

# Técnicas de análisis Multivariado

**Alumnos:***Huertas Quispe, Anthony Enrique**Torres Salinas, Karina Hesi**Córdova Proleón, Christian Therius***Cod:** 20173728**Cod:** 20164111**Cod:** 20173970**Semestre:** 2017-II**Tema:** Lista 2

PROF. ENVER TARAZONA



Pontificia Universidad Católica del Perú  
Escuela de Posgrado  
Maestría en Estadística

## Pregunta 1 (4 puntos)

El conjunto de datos del archivo `wisc.bc_data.csv` contiene información para 569 instancias relacionadas a biopsias por sospecha de cáncer de mama, cada una de ellas conteniendo 32 atributos. Uno de los atributos es un número de identificación, otro es el diagnóstico de cáncer y los 30 atributos restantes están relacionados a mediciones numéricas realizadas en laboratorio. El diagnóstico fue codificado como M que indica que es maligno o B que indica que es benigno. El objetivo del análisis es predecir si un paciente tiene cáncer basado en los resultados del laboratorio. Compare varios modelos de clasificación (logístico, LDA, QDA y RDA) y encuentre aquel que tenga mejor desempeño para el conjunto de datos de entrenamiento. Utilice diversos criterios de comparación la exactitud total, índice Kappa de Cohen, la sensibilidad, especificidad y el AUC (área bajo la curva ROC). Presente una gráfica comparativa que muestre las curvas ROC para los clasificadores usados.

**Desarrollo.** Se estableció en Anexo 1, como categoría de referencia al diagnóstico de cáncer “maligno”. La base de datos “sf wisc.bc\_data.csv” que cuenta con 569 registros se dividió en 2 base de datos. La primera parte tendría el 80 % de la base de datos original, esto nos servirá como entrenamiento para la construcción del modelo; mientras que el otro 20 % de la base de datos original la usaremos para la validación del modelo.

Indicadores	Logístico	LDA	QDA	RDA
Exactitud total	0.930	0.947	0.965	0.956
Kappa	0.854	0.889	0.927	0.908
Sensibilidad	0.913	0.976	0.957	0.977
Especificidad	0.941	0.931	0.971	0.944
AUC	0.927	0.938	0.964	0.949

Del cuadro resumen, se tiene que el modelo con el análisis discriminante cuadrático es el que tiene una mejor precisión o tasa de aciertos (96.5 %), mientras que con el modelo logístico se consigue la menor precisión (93.0 %). En cuanto al índice Kappa de Cohen, se tiene una mejor concordancia (casi perfecta) utilizando el modelo discriminante cuadrático en comparación a los otros tres modelos, pues se tiene un índice de 0.927.

La sensibilidad con mayor valor se consigue usando el modelo discriminante lineal (0.976) y el modelo discriminante regularizado (0.977). La especificidad más grande se obtendría usando el modelo discriminante cuadrático (0.971), lo que significa que el 97.1 % son correctamente diagnosticados como cáncer “benigno”. La curva ROC que presenta una mayor área es también utilizando el modelo discriminante cuadrático (0.964).

En conclusión, el mejor modelo que presenta un mejor desempeño con la base de datos para la evaluación del modelo, sería empleando el discriminante cuadrático.

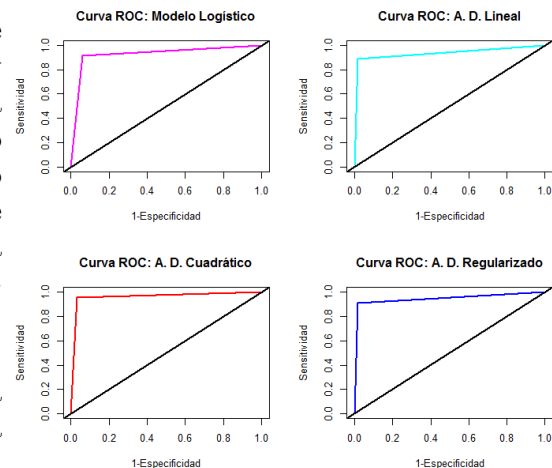


Figura 1: Curvas ROC

## Pregunta 2 (4 puntos)

Wolford y Hollingsworth (1974) estaban interesados en las confusiones que se producen cuando una persona intenta identificar las letras del abecedario visualizadas solo durante algunos milisegundos. Se construyó una matriz de confusión que muestra la frecuencia con la que cada letra de estímulo se llamó erróneamente por otra. Una sección de esta matriz se muestra en la tabla a continuación.

Letter	C	D	G	H	M	N	Q	W
C	-							
D	5	-						
G	12	2	-					
H	2	4	3	-				
M	2	3	2	19	-			
N	2	4	1	18	16	-		
Q	9	20	9	1	2	8	-	
W	1	5	2	5	18	13	4	-

Realice un análisis de Escalamiento Multidimensional e interprete los resultados. ¿Qué se puede concluir en relación a la confusión entre las letras?

**Desarrollo.** Para ver que podemos concluir en relación a la confusión entre las letras, hay que encontrar una matriz de similitud entre los elementos a partir de la matriz de confusión, para luego construir una matriz de distancias de disimilitudes  $d_{ij}$ .

Entonces de las similitudes  $s_{ij}$ , escogemos una similitud máxima  $c \geq \max\{s_{ij}\}$ , tal que la matriz de distancias de similitudes tendría los elementos:

$$d_{ij} = \begin{cases} c - s_{ij} & \text{si } i \neq j, \\ 0 & \text{si } i = j \end{cases}$$

Dado que las disimilitudes dependen del valor de  $c$  elegido el escalamiento multidimensional no métrico sería el más apropiado. Sea la matriz de confusión:

Listing 1: Matriz de confusión.

---

```
1 letras=as.dist(as.matrix(letras))
2 letras
```

---

```

      C  D  G  H  M  N  Q
D  5
G 12  2
H  2  4  3
M  2  3  2 19
N  2  4  1 18 16
Q  9 20  9  1  2  8
W  1  5  2  5 18 13  4
```

Figura 2: Matriz de confusión

Escogemos  $c = 21$  que es el máximo de los  $s_{ij} + 1$ , entonces tenemos nuestra matriz de distancias de disimilitudes:

Listing 2: Matriz de distancias de disimilaridades.

---

```
1 letras=21-letras
2 letras
```

---

```

      C  D  G  H  M  N  Q
D 16
G  9 19
H 19 17 18
M 19 18 19  2
N 19 17 20  3  5
Q 12  1 12 20 19 13
W 20 16 19 16  3  8 17

```

Figura 3: Matriz de distancias de disimilaridades

Usamos la función `smacofsym` porque estamos usando una matriz simétrica, el número de dimensiones en el cual queremos representar la distancia entre las letras es 2, entonces de los resultados podemos observar que hemos necesitado 30 iteraciones para generar la estimación final. La medida de estrés es de 0.06, en general podríamos decir que tenemos una representación buena y eso es adecuado.

Listing 3: Matriz de distancias de disimilaridades.

---

```
1 res=smacofSym(letras,ndim=2, type = "ordinal")
2 res
3 summary(res)
```

---

<pre> Call: smacofsym(delta = letras, ndim = 2, type = "ordinal")  Model: Symmetric SMACOF Number of objects: 8 Stress-1 value: 0.06 Number of iterations: 30 </pre>	<pre> Configurations:       D1      D2 C -0.7234 -0.2231 D -0.3098  0.5664 G -0.6208 -0.5749 H  0.5748 -0.3876 M  0.6656 -0.0829 N  0.4167  0.0494 Q -0.4811  0.2942 W  0.4781  0.3584 </pre>
<pre> Stress per point (in %):       C      D      G      H      M      N      Q      W 2.50  5.76 15.36 15.38 12.94 23.24  7.55 17.27 </pre>	

Figura 4: Cuadro 1.

Del cuadro anterior observamos que los puntos de estrés son altos, siendo la letra N la que está peor representada y la que tiene menor estrés es la letra C aunque no es cercano a cero.

Evaluando los resultados de manera más general hemos visto que la medida de estrés es regularmente baja con el cuál podríamos tener resultados con el que podríamos estar satisfechos.

Del siguiente mapa perceptual observamos que hay cercanía entre las letras D y Q, así como las letras C y G, otro grupo que podría presentar ciertas características similares son las letras W, N, M y H. La dimensión 1 si bien no tiene una explicación ayuda a diferenciar un poco más claro los grupos de letras (W, N, M y H), (D, Q) y (C, G).

Para sustentar que los resultados sean coherentes observamos la siguiente gráfica de Shepard, donde la idea es observar una función que sea no decreciente, observamos la gráfica vemos que

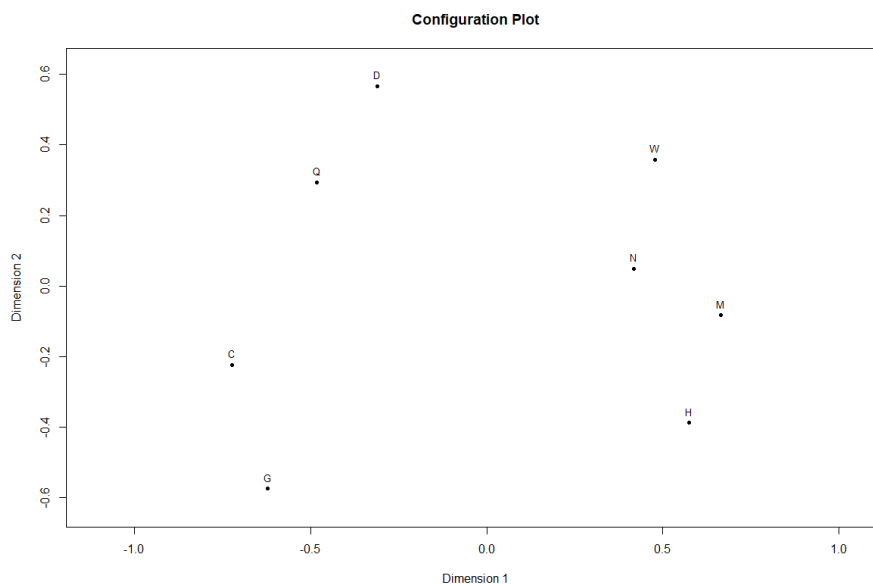


Figura 5: Configuración.

tenemos una función no decreciente y eso nos ayuda a poder evaluar que estamos teniendo resultados adecuados para representar a las letras en 2 dimensiones.

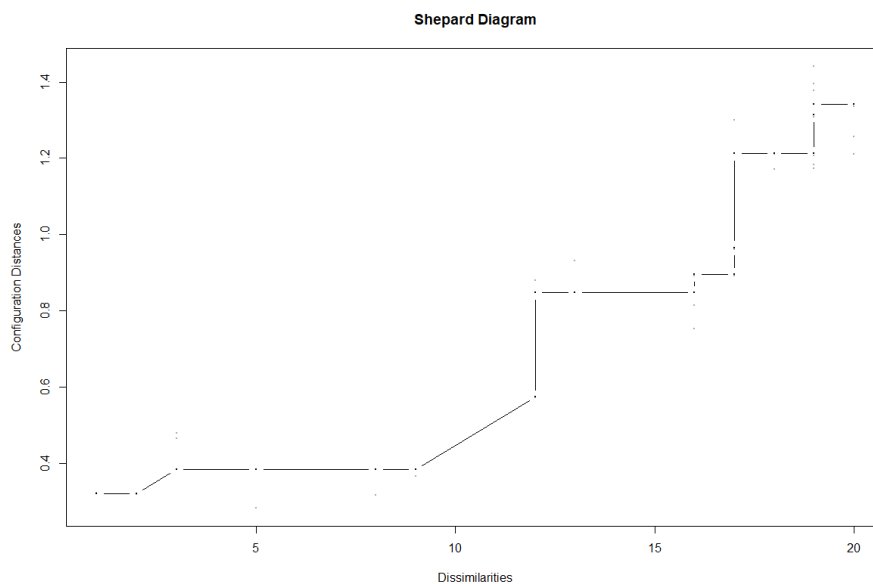


Figura 6: Diagrama Shepard.

