



## **Reconocimiento de patrones**

### **Tarea 6**

Clustering 3

**Alumno:**

Pérez Rodríguez Raúl Francisco

**Octubre 2017**

Realice la identificación de clases por k-means de los datos de los sobrevivientes del Titanic.

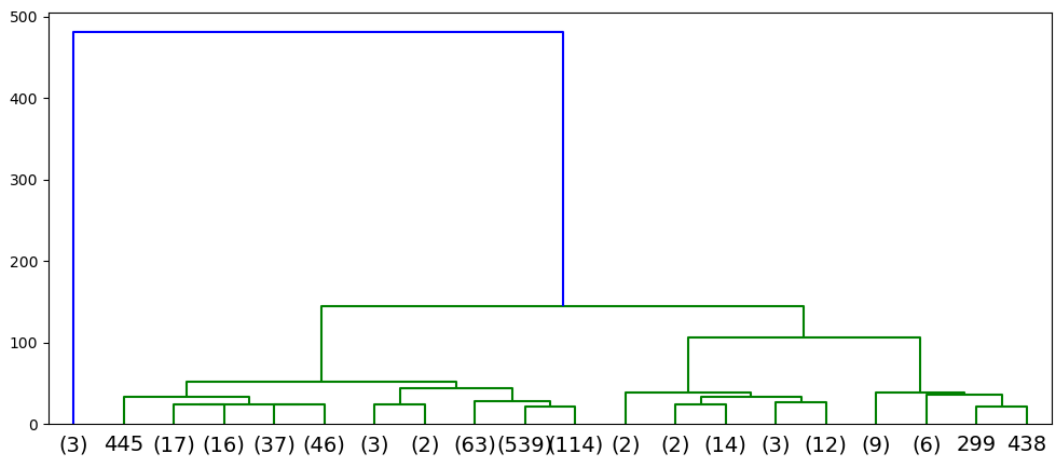
Dendrogramas

El conjunto de datos iniciales del Titanic, contienen 10 variables con 891 instancias. Eliminando las variables no numéricas, quedan las variables Survived ,Pclass, Age, SibSp, Parch, Fare.

Data Dictionary

Variable	Definition	Key
survival	Survival	0 = No, 1 = Yes
pclass	Ticket class	1 = 1st, 2 = 2nd, 3 = 3rd
sex	Sex	
Age	Age in years	
sibsp	# of siblings / spouses aboard the Titanic	
parch	# of parents / children aboard the Titanic	
ticket	Ticket number	
fare	Passenger fare	
cabin	Cabin number	
embarked	Port of Embarkation	C = Cherbourg, Q = Queenstown, S = Southampton

Utilizando la técnica de dendrograma con el método de centroid y truncando para obtener los últimos 20 clusters obtenemos la siguiente gráfica.



Analizando los últimos 20 clusters del dendrograma se pueden observar que puede haber entre 3 y 5 grupos.

## K-Means

Utilizando el dendrograma como una primera prueba, se observaron la formación de entre 3 y 5 grupos de datos. Con estos resultados se utilizara el k-means con k igual a 3, 4 y 5 para observar la formación de los grupos con una selección aleatoria de los clusters iniciales.

### Con k igual a 3

Realizando 3 corridas con el k-means con k igual a 3, se obtuvieron los mismos clusters

```
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6> python .\main.py

Cluster 0: 142 miembros.
[ 0.67  1.25 34.96 0.97  0.48 83.39]

Cluster 1: 729 miembros.
[ 0.32  2.55 28.64 0.43  0.34 15.45]

Cluster 2: 20 miembros.
[ 0.7   1.   31.02 0.75   1.05 279.31]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6> python .\main.py

Cluster 0: 142 miembros.
[ 0.67  1.25 34.96 0.97  0.48 83.39]

Cluster 1: 729 miembros.
[ 0.32  2.55 28.64 0.43  0.34 15.45]

Cluster 2: 20 miembros.
[ 0.7   1.   31.02 0.75   1.05 279.31]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6> python .\main.py

Cluster 0: 20 miembros.
[ 0.7   1.   31.02 0.75   1.05 279.31]

Cluster 1: 142 miembros.
[ 0.67  1.25 34.96 0.97  0.48 83.39]

Cluster 2: 729 miembros.
[ 0.32  2.55 28.64 0.43  0.34 15.45]
```

## Con k igual a 4

Realizando 3 corridas con el k-means con k igual a 4, se obtuvieron clusters diferentes excepto el cluster con 20 elementos, el cual se repite con k igual a 3.

```
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 275 miembros.
[ 0.36  2.73 17.02  0.74  0.52 15.87]
Cluster 1: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05 279.31]
Cluster 2: 141 miembros.
[ 0.67  1.25 34.71  0.98  0.48 83.63]
Cluster 3: 455 miembros.
[ 0.29  2.44 35.75  0.24  0.24 15.27]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 33 miembros.
[ 0.76  1.   32.18  0.61  0.79 131.11]
Cluster 1: 712 miembros.
[ 0.32  2.56 28.65  0.4   0.31 14.77]
Cluster 2: 126 miembros.
[ 0.59  1.41 34.76  1.17  0.56 65.6 ]
Cluster 3: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05 279.31]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 253 miembros.
[ 0.38  2.7  16.3   0.77  0.56 16.79]
Cluster 1: 142 miembros.
[ 0.67  1.25 34.96  0.97  0.48 83.39]
Cluster 2: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05 279.31]
Cluster 3: 476 miembros.
[ 0.29  2.47 35.19  0.25  0.23 14.75]
```

## Con k igual a 5

Realizando 3 corridas con el k-means con k igual a 5, se obtuvieron clusters diferentes excepto el cluster con 20 elementos, el cual se repite con k igual a 3. Pero a comparación con k igual a 4, en este caso el número de miembros por cluster no es tan diferente.

```
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 95 miembros.
[ 0.55  2.55  9.16  1.74  1.21  27.89]
Cluster 1: 97 miembros.
[ 0.68  1.21  34.12  1.16  0.59  96.5 ]
Cluster 2: 172 miembros.
[ 0.46  1.58  41.07  0.41  0.51  34.69]
Cluster 3: 507 miembros.
[ 0.26  2.78  28.79  0.2   0.12  10.12]
Cluster 4: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05  279.31]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 98 miembros.
[ 0.68  1.2  34.01  1.15  0.59  96.16]
Cluster 1: 94 miembros.
[ 0.56  2.52  9.14  1.72  1.24  28.37]
Cluster 2: 506 miembros.
[ 0.26  2.78  28.6   0.21  0.12  10.13]
Cluster 3: 173 miembros.
[ 0.45  1.59  41.5   0.4   0.49  34.05]
Cluster 4: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05  279.31]
PS C:\Users\Raul Perez\Documents\GitHub\Reconocimiento-de-patrones\06-Clustering-3\Tarea 6>
Cluster 0: 412 miembros.
[ 0.25  2.75  32.31  0.18  0.1   10.18]
Cluster 1: 191 miembros.
[ 0.41  2.73  14.39  0.84  0.62  16.71]
Cluster 2: 95 miembros.
[ 0.68  1.2  34.28  1.17  0.58  97.15]
Cluster 3: 173 miembros.
[ 0.46  1.55  37.72  0.6   0.59  37.53]
Cluster 4: 20 miembros.
[ 0.7   1.   31.02  0.75  1.05  279.31]
```

## Conclusión

Analizando los resultados al utilizar la técnica de k-means con los valores de k iguales a 3, 4 y 5. Se concluye que es más probable que haya 3 grupos de datos en el conjunto de datos.