

Parte 1

Aqui importaremos librerias

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

from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, plot_confusion_matrix

from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error

%matplotlib notebook
```

```
In [2]: ## Se agrega el DataFrame que se trabajará en la actividad
df = pd.read_csv('celulares.csv', sep=',')
## Se desprende que la cantidad de filas y columnas que tenemos en dicha DataFrame = df son 2200 rows y 22 cols
df.info()
## mediante el comando df.dtypes() se eliminaron todas las filas que contuvieran valores nulos del data frame,
## esto debido a que no se le puede atribuir un valor en específico y sin ello, nosotros no podemos ver como va
## Por otro lado, por ejemplo, en caso de que los componentes que tienen valores NaN sean de gama baja o muy bi
## esto podría aumentar el presupuesto para los demás componentes, con lo cual, se disminuiría el dataframe.
df
```

```
Out[2]: Unnamed:  battery_power  blue  clock_speed  dual_sim  fc  four_g  int_memory  m_dep  mobile_wt  ...  px_height  px_width  ram
0  1  16630  1.0  0.5  1.0  13.0  1.0  270  0.3  1690  ...  8310  14390  20840
1  0  16630  1.0  0.5  1.0  13.0  0.0  400  0.0  1470  ...  9510  15450  13860
2  2  13230  1.0  2.5  1.0  10.0  1.0  280  0.2  1310  ...  1620  6190  18920
3  3  10990  0.0  0.5  0.0  13.0  1.0  610  0.3  1460  ...  3930  10960  16990
4  4  15630  0.0  1.7  1.0  10.0  0.0  160  0.1  1510  ...  4100  5720  39220
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
2195 2195 18870 0.0 2.3 0.0 0.0 1.0 9.0 0.1 1910  ... 7120 14420 9900
2196 2196 5410 1.0 2.3 0.0 4.0 0.0 510 0.4 2000  ... 10120 12260 4030
2197 2197 18140 0.0 1.4 1.0 1.0 1.0 9.0 0.4 1410  ... 7560 7860 35560
2198 2198 10270 1.0 2.2 0.0 0.0 0.0 630 0.8 1020  ... 1520 7140 17320
2199 2199 18870 1.0 1.9 0.0 4.0 1.0 62.0 0.5 830  ... 4190 7360 27570
2000 rows x 22 columns
```

```
In [3]: df[['battery_power', 'pc', 'px_width', 'ram', 'three_g', 'wifi', 'price_range']].groupby('price_range').describe()
```