En el siguiente gráfico de residuos, esperamos que haya cierta relación lineal, del gráfico se observa que hay una tendencia lineal entre las disimilaridades y las distancias que tenemos, asimismo no se observan outliers.

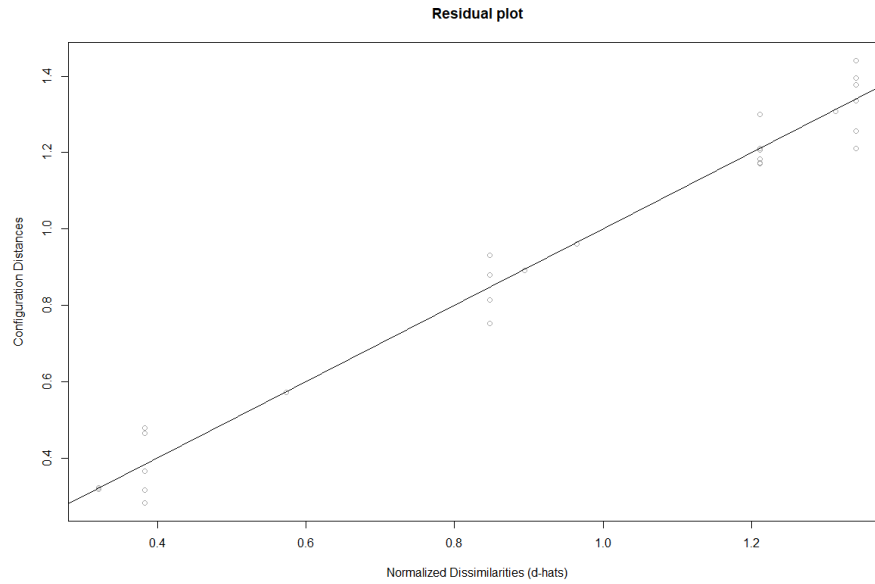


Figura 7: Gráfica de residuales.

Para evaluar que tan asociadas están esas 2 medidas, se sacó la correlación al cuadrado saliendo **0.9758**, un valor muy cercano a 1.

En el siguiente gráfico podemos observar el estrés individual para cada uno de las letras, ya habíamos observado antes que el que se encuentra mejor representado es la letra C y la letra N es el que tiene mayor medida de estrés en comparación a todos los demás.

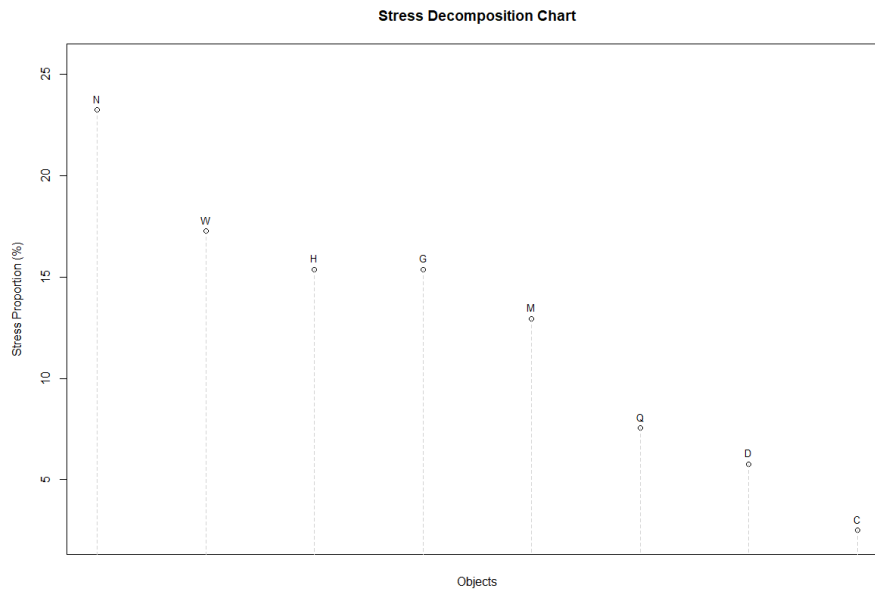


Figura 8: Gráfico de Estrés.

Con el gráfico de burbuja tenemos la representación del mapa perceptual y en simultáneo cuál es la medida de estrés individual que tenemos para cada una de las letras:

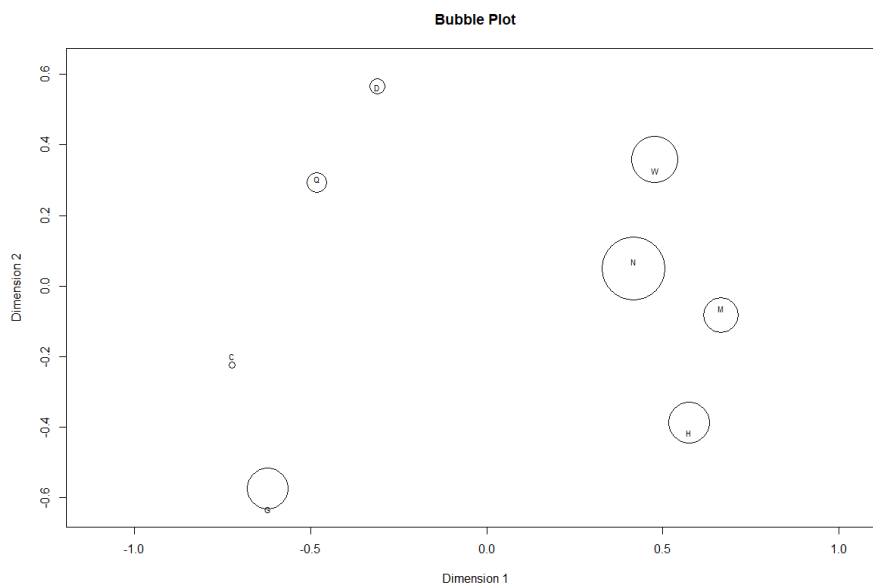


Figura 9: Gráfico de burbuja.

Usando un agrupamiento jerárquico aglomerativo con enlace promedio, se forman los siguientes cluster: (W, H, M, N) , (D, Q) y (C, G).

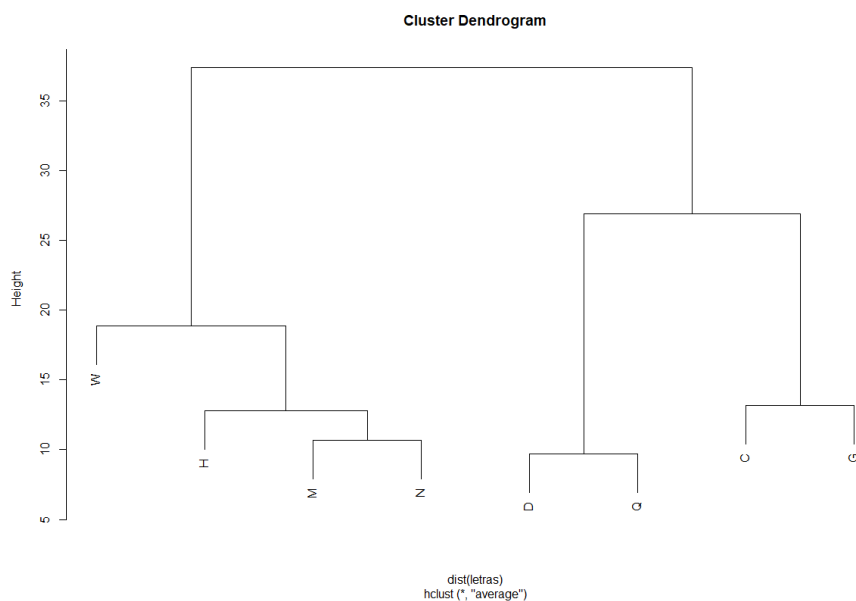


Figura 10: Dendograma.

Lo que se concluye es que dado que las personas tienen poco tiempo para visualizar las letras del abecedario, tienden a confundirlas de acuerdo a ciertas características y formas de éstas, por ejemplo (D y Q) son letras cerradas y parecieran redondas, las letras (C y G) son semi cerradas y el último grupo (W, H, M y N) tienen trazos marcados tanto verticales como diagonales unidas por ángulos.

### Pregunta 3 (8 puntos)

Aplique al menos 2 de las metodologías presentadas en el curso (o una distinta) a un conjunto de datos reales.

**Desarrollo.** Se hará uso de una base de datos original desarrollada por *School of Information and Computer Science, University of California*. Esta base de datos cuenta con seis variables biomecánicas (explicativas) medidas en pacientes ortópedicos con el fin de categorizar si un paciente presenta alguna malformación (anormal) o no la presente (normal); :

- incid.pelvica: Incidencia pélvica.
- incl.pelvica: Inclinação pélvica.
- angulo.lumbar: Ángulo de lordosis lumbar.
- pend.sacral: Ángulo lumbar.
- rad.pelvico: Radio pélvico.
- grado.spond: Grado de spondylolisthesis.
- clase: Anormal y Normal.
- subclase: Anormal (Hernia, Spondylolisthesis) y Normal.

### Clustering

El objetivo será agrupar los datos mediante un lenguaje no supervisado, de tal forma que lo obtenido infiera los grupos que originalmente son presentados es decir Normal y Anormal.

Primero diseñaremos una primera base haciendo uso de todos las variables, adicionando otras más, para luego tomar el respectivo subconjunto de variables con las cuales trabajaremos en esta metodología.

Listing 4: Base de datos.

---

```
1 head(datos)
```

---

	incid.pelvica	incl.pelvica	angulo.lumbar	pend.sacral	rad.pelvico	grado.spond	clase
1	63.02782	22.552586	39.60912	40.47523	98.67292	-0.254400	Abnormal
2	39.05695	10.060991	25.01538	28.99596	114.40543	4.564259	Abnormal
3	68.83202	22.218482	50.09219	46.61354	105.98514	-3.530317	Abnormal
4	69.29701	24.652878	44.31124	44.64413	101.86850	11.211523	Abnormal
5	49.71286	9.652075	28.31741	40.06078	108.16872	7.918501	Abnormal
6	40.25020	13.921907	25.12495	26.32829	130.32787	2.230652	Abnormal
	subclase	clase.grupo	subclase.grupo				
1	Hernia	1	1				
2	Hernia	1	1				
3	Hernia	1	1				
4	Hernia	1	1				
5	Hernia	1	1				
6	Hernia	1	1				

Figura 11: Datos Generales; Clase.grupo: 2 = Normal, 1 = Anormal; Clase.grupo: 2 = Normal, 1 = Anormal Hernia, 3 = Anormal Spondylolisthesis



Se desarrollará un resumen descriptivo de las variables explicativas, y un análisis de caja usando los grupos de clase indicados en la base.

incid.pelvica	incl.pelvica	angulo.lumbar	pend.sacral	rad.pelvico
Min. : 26.15	Min. : -6.555	Min. : 14.00	Min. : 13.37	Min. : 70.08
1st Qu.: 46.43	1st Qu.: 10.667	1st Qu.: 37.00	1st Qu.: 33.35	1st Qu.: 110.71
Median : 58.69	Median : 16.358	Median : 49.56	Median : 42.40	Median : 118.27
Mean : 60.50	Mean : 17.543	Mean : 51.93	Mean : 42.95	Mean : 117.92
3rd Qu.: 72.88	3rd Qu.: 22.120	3rd Qu.: 63.00	3rd Qu.: 52.70	3rd Qu.: 125.47
Max. : 129.83	Max. : 49.432	Max. : 125.74	Max. : 121.43	Max. : 163.07
grado.spond	clase			
Min. : -11.058	Abnormal: 210			
1st Qu.: 1.604	Normal : 100			
Median : 11.768				
Mean : 26.297				
3rd Qu.: 41.287				
Max. : 418.543				

