

▼ Exercici 1

Parteix el conjunt de dades adjunt en train i test. Estudia els dos conjunts per separat, a nivell descriptiu.

També adjunt trobaràs una descripció de les diferents variables del dataset.

<https://www.kaggle.com/code/prasadperera/the-boston-housing-dataset>

```

1 # Tratamiento de datos
2 # =====
3 import pandas as pd
4 import numpy as np
5
6 # Gráficos
7 # =====
8 import matplotlib.pyplot as plt
9 from matplotlib import style
10 import seaborn as sns
11
12 # Preprocesado y análisis
13 # =====
14 import statsmodels.api as sm
15 from scipy import stats
16
17 # Configuración matplotlib
18 # =====
19 plt.style.use('ggplot')
20
21 # Configuración warnings
22 # =====
23 import warnings
24 warnings.filterwarnings('ignore')

```

```

1 # Activo Google Drive
2
3 from google.colab import drive
4 drive.mount('/content/drive')

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.r



```

1 # Abro el fichero
2

```

```

3 #path='https://drive.google.com/file/d/1j8WdhTxMpngrUigfcfrJynng-9JHo9x
4 path = ('/content/drive/MyDrive/01_COLAB/housing data.csv')
5 #data = pd.read_csv(file,
6
7 df= pd.read_csv(path, sep=',', encoding="latin-1")
8 df.shape
9
10 nRow, nCol = df.shape
11 print(f'Hay {nRow} filas con {nCol} columnas')
12 df1=df.copy()
13 print('\nImprimo el primer registro, solo para ver como es:\n')
14 df1.iloc[0]

```

Hay 505 filas con 14 columnas

Imprimo el primer registro, solo para ver como es:

```

0.00632      0.02731
18.00        0.00000
2.310        7.07000
0            0.00000
0.5380       0.46900
6.5750       6.42100
65.20        78.90000
4.0900       4.96710
1            2.00000
296.0        242.00000
15.30        17.80000
396.90       396.90000
4.98         9.14000
24.00        21.60000
Name: 0, dtype: float64

```

Veo que es un fichero sin nombre en las columnas. Utilizo el fichero de TXT e identifico los nombres de las columnas.

```

1 columnas = list(df1.columns[:-1])
2 #columnas[0][0]
3 columnas

```

```

['0.00632',
'18.00',
'2.310',
'0',
'0.5380',
'6.5750',
'65.20',
'4.0900',
'1',
'296.0',
'15.30',
'396.90',
'4.98']

```

```

1 # Asigno una lista con los nombres de las columnas
2 Nombre_Columnas = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV']
3
4 df1.columns = [Nombre_Columnas]
5 df1.iloc[0]

```

```

CRIM      0.02731
ZN        0.00000
INDUS     7.07000
CHAS      0.00000
NOX       0.46900
RM        6.42100
AGE       78.90000
DIS       4.96710
RAD       2.00000
TAX      242.00000
PTRATIO   17.80000
B        396.90000
LSTAT     9.14000
MEDV     21.60000
Name: 0, dtype: float64

```

```

1 # Analizo los datos con una estadística descriptiva y confirmo que tengo los datos
2 print(df1.describe())

```

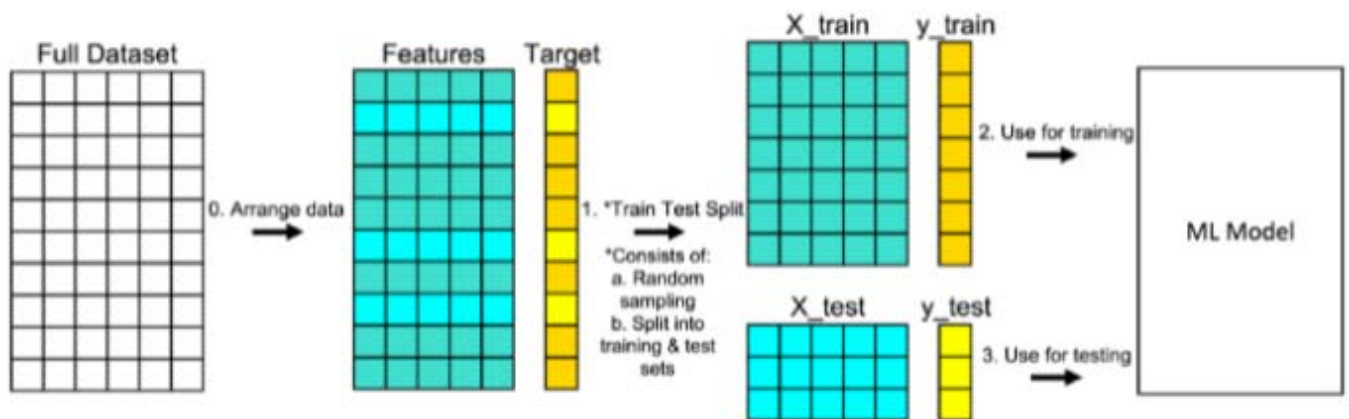
	CRIM	ZN	INDUS	CHAS	NOX	RM
count	505.000000	505.000000	505.000000	505.000000	505.000000	505.000000
mean	3.620667	11.350495	11.154257	0.069307	0.554728	6.284059
std	8.608572	23.343704	6.855868	0.254227	0.115990	0.703195
min	0.009060	0.000000	0.460000	0.000000	0.385000	3.561000
25%	0.082210	0.000000	5.190000	0.000000	0.449000	5.885000
50%	0.259150	0.000000	9.690000	0.000000	0.538000	6.208000
75%	3.678220	12.500000	18.100000	0.000000	0.624000	6.625000
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000

	AGE	DIS	RAD	TAX	PTRATIO	B
count	505.000000	505.000000	505.000000	505.000000	505.000000	505.000000
mean	68.581584	3.794459	9.566337	408.459406	18.461782	356.594376
std	28.176371	2.107757	8.707553	168.629992	2.162520	91.367787
min	2.900000	1.129600	1.000000	187.000000	12.600000	0.320000
25%	45.000000	2.100000	4.000000	279.000000	17.400000	375.330000
50%	77.700000	3.199200	5.000000	330.000000	19.100000	391.430000
75%	94.100000	5.211900	24.000000	666.000000	20.200000	396.210000
max	100.000000	12.126500	24.000000	711.000000	22.000000	396.900000

	LSTAT	MEDV
count	505.000000	505.000000
mean	12.668257	22.529901
std	7.139950	9.205991
min	1.730000	5.000000
25%	7.010000	17.000000
50%	11.380000	21.200000
75%	16.960000	25.000000
max	37.970000	50.000000