```
Out[3]: battery_power  pc  ...  three_g
price_range count mean 25% 75% max count mean 75% max count mean std min
0.0 5000 1126868 438614528 5010 84300 12060 159625 19960 5000 9374 ... 1.0 1.0 5000 0.04 0.500485 0.0
1.0 5000 1228320 452863065 5010 81650 12195 163575 19980 5000 10018 ... 1.0 1.0 5000 0.04 0.500485 0.0
2.0 5000 1379384 414992261 5030 103475 14495 173330 19940 5000 10150 ... 1.0 1.0 5000 0.04 0.499924 0.0
3.0 5000 1379384 414992261 5030 103475 14495 173330 19940 5000 10150 ... 1.0 1.0 5000 0.04 0.499924 0.0
4 rows x 18 columns
```

No influyen en la predicción

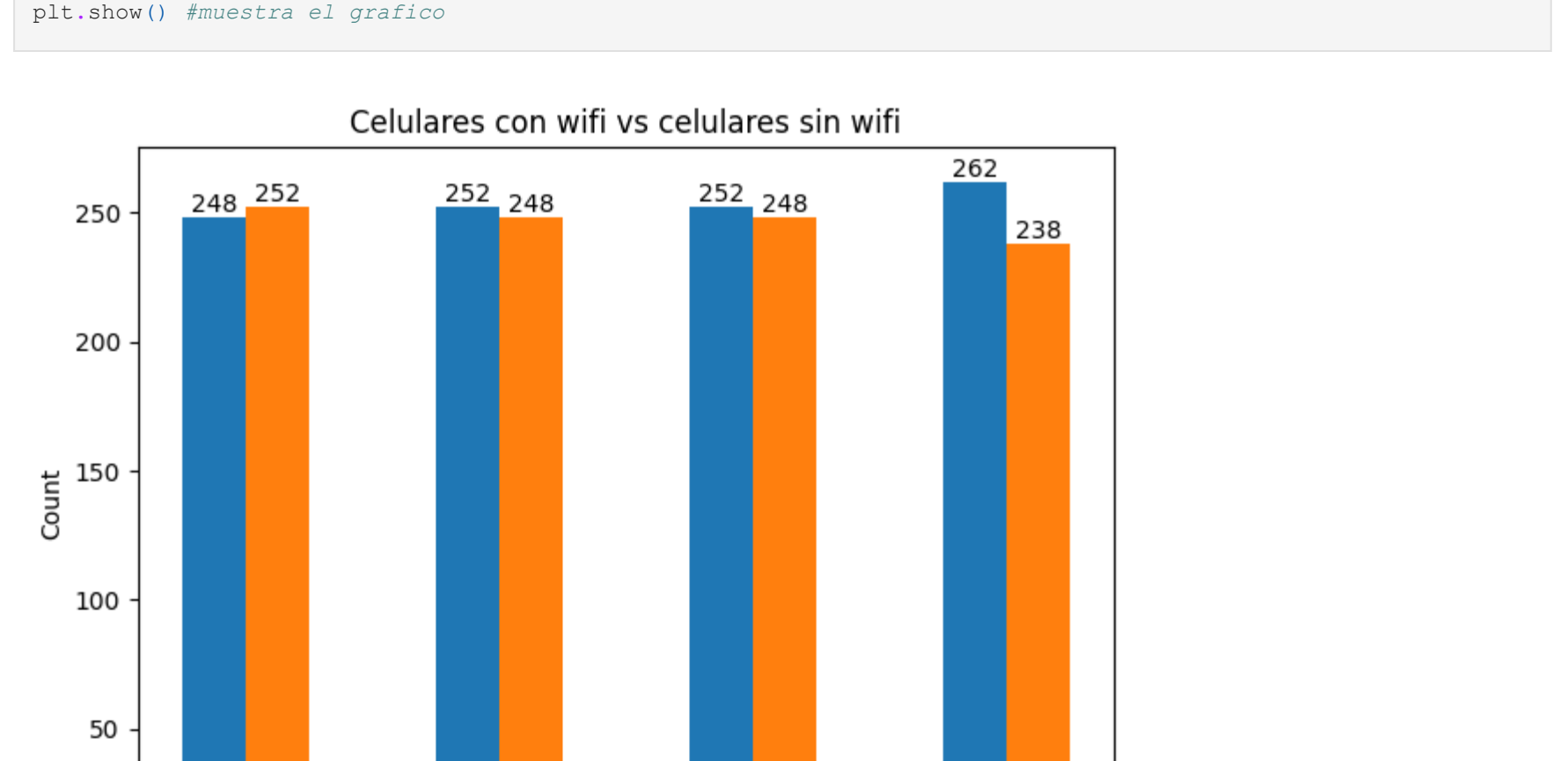
- unnamed -> No se toma en cuenta debido a que no aporta ninguna diferencia al momento de clasificar y poder hacer la selección.
- blue -> No se toma en cuenta debido a que los promedios generan ciertas inconsistencias en cada rango de precio, es decir, el promedio de la cámara con respecto al rango de precio 1.0 es el mismo que el promedio de celulares con bluetooth en el rango de precio 2.0. Por otro lado, vemos que al menos el 25% de los celulares de todos los rangos de precios poseen bluetooth, por lo que tener o no bluetooth no parece ser una métrica confiable para predecir el rango de precio.
- clock_speed -> No se toma en cuenta, puesto que, las mayores velocidades de ejecución de tareas por parte del procesador, los tienen aquellos ubicados en el rango de precio 0.0. Por otro lado, podemos ver que el promedio es inconsistente, pues, no es descendente ni ascendente en cuanto al rango de precio.
- dual_sim -> No se toma en cuenta debido a la inconsistencia que presentan los promedios en el rango de precios 1.0 y 2.0, esto, a causa de que hay más celulares que poseen dual_sim en el rango de precio 1.0 que en el rango de precio 2.0.
- fc -> Esta característica no se toma en cuenta debido a que el promedio no es consistente con la tendencia que presenta el precio. Además, tenemos que en el rango de precio 0.0, hay un celular con 19 megapíxeles en la cámara frontal, el cual, es un máximo global, pasando así, a los celulares de mayor rango de precio.
- four_g -> Se elimina debido a la falta de congruencia en los promedios (lo cual, se observa al pasar del rango de precio 1.0 al rango de precio 2.0).
- m_dep -> No se toma en cuenta debido a que estrictamente es inconsistente con la tendencia de los rangos de precios.
- mobile_wt -> Tampoco sera tomado en cuenta, pues, el rango de precio 2.0 rompe con la tendencia general que había en el peso de los dispositivos, por lo que, los datos dañarían la capacidad de predicción.
- m_cores -> Estos no se tomarán en cuenta debido a que presentan una distribución prácticamente al azar, por lo que, solo dañarían la predicción.
- px_height -> Muestra inconsistencia con la tendencia de los precios, por lo cual, no se toma en cuenta.
- px_w -> No se toma en cuenta, pues, según parece, el alto de la pantalla disminuye a medida que aumenta el precio del celular, desde el rango 0.0 al rango 3.0, sin embargo, vemos que los datos dañarían el rango de precio 3.0 rompen con esta tendencia, por lo cual, podría conducir a errores.
- screen_time -> No se toma en cuenta debido a que presenta incongruencia en el rango de precio 1.0.
- talk_time -> No se tomara como característica para predecir debido a que presenta una distribución relativamente al azar.
- touch_screen -> No se tomara como característica pues, no parece tener una tendencia general, por lo cual, no aportaría mucha información valiosa para la predicción.

Dejar

- battery_power -> Este dato se tomara en cuenta debido a que, el promedio parece tener cierta congruencia con los rangos de precios, es decir, a medida que aumenta el poder de la batería, tiende a aumentar el precio del celular. Esto se puede ver reflejado en el promedio, en el 50% y en el 75% del describe.
- pc -> Este dato sera tomado en cuenta, pues, posee una distribución uniforme y acorde a la tendencia que presentan los precios.
- px_width -> Se tomara en cuenta, pues, debido a que presenta cierta incongruencia con la tendencia de los precios, es mínima, por lo cual, en general, aportará información valiosa para la predicción.
- ram -> La ram es altamente congruente con la tendencia de los precios, por lo cual, se utilizara como característica para predecir los rangos de precios.
- three_g -> Este se utilizara como característica para predecir los rangos de precios, debido a que a pesar de presentar cierta incongruencia en el rango de precio 3.0, en general va acorde con la tendencia de estos.
- wifi -> A pesar de presentar un promedio similar en el rango de precio 1.0 y en el rango de precio 2.0, en general, presenta una tendencia al aumento junto a los rangos de precios.
- price_range -> sera el criterio en base al cual se hará la clasificación.
- int_memory -> Este se dejara, pues, se usara como criterio base para la regresión.

