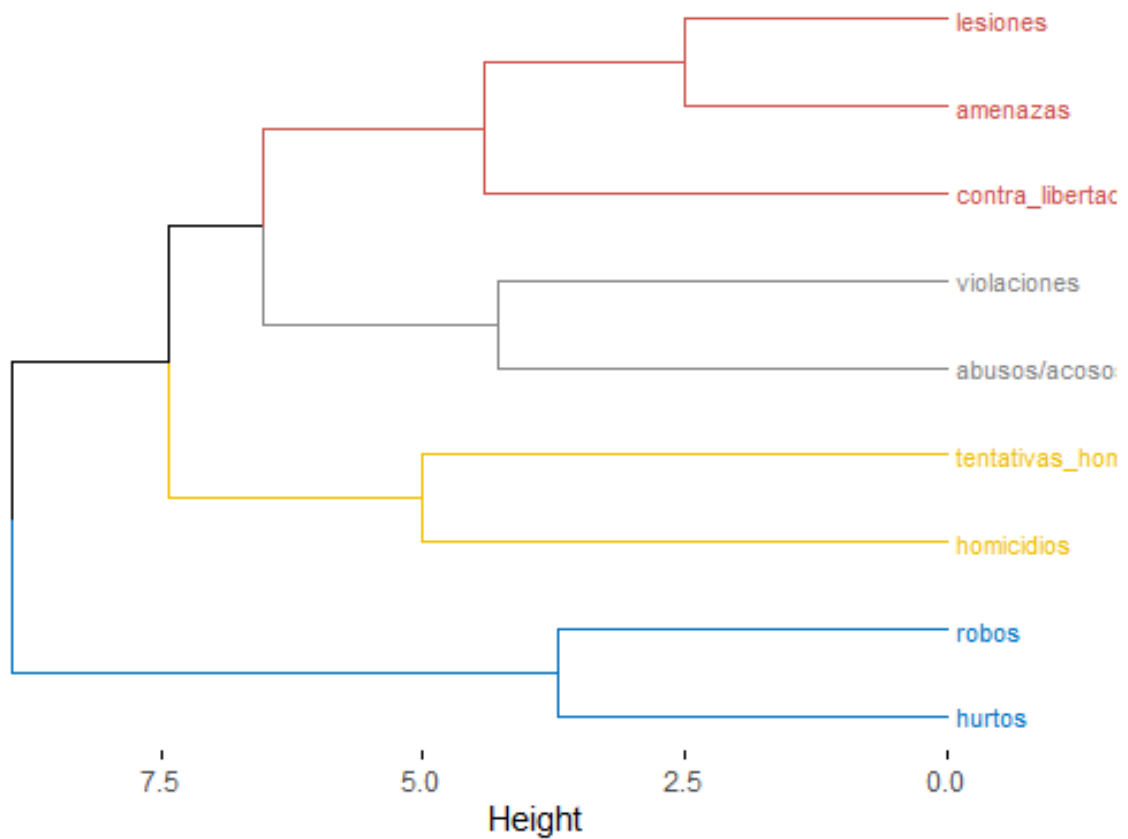


Fig. 15: Dendrograma por crímenes

3 Conclusión

4 Anexo

4.1 Código de limpieza en R

```
1 library(readr)
2 library(tidyverse)
3 library(reshape2)
4
5 #funcion para normalizar nombres de columnas
6 dbSafeNames = function(names) {
7   names = gsub('[^a-z0-9]+', '_', tolower(names))
8   names = make.names(names, unique = TRUE, allow_ = TRUE)
9   names = gsub('.', '_', names, fixed = TRUE)
10 }
11
12
13 ### Leer dataset
14 data <- read_delim("snic-provincias.csv", ";", escape_double =
15   FALSE,
16   locale = locale(encoding = "ISO-8859-1"), trim_
17   ws = TRUE)
18 colnames(data) = dbSafeNames(colnames(data))
19
20 habitantesArg <- read_csv("habitantesArg.csv", col_names = FALSE)
21 ###
22
23 ### Procesar datos
24 str(data)
25 str(habitantesArg)
26
27 # Filtrar por a o y por delito de interes
28 data <- filter(data, data$anio == 2019,
29   data$codigo_delito_snic_id %in% c(1, 2, 5, 10, 11,
30   13, 14, 15, 19))
31
32 # Quedarnos con columnas de valor(prov, delito, casos)
33 data <- data[,c(3,5,6)]
34
35 # Pasar los vos valores de nombre de delito y canditad como columna
36   -valor
37 data <- dcast(data, provincia_nombre ~ codigo_delito_snic_nombre,
38   value.var = "cantidad_hechos" )
39
40 # Unir datos de habitantes con datos criminales
41 data <- left_join(data, habitantesArg,
42   by= c("provincia_nombre" = "X1"))
43
44 # Convertir la provincia en nombre de columna
45 data <- data %>%
46   remove_rownames %>%