Figura 12: Resumen descriptivo

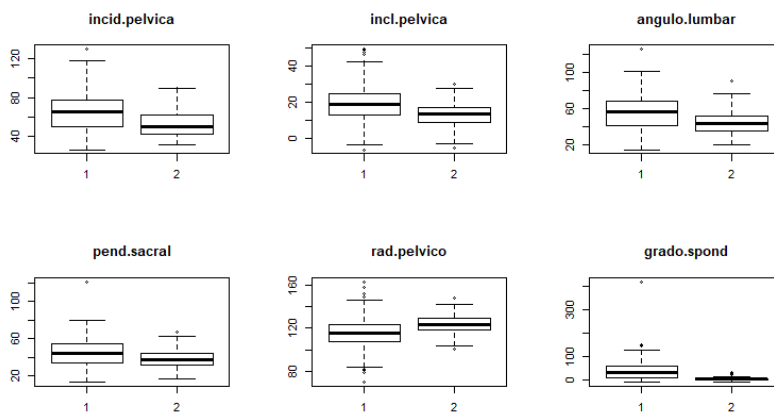


Figura 13: Diagrama de Caja - Grupos de clase originales

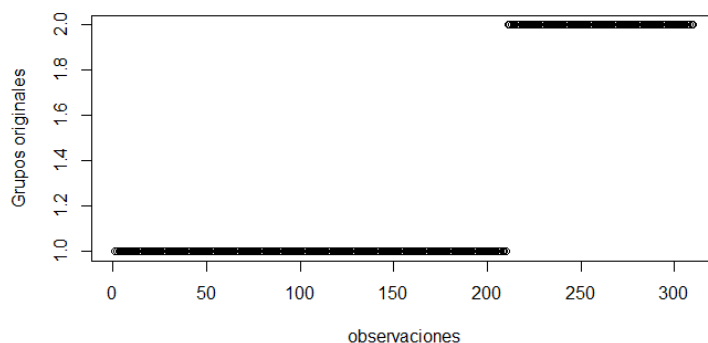


Figura 14: Observaciones vs Grupo de Clase

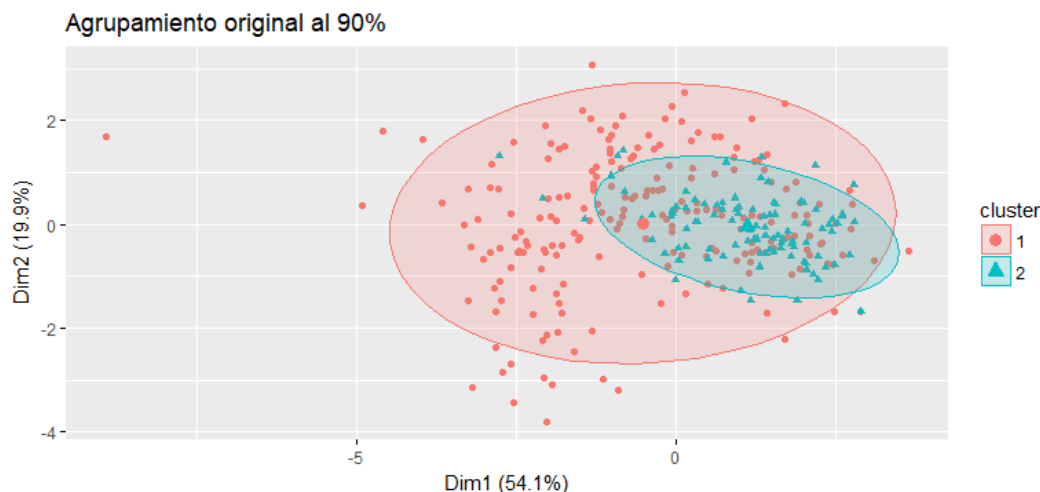


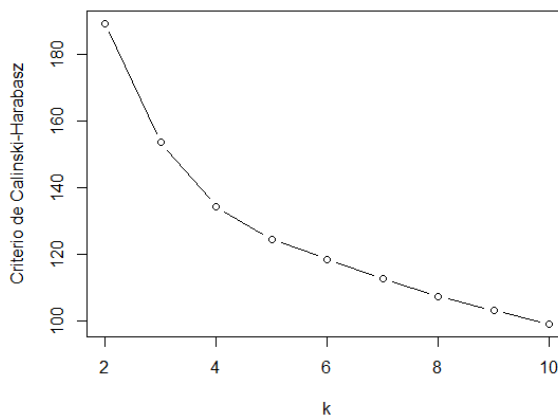
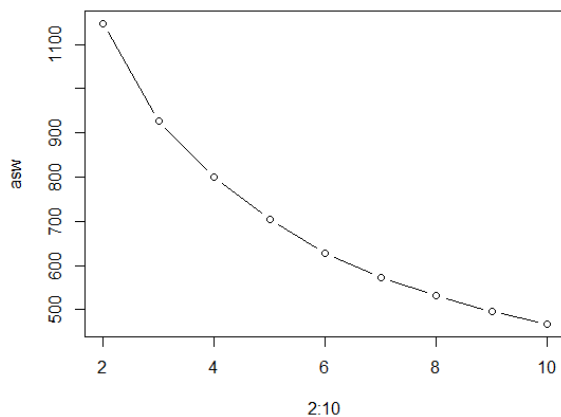
Figura 15: Agrupamiento por clases originales en un 90 %

Los gráficos 13, 14 y 15 nos representan los datos respecto a sus grupos originales, donde 1 = Anormal y 2 = Normal. El fin del uso de esta metodología es modelar los datos para que se agrupen mediante la información procedente de sus variables explicativas sin variables respuesta; y observar que efectivamente el ajuste por clustering nos lleva a un aproximado a la realidad. Esto con el objetivo de poder usar un modelo para agrupar datos nuevos.

Como hemos observado solo estamos analizando sobre la clase y no la subclase, es decir con 2 agrupamientos. Analicemos si en efecto, tiene sentido el uso de dos conglomerados mediante distintos métodos.

Listing 5: kmeans.

```
1 SC.cluster(datos1,1)
2 silueta(datos1,1)
3 kmeansruns(scale(datos1),criterion="asw")
4 Cal.Har(datos1,1)
5 kmeansruns(scale(datos1),criterion="ch")
```



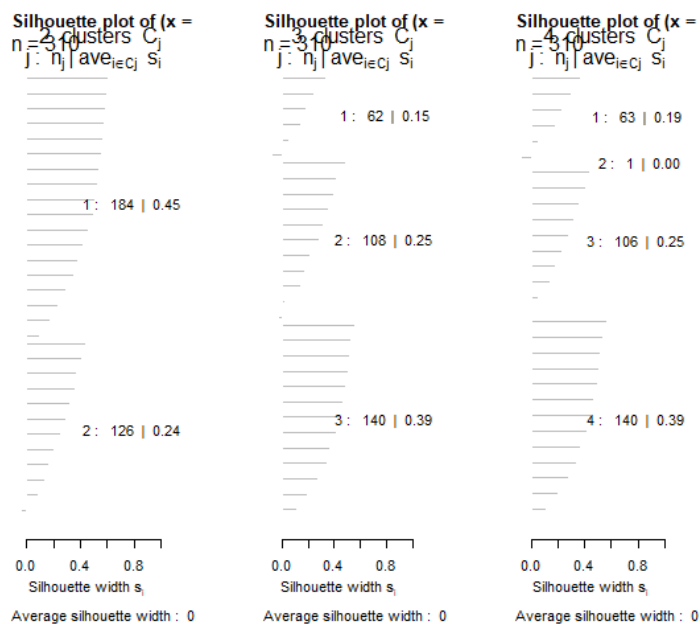


Figura 16: kmeans - Silueta

Listing 6: PAM.

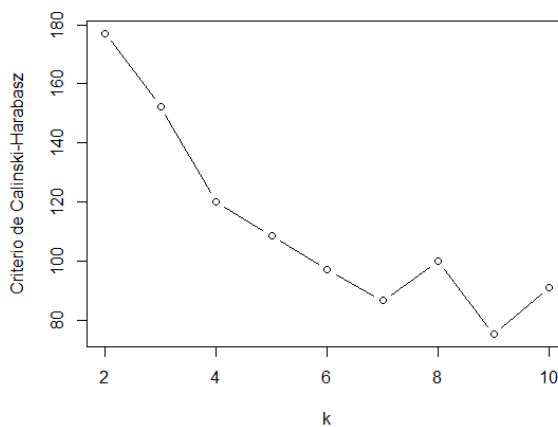
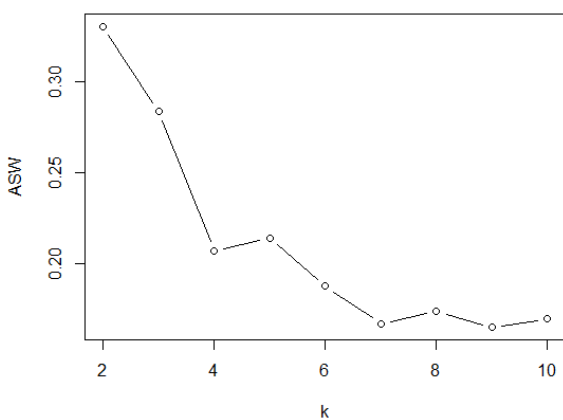
---

```

1 SC.cluster(datos1,2)
2 silueta(datos1,2)
3 pamk(scale(datos1),criterion="asw")
4 Cal.Har(datos1,2)
5 pamk(scale(datos1),criterion="ch")

```

---



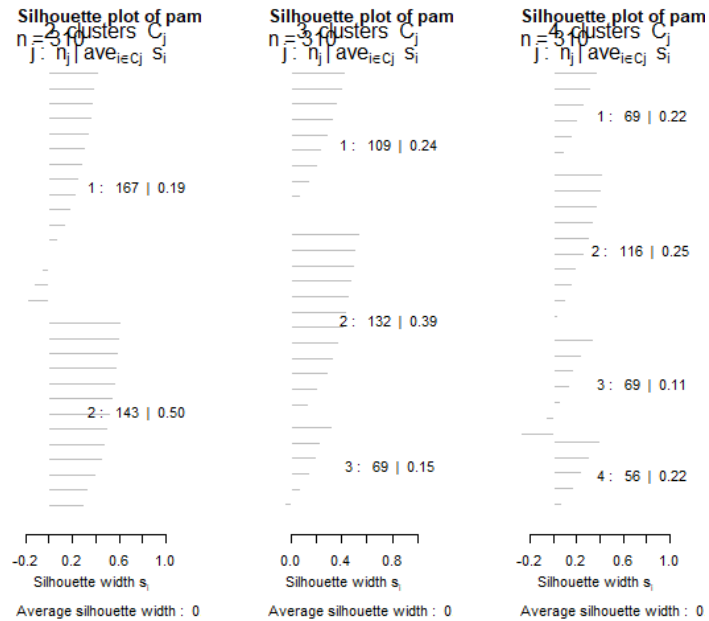


Figura 17: PAM - Silueta

Listing 7: Clara.