```
In [4]: seleccion_df=df[['battery_power', 'pc', 'px_width', 'int_memory', 'ram', 'three_g', 'wifi', 'price_range']]
seleccion_df
seleccion_dfint=df[['battery_power', 'pc', 'px_width', 'ram', 'three_g', 'wifi', 'int_memory']]
```

```
In [7]: ax = sns.boxplot(x='price_range', y='battery_power', data=seleccion_df)
ax = plt.xticks(rotation=45, labels='precios', fontweight='bold', fontsize=12)
```



En este grafico se relaciona el poder de batería con el rango de precio. Cada cajita tiene una línea, la cual es el promedio de poder de batería en cada rango de precio, por lo tanto, podemos ver que el promedio aproximado de poder de batería en los rangos de precios 0.0, 1.0, 2.0 y 3.0, son 1050, 1200, 1210 y 1450, respectivamente. En consecuencia a los datos dados anteriormente, tenemos que hay una correlación positiva entre ambas características, es decir, a medida que aumenta el precio del celular, aumenta así mismo el poder de batería.

```
In [5]: seleccion_df1=seleccion_df[seleccion_df['wifi']==0.0]
seleccion_df1.groupby('price_range').count()
```

```
Out[5]: battery_power  pc  px_width  int_memory  ram  three_g  wifi
price_range count mean 25% 75% max count mean 75% max count mean std min
0.0 252 252 252 252 252 252 252
1.0 248 248 248 248 248 248 248
2.0 248 248 248 248 248 248 248
3.0 238 238 238 238 238 238 238
```

```
In [6]: seleccion_df2=seleccion_df[seleccion_df['wifi']==1.0]
seleccion_df2.groupby('price_range').count()
```

```
Out[6]: battery_power  pc  px_width  int_memory  ram  three_g  wifi
price_range count mean 25% 75% max count mean 75% max count mean std min
0.0 248 248 248 248 248 248 248
1.0 252 252 252 252 252 252 252
2.0 252 252 252 252 252 252 252
3.0 262 262 262 262 262 262 262
```

