

Informe del Cuarto Mini Proyecto de Estadística. Curso 2019-2020

Daniel Alberto García Pérez
Grupo C412

D.GARCIA@ESTUDIANTES.MATCOM.UH.CU

Leonel Alejandro García López
Grupo C412

L.GARCIA3@ESTUDIANTES.MATCOM.UH.CU

Roberto Marti Cedeño
Grupo C412

R.MARTI@ESTUDIANTES.MATCOM.UH.CU

Tutor(es):

Msc. Dalia Diaz Sistachs, *Facultad de Matemática y Computación, Universidad de La Habana*

Tema: Estadística, Técnicas de Clasificación.

1. Ejercicios

Mostramos a continuación los enunciados de los ejercicios propuestos. Todos los códigos relacionados con la respuesta a los ejercicios propuestos se encuentran dentro de la carpeta Code adjunta al proyecto.

1.1 Ejercicio 1

Trabajando para encontrar una sintomatología estándar para comportamientos de anorexia y bulimia se realice un estudio de 55 adolescentes con desórdenes alimenticios conocidos. En cada observación las pacientes fueron valoradas en síntomas diferentes.

1. Realice un análisis de componentes principales. Calcule las componentes, interprete sus valores.
2. Agrupe los datos utilizando técnicas de clúster jerárquico.
3. Realice un agrupamiento utilizando el algoritmo k-means.
4. Diga si es posible construir un árbol de clasificación usando CART tomando como variable respuesta el diagnostico. De ser posible constrúyalo. De no ser posible existe otra variable que se pueda usar como dependiente. Explique su respuesta.

2. Solución

Comenzaremos analizando las variables presentes en los datos. Debido a la gran cantidad de variables presentes en la muestra, y al tamaño de sus gráficas pasaremos directamente a analizar la matriz de correlación de forma gráfica.

Como podemos apreciar en el gráfico (Figura 1), la matriz no se encuentra altamente correlacionada. Por lo que podemos decir que la mayoría de sus variables son independientes. Pasamos entonces al análisis de sus componentes.

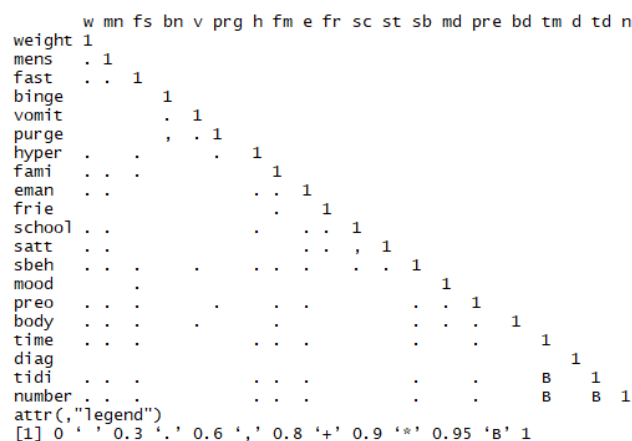


Figure 1: Matriz de correlación variante gráfica.

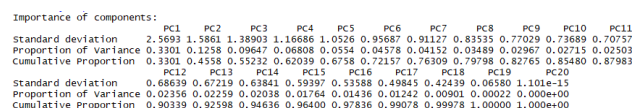


Figure 2: Importancia de las componentes.

Siguiendo el criterio de Kaiser en la importancia de los componentes (Figura 2), podemos llegar a la conclusión de que las componentes principales son las 5 primeras. Sin embargo, incumplimos el criterio del porcentaje dado que hasta la 5ta componente acumulamos solo un 68% de los datos aproximadamente. Se acordó entonces incluir una sexta componente dado que su valor propio es aproximadamente 1 y aglomeramos mas del 70% de los datos de la muestra.

Entonces nuestra primera componente se caracteriza por pacientes que presentan un alto valor de peso corporal, menstruación, que suelen restringirse el consumo de comida (fasting), son hiperactivos, mantiene buenas relaciones familiares, son independientes, poseen diver-

sos amigos, tienen un buen desempeño en la escuela o trabajo, poseen una actitud y comportamiento sexual positivo, también un buen estado de ánimo así como se preocupan por su percepción física, su peso y alimentos que ingieren. Por ultimo, su entrevista fue larga, interactuaron bastante durante la misma y se caracterizan por numero de paciente elevado.

Nuestra segunda componente se caracteriza por pacientes que no suelen darse atracones de comida, no suelen vomitar ni suelen provocarse el vómito. Poseen pocos amigos. Poseen una baja actitud y comportamiento sexual. Por otro lado, su tiempo de entrevista fue largo e interactuaron bastante durante la misma, su diagnóstico fue elevado y suelen pertenecer al ultimo grupo de pacientes.

La tercera componente se caracteriza por pacientes que, no se atracan de comida, no suelen provocarse pérdidas de comida y cuya entrevista fue de poca duración. Por otro lado, poseen buenas relaciones de amistad, un alto desempeño escolar o laboral y poseen una conducta sexual alta.

La cuarta componente se caracteriza por pacientes que, poseen poca menstruación, se limitan a la hora de comer, suelen estar deprimidos, no se preocupan por la comida ni por su peso, suelen ser delgados, y su diagnóstico es bajo. Por otro lado, son hiperactivos, son independientes, poseen un buen desempeño escolar o laboral, presentan una actitud sexual positiva, su entrevista fue de larga duración e interactuaron bastante en la misma. Son principalmente del ultimo grupo de pacientes.