47   column_to_rownames(var="provincia_nombre")
48
49 # Renombrar columnas
```

```

47 colnames(data) <- c("amenazas","homicidios", "tentativas_homicidios",
48   "hurto", "lesiones",
49   "abusos/acoso", "contra_libertad", "robos", "
50   violaciones", "totHabitantes")
51 # Obtener la tasa dividiendo por los habitantes x/c 100.000
52   habitantes
53 for (i in 1:24){
54   for (j in 1:9){
55     data[i,j] = data[i,j] / data[i,10] * 100000
56   }
57 }
58 # Eliminamos el total de habitantes ya que no nos es mas util
59 data <- data[,1:9]
60 ###
61
62 ### Guardamos el archivo para leerlo en el script de clustering
63 write.csv(data, "tasasCriminales.csv")
64 ###
65
66 rm(data, habitantesArg, i, j, dbSafeNames)

```

4.2 Código de clustering en R

```

1 #####
2 # Agrupamiento de provincias a partir de estadísticas criminales
3   en Argentina
4 #
5 # Grupo B: Benitez, Garcia, Rodriguez, Rechimon
6 # fecha de creacion: 12/11/2020
7 # actualizacion: 16/11/2020 - comentarios
8 #
9 #####
10
11 ### Bibliotecas
12 library(readr)
13 library(tidyverse) # select, pipes, gather
14 library(ggplot2)
15 library(corrplot) # corrplot
16 library(psych)
17 library(plotly)
18 library(dplyr)
19 library(cluster)
20 library(factoextra)
21 library(pheatmap) # dendograma doble con mapa de calor
22 library(clValid) # comparar metodos de agrupamiento
23 ###
24
25 ### Leer data
26 # Tasa: casos cada 100k de habitantes por provincia
27 tasasCriminales <- read_csv("tasasCriminales.csv")
28 head(tasasCriminales)
29 ###

