

Autor: Juan Manuel Rodríguez Gómez

Asignatura: Estadística Multivariante (Prácticas)

Tarea Voluntaria 1 (Glaucoma DB)

1. Lectura del Conjunto de Datos

Se han eliminado en Excel los valores no numéricos de las celdas del conjunto de datos

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: # Lectura del conjunto de datos
df = pd.read_csv("C:/Users/LENOVO/Desktop/glaucoma.csv", delimiter=";", header=0)

In [3]: # Mostramos el DataFrame
df

Out[3]:
```

OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
0	0.0	0.0	112	4.0	1.5	174.0	NaN	31	0	0.0	0	3.0	0
1	1.0	NaN	108	4.0	1.2	128.0	1.0	29	23	19.0	24	2.0	4
2	0.0	1.0	123	4.0	1.1	133.0	1.0	36	30	30.0	30	1.0	4
3	1.0	2.0	131	4.0	1.5	191.0	1.0	14	0	21.0	14	1.0	0
4	0.0	2.0	156	4.0	1.2	182.0	1.0	14	0	16.0	17	1.0	0
...
116	NaN	16.0	102	4.0	1.4	141.0	NaN	23	0	0.0	0	1.0	0
117	NaN	16.0	107	4.0	1.4	149.0	NaN	26	0	0.0	0	0.0	0
118	1.0	1.0	140	4.0	NaN	211.0	NaN	21	0	0.0	0	2.0	0
119	0.0	1.0	198	45019.0	NaN	235.0	NaN	17	0	0.0	0	0.0	0
120	1.0	1.0	135	4.0	1.5	219.0	NaN	22	0	0.0	0	0.0	0

121 rows x 19 columns

```
In [4]: # Breve información del DataFrame
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121 entries, 0 to 120
Data columns (total 19 columns):
# Column Non-Null Count Dtype
---  ---
0 OJO 117 non-null float64
1 TIPO_GLAUCOMA 119 non-null float64
2 N_IMPACTOS 121 non-null int64
3 CUADRANTES 118 non-null float64
4 ENERGIA_IMPACTO 107 non-null float64
5 ENERGIA_TOTAL 121 non-null float64
6 CIRUJIA_PREVIA 85 non-null float64
7 PIO_PRE_SLT 121 non-null int64
8 PIO_1_SEMANA 121 non-null float64
9 PIO_1_MES 116 non-null float64
10 PIO_3_MES 121 non-null int64
11 FARMACOS_PRE 117 non-null float64
12 FARMACOS_1_MES 121 non-null int64
13 FARMACOS_3_MES 121 non-null int64
14 DOLOR 61 non-null float64
15 SEXO 60 non-null float64
16 EDAD 121 non-null int64
17 PIO_NORMAL 116 non-null float64
18 PIO_NORMAL_CAT 121 non-null int64
dtypes: float64(11), int64(8)
memory usage: 18.1 KB
```

2. Reemplazamiento de Valores Nulos del Conjunto de Datos

Vamos a sustituir los valores nulos por la mediana de la columna correspondiente

```
In [5]: # Vemos las filas del DataFrame que contienen valores nulos
df.isna().any()

Out[5]:
```

OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
0	0.0	True	True	True	True	True	True	True	True	True	True	True	True
1	1.0	NaN	True	True	True	True	True	True	True	True	True	True	True
2	0.0	1.0	True	True	True	True	True	True	True	True	True	True	True
3	1.0	2.0	True	True	True	True	True	True	True	True	True	True	True
4	0.0	2.0	True	True	True	True	True	True	True	True	True	True	True
...
116	NaN	16.0	102	4.0	1.4	141.0	NaN	23	0	0.0	0	1.0	0
117	NaN	16.0	107	4.0	1.4	149.0	NaN	26	0	0.0	0	0.0	0
118	1.0	1.0	140	4.0	NaN	211.0	NaN	21	0	0.0	0	2.0	0
119	0.0	1.0	198	45019.0	NaN	235.0	NaN	17	0	0.0	0	0.0	0
120	1.0	1.0	135	4.0	1.5	219.0	NaN	22	0	0.0	0	0.0	0

64 rows x 19 columns