```
1 # Verifico que no hay nulos y los tipos de los datos
2 df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 505 entries, 0 to 504
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype  
---  -
 0   (CRIM,)               505 non-null   float64
 1   (ZN,)                 505 non-null   float64
 2   (INDUS,)              505 non-null   float64
 3   (CHAS,)               505 non-null   int64  
 4   (NOX,)                505 non-null   float64
 5   (RM,)                 505 non-null   float64
 6   (AGE,)                505 non-null   float64
 7   (DIS,)                505 non-null   float64
 8   (RAD,)                505 non-null   int64  
 9   (TAX,)                505 non-null   float64
10  (PTRATIO,)            505 non-null   float64
11  (B,)                  505 non-null   float64
12  (LSTAT,)              505 non-null   float64
13  (MEDV,)               505 non-null   float64
dtypes: float64(12), int64(2)
memory usage: 55.4 KB
```



Vamos a Dividir la base de datos en Train y en Test

```
1 !pip install scikit-learn
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: numpy>=1.14.6 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: scipy>=1.1.0 in /usr/local/lib/python3.7/dist-packages
```

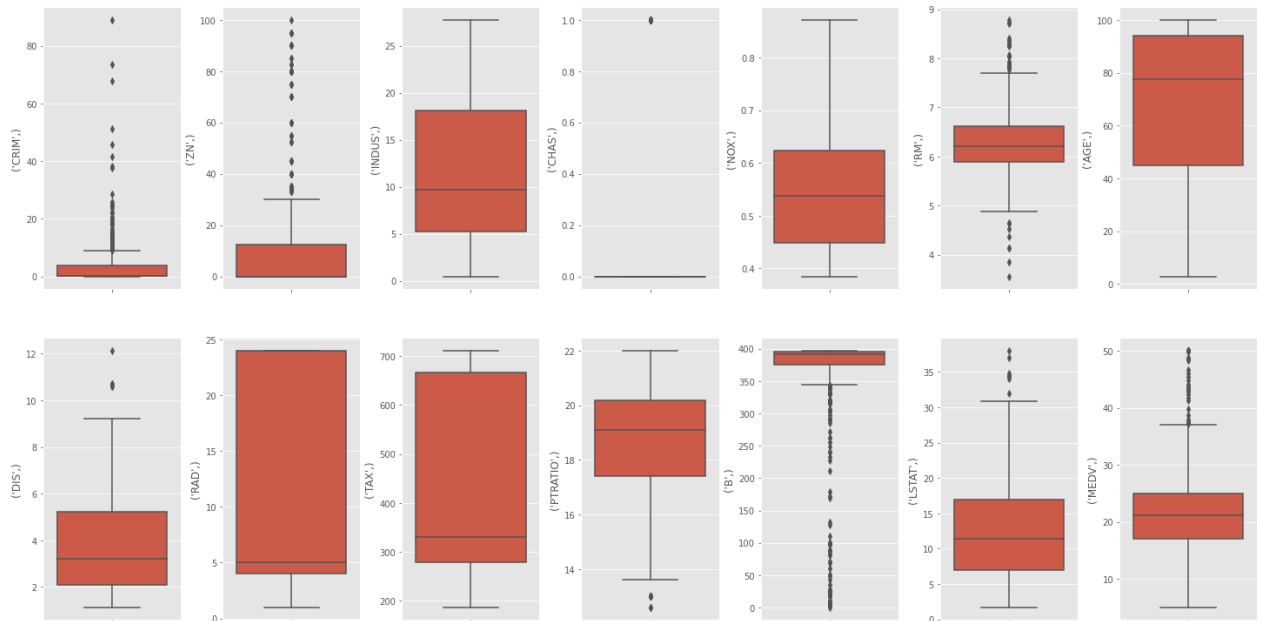
```
1
2 from sklearn.model_selection import train_test_split
3
```

```
4 features = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX']
5 #features = columnas
6 x = df1.loc[:, features]
7 y = df1.loc[:, ['MEDV']]
8
9 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size =0.2,
10                                                random_state = 0)
11
12 print("xtrain shape : ", xtrain.shape)
13 print("xtest shape : ", xtest.shape)
14 print("ytrain shape : ", ytrain.shape)
15 print("ytest shape : ", ytest.shape)

xtrain shape : (404, 5)
xtest shape : (101, 5)
ytrain shape : (404, 1)
ytest shape : (101, 1)
```

Vamos a ver gráficamente como es cada base de datos (train y test)

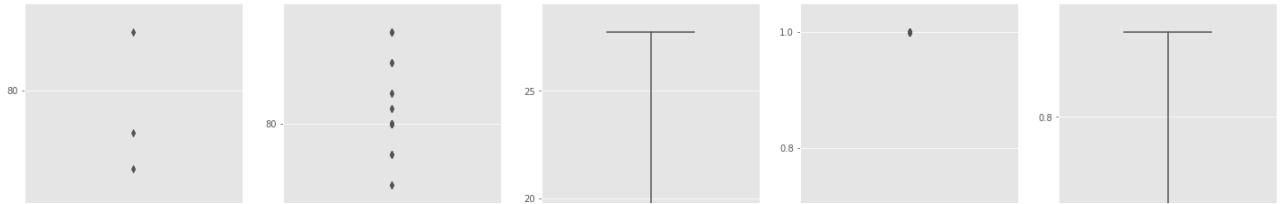
```
1 # Dibujamos todas las columnas en su forma original
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 from scipy import stats
5
6 fig, axs = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
7 index = 0
8
9 axs = axs.flatten()
10
11 for k,v in df1.items():
12     sns.boxplot(y=k, data=df1, ax=axs[index])
13     index += 1
14 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)
```



```

1 # Ahora dibujo los datos de TRAIN de las columnas 'CRIM','ZN','INDUS','L
2
3 fig, axs = plt.subplots(ncols=5, nrows=1, figsize=(20, 10))
4 index = 0
5 axs = axs.flatten()
6 for k,v in xtrain.items():
7     sns.boxplot(y=k, data=xtrain, ax=axs[index])
8     index += 1
9 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)