```
In [18]: labels = ['0.0', '1.0', '2.0', '3.0'] # Ponemos los rangos de precios en las etiquetas del grafico, las cuales
# Si = 252, 248, 248, 238 # No = 252, 248, 248, 238 # Se establecen los valores de celulares que no tienen wifi
Now = [252, 248, 248, 238] # Se establecen los valores de celulares que no tienen wifi
```

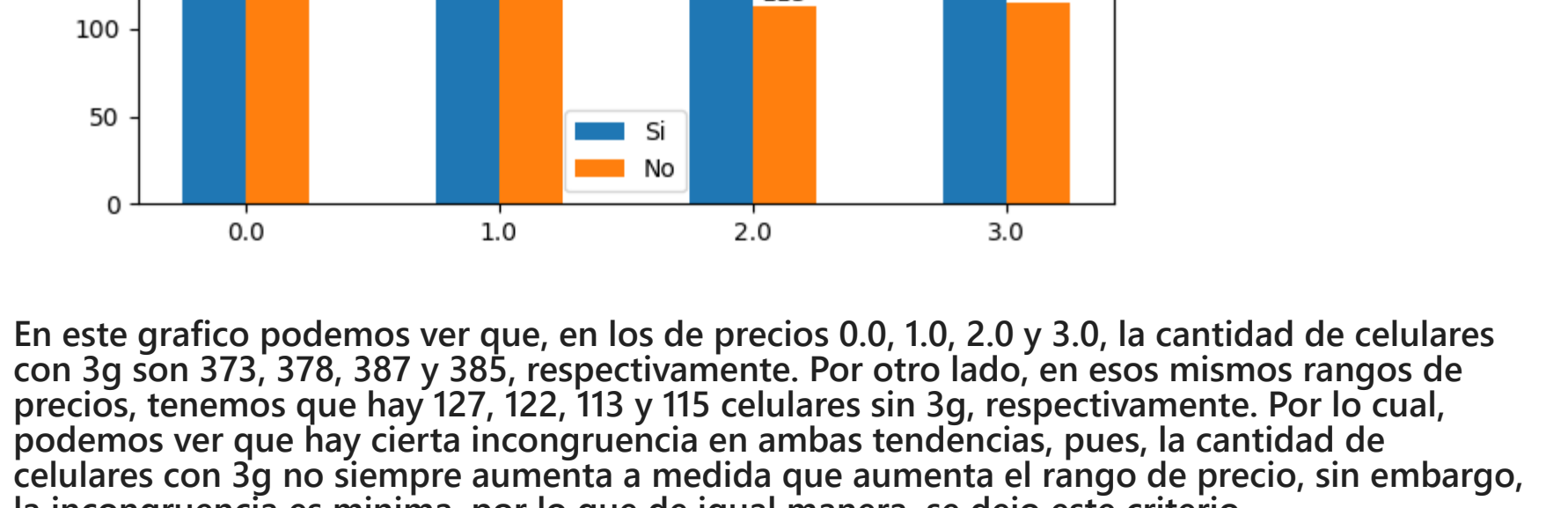
```
x = np.arange(len(labels)) # the label locations
width = 0.25 # the width of the bars
```

```
fig, ax = plt.subplots() # Hace posible que haya dos barras para poder comparar el si y el no
rects1 = ax.bar(x + width/2, Si, width, label='Si') # Genera la barra que se utilizara para los que tienen wifi
rects2 = ax.bar(x - width/2, No, width, label='No') # Genera la barra que se utilizara para los que no tienen wifi
```

```
# Add some text for labels, title and custom x-axis tick labels, etc.
ax.set_ylabel('Count') # Etiqueta eje y
ax.set_title('Celulares con wifi vs celulares sin wifi') # Titulo del grafico #
ax.set_xticks(x) # Cuadrada eje x con las etiquetas
ax.set_xticklabels(labels) # Se ponen las etiquetas en el eje x
ax.legend() # Etiqueta de los dos grupos
```

```
ax.bar_label(rects1, padding=0) # Cantidad de celulares que tienen wifi
ax.bar_label(rects2, padding=0) # Cantidad de celulares que no tienen wifi
```

```
fig.tight_layout() # Cuadrada el cartelito de acuerdo a la ubicación del grafico
plt.show() # Muestra el grafico
```



Este grafico nos muestra como se distribuyen la cantidad de celulares que si tienen wifi y los que no tienen wifi segun el rango de precio. Es decir, en los rangos de precio 0.0, 1.0, 2.0 y 3.0 hay 248, 252, 248 y 238 celulares que tienen wifi, respectivamente. Por otro lado, en los rangos de precio 0.0, 1.0, 2.0, 3.0 hay 252, 248, 248 y 238 celulares que no tienen wifi, respectivamente. Debido a lo anteriormente mencionado, podemos ver que la cantidad de celulares con wifi es congruente con los rangos de precios, pues, a medida que aumenta el rango de precio, tiende a aumentar la cantidad de celulares con wifi.

Idea de grafico extraida de

https://gallery-plotlib.org/stable/gallery/lines_and_markers/barchart.html#sph-glr-gallery-lines-bars-and-markers-barchart-py