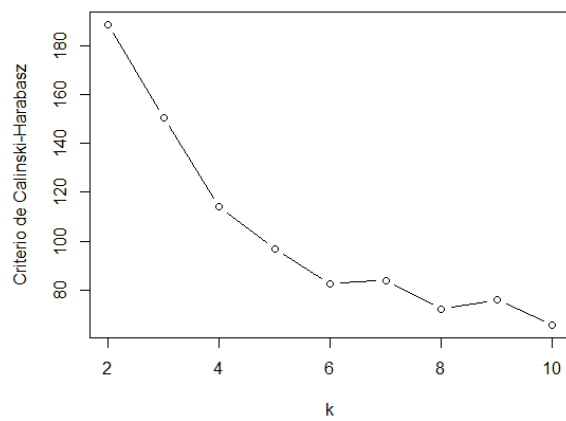
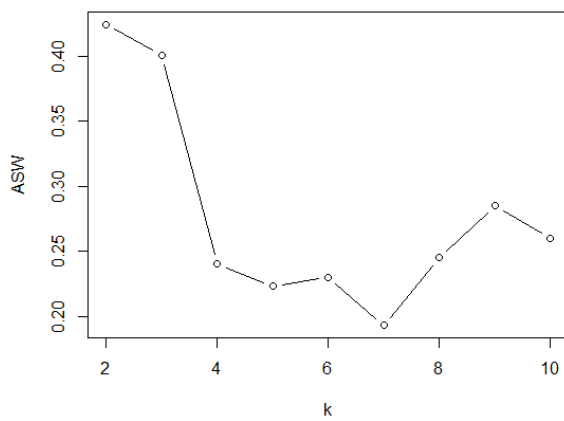
---

```

1 SC.cluster(datos1,3)
2 silueta(datos1,3)
3 pamk(scale(datos1),criterion="asw",usepam=FALSE)
4 Cal.Har(datos1,3)
5 pamk(scale(datos1),criterion="ch",usepam=FALSE)

```

---



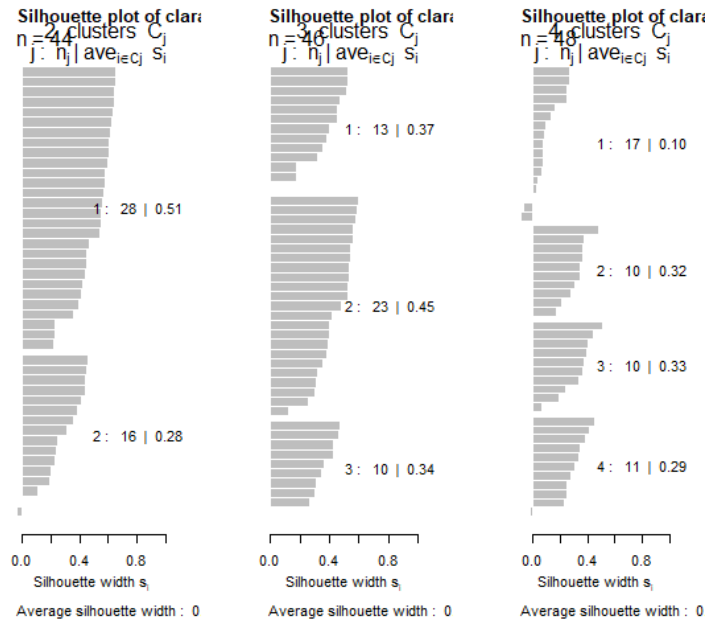


Figura 18: PAM - Silueta

Listing 8: Fanny.

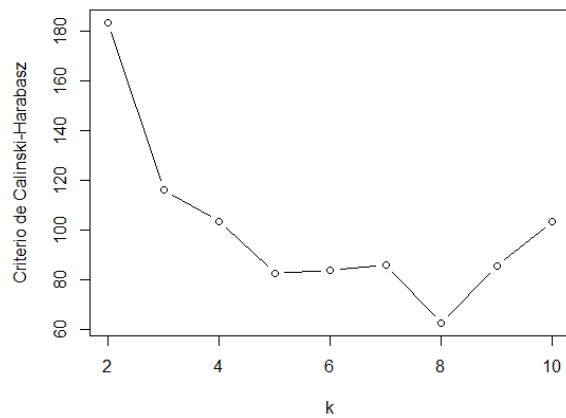
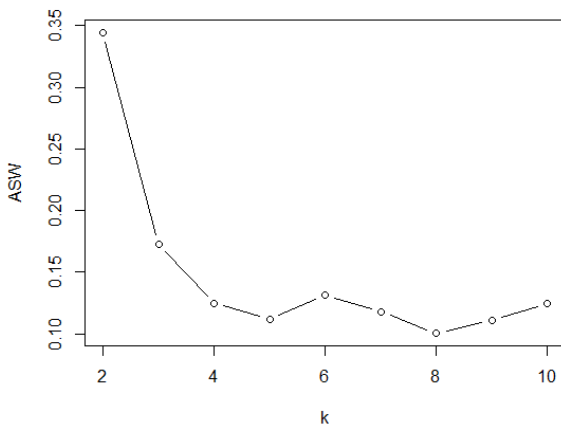
---

```

1 SC.cluster(datos1,4)
2 silueta(datos1,4)
3 Cal.Har(datos1,4)

```

---



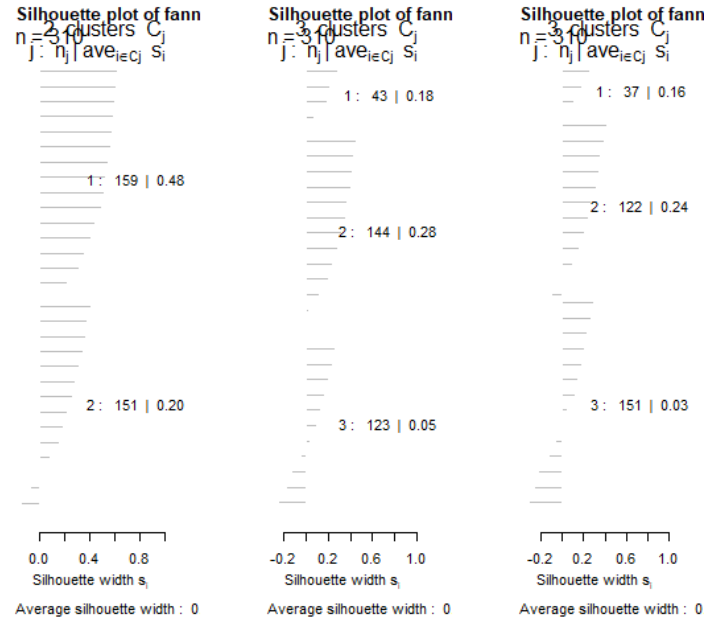


Figura 19: PAM - Silueta

Después de haber obtenido todos los gráficos, en efecto parece ser eficiente el uso de dos conglomerados, tanto por el índice de mayor valor generado en la silueta de cluster 2 como por los criterios de suma de cuadrados y Criterio de Calinski-Harabasz.

A continuación, Se observarán como se han agrupado los datos mediante los métodos.

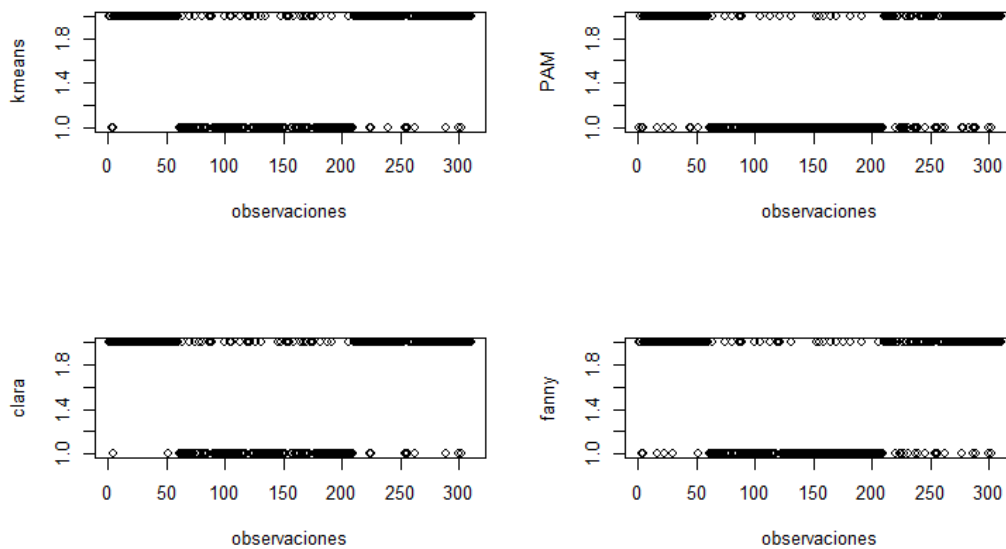


Figura 20: Agrupamiento

Como se observa, en relación al agrupamiento real (figura 14), la tendencia parece darse. Usaremos a continuación métodos de cluster jerárquicos diseñados por arboles de decisión.

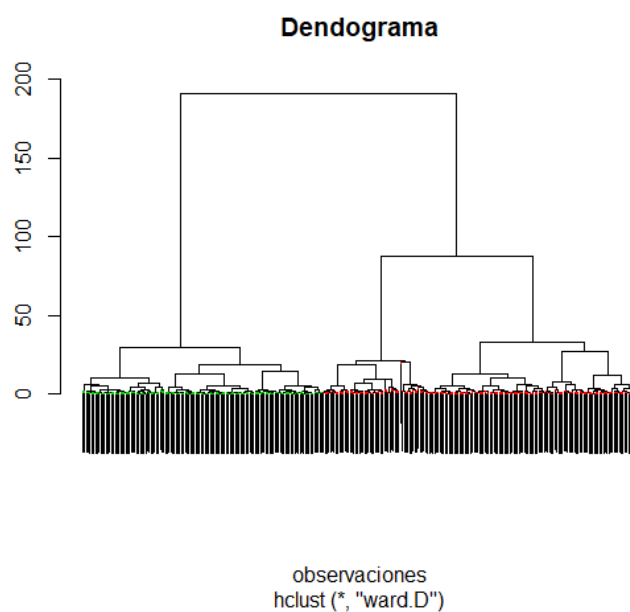


Figura 21: Dendograma HIERARCHICAL

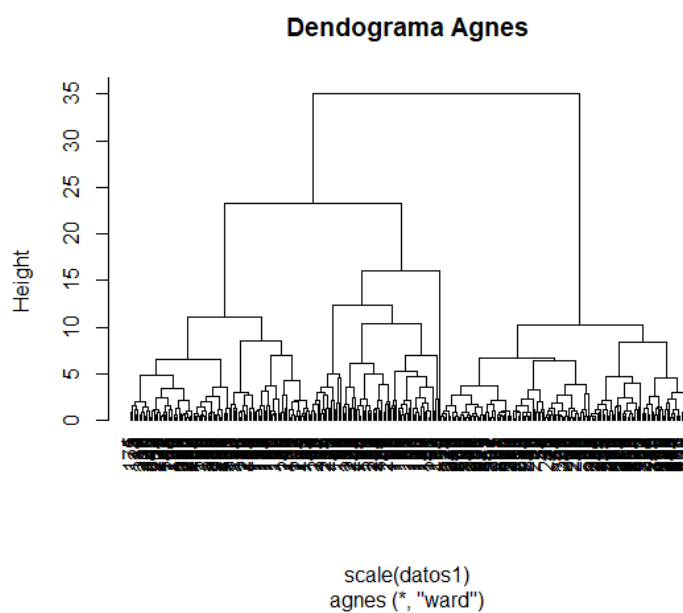


Figura 22: Dendograma Agnes

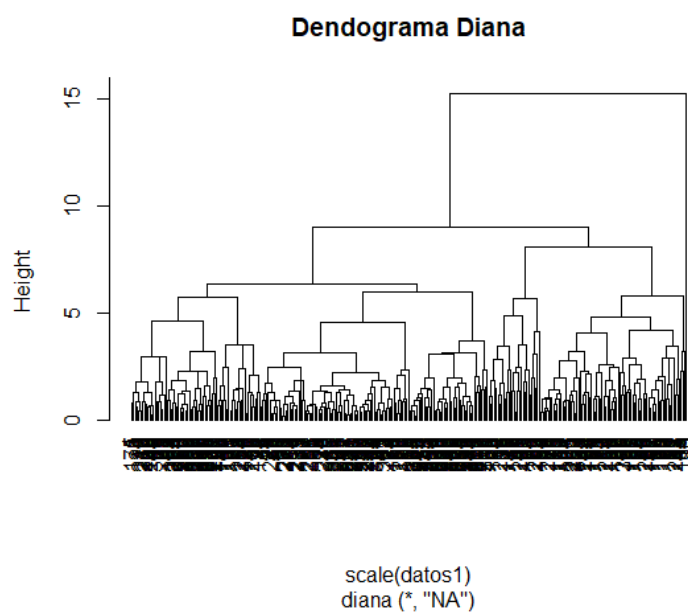


Figura 23: Dendrograma Diana

Analicemos gráficamente, el agrupamiento de los datos para observar mejor la tendencia en la agrupación.