```
In [7]: # Copiamos el conjunto de datos para no alterar el original
df_copia = df.copy()

In [8]: # Reemplazamos los valores nulos de cada columna con la mediana correspondiente
# de cada atributo usando la clase Imputer de Sklearn
from sklearn.impute import SimpleImputer

imputer = SimpleImputer(strategy="median")

In [9]: # Se le proporcionan los atributos numéricos (en este caso, todos los atributos del dataframe son numéricos, pero
# hay que tener cuidado ya que la clase imputer no admite valores categóricos, solo valores numéricos)
# por lo que calculo los valores
imputer.fit(df_copia)

Out[9]: SimpleImputer(strategy='median')

In [10]: # Reemplazamos los valores nulos
df_copia_nonan = imputer.transform(df_copia)

In [11]: # Transformamos el resultado a un DataFrame de Pandas
df_copia = pd.DataFrame(df_copia_nonan, columns=df.columns)

In [12]: df_copia.head(10)
```

```
Out[12]:
```

OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
0	0.0	0.0	112.0	4.0	1.5	174.0	NaN	31	0	0.0	0	3.0	0
1	1.0	4.0	108.0	4.0	1.2	128.0	1.0	29	23.0	19.0	24.0	3.0	4.0
2	0.0	1.0	123.0	4.0	1.1	133.0	1.0	36	30.0	30.0	30.0	1.0	4.0
3	1.0	2.0	131.0	4.0	1.5	191.0	1.0	14	0	21.0	14.0	1.0	0
4	0.0	2.0	156.0	4.0	1.2	182.0	1.0	14	0	16.0	17.0	1.0	0
...
116	NaN	16.0	102	4.0	1.4	141.0	NaN	23	0	0.0	0	1.0	0
117	NaN	16.0	107	4.0	1.4	149.0	NaN	26	0	0.0	0	0.0	0
118	1.0	1.0	140	4.0	NaN	211.0	NaN	21	0	0.0	0	2.0	0
119	0.0	1.0	198	45019.0	NaN	235.0	NaN	17	0	0.0	0	0.0	0
120	1.0	1.0	135	4.0	1.5	219.0	NaN	22	0	0.0	0	0.0	0

```
In [13]: # Comprobamos que el DataFrame ya no tiene valores nulos
df_copia.isna().any()

Out[13]:
```

OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
0	0.0	False	False	False	False	False	False	False	False	False	False	False	False
1	1.0	False	False	False	False	False	False	False	False	False	False	False	False
2	0.0	False	False	False	False	False	False	False	False	False	False	False	False
3	1.0	False	False	False	False	False	False	False	False	False	False	False	False
4	0.0	False	False	False	False	False	False	False	False	False	False	False	False
...
116	NaN	False	False	False	False	False	False	False	False	False	False	False	False
117	NaN	False	False	False	False	False	False	False	False	False	False	False	False
118	1.0	False	False	False	False	False	False	False	False	False	False	False	False
119	0.0	False	False	False	False	False	False	False	False	False	False	False	False
120	1.0	False	False	False	False	False	False	False	False	False	False	False	False

3. Reemplazamiento de Valores Atípicos del Conjunto de Datos

Usamos el método intercuartílico (IQR) para detectar los outliers del dataframe

```
In [14]: def detectar_outliers(df):
"""
Detecta y devuelve los valores atípicos en cada columna de un DataFrame utilizando el método IQR.
"""
outliers_dict = {}

for columna in df.columns:
    # Calcular el IQR (rango intercuartílico) para la columna actual
    Q1 = df[columna].quantile(0.25)
    Q3 = df[columna].quantile(0.75)
    IQR = Q3 - Q1

    # Calcular los límites para identificar outliers en la columna actual
    limite_inferior = Q1 - 1.5 * IQR
    limite_superior = Q3 + 1.5 * IQR

    # Identificar los valores atípicos en la columna actual
    outliers = df[(df[columna] < limite_inferior) | (df[columna] > limite_superior)][columna]

    # Almacenar los valores atípicos en el diccionario de resultados
    outliers_dict[columna] = outliers

return outliers_dict

In [15]: outliers_df_copia = detectar_outliers(df_copia)
print(outliers_df_copia)
```