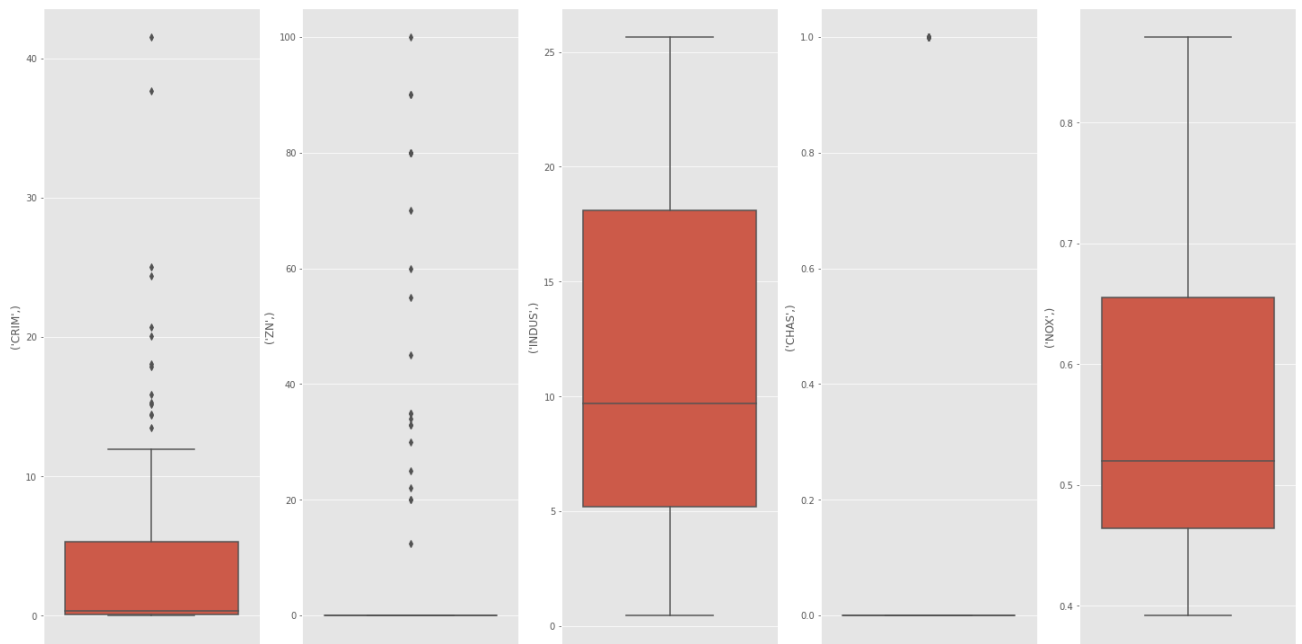
```



```

1 fig, axs = plt.subplots(ncols=5, nrows=1, figsize=(20, 10))
2 index = 0
3 axs = axs.flatten()
4 for k,v in xtest.items():
5     sns.boxplot(y=k, data=xtest, ax=axs[index])
6     index += 1
7 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)

```



Conclusion del split:

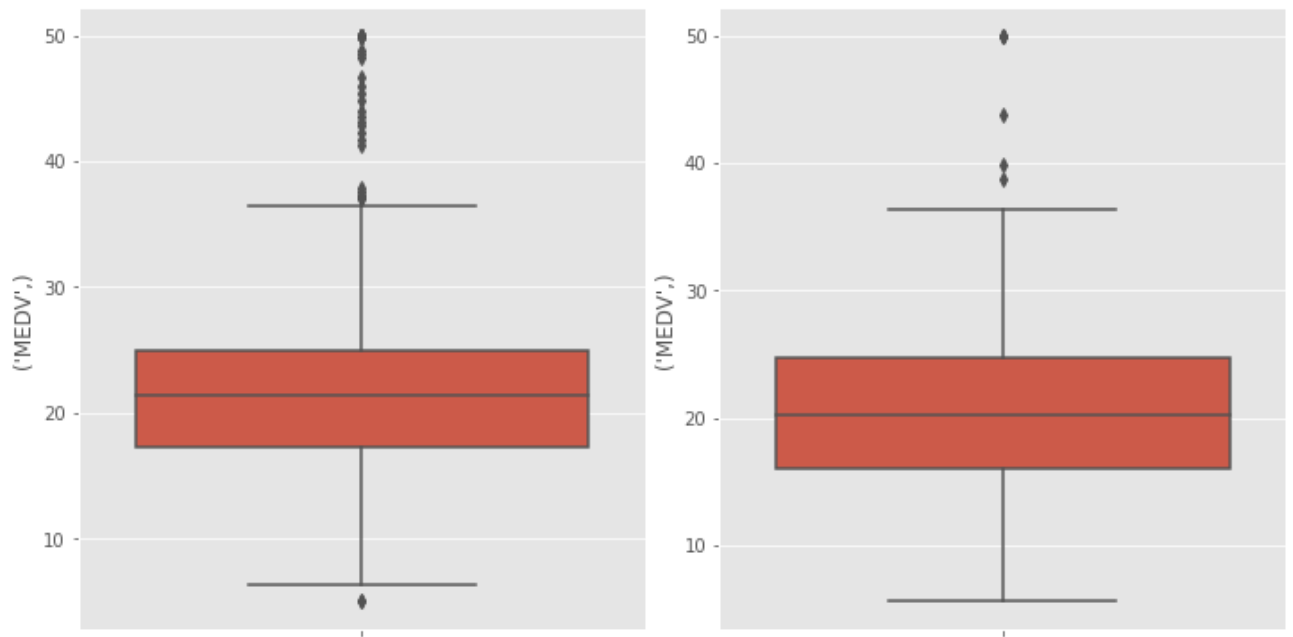
Veo que las formas son muy semejantes por lo que considero que la división que hace es aleatoria y los datos son representativos de la base de trabajo original.

Compruebo graficamente (boxplot y aproximacion a curva normal) con la variable de salida

```

1 #Dibujo ytrain e ytest
2
3 fig, axs = plt.subplots(ncols=2, nrows=1, figsize=(10, 5))
4 index = 0
5 axs = axs.flatten()
6
7 for k,v in ytrain.items():
8     sns.boxplot(y=k, data=ytrain, ax=axs[0])
9     sns.boxplot(y=k, data=ytest, ax=axs[1])
10    index += 1
11 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)