```
In [9]: seleccion_df1=seleccion_df[seleccion_df['three_g']==1.0]
seleccion_df1.groupby('price_range').count()
```

```
Out[9]: battery_power  pc  px_width  int_memory  ram  three_g  wifi
price_range count mean 25% 75% max count mean 75% max count mean std min
0.0 373 373 373 373 373 373 373
1.0 378 378 378 378 378 378 378
2.0 387 387 387 387 387 387 387
3.0 385 385 385 385 385 385 385
```

```
In [10]: seleccion_df2=seleccion_df[seleccion_df['three_g']==0.0]
seleccion_df2.groupby('price_range').count()
```

```
Out[10]: battery_power  pc  px_width  int_memory  ram  three_g  wifi
price_range count mean 25% 75% max count mean 75% max count mean std min
0.0 127 127 127 127 127 127 127
1.0 122 122 122 122 122 122 122
2.0 113 113 113 113 113 113 113
3.0 115 115 115 115 115 115 115
```

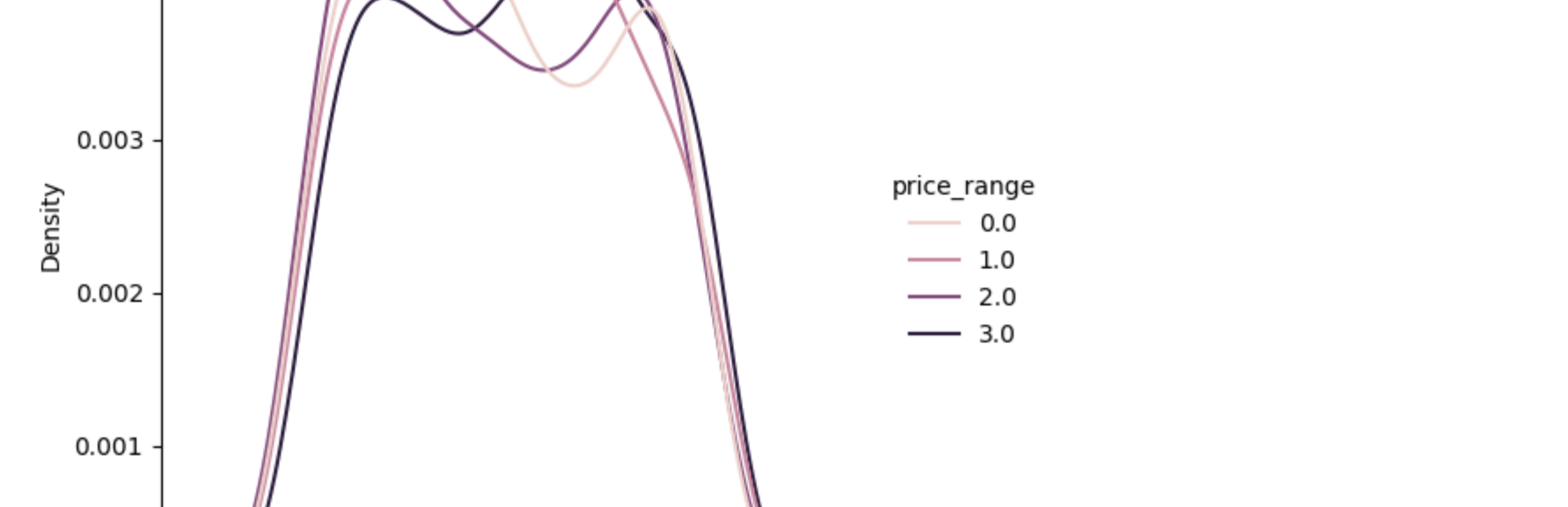
```
In [14]: import matplotlib.pyplot as plt
import numpy as np
labels = ['0.0', '1.0', '2.0', '3.0']
Now = [373, 378, 387, 385]
```

```
x5 = np.arange(len(labels)) # the label locations
width = 0.25 # the width of the bars
```