```
{'OJO': Series([], Name: OJO, dtype: float64), 'TIPO_GLAUCOMA': 89    14.0
103    15.0
104    15.0
105    15.0
111    16.0
112    16.0
113    16.0
114    16.0
115    16.0
116    16.0
117    16.0
Name: TIPO_GLAUCOMA, dtype: float64, 'N_IMPACTOS': 14    246.0
24    262.0
25    0.0
68    246.0
101    0.0
119    198.0
Name: N_IMPACTOS, dtype: float64, 'CUADRANTES': 20    3.0
25    0.0
26    3.5
29    3.5
32    3.0
33    3.0
25    3.0
37    3.0
83    2.0
90    3.0
91    3.0
92    3.0
93    3.0
101    0.0
119    45019.0
Name: CUADRANTES, dtype: float64, 'ENERGIA_IMPACTO': 25    0.0
75    2.4
Name: ENERGIA_IMPACTO, dtype: float64, 'ENERGIA_TOTAL': 7    361.0
14    387.0
17    312.0
19    326.0
24    388.0
25    0.0
68    387.0
71    0.0
74    0.0
101    0.0
Name: ENERGIA_TOTAL, dtype: float64, 'CIRUJIA_PREVIA': 5    0.0
6    0.0
7    0.0
9    0.0
10    0.0
14    0.0
18    0.0
20    0.0
35    0.0
36    0.0
37    0.0
41    0.0
45    0.0
46    0.0
Name: CIRUJIA_PREVIA, dtype: float64, 'PIO_PRE_SLT': 15    46.0
25    0.0
30    46.0
25    29.0
74    0.0
Name: PIO_PRE_SLT, dtype: float64, 'PIO_1_SEMANA': Series([], Name: PIO_1_SEMANA, dtype: float64), 'PIO_1_MES': Series([], Name: PIO_1_MES, dtype: float64), 'PIO_3_MES': 46    43.0
Name: PIO_3_MES, dtype: float64, 'FARMACOS_PRE': Series([], Name: FARMACOS_PRE, dtype: float64), 'FARMACOS_1_MES': Series([], Name: FARMACOS_1_MES, dtype: float64), 'FARMACOS_3_MES': 1    4.0
2    4.0
5    3.0
6    3.0
8    2.0
10    2.0
14    1.0
18    4.0
19    2.0
20    3.0
27    2.0
30    4.0
31    2.0
32    2.0
34    2.0
35    2.0
48    2.0
49    2.0
51    1.0
52    1.0
Name: FARMACOS_3_MES, dtype: float64, 'DOLOR': 0    0.0
7    0.0
8    0.0
Name: DOLOR, dtype: float64, 'SEXO': 1    0.0
2    0.0
7    0.0
9    0.0
13    0.0
18    0.0
21    0.0
22    0.0
27    0.0
28    0.0
31    0.0
32    0.0
33    0.0
34    0.0
35    0.0
37    0.0
44    0.0
45    0.0
46    0.0
47    0.0
48    0.0
49    0.0
53    0.0
54    0.0
58    0.0
59    0.0
Name: SEXO, dtype: float64, 'EDAD': Series([], Name: EDAD, dtype: float64), 'PIO_NORMAL': Series([], Name: PIO_NORMAL, dtype: float64), 'PIO_NORMAL_CAT': Series([], Name: PIO_NORMAL_CAT, dtype: float64)}
```

El outlier más a destacar es el de la fila 119 de la columna "CUADRANTES", luego, solo vamos a sustituir dicho outlier por la mediana de la columna correspondiente

```
In [16]: df_copia.at[119, "CUADRANTES"] = df_copia["CUADRANTES"].median()