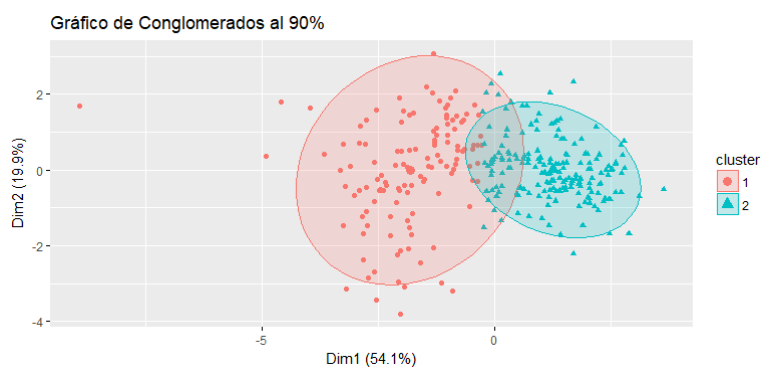


Figura 24: kmeans - Agrupamiento por clases en un 90 %



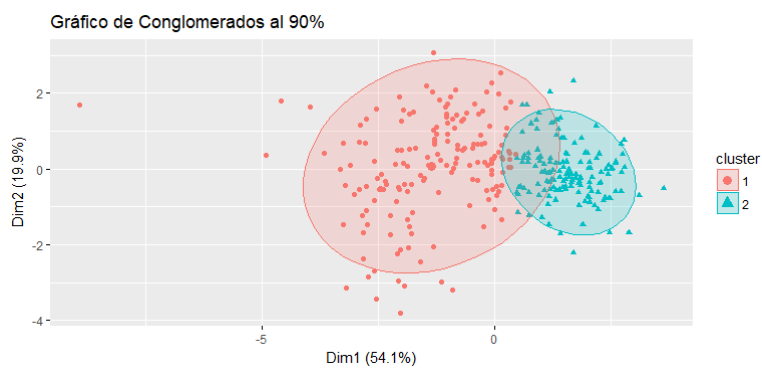


Figura 25: PAM - Agrupamiento por clases en un 90 %

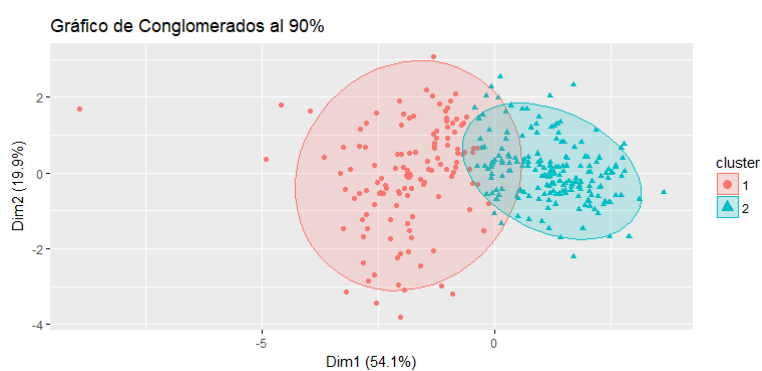


Figura 26: Clara - Agrupamiento por clases en un 90 %

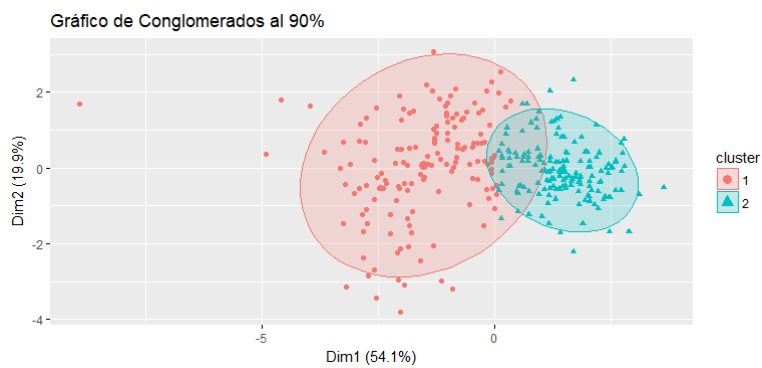


Figura 27: Fanny - Agrupamiento por clases en un 90 %

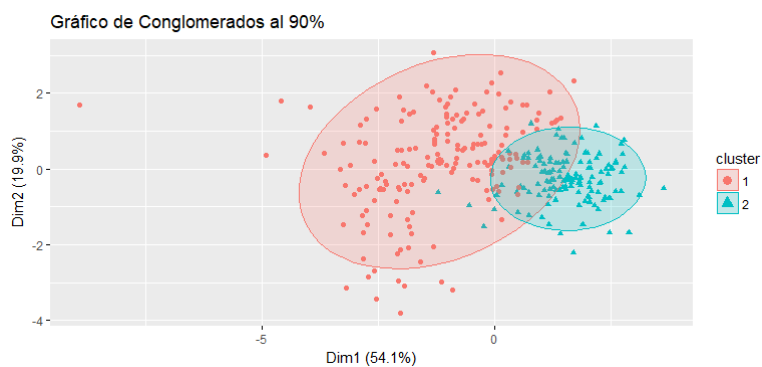


Figura 28: HIERARCHICAL - Agrupamiento por clases en un 90 %

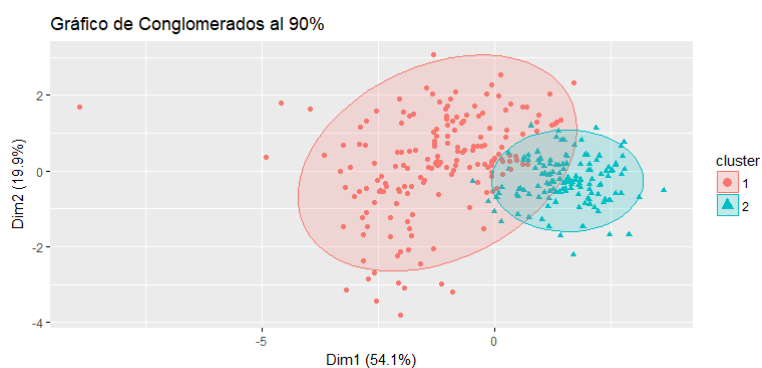


Figura 29: Agnes - Agrupamiento por clases en un 90 %

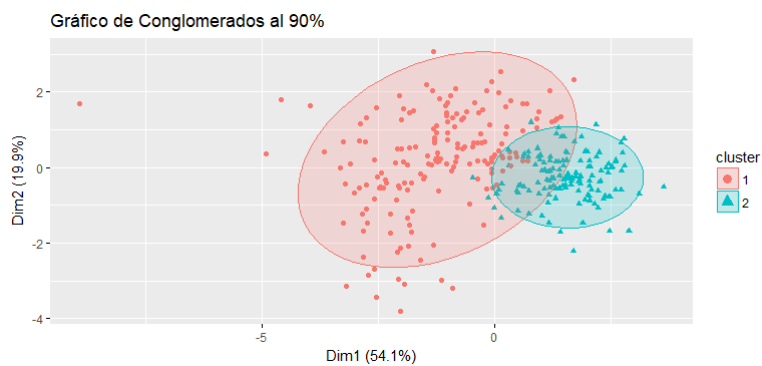


Figura 30: Diana - Agrupamiento por clases en un 90 %

En efecto, se sigue un agrupamiento acorde a lo establecido por la base de datos original, sin embargo es necesario realizar un análisis, con mejor resumen por diagrama de caja, de estos nuevos agrupamiento para corroborar que el modelo trabaja de forma adecuada.