```
fig, ax = plt.subplots()
rects1 = ax.bar(x5 + width/2, Si, width, label='Si')
rects2 = ax.bar(x5 - width/2, No, width, label='No')
```

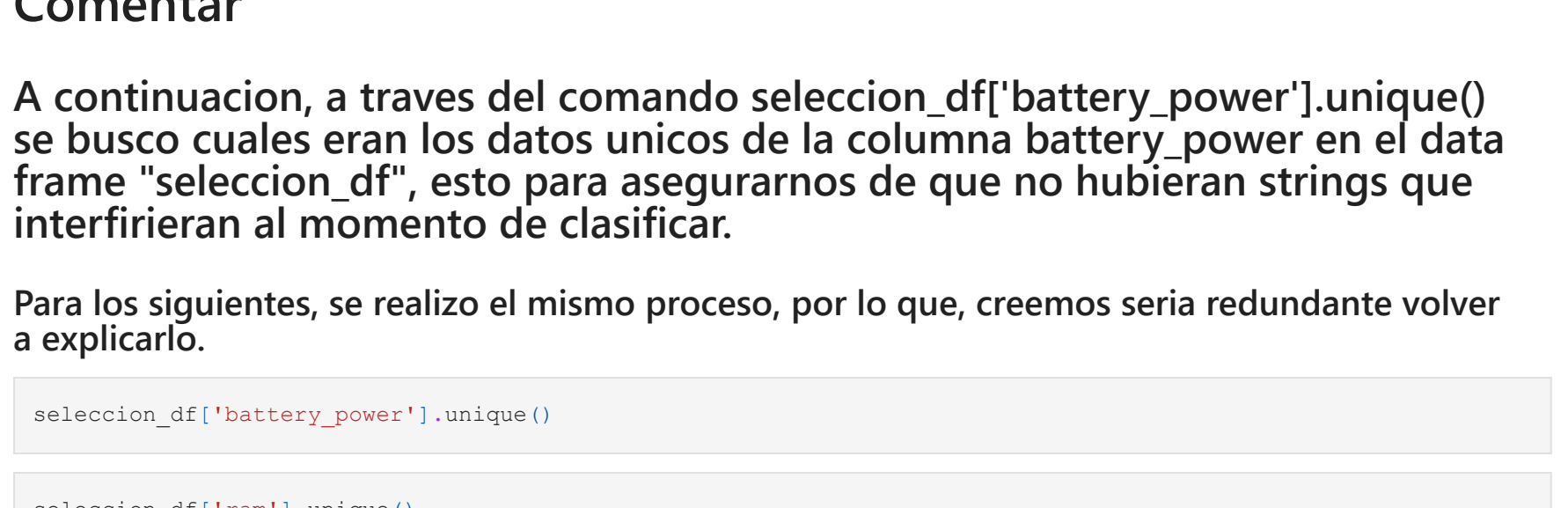
```
# Add some text for labels, title and custom x-axis tick labels, etc.
ax.set_ylabel('Count')
ax.set_title('Celulares con 3G vs celulares sin 3G')
ax.set_xticks(x5)
ax.set_xticklabels(labels)
ax.legend()
```