In [17]: df_copia

Out[17]:
```

OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
0	0.0	0.0	112.0	4.0	1.5	174.0	1.0	31.0	0.0	0.0	0.0	3.0	0.0
1	1.0	4.0	108.0	4.0	1.2	128.0	1.0	29.0	23.0	19.0	24.0	3.0	4.0
2	0.0	1.0	123.0	4.0	1.1	133.0	1.0	36.0	30.0	30.0	30.0	1.0	4.0
3	1.0	2.0	131.0	4.0	1.5	191.0	1.0	14.0	0.0	21.0	14.0	1.0	0.0
4	0.0	2.0	156.0	4.0	1.2	182.0	1.0	14.0	0	16.0	17.0	1.0	0
...
116	0.0	16.0	102.0	4.0	1.4	141.0	1.0	23.0	0	0.0	0	1.0	0.0
117	0.0	16.0	107.0	4.0	1.4	149.0	1.0	26.0	0	0.0	0	0.0	0.0
118	1.0	1.0	140.0	4.0	1.5	211.0	1.0	21.0	0	0.0	0	2.0	0.0
119	0.0	1.0	198.0	4.0	1.5	235.0	1.0	17.0	0	0.0	0	0.0	0.0
120	1.0	1.0	135.0	4.0	1.5	219.0	1.0	22.0	0	0.0	0	0.0	0.0

121 rows x 19 columns

4. Matriz de Correlación

Una vez tratados todos los datos del dataframe, visualizamos su matriz de correlación para poder sacar alguna conclusión

```
In [18]: # Matriz de Correlación
matriz_corr = df_copia.corr()

matriz_corr

Out[18]:
```

	OJO	TIPO_GLAUCOMA	N_IMPACTOS	CUADRANTES	ENERGIA_IMPACTO	ENERGIA_TOTAL	CIRUJIA_PREVIA	PIO_PRE_SLT	PIO_1_SEMANA	PIO_1_MES	PIO_3_MES	FARMACOS_PRE	FARMACOS_1_MES	FARMACOS_3_MES
OJO	1.000000	-0.167514	0.091126	0.011817	0.118807	0.061519	-0.077122	0.127733	0.071269	-0.036483	-0.014383	-0.027256	0.0333	
TIPO_GLAUCOMA	-0.167514	1.000000	-0.084669	0.146277	-0.048604	-0.060582	0.194778	-0.127291	-0.095710	-0.260416	-0.267442	-0.076669	-0.3334	
N_IMPACTOS	0.091126	-0.084669	1.000000	0.489310	0.294747	0.719566	-0.259216	0.200045	-0.006277	0.157316	0.236205	0.049481	0.1112	
CUADRANTES	0.011817	0.146277	0.489310	1.000000	0.218216	0.416060	0.030013	0.131072	0.086378	0.084690	0.042076	-0.009880	0.0405	
ENERGIA_IMPACTO	0.118807	-0.048604	0.294747	0.218216	1.000000	0.623847	-0.259178	0.205444	0.071358	0.085472	0.080678	0.115866	-0.0515	
ENERGIA_TOTAL	0.061519	-0.060582	0.719566	0.416060	0.623847	1.000000	-0.235751	0.204123	0.062095	0.155687	0.151153	0.115815	0.0161	
CIRUJIA_PREVIA	-0.077122	0.194778	-0.259216	0.030013	-0.259178	-0.235751	1.000000	-0.204823	-0.062799	-0.336719	-0.401157	-0.249314	-0.3917	
PIO_PRE_SLT	0.127733	-0.127291	0.200045	0.131072	0.205444	0.204123	-0.204823	1.000000	0.166887	0.329483	0.361240	-0.041010	0.2116	
PIO_1_SEMANA	0.071269	-0.095710	-0.006277	0.086378	0.071358	0.062095	-0.336719	0.329483	1.000000	0.645000	0.401204	0.222228	0.4597	
PIO_1_MES	0.071269	-0.095710	-0.006277	0.086378	0.071358	0.062095	-0.336719	0.329483	1.000000	0.645000	0.401204	0.222228	0.4597	
PIO_3_MES	-0.036483	-0.260416	0.157316	0.084690	0.085472	0.155687	-0.401157	0.361240	0.401204	1.000000	0.655155	0.126318	0.4597	
FARMACOS_PRE	-0.027256	-0.076669	0.049481	-0.009880	0.115866	0.151153	-0.249314	-0.041010	0.222228	0.222228	1.000000	0.126318	0.4597	
FARMACOS_1_MES	0.0333	-0.334635	0.112343	0.045985	-0.051138	0.016888	-0.336719	0.329483	0.459320	0.557144	0.494319	0.578539	1.0000	
FARMACOS_														