```

```
30 ### Procesamiento de datos
31 # X1 como rownames
32 tasasCriminales <- tasasCriminales %>%
33   remove_rownames %>%
34   column_to_rownames(var="X1")
35
36 # Visualizamos
37 X11()
38 boxplot(tasasCriminales,
39         names = colnames(tasasCriminales), las=2,
40         xlab = "Provincias", ylab = "Tasas")
41
42 # Correlaciones
43 corrpplot.mixed(cor(tasasCriminales), tl.pos = "lt")
44
45 # Estandarizamos las variables
46 scaled.tasas <- as.data.frame(scale(tasasCriminales))
47
48 # Analisis de Componentes Principales
49 KMO(tasasCriminales) # Mayor es mejor, va de 0 a 1, msa = medida de
   adecuacion del muestreo
50 ###
51
52 ### Clustering de Provincias
53 # Calculo de distancias
54 # metodos:
55 # "euclidean", "maximum", "manhattan", "canberra", "binary",
56 # "pearson", "spearman", "kendall", "minkowski"
57 dist.eucl <- dist(scaled.tasas, method = "euclidean")
58 fviz_dist(dist.eucl)
59
60 dist.maximum <- dist(scaled.tasas, method = "maximum")
61 fviz_dist(dist.maximum)
62
63 dist.manh <- dist(scaled.tasas, method = "manhattan")
64 fviz_dist(dist.manh)
65
66 dist.mink <- get_dist(scaled.tasas, method = "minkowski")
67 fviz_dist(dist.mink)
68
69 # Numero optimo de clusters
70 fviz_nbclust(scaled.tasas, kmeans, method = "wss") # 3 < k < 6
71 fviz_nbclust(scaled.tasas, kmeans, method = "silhouette") # 3 o 5
72
73 # Numero de clusters y metodo de agrupamiento optimo
74 comparacion <- clValid(
75   obj      = scaled.tasas,
76   nClust   = 3:6, # en promedio grupos de a 4 provincias
77   clMethods = c("hierarchial", "kmeans", "pam"),
78   validation = c("stability", "internal"))
79 summary(comparacion)
80
81 # Guardamos los grupos optimos
82 k <- 5
83
84 set.seed(124) # experimentos replicables
```

```

85
86 # Clustering jerarquico
87 # metodos:
88 # "ward.D", "ward.D2", "single", "complete", "average" (= UPGMA),
89 # "mcquitty" (= WPGMA), "median" (= WPGMC) or "centroid" (= UPGMC).
90 hc.res <- hclust(dist.eucl, method = "ward.D2")
91 fviz_dend(hc.res, cex = 0.6, k, palette = "jco", horiz = F)
92
93 # k-means clustering
94 km.res <- kmeans(scaled.tasas, k, nstart = 25)
95 fviz_cluster(km.res, data = scaled.tasas, palette = "jco",
96               ggtheme = theme_minimal())
97
98 # pam clustering
99 pam.res <- pam(x = scaled.tasas, k = k, metric = "euclid")
100 fviz_cluster(pam.res, data = scaled.tasas, palette = "jco",
101              ggtheme = theme_minimal())
102
103
104 # heatmap
105 pheatmap(mat = scaled.tasas, scale = "none",
106           clustering_distance_rows = "euclidean",
107           clustering_distance_cols = "euclidean",
108           clustering_method = "ward.D2",
109           cutree_rows = k, fontsize = 8)
110 ###
111
112 ### Guardar Datos
113 # Asignar variable de clase y de provincia
114 tasasCriminales$clase <- as.factor(cutree(hc.res, k = 5))
115 write.csv(tasasCriminales, "tasasCriminalesClasificadas.csv")
116
117 # Pasar a formato vertical
118 data_long <- gather(tasasCriminales, crimen, tasa, 1:9, factor_key=
119                     TRUE)
120
121 # Plotear grupos
122 acp <- prcomp(tasasCriminales[,1:9])
123 princomp <- as.data.frame(acp$x)
124 princomp$clase <- tasasCriminales[,10]
125
126 p <- plot_ly(x=princomp$PC1,y=princomp$PC2,text=rownames(princomp),
127              mode="markers",color = princomp$clase,marker=list(size
128                          =11))
129 p <- layout(p,title="PCA Clusters from Hierachical Clustering",
130             xaxis=list(title="PC1"),
131             yaxis=list(title="PC2"))
132
133 p
134
135 ggplot(data_long, aes(x = crimen, y = tasa, group=clase, colour =
136                       clase)) +
137   stat_summary(fun = mean, geom="pointrange", size = 1)+
138   stat_summary(geom="line")
139 ###

```

```
138
139 ### Clustering de Crimenos
140 # Transpuesta para agrupar por crimenos
141 transpuesta <- as.data.frame(t(scaled.tasas))
142
143 # Calculo de distancias
144 dist.eucl2 <- dist(transpuesta, method = "euclidean")
145
146 # Guardamos los grupos optimos
147 k2 <- 4
148
149 set.seed(200) # experimentos replicables
150
151 # Clustering jerarquico
152 hc.res2 <- hclust(dist.eucl2, method = "ward.D2")
153 fviz_dend(hc.res2, cex = 0.6, k2, palette = "jco", horiz = T)
154
155 # k-means clustering
156 km.res2 <- kmeans(transpuesta, k2, nstart = 25)
157 fviz_cluster(km.res2, data = transpuesta, palette = "jco",
158               ggtheme = theme_minimal())
159
160 # pam clustering
161 pam.res2 <- pam(x = transpuesta, k = k2, metric = "euclid")
162 fviz_cluster(pam.res2, data = transpuesta, palette = "jco",
163               ggtheme = theme_minimal())
164 ###
```