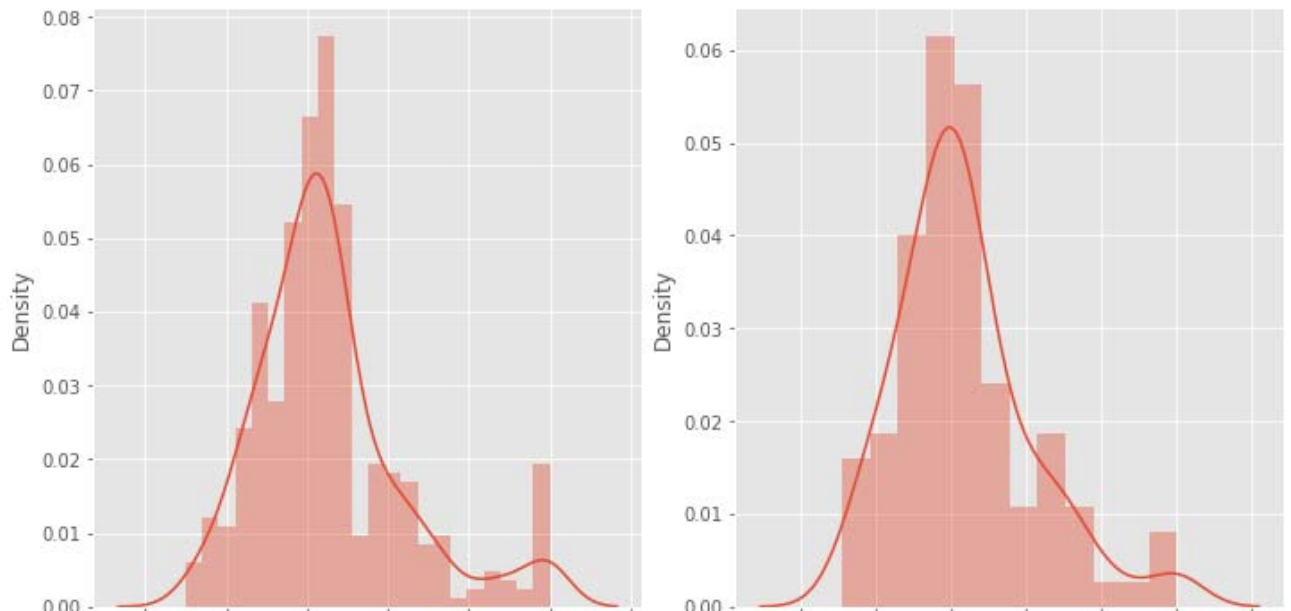
```



```

1
2 fig, axs = plt.subplots(ncols=2, nrows=1, figsize=(10, 5))
3 index = 0
4 axs = axs.flatten()
5 #for k,v in ytrain.items():
6     #sns.distplot(v, ax=axs[index])
7 sns.distplot(ytrain['MEDV'], ax=axs[0])
8 sns.distplot(ytest['MEDV'], ax=axs[1])
9     #sns.distplot(v, ax=axs[index])
10
11     #index += 1
12 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)

```

Veo que visualmente las dos graficas de los datos split son muy parecidas, por lo que deduzco que lo ha hecho correctamente.

▼ - Exercici 2

Aplica algun procés de transformació (estandarditzar les dades numèriques, crear columnes dummies, polinomis...).

```
1 from sklearn.preprocessing import MinMaxScaler
2 # define data
3
4 print('Media de CRIM:',xtrain['CRIM'].mean())
5 # define min max scaler
6 scaler = MinMaxScaler()
7 # transform data
8 scaled = scaler.fit_transform(xtrain['CRIM'])
9 print('Media de CRIM transformado ',scaled.mean())
```

```
Media de CRIM: CRIM      3.446082
dtype: float64
Media de CRIM transformado  0.03863248924499451
```

para que sirve MinMaxScaler:

Transforma las características escalando cada una de ellas a un rango determinado. Este estimador escala y traduce cada característica individualmente de manera que se encuentre en el rango dado en el conjunto de entrenamiento, por ejemplo, entre cero y uno.

```
1 # Normalizo la columna CRIM
2
```

```

3 from sklearn.preprocessing import Normalizer
4 scaler = Normalizer().fit(xtrain['CRIM'])
5 normalized_CRIM = scaler.transform(xtrain['CRIM'])
6
7 print(normalized_CRIM[3])

[1.]

```

Para que sirve normalizar:

Normalizar las muestras individualmente a la norma de la unidad. Cada muestra (es decir, cada fila de la matriz de datos) con al menos un componente no nulo se reescala independientemente de las demás muestras para que su norma (l_1 , l_2 o inf) sea igual a uno.

Se utiliza en el machine learning, porque algunos valores de características difieren de otros varias veces. Las características **con valores más altos** dominarán el proceso de aprendizaje.

```

1 # BINARIZACION
2
3 from sklearn.preprocessing import Binarizer
4
5
6 binarizer = Binarizer(threshold=0.2).fit(xtrain['CRIM'])
7 binary_CRIM = binarizer.transform(xtrain['CRIM'])
8

```

Para que sirve binarizar:

Binarizar los datos (establecer los valores de las características en 0 o 1) según un umbral. en este caso es 0.2

Los valores superiores al mapa de umbral a 1, mientras que los valores inferiores o iguales al mapa de umbral a 0. Con el umbral por defecto de 0, sólo los valores positivos se asignan a 1.

La binarización es una operación común en los datos de recuento de texto en la que el analista puede decidir considerar únicamente la presencia o ausencia de una característica en lugar de un número cuantificado de ocurrencias, por ejemplo.