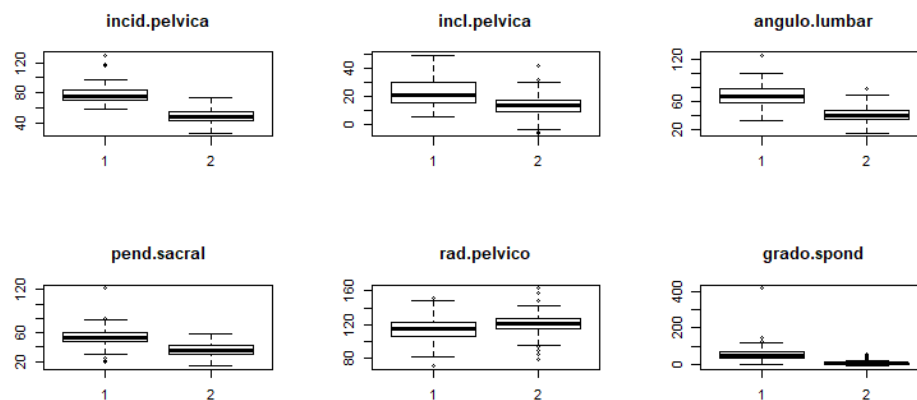


Figura 31: kmeans

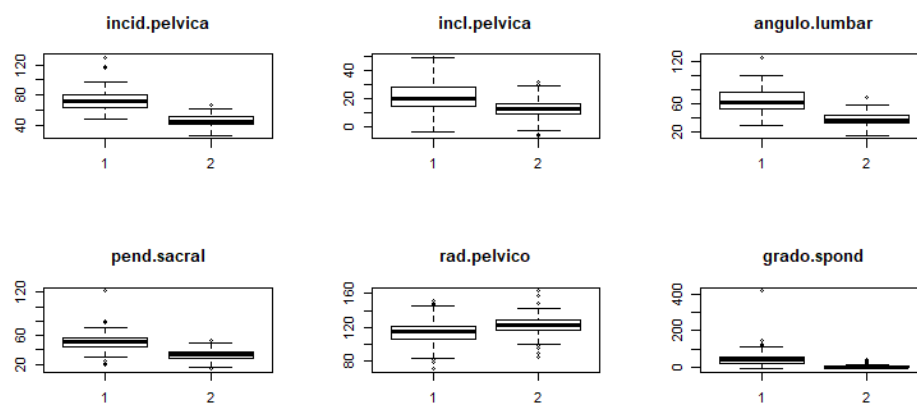


Figura 32: PAM

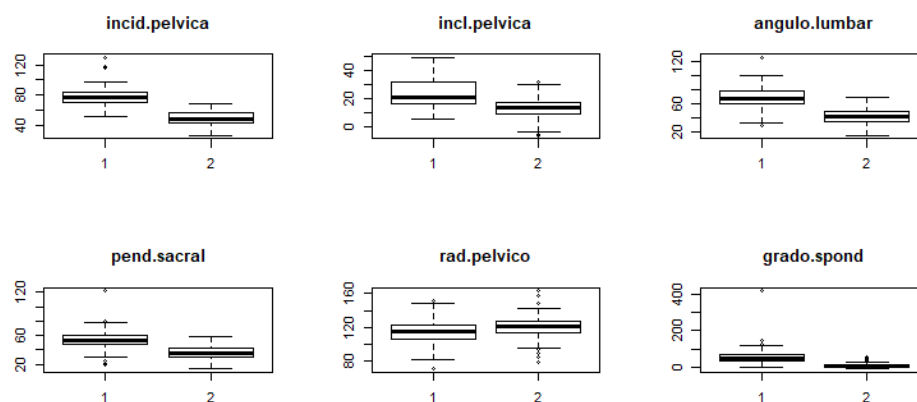


Figura 33: Clara

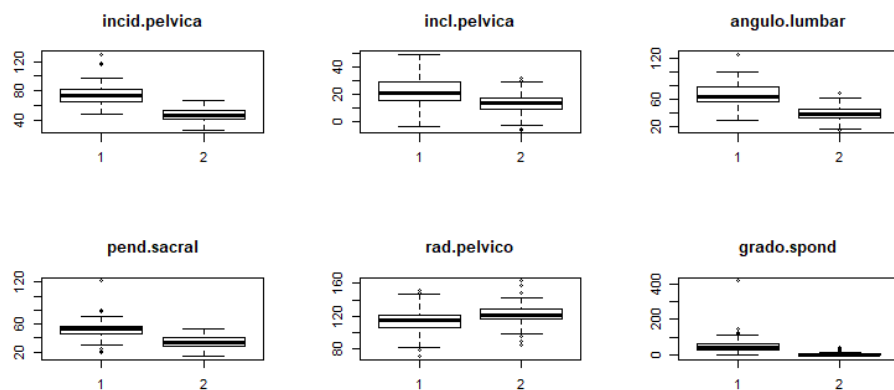


Figura 34: Fanny

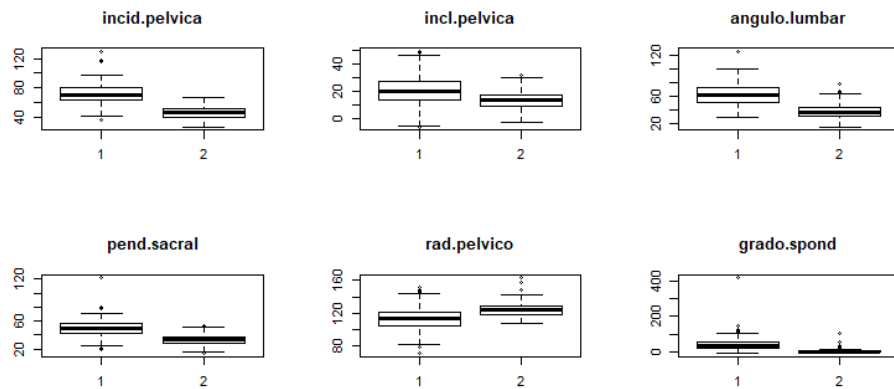


Figura 35: Hierarchical

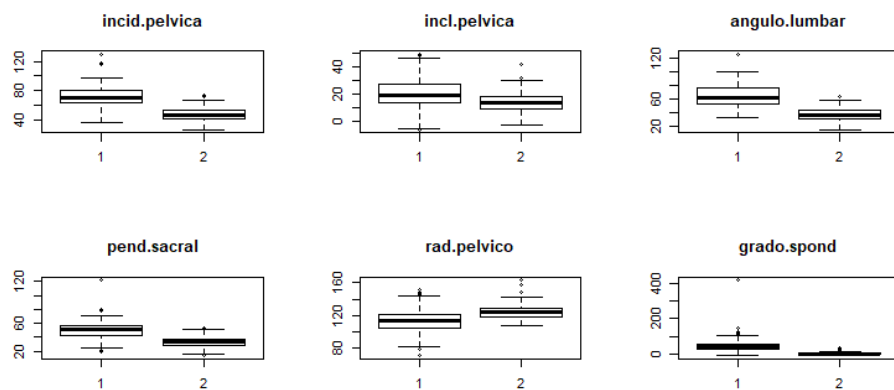


Figura 36: Agnes

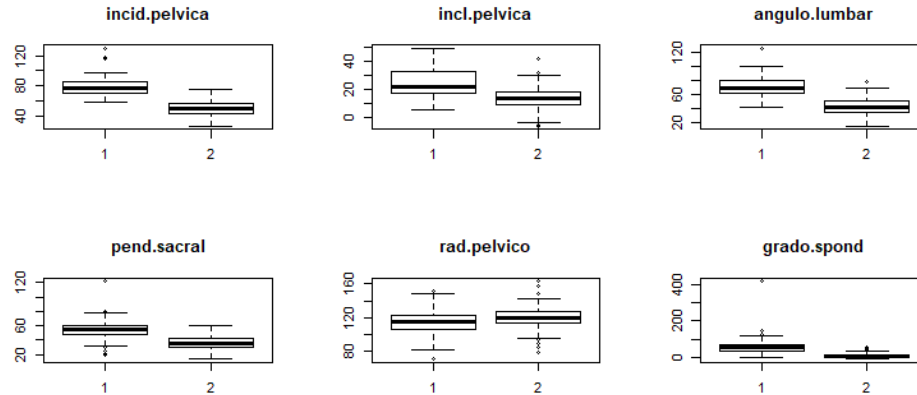


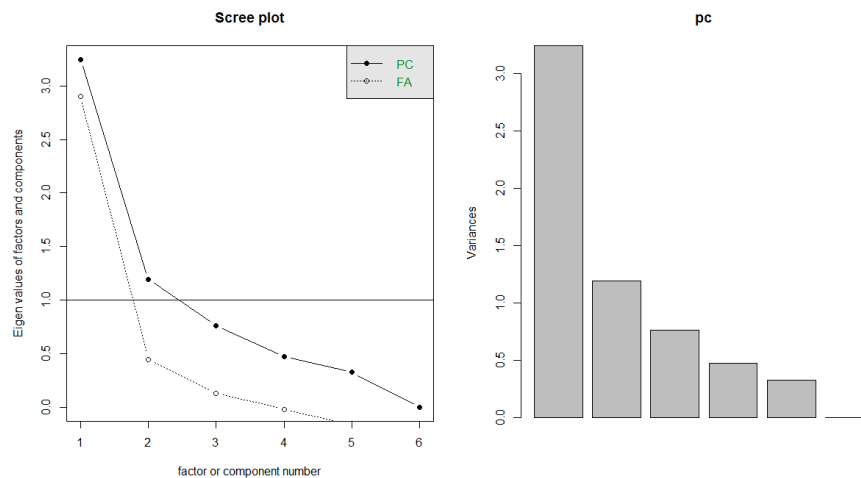
Figura 37: Diana

Evaluando el diagrama de cada de la Figura 13. En conclusión se ha efectuado un correcto agrupamiento por parte de los modelos.

### Análisis de Componentes Principales (PCA)

Esta metodología buscará relacionar las distintas variables de la base de datos, con el objetivo de analizar correspondencias entre ellas.

El gráfico de Sedimentación nos ayudará a elegir el número de componente que retengan la mayor información posible. En el siguiente gráfico (línea "PC"), se observa que a partir del tercer componente, la pendiente formada por los autovalores se estabiliza, por lo que se decide usar solo dos componentes, de otra manera, si tomamos en cuenta la regla de Kaiser que nos sugiere tomar tantos componentes como autovalores mayores a la unidad, optaríamos también usar dos componentes. Usar las dos primeras componentes, implica también que ambas expliquen el 86.8 % de la varianza total. La tercera componente y cuarta componente solo aportarían el 8.9 % y 4.3 % de la varianza total.



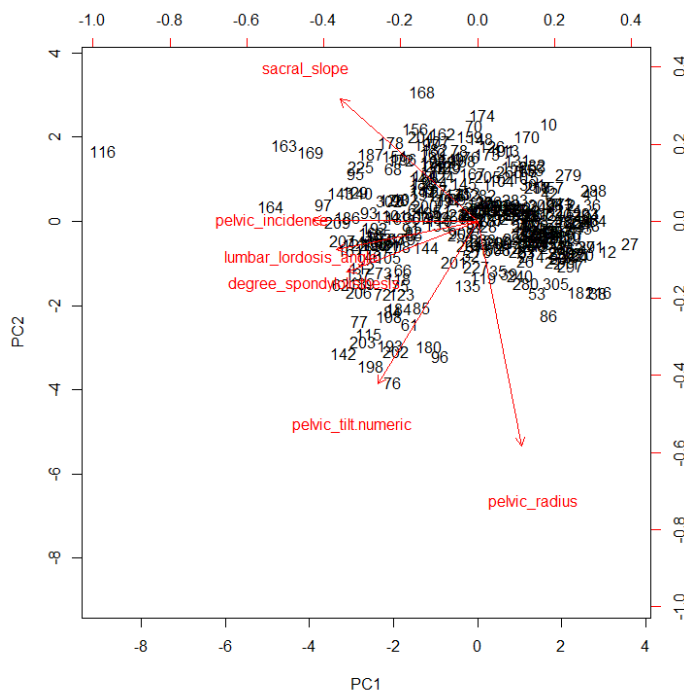
Importance of components:

	PC1	PC2	PC3	PC4	PC5	PC6
Standard deviation	1.802	1.093	0.872	0.6874	0.5710	1.94e-10
Proportion of Variance	0.541	0.199	0.127	0.0788	0.0543	0.00e+00
Cumulative Proportion	0.541	0.740	0.867	0.9457	1.0000	1.00e+00

Ahora bien, si calculamos las cargas de las dos primeras componentes, se observa que las variables `pelvic_incidence`, `lumbar_lordosis_angle`, `degree_spondylolisthesis` y `sacral_slope` pueden ser explicados en la primera componente, mientras que las variables `pelvic_tilt.numeric` y `pelvic_radius` pueden ser explicadas en la segunda componente, pues éstas tienen un mayor peso (en valor absoluto) respecto solo a las componentes mencionadas.

	PC1	PC2	PC3	PC4	PC5	PC6
<code>pelvic_incidence</code>	-0.53514	0.0021937	-0.096069	0.1027990	-0.42346	-7.1729e-01
<code>pelvic_tilt.numeric</code>	-0.32358	-0.5275454	-0.648701	0.0064412	-0.15056	4.1649e-01
<code>lumbar_lordosis_angle</code>	-0.45797	-0.0928751	0.152338	0.5480997	0.67677	1.4568e-11
<code>sacral_slope</code>	-0.44591	0.3961573	0.360313	0.1272009	-0.43150	5.5860e-01
<code>pelvic_radius</code>	0.14350	-0.7277556	0.585991	0.1742130	-0.27576	3.4466e-12
<code>degree_spondylolisthesis</code>	-0.42398	-0.1627769	0.271184	-0.8015281	0.27885	-8.3036e-12

En el gráfico biplot, se aprecia que un buen grupo de observaciones que se contraponen o es explicado de manera inversa con las variables de `pelvic_incidence`, `lumbar_lordosis_angle` y `degree_spondylolisthesis`, en tanto, las variables `sacral_slope` y `pelvic_radius` se relacionan de manera inversa.