La quinta componente se caracteriza por pacientes que, no se atracan de comida, tienen relaciones familiares pobres, dependen de su familia, tienen relaciones de amistad pobres, y su diagnóstico es bajo. Por otro lado, presentan un alto peso corporal, se restringen la ingestión de alimentos y poseen un alto autoestima.

La ultima componente, se caracteriza por pacientes que no suelen vomitar, no son hiperactivos, de complejión delgada, que se atracan de comida, poseen un buen desempeño laboral o educativo, presentan una actitud sexual positiva y un estado de ánimo elevado.

Siguiendo los resultados obtenidos del análisis de las componentes principales, realizamos un análisis de tipo clúster (Figuras 3 y 4). Podemos destacar lo acertado de considerar 6 componentes principales en lugar de 5 dado que, en el caso de 5 componentes se pierde información sobre una rama considerable de la jerarquía.

Después de la ejecución del algoritmo kmeans, con 6 componentes como hemos mantenido hasta ahora, (Figura 5) obtenemos una similitud por componentes, de un 44% resultado directo de la desproporción presente en la 2da de las componentes, lo que significa que en los datos propuestos abundan los pacientes con las características propias de esa componente. Sucesivos intentos de corridas del algoritmo con mas componentes incrementan dicha similitud, pero incumplen con los análisis determinados anteriormente, este comportamiento se debe posiblemente a que las observaciones de las variables aleatorias no superan un rango determinado entre 1 y 4.

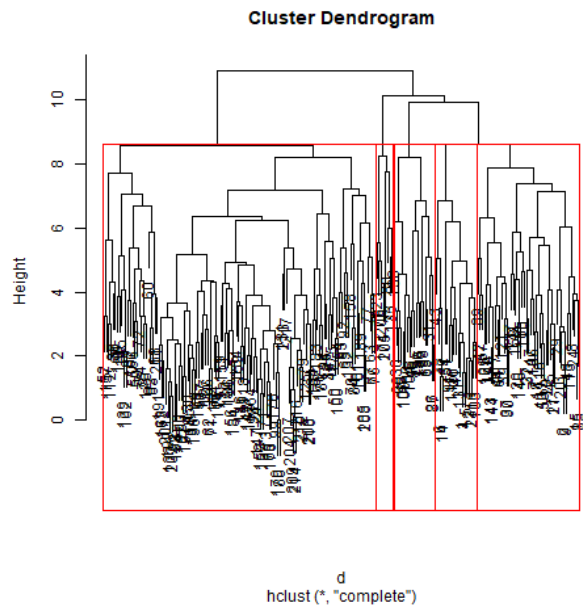


Figure 3: Clúster jerárquico de 5 componentes.

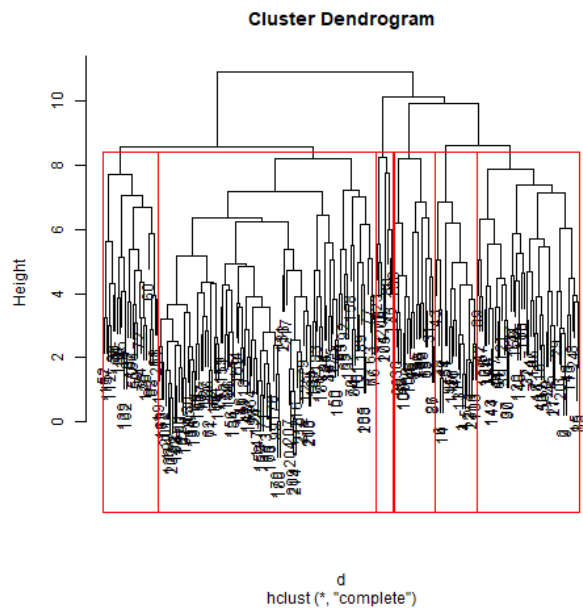


Figure 4: Clúster jerárquico de 6 componentes.

[illegible]

Figure 5: Resultado del algoritmo Kmeans con 6 componentes.

Finalmente se realizo el análisis del árbol de clasificación tomando como variable el diagnóstico de los pacientes. (Figura 6)

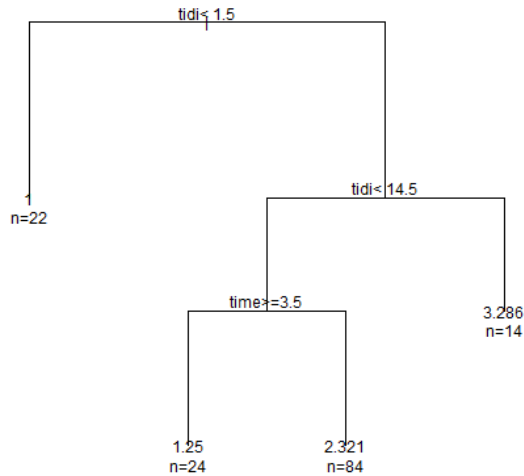


Figure 6: Árbol de Decisión

Como resultado se obtuvo que para la predicción del diagnóstico empleando el árbol de decisión se toman en cuenta 2 variables, tidi y el tiempo de la entrevista. Si el paciente presenta un valor de tidi inferior a 1,5 podemos predecir que posee un diagnóstico de tipo 1. En cambio, si el paciente posee un valor de tidi superior a 14,5 su diagnóstico es de tipo 3. Finalmente si el tiempo de entrevista del paciente es superior a 3.5 se le diagnostica el trastorno alimenticio 2, eoc, su diagnóstico es de tipo 1.