```
class sklearn.dummy.DummyRegressor(*, strategy='mean', constant=None, quantile=None)
```

<https://scikit-learn.org/stable/modules/generated/sklearn.dummy.DummyRegressor.html>

```

1 #Generating Polynomial Features
2
3 from sklearn.preprocessing import PolynomialFeatures
4 poly = PolynomialFeatures(2)

```

```

5 polynomial_CRIMS = poly.fit_transform(xtrain['CRIM'])
6 print('entrada CRIM:\n..... 3 primeras filas')
7 print( xtrain['CRIM'][:3])
8 print ('\\nColumnas nuevas, transformadas: ')
9 print(polynomial_CRIMS[:3])

```

```

entrada CRIM:
..... 3 primeras filas
      CRIM
261  0.52014
71   0.09164
479  5.82401

```

```

Columnas nuevas, transformadas:
[[1.00000000e+00 5.20140000e-01 2.70545620e-01]
 [1.00000000e+00 9.16400000e-02 8.39788960e-03]
 [1.00000000e+00 5.82401000e+00 3.39190925e+01]]

```

```
1 xtrain['CRIM'][:3]
```

	CRIM
261	0.52014
71	0.09164
479	5.82401

Conclusion de PolynomialFeatures:

Cuando aplicamos la transformacion polinómica de tipo 2 a una sola columna (CRIM) lo que estamos haciendo es las siguientes operacion, obtenemos otras 3 columnas nuevas con el valor introducido:

$$X \Rightarrow 1, X, X^2$$

por ejemplo 0.052014 --> 1, 0.052014, 0.052014² --> [1, 0.520, 0.270]

El problema es que todas las columnas las veremos como correlacionadas

Y si hubieramos introducido 2 columnas entonces hubiera creados estas columnas nuevas:

$$X_1, X_2 = 1, X_1, X_2, X_1^2, X_1 * X_2, X_2^2$$

Esto nos va ser muy útil cuando calculemos los residuos para ver como se acomoda una ecuación.

Dibujo todas las transformaciones:

```

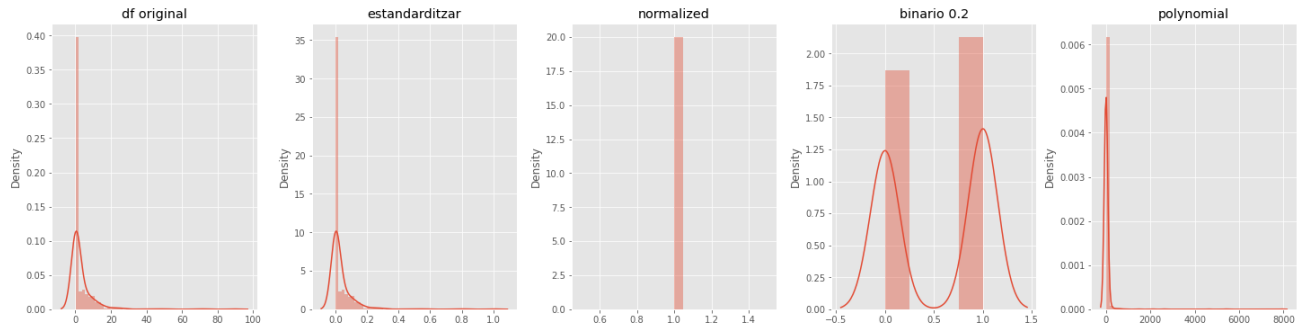
1 fig, axs = plt.subplots(ncols=5, nrows=1, figsize=(20, 5))
2 index = 0
3 axs = axs.flatten()
4 #for k,v in ytrain.items():
5     #sns.distplot(v, ax=axs[index])

```

```

6 sns.distplot(xtrain['CRIM'], ax=axes[0]).set(title='df original')
7 sns.distplot(scaled, ax=axes[1]).set(title='estandarditzar')
8 sns.distplot(normalized_CRIM, ax=axes[2]).set(title='normalized')
9 sns.distplot(binary_CRIM, ax=axes[3]).set(title='binario 0.2')
10 sns.distplot(polynomial_CRIMS, ax=axes[4]).set(title='polynomial')
11 plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)

```



▼ Exercici 3

Resumeix les noves columnes generades de manera estadística i gràfica

Reflexión sobre las transformaciones desde un punto de vista operativo:

Es importante entender para qué queremos las transformaciones aunque el valor que obtengamos de la estadística descriptiva nos pueda hacer pensar que no es lo mismo.

Vamos a conseguir que todos los valores de todas las columnas puedan participar a la hora de crear un modelo, porque si hay valores de magnitud muy altos respecto a otros muy pequeños (ejemplo mm v Km) hace que los pequeños cuando no se transforman difícilmente entran a participar del modelo.

Otra gran ventaja es la operatividad en el cálculo que al estar transformados es mucho más rápida.

También nos permite romper la colinearidad de los datos en diferentes columnas porque por ejemplo $X \rightarrow X^2 \rightarrow X^n$ la hipótesis de correlación nos saldrá baja y querrá decir que $X \rightarrow X^2 \rightarrow X^n$ están correlacionados

```

1 binary_CRIM = pd.DataFrame(binary_CRIM)
2 binary_CRIM[:3]

```

	0
0	1.0
1	0.0
2	1.0

1 # Analizo los datos con estadística descriptiva:

2

3 print(binary_CRIM.describe())

```

count    404.000000
mean      0.532178
std       0.499582
min       0.000000
25%      0.000000
50%      1.000000
75%      1.000000
max       1.000000

```

1 normalized_CRIM = pd.DataFrame(normalized_CRIM)

2 print(normalized_CRIM.describe())

```

count    404.0
mean      1.0
std       0.0
min       1.0
25%      1.0
50%      1.0
75%      1.0
max       1.0

```

1 polynomial_CRIMS =pd.DataFrame(polynomial_CRIMS)

2 print(polynomial_CRIMS.describe())

```

count    0      1      2
count    404.0  404.000000  404.000000
mean      1.0    3.446082   88.744553
std       0.0    8.778372  561.399813
min       1.0    0.009060   0.000082
25%      1.0    0.079710   0.006354
50%      1.0    0.243125   0.059114
75%      1.0    3.202962  10.263617
max       1.0   88.976200  7916.764166

```

▼ Visualización gráfica de las transformaciones:

Voy a ver cual es el impacto de las principales transformaciones y como le afectan los outliers

```

1 features = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX']
2
3 x = df1.loc[:, features]
4 x= x['CRIM']
5 y = df1.loc[:, ['MEDV']]

1 from sklearn.model_selection import train_test_split
2
3 X_full, xtest, y_ful, ytest = train_test_split(x, y, test_size =0.2,
4                                             random_state = 0)
5
6 print("X_full shape : ", X_full.shape)
7 print("xtest shape : ", xtest.shape)
8 print("y_ful shape : ", y_ful.shape)
9 print("ytest shape : ", ytest.shape)