## Regresión Logística:

Se hará una evaluación por modelo logístico. Asumiendo clases donde 0 = anormal y 1=normal.

```
call:
  glm(formula = clase_grupo ~ pelvic_incidence + pelvic_tilt.numeric +
      lumbar_lordosis_angle + sacral_slope + pelvic_radius + degree_spondylolisthesis
      family = binomial, data = datos)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.2678  -0.3639  -0.0289   0.4081   2.7317

Coefficients:
      Estimate Std. Error z value Pr(>|z|)
(Intercept)    -1.530e+01  3.315e+00  -4.615 3.93e-06 ***
pelvic_incidence  2.517e+07  4.017e+07   0.627  0.531
pelvic_tilt.numeric -2.517e+07  4.017e+07  -0.627  0.531
lumbar_lordosis_angle  1.794e-02  2.290e-02   0.784  0.433
sacral_slope    -2.517e+07  4.017e+07  -0.627  0.531
pelvic_radius     1.077e-01  2.318e-02   4.645 3.39e-06 ***
degree_spondylolisthesis -1.693e-01  2.335e-02  -7.248 4.23e-13 ***
---
signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 389.86 on 309 degrees of freedom
Residual deviance: 177.87 on 303 degrees of freedom
AIC: 191.87

Number of Fisher Scoring iterations: 8
```

Figura 38: Estimación del modelo

```
      ytrue
ypred  0  1
0 196  28
1  14  72
> testerr <- mean(ypred!=ytrue)
> testerr
[1] 0.1354839
> (S <- mc[2,2]/sum(mc[,2]))
[1] 0.72
> (E <- mc[1,1]/sum(mc[,1]))
[1] 0.9333333
```

Figura 39: Matriz de confusión (error de clasificación)

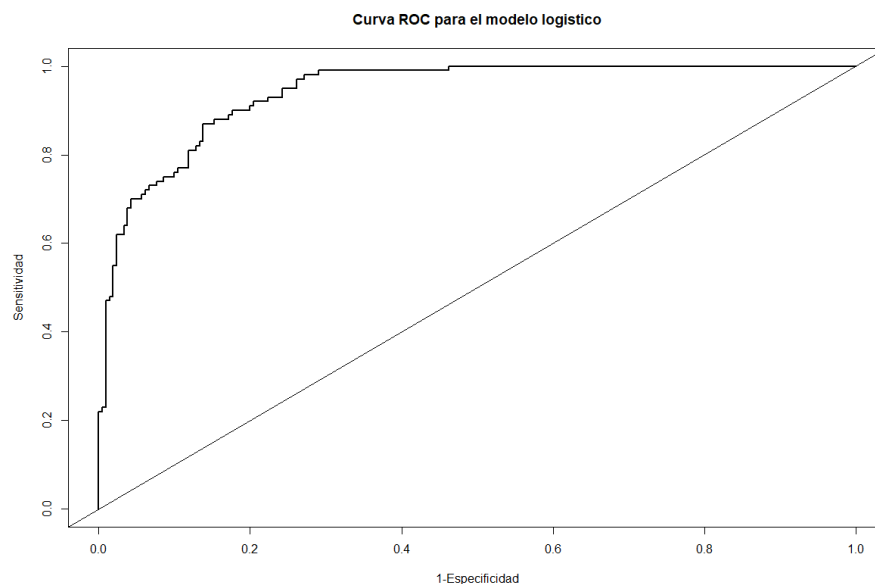


Figura 40: CURVA ROC

	incid.pelvica	incl.pelvica	angulo.lumbar	pend.sacral	rad.pelvico	grado.spond	clase
1	63.02782	22.552586	39.60912	40.47523	98.67292	-0.254400	Abnormal
2	39.05695	10.060991	25.01538	28.99596	114.40543	4.564259	Abnormal
3	68.83202	22.218482	50.09219	46.61354	105.98514	-3.530317	Abnormal
4	69.29701	24.652878	44.31124	44.64413	101.86850	11.211523	Abnormal
5	49.71286	9.652075	28.31741	40.06078	108.16872	7.918501	Abnormal
6	40.25020	13.921907	25.12495	26.32829	130.32787	2.230652	Abnormal

	subclase	clase.grupo	subclase.grupo	clas.grupo	ypred	yprob
1	Hernia	0	1	0	1	0.16443417
2	Hernia	0	1	0	1	0.22865463
3	Hernia	0	1	0	1	0.62861031
4	Hernia	0	1	0	1	0.04161865
5	Hernia	0	1	0	1	0.22036960
6	Hernia	0	1	0	1	0.52556338

Figura 41: Base de datos con categorías con las probabilidades y categorías predichas



## Anexo 1

Listing 9: Código R. Pregunta 1

---

```

1 rm(list=ls(all=TRUE))
2
3 ### Instalar y cargar paquetes
4
5 library(MASS)
6 library(rpart)
7 library(caret)
8 library(klaR)
9 library(pROC)
10
11 ### Cargar base de datos: wisc_bc_data.csv
12
13 # datos = read.csv(https://raw.githubusercontent.com/joshishwetha/dsx-spark/master/
14 #   data.csv)[-1]
15 # datos = read.csv("http://www3.nd.edu/~steve/computing_with_data/Data/wisc_bc_data.
16 #   csv")[-1]
17 datos = read.csv("wisc_bc_data.csv")[-1]
18 datos$diagnosis = ifelse(datos$diagnosis=='M',1,0)
19 datos$diagnosis = as.factor(datos$diagnosis)
20
21 ### Base de datos de entrenamiento y de test
22 N = nrow(datos)
23 p = 0.8
24 n = round(N*p,0)
25
26 set.seed(1234)
27 indic = sample(N,n,replace=FALSE)
28 datos.train = datos[indic,]
29 datos.test = datos[-indic,]
30
31 #-----#
32 # MODELO LOGISTICO #
33 #-----#
34
35 ## Seleccion de Variables
36 model.glm = glm(diagnosis ~ ., family=binomial, data=datos.train)
37 model.step = stepAIC(model.glm)
38 model.step
39 # X1 = datos.test[c(2,6,7,9,11,12,14,17,19,23,24,25,26,28)] #Variables elegidas
40 #   segun stepAIC
41 X1 = datos.test[attr(model.step$terms,"term.labels")]
42 head(X1)
43
44 ## Matriz de confusion
45 pi.glm = predict(model.step, type="response", newdata=X1)
46 cut = 0.5
47 y.glm = as.numeric(pi.glm >= cut)
48 y.true = datos.test$diagnosis
49 mc.glm = confusionMatrix(y.true, y.glm, positive="1")
50 roc.glm = roc(response=y.true, predictor=y.glm)
51 area.glm = roc.glm$auc
52

```

```

53 #-----#
54 # ANALISIS DISCRIMINANTE LINEAL (LDA) #
55 #-----#
56
57 ## Seleccion de Variables
58 wilks.lda = greedy.wilks(diagnosis ~ ., data=datos.train)
59 model.lda = lda(wilks.lda$formula, data=datos.train)
60 model.lda
61 X2 = datos.test[attr(model.lda$means,"dimnames")[[2]]] #Variables elegidas segun
    greedy.wilks
62 head(X2)
63
64 ## Matriz de confusion
65 y.lda = predict(model.lda, newdata=X2)$class
66 y.true = datos.test$diagnosis
67 mc.lda = confusionMatrix(y.true, y.lda, positive="1")
68 roc.lda = roc(response=y.true, predictor=as.numeric(y.lda))
69 area.lda = roc.lda$auc
70
71
72 #-----#
73 # ANALISIS DISCRIMINANTE CUADRATICO (QDA) #
74 #-----#
75
76 ## Seleccion de Variables
77 wilks.qda = greedy.wilks(diagnosis ~ ., data=datos.train)
78 model.qda = qda(wilks.qda$formula, data=datos.train)
79 model.qda
80 X3 = datos.test[attr(model.qda$means,"dimnames")[[2]]] #Variables elegidas segun
    greedy.wilks
81 head(X3)
82
83 ## Matriz de confusion
84 y.qda = predict(model.qda, newdata=X3)$class
85 y.true = datos.test$diagnosis
86 mc.qda = confusionMatrix(y.true, y.qda, positive="1")
87 roc.qda = roc(response=y.true, predictor=as.numeric(y.qda))
88 area.qda = roc.qda$auc
89
90
91 #-----#
92 # ANALISIS DISCRIMINANTE REGULARIZADO (RDA) #
93 #-----#
94
95 ## Seleccion de Variables
96 wilks.rda = greedy.wilks(diagnosis ~ ., data=datos.train)
97 model.rda = rda(wilks.rda$formula, data=datos.train)
98 model.rda
99 X4 = datos.test[attr(model.rda$means,"dimnames")[[1]]] #Variables elegidas segun
    greedy.wilks
100 head(X4)
101
102 ## Matriz de confusion
103 y.rda = predict(model.rda, newdata=X4)$class
104 y.true = datos.test$diagnosis
105 mc.rda = confusionMatrix(y.true, y.rda, positive="1")
106 roc.rda = roc(response=y.true, predictor=as.numeric(y.rda))
107 area.rda = roc.rda$auc

```

```

108
109
110 #-----#
111 # RESUMEN #
112 #-----#
113
114 GLM = c(mc.glm$overall[1:2], mc.glm$byClass[1:2], AUC=area.glm)
115 LDA = c(mc.lda$overall[1:2], mc.lda$byClass[1:2], AUC=area.lda)
116 QDA = c(mc.qda$overall[1:2], mc.qda$byClass[1:2], AUC=area.qda)
117 RDA = c(mc.rda$overall[1:2], mc.rda$byClass[1:2], AUC=area.rda)
118 tabla = round(cbind(GLM, LDA, QDA, RDA),3)
119 tabla
120
121
122 #-----#
123 # Grafica de Curvas ROC #
124 #-----#
125
126 par(mfrow=c(2,2))
127 plot(1-roc.glm$specificities, roc.glm$sensitivities, type="l", lwd=2, col="magenta",
128      ylab="Sensitividad", xlab="1-Especificidad",main="Curva ROC: Modelo Logistico")
129      abline(a=0, b=1, lwd=2)
130 plot(1-roc.lda$specificities, roc.lda$sensitivities, type="l", lwd=2, col="cyan",
131      ylab="Sensitividad", xlab="1-Especificidad",main="Curva ROC: A. D. Lineal")
132      abline(a=0, b=1, lwd=2)
133 plot(1-roc.qda$specificities, roc.qda$sensitivities, type="l", lwd=2, col="red",
134      ylab="Sensitividad", xlab="1-Especificidad",main="Curva ROC: A. D. Cuadratico")
135      abline(a=0, b=1, lwd=2)
136 plot(1-roc.rda$specificities, roc.rda$sensitivities, type="l", lwd=2, col="blue",
137      ylab="Sensitividad", xlab="1-Especificidad",main="Curva ROC: A. D. Regularizado")
138      abline(a=0, b=1, lwd=2)

```

---

## Anexo 2

Listing 10: Código R. Pregunta 2

---

```
1 ## Pregunta 2
2
3 rm(list=ls())
4
5 library(MASS)
6 library(smacof)
7 library(cluster)
8
9 letras=read.csv(file.choose(), sep=";")
10 filas=letras[,1]
11 letras=letras[,-1]
12 rownames(letras)=filas
13 letras=as.dist(as.matrix(letras))
14 letras=21-letras
15
16
17 # Escalamiento no metrico
18
19
20 res=smacofSym(letras, ndim=2, type = "ordinal")
21 res
22 summary(res)
23
24 plot(res, plot.type="confplot")
25 plot(res, plot.type="Shepard")
26 plot(res, plot.type="resplot")
27 (cor(res$confdist, res$dhat))^2
28 plot(res, plot.type="stressplot")
29 plot(res, plot.type="bubbleplot")
30
31
32 ## Cluster jerarquico
33
34 a=hclust(dist(letras), method="average")
35 plot(a)
```

---

## Anexo 3

Listing 11: Código R. Pregunta 3.