```
ax1.bar_label(rects1, padding=0)
ax1.bar_label(rects2, padding=0)
fig.tight_layout()
plt.show()
```



En este grafico podemos ver que, en los rangos de precios 0.0, 1.0, 2.0 y 3.0, la cantidad de celulares con 3g son 373, 378, 387 y 385, respectivamente. Por otro lado, en esos mismos rangos de precios, tenemos que hay 127, 122, 113 y 115 celulares sin 3g, respectivamente. Por lo cual, podemos ver que hay cierta incongruencia en ambas tendencias, pues, la cantidad de celulares con 3g no siempre aumenta a medida que aumenta el precio, sin embargo, la incongruencia es minima, por lo que de igual manera, se dejo este criterio.

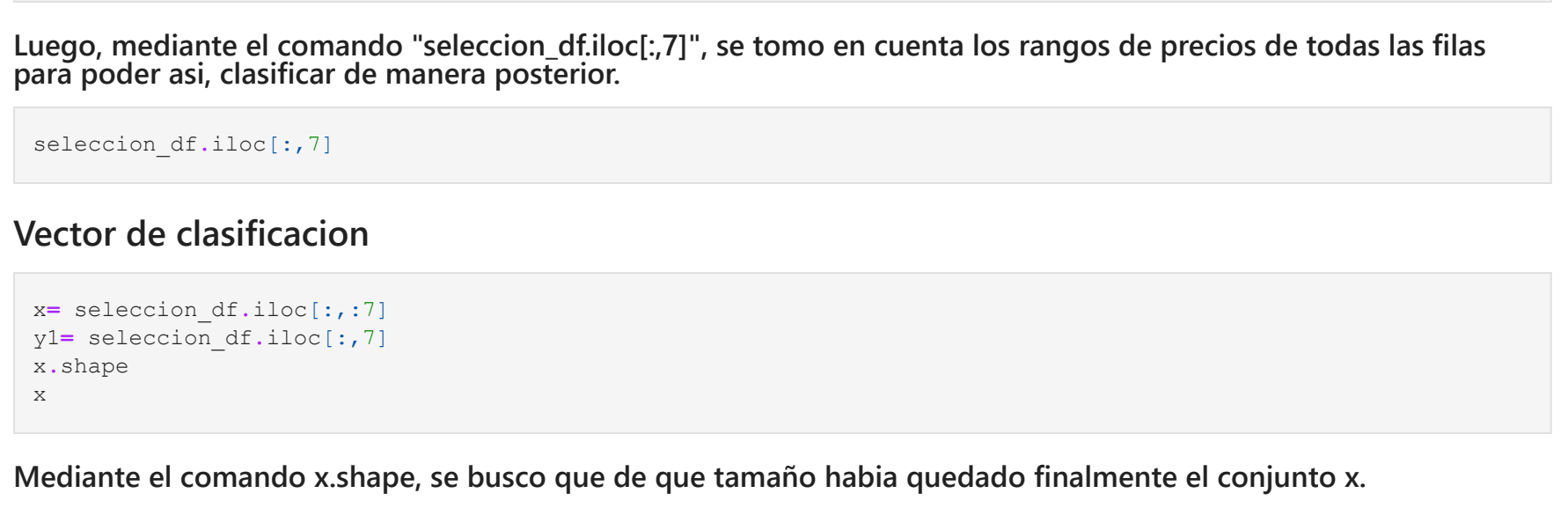
```
In [12]: sns.displot(data=seleccion_df, x='px_width', hue='price_range', kind='kde')
```



Como se logra ver en este grafico, en primer lugar, en el rango 0.0, hay una gran cantidad de celulares con aproximadamente 750 y 1300 pixeles de ancho de cámara. En segundo lugar, tenemos que en el rango 1.0, hay una gran cantidad de celulares con 750, 1200 y 1750 pixeles de ancho de cámara. Luego, en el rango 2.0, tenemos que hay una gran cantidad de celulares con aproximadamente 750 y 1600 pixeles de ancho. Finalmente, tenemos que en el rango de precio 3.0 hay una gran cantidad de celulares con aproximadamente, 1800 pixeles de ancho en el rango de precio.

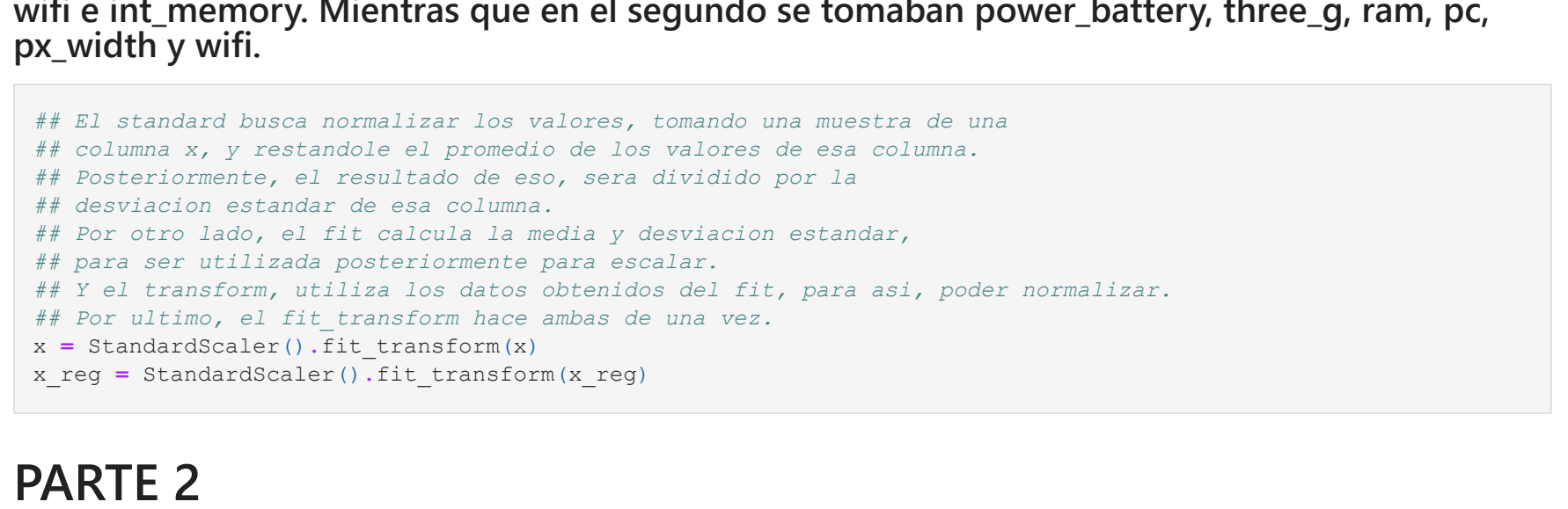
Por otro lado, podemos ver los datos en general parecen seguir cierta tendencia, en la cual, al aumentar el precio, en general aumenta el ancho de resolucion en pixeles.