X_full shape : (404, 1)
xtest shape : (101, 1)
y_ful shape : (404, 1)
ytest shape : (101, 1)

```

https://runebook.dev/es/docs/scikit_learn/auto_examples/preprocessing/plot_all_scaling

Aprovecharé esta parte para dar mejor visualizaciona los datos. <https://github.com/puchee99?tab=repositories>

```

1 # Aquí tengo las distribuciones que voy a trabajar.
2
3 distributions = [
4     ('Datos sin escalar', X),
5     ('Datos después de escalar',
6      StandardScaler().fit_transform(X)),
7     ('Datos despues de aplicar min-max scalado',
8      MinMaxScaler().fit_transform(X)),
9     ('Datos después de de la transformacion de cuartiles (uniforme)',
10      QuantileTransformer(output_distribution='uniform')
11      .fit_transform(X)),
12     ('Datos después de la transformacion (gaussiana)',
13      QuantileTransformer(output_distribution='normal')
14      .fit_transform(X)),
15 ]
16

```

```

1 #@title
2
3

```



```
4 dataset = fetch_california_housing()
5
6 X_full, y_full = dataset.data, dataset.target
7
8 #
9 X = X_full[:, [0, 2]]
10
11
12 # escalar la salida entre 0 y 1 para la barra de colores
13 y = minmax_scale(y_full)
14
15 # plasma no existe en matplotlib <1,5
16 cmap = getattr(cm, 'plasma_r', cm.hot_r)
17
18 #
19 def create_axes(title, figsize=(16, 6)):
20     fig = plt.figure(figsize=figsize)
21     fig.suptitle(title)
22
23     # definir el eje para el primer gráfico
24     left, width = 0.1, 0.22
25     bottom, height = 0.1, 0.7
26     bottom_h = height + 0.15
27     left_h = left + width + 0.02
28
29     rect_scatter = [left, bottom, width, height]
30     rect_histx = [left, bottom_h, width, 0.1]
31     rect_histy = [left_h, bottom, 0.05, height]
32
33     ax_scatter = plt.axes(rect_scatter)
34     ax_histx = plt.axes(rect_histx)
35     ax_histy = plt.axes(rect_histy)
36
37     # definir el eje para el gráfico ampliado
38     left = width + left + 0.2
39     left_h = left + width + 0.02
40
41     rect_scatter = [left, bottom, width, height]
42     rect_histx = [left, bottom_h, width, 0.1]
43     rect_histy = [left_h, bottom, 0.05, height]
44
45     ax_scatter_zoom = plt.axes(rect_scatter)
46     ax_histx_zoom = plt.axes(rect_histx)
47     ax_histy_zoom = plt.axes(rect_histy)
48
49     # definir el eje de la barra de colores
50     left, width = width + left + 0.13, 0.01
```

```

51
52 rect_colorbar = [left, bottom, width, height]
53 ax_colorbar = plt.axes(rect_colorbar)
54
55 return ((ax_scatter, ax_histy, ax_histx),
56         (ax_scatter_zoom, ax_histy_zoom, ax_histx_zoom),
57         ax_colorbar)
58
59 # Dibujar las distribuciones:
60
61 def plot_distribution(axes, X, y, hist_nbins=50, title="",
62                     x0_label="", x1_label=""):
63     ax, hist_X1, hist_X0 = axes
64
65     ax.set_title(title)
66     ax.set_xlabel(x0_label)
67     ax.set_ylabel(x1_label)
68
69     # El diagrama de dispersión
70     colors = cmap(y)
71     ax.scatter(X[:, 0], X[:, 1], alpha=0.5, marker='o', s=5, lw=0, c=co
72
73     # Retirar la parte superior y la columna derecha por estética
74     # hacer un buen diseño de ejes
75     ax.spines['top'].set_visible(False)
76     ax.spines['right'].set_visible(False)
77     ax.get_xaxis().tick_bottom()
78     ax.get_yaxis().tick_left()
79     ax.spines['left'].set_position(('outward', 10))
80     ax.spines['bottom'].set_position(('outward', 10))
81
82     # Histograma para el eje X1 (función 5)
83     hist_X1.set_ylim(ax.get_ylim())
84     hist_X1.hist(X[:, 1], bins=hist_nbins, orientation='horizontal',
85                 color='grey', ec='grey')
86     hist_X1.axis('off')
87
88     # Histograma para el eje X0 (función 0)
89     hist_X0.set_xlim(ax.get_xlim())
90     hist_X0.hist(X[:, 0], bins=hist_nbins, orientation='vertical',
91                 color='grey', ec='grey')
92     hist_X0.axis('off')

```

1 #creamos los plots para ver como varia con cada transformación:

2

3