---

```

1 #####
2 ##### TRABAJO GRUPAL — PREG. 3 — ANALISIS DE DATOS
3 #####
4 ##### BASE DE DATOS: CARACTERISTICAS BIOMECANICAS DE PACIENTES ORTOPEDICOS
5 #####
6
7 library(cluster)
8 library(fpc)
9 library(NbClust)
10 library(clValid)
11 library(dendextend)
12 library(circlize)
13 library(sparcl)
14
15 ## LECTURA DE DATOS
16 datos <- read.csv("C:../ortopedico.csv")
17 head(datos)
18 sapply(datos,class)
19
20 attach(datos)
21
22 ## ADICION DE GRUPOS ORIGINALES
23 dim.class=c(summary(class))
24 datos$clase.grupo=c(rep(1,dim.class[1]),rep(2,dim.class[2]))
25
26 dim.sub=c(summary(sub.class))
27 datos$subclase.grupo=c(rep(1,dim.sub[1]),rep(3,dim.sub[3]),rep(2,dim.sub[2]))
28
29 colnames(datos) = c("incid.pelvica","incl.pelvica","angulo.lumbar","pend.sacral",
30 "rad.pelvico","grado.spond","clase","subclase",
31 colnames(datos[,9:10]))
32
33 head(datos)
34
35 ## ANALISIS EXPLORATORIO
36 dim(datos)
37 summary(datos[,1:7])
38
39 par(mfrow=c(2,3))
40 for (i in 1:6) {
41   boxplot(datos[,i]~datos$clase.grupo, main = names(datos[i]), type="l")
42 }
43
44 par(mfrow=c(1,1))
45 plot(datos$clase.grupo,ylab= "Grupos originales",xlab = "observaciones")
46
47 ## DATOS + GRUPO CLASS
48 datos1 <- datos[,1:6]
49
50 par(mfrow=c(1,1))
51 fviz_cluster(list(data = datos1,cluster=datos[,9]),
52 ellipse.type="norm",ellipse.level=0.9,
53 geom = "point",
54 main = "Agrupamiento original al 90%")
55

```

```

56 #####
57 ##### PRIMERA METODOLOGIA: CLUSTERING
58 #####
59
60 ### PASO 1. DEFINICION DE FUNCIONES (ALTERNATIVO)
61
62 # 1: kmeans, 2: PAM, 3: clara, 4: fanny
63 metodo <- function(datos1,h,nm){
64   set.seed(2010)
65   if(nm == 2){ pam(scale(datos1),h)
66 }else{if(nm == 3){ clara(scale(datos1),h)
67 }else{if(nm == 4){ fanny(scale(datos1),h,maxit=5000)
68 }else{ kmeans(scale(datos1),h,nstart = 100)}}}
69 }
70
71 # SUMA DE CUADRADOS DENTRO DE CADA CLUSTER
72 SC.cluster <- function(datos1, nm){
73   asw<-numeric()
74   if(nm ==1){for(h in 2:10){asw[h-1]=metodo(datos1,h,nm)$tot.withinss}
75   plot(2:10,asw,type="b")
76 }else{for(h in 2:10){asw[h-1]=metodo(datos1,h,nm)$silinfo$avg.width}
77   plot(2:10,asw,type="b",xlab="k",ylab="ASW")
78 }
79 }
80
81 # SILUETA
82 silueta <- function(datos1, nm){
83   diss.datos=daisy(scale(datos1))
84   par(mfrow=c(1,3))
85   if(nm ==1){for(h in 2:4){
86     plot(silhouette(metodo(datos1,h,nm)$cluster,diss.datos))}
87 }else{for(h in 2:4){plot(metodo(datos1,h,nm),which.plots=2)}}
88 }
89
90 # CRITERIO DE CALINSKI — HARABASZ
91 Cal.Har <- function(datos1, nm){
92   ch<-numeric()
93   par(mfrow=c(1,1))
94   if(nm ==1){for(h in 2:10){
95     ch[h-1] = calinhara(scale(datos1),metodo(datos1,h,nm)$cluster)}
96 }else{for(h in 2:10){
97   ch[h-1] = calinhara(scale(datos1),metodo(datos1,h,nm)$clustering)}}
98   plot(2:10,ch,type="b",xlab="k", ylab="Criterio de Calinski–Harabasz")
99 }
100
101 # GRAFICAS CLUSTER
102 plot.cluster <- function(datos1,clus){
103   par(mfrow=c(1,1))
104   fviz_cluster(list(data = datos1,cluster=clus),
105     ellipse.type="norm",ellipse.level=0.9,
106     geom = "point",
107     main = "Grafico de Conglomerados al 90%")
108 }
109
110 ### PASO 2. DETERMINANDO NUMERO DE CONGLOMERADOS
111 res.cluster=list()
112
113 ## KMEANS

```

```

114 SC.cluster(datos1,1)
115
116 silueta(datos1,1)
117 kmeansruns(scale(datos1),criterion="asw")
118
119 Cal.Har(datos1,1)
120 kmeansruns(scale(datos1),criterion="ch")
121
122 set.seed(2000)
123 res.cluster[[1]]=kmeans(scale(datos1),2)$cluster
124
125 par(mfrow=c(1,1))
126 plot.cluster(datos1,res.cluster[[1]])
127
128 ## PAM
129 SC.cluster(datos1,2)
130
131 silueta(datos1,2)
132 pamk(scale(datos1),criterion="asw")
133
134 Cal.Har(datos1,2)
135 pamk(scale(datos1),criterion="ch")
136
137 set.seed(2000)
138 res.cluster[[2]]=pam(scale(datos1),2)$clustering
139
140 par(mfrow=c(1,1))
141 plot.cluster(datos1,res.cluster[[2]])
142
143 ## CLARA
144 SC.cluster(datos1,3)
145
146 silueta(datos1,3)
147 pamk(scale(datos1),criterion="asw",usepam=FALSE)
148
149 Cal.Har(datos1,3)
150 pamk(scale(datos1),criterion="ch",usepam=FALSE)
151
152 set.seed(2000)
153 res.cluster[[3]]=clara(scale(datos1),2)$clustering
154 for (i in 1:dim(datos)[1]) {
155   if(res.cluster[[3]][i] == 1){
156     res.cluster[[3]][i]=2
157   }else{res.cluster[[3]][i] = 1
158   }
159 }
160
161 par(mfrow=c(1,1))
162 plot.cluster(datos1,res.cluster[[3]])
163
164 ## FANNY
165 SC.cluster(datos1,4)
166 silueta(datos1,4)
167 Cal.Har(datos1,4)
168
169 set.seed(2000)
170 res.cluster[[4]]=fanny(scale(datos1),2)$clustering
171 for (i in 1:dim(datos)[1]) {
172   if(res.cluster[[4]][i] == 1){

```

```

173 res.cluster[[4]][i]=2
174 }else{res.cluster[[4]][i] = 1
175 }
176 }
177
178 par(mfrow=c(1,1))
179 plot(cluster(datos1,res.cluster[[4]]))
180
181 #table(datos[,9],res$cluster)
182
183 par(mfrow=c(2,2))
184 plot(res.cluster[[1]],xlab = "observaciones",ylab = "kmeans")
185 plot(res.cluster[[2]],xlab = "observaciones",ylab = "PAM")
186 plot(res.cluster[[3]],xlab = "observaciones",ylab = "clara")
187 plot(res.cluster[[4]],xlab = "observaciones",ylab = "fanny")
188
189 ## CLUSTER JERARQUICO: HIERARCHICAL
190 set.seed(2000)
191 hc <- hclust(dist(scale(datos1)),method="ward.D")
192 res.cluster[[5]] = cutree(hc , 2)
193
194 dend <- as.dendrogram(hc) %>%
195 color_branches(k=2) %>%
196 color_labels
197
198 par(mar = rep(0,4))
199 circlize_dendrogram(dend, labels_track_height = 0.3,
200 dend_track_height = .6,
201 main = "Dendograma")
202
203 par(mfrow=c(1,1))
204 ColorDendrogram(hc, y = res.cluster[[5]],
205 labels = names(res.cluster[[5]]),
206 main = "Dendograma",
207 xlab= "observaciones",
208 branchlength = 0.4)
209
210 par(mfrow=c(1,1))
211 plot(cluster(datos1,res.cluster[[5]]))
212
213 ## CLUSTER JERARQUICO: AGNES
214 agnes.single=agnes(scale(datos1),method="single")
215 agnes.single$ac
216 par(mfrow=c(1,1))
217 pltree(agnes.single,cex=1,hang=-1, main = "Dendograma Agnes")
218
219 agnes.ward=agnes(scale(datos1),method="ward")
220 agnes.ward$ac
221 par(mfrow=c(1,1))
222 pltree(agnes.ward,cex=1,hang=-1, main = "Dendograma Agnes")
223
224 diss.datos=daisy(scale(datos1))
225
226 agnes.ward=as.hclust(agnes.ward)
227 par(mfrow=c(1,3))
228 for(h in 2:4){
229 res=cutree(agnes.ward,k=h)
230 plot(silhouette(res,diss.datos))
231 }

```



```

232
233 set.seed(2000)
234 res.cluster[[6]]=cutree(agnes.ward,k=2)
235
236 par(mfrow=c(1,1))
237 plot.cluster(datos1,res.cluster[[6]])
238
239 ## DIANA
240 diana.metodo=diana(scale(datos1))
241
242 par(mfrow=c(1,1))
243 pltree(diana.metodo,cex=1,hang=-1, main = "Dendograma Diana")
244
245 diana.metodo=as.hclust(diana.metodo)
246 par(mfrow=c(1,3))
247 for(h in 2:4){
248   res=cutree(diana.metodo,k=h+1)
249   plot(silhouette(res,diss.datos))
250 }
251
252 set.seed(2000)
253 res.cluster[[7]]=cutree(diana.metodo,k=3)
254 for (i in 1:dim(datos)[1]) {
255   if(res.cluster[[7]][i] == 1){
256     res.cluster[[7]][i]=2
257   }else{res.cluster[[7]][i] = 1
258   }
259 }
260
261 par(mfrow=c(1,1))
262 plot.cluster(datos1,res.cluster[[7]])
263
264 par(mfrow=c(1,2))
265 plot(res.cluster[[5]],xlab = "observaciones",ylab = "Hierarchical")
266 plot(res.cluster[[6]],xlab = "observaciones",ylab = "Agnes: Wald")
267
268 par(mfrow=c(1,1))
269 plot(res.cluster[[7]],xlab = "observaciones",ylab = "Diana")
270
271 ### ————— PASO 3. PERFILADO Y CARACTERIZACION DE CLUSTERS
272
273 ## ADICION DE CLUSTER A LA BASE
274 datos.new <- datos[,c(1:6,9)]
275 for (i in 1:7) {
276   datos.new <- cbind(datos.new,res.cluster[[i]])
277 }
278
279 colnames(datos.new)<-c(colnames(datos.new[,1:7]),"cluster.km","cluster.PAM",
280 "cluster.clara","cluster.fanny","cluster.h",
281 "cluster.agnes","cluster.diana")
282 head(datos.new)
283
284 # TABLA DE MEDIAS
285 med=list()
286 for (i in 1:7) {
287   med[[i]] = aggregate(x = datos.new[,1:6],
288   by = list(datos.new[,i+7]),
289   FUN = mean)
290 }

```

```

291 med[[1]]
292
293 # ANALISIS EXPLORATIVO POR CLUSTER
294 for (j in 1:7) {
295   par(mfrow=c(2,3))
296   for (i in 1:6) {
297     boxplot(datos.new[,i]~datos.new[,7+j], main=names(datos.new[i]), type="l")
298   }
299 }
300
301 ## OTRAS TECNICAS
302 clmethods <- c("hierarchical","kmeans","pam","agnes","diana")
303
304 # Medidas de validacion interna
305 intern <- clValid(scale(datos1), nClust = 2:10,
306 clMethods = clmethods, validation = "internal")
307 summary(intern)
308
309 par(mfrow=c(1,1))
310 plot(intern)
311
312 # Medidas de estabilidad
313 stab <- clValid(scale(datos1), nClust = 2:10, clMethods = clmethods,
314 validation = "stability")
315 summary(stab)
316
317 # Mostrar solo scores optimos
318 optimalScores(intern)
319 optimalScores(stab)
320
321
322 ##### SEGUNDO METODO PCA
323
324 ## Cargar base de datos: ortopedico.csv
325 datos = read.csv("C:../ortopedico.csv")[,-7]
326 head(datos)
327
328 ## Grafico de Sedimentacion
329 library(psych)
330 par(mfrow=c(1,2))
331 scree(datos)
332 pc = prcomp(x=datos, scale=TRUE, center=TRUE)
333 plot(pc)
334 par(mfrow=c(1,1))
335
336 ## Proporcion de varianza explicada
337 summary(pc)
338
339 ## Carga de las componentes
340 pc$rotation
341
342 ## Grafico biplot
343 biplot(pc, scale=0)

```

---