```

4 #@title
5
6 def make_plot(item_idx):
7     titulo, X = distributions[item_idx]
8     print('--> ', pd.DataFrame(X[0:3])) # Para ver que datos saca
9
10    ax_zoom_out, ax_zoom_in, ax_colorbar = create_axes(titulo)
11    axarr = (ax_zoom_out, ax_zoom_in)
12    plot_distribution(axarr[0], X, y, hist_nbins=200,
13                    x0_label="Media de Ingresos",
14                    x1_label="Cantidad de edificios_1",
15                    title="Todos los datos")
16
17    # zoom-in... Para quitar los outliers
18    zoom_in_percentile_range = (0, 99)
19    cutoffs_X0 = np.percentile(X[:, 0], zoom_in_percentile_range)
20    cutoffs_X1 = np.percentile(X[:, 1], zoom_in_percentile_range)
21
22    non_outliers_mask = (
23        np.all(X > [cutoffs_X0[0], cutoffs_X1[0]], axis=1) &
24        np.all(X < [cutoffs_X0[1], cutoffs_X1[1]], axis=1))
25    plot_distribution(axarr[1], X[non_outliers_mask], y[non_outliers_ma
26                    hist_nbins=50,
27                    x0_label="Media de Ingresos",
28                    x1_label="Cantidad de edificios",
29                    title="Sin outliers")
30
31 #Dibujamos la escal de colores de y
32    norm = mpl.colors.Normalize(y_full.min(), y_full.max())
33    mpl.colorbar.ColorbarBase(ax_colorbar, cmap=cmap,
34                              norm=norm, orientation='vertical',
35                              label='Color con la escala de y')

```

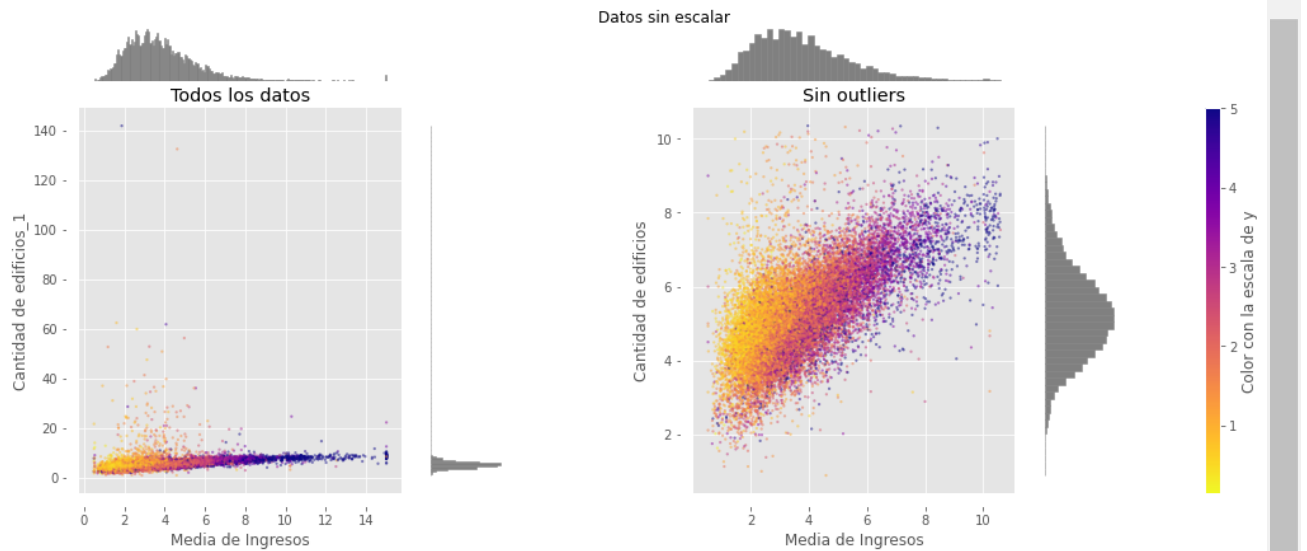
Dibujo todas las transformaciones y veo la importancia en
 el posible impacto de los outliers en el calculo del modelo:

[Mostrar código](#)

```

0 8.3252 6.984127
1 8.3014 6.238137
2 7.2574 8.288136

```



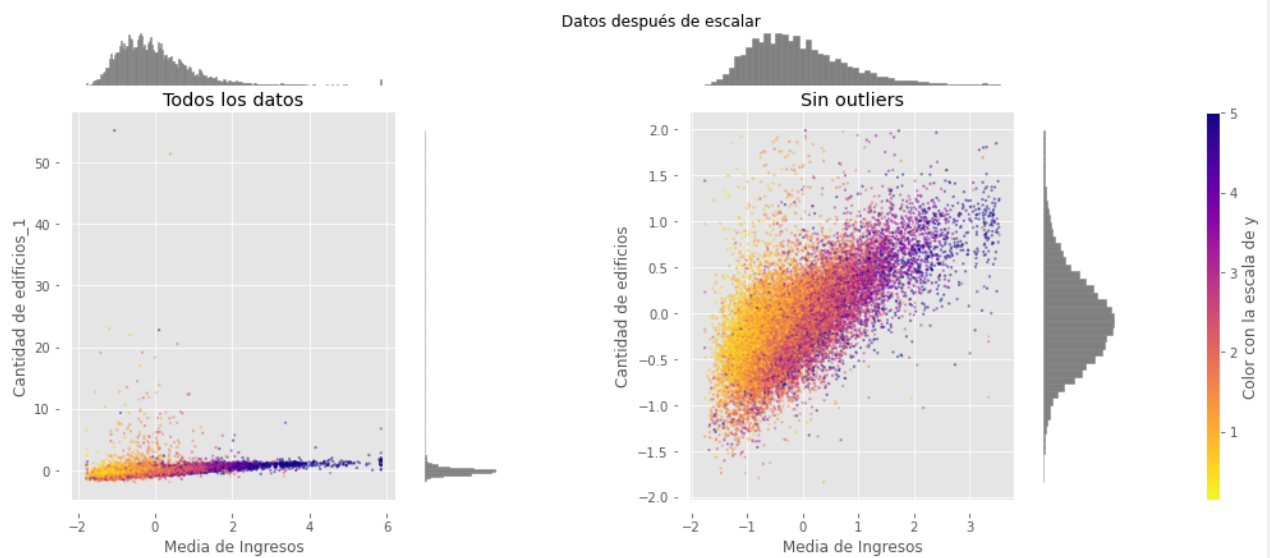
1 estudio: Datos después de escalar

```
-->      0      1
```

```

0 2.344766 0.628559
1 2.332238 0.327041
2 1.782699 1.155620

```



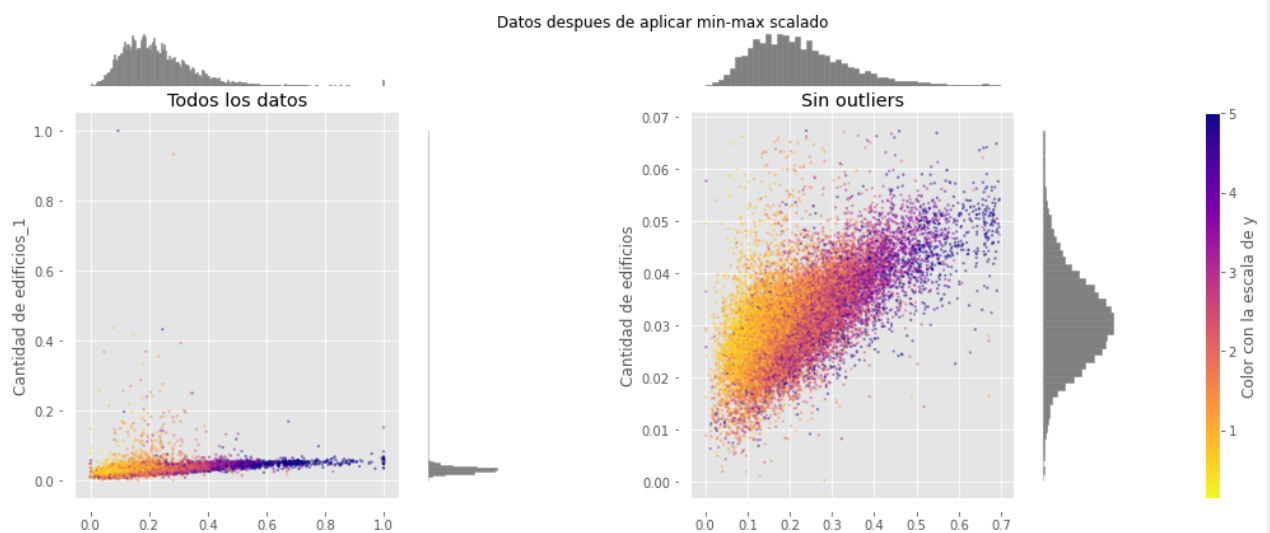
2 estudio: Datos despues de aplicar min-max scalado

```
-->      0      1
```

```

0 0.539668 0.043512
1 0.538027 0.038224
2 0.466028 0.052756

```

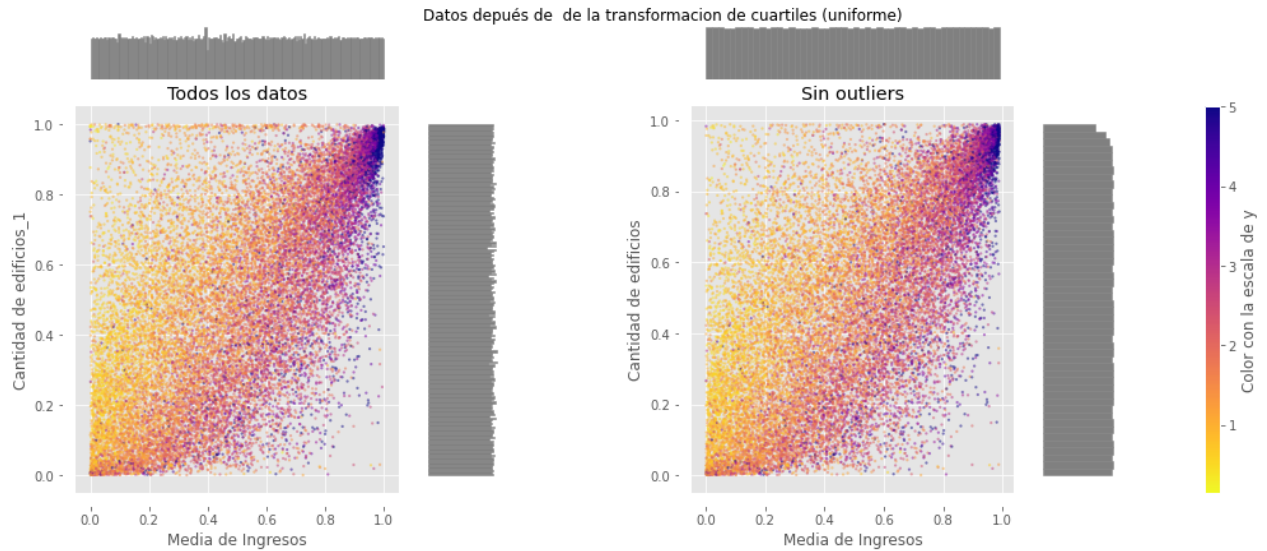


Media de Ingresos

Media de Ingresos

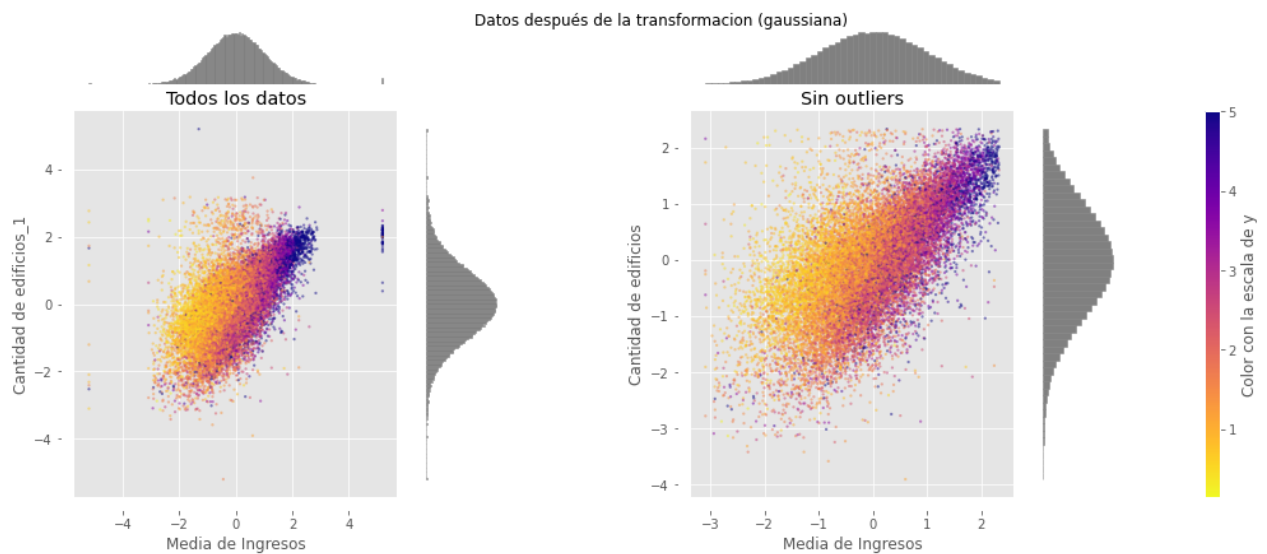
3 estudio: Datos después de de la transformacion de cuartiles (uniforme)

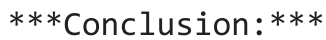
```
-->      0      1
0  0.972520  0.902275
1  0.972217  0.794595
2  0.948467  0.973393
```



4 estudio: Datos después de la transformacion (gaussiana)

```
-->      0      1
0  1.919185  1.294622
1  1.914418  0.822470
2  1.630171  1.933184
```





Conclusion:

Como hemos visto en los ejercicios transformar los datos será algo muy importante para hacer un buen modelo, evitar colinealidad, evitar la influencia solo de valores grandes, incrementar velocidad de cálculo, etc