```
In [16]: graph1=sns.histplot(data=seleccion_df, x='pc', hue='price_range', multiple='dodge', shrink=.8, legend=True)
plt.show()
```



La presente grafica, enfrenta la cantidad celulares versus la cantidad de pixeles que posee su cámara principal, distinguiendo la su rango de precio de 0.0 la mayoría de los celulares se encuentra por debajo de encontrando sus picos por debajo de los 20 pixeles. En el rango de precio 1.0 encontramos como se mantiene, y tiende a tener una leve subida en cuanto se va aumentando su cantidad de pixeles. El rango de precio 2.0 es muy parecido al anterior no obstante tiene la tendencia subir con aun mas cantidad de pixeles. Por ultimo encontramos el rango de precio 3.0 que en su mayoría encontramos que la cámara principal tiene 20 pixeles y es aquí donde se encuentra la mayoría de celulares.

```
In [17]: sns.displot(data=seleccion_df, x='int_memory', hue='price_range', kind='kde')
```



En la grafica anterior podemos ver como la memoria interna de un celular va a variar segun el rango de precio que tenga, en primera instancia se puede observar que en el rango de celulares correspondian a un rango de precio 1.0 cuando realmente eran de un rango de precio 2.0. Luego, cada vez se disminuye más llegando a una baja cada vez que aumentamos la memoria. Luego tenemos que la cantidad de memoria interna en los celulares con rangos de precio 1.0 se mantiene entre los 16 y 32GB, para posteriormente tenga un descenso. Llegando así al precio de rango 2.0 en donde se pueden apreciar dos picos importantes en primer lugar en los de 16GB y luego llegando a los 60 GB siendo este ultimo pico, mucho menor que el anterior caso donde el promedio era de 32GB. Por otro lado, tenemos que el fit transform, nos ayuda a ajustar y transformar los datos de precio 2.0 y realmente eran de un rango de precio 1.0.

```
sns.set_theme(style='whitegrid') sns.stripplot(data=seleccion_df, x='ram', y='price_range') plt.xticks(rotation=90)
```

Comentar

A continuación, a traves del comando seleccion_df['battery_power'].unique() se busco cuales eran los datos unicos de la columna battery_power en el data frame "seleccion_df", esto para asegurarnos de que no hubieran strings que interfirieran al momento de clasificar.

Para los siguientes, se realizo el mismo comando, por lo que, nuevamente seria redundante volver a explicarlo.

```
In [ ]: seleccion_df['battery_power'].unique()
```

```
In [ ]: seleccion_df['ram'].unique()
```

```
In [ ]: seleccion_df['px_width'].unique()
```

```
In [ ]: seleccion_df['pc'].unique()
```

```
In [ ]: seleccion_df['int_memory'].unique()
```

```
In [ ]: seleccion_df['three_g'].unique()
```

```
In [ ]: seleccion_df['wifi'].unique()
```

```
In [ ]: seleccion_df
```

Se crea un nuevo dataframe, esta vez sin tomar en cuenta el price_range.

```
In [ ]: seleccion_dfmt = seleccion_df.drop('price_range', axis=1)
```

Como buscamos modificar cuales son los rangos de precios de cada celular segun sus especificaciones, mediante el comando "seleccion_dfmt.iloc[:,7]" se separaron las características de los celulares del rango de precio.

```
In [ ]: seleccion_dfmt.iloc[:,7]
```

Luego, mediante el comando "seleccion_dfmt.iloc[:,7]", se tomo en cuenta los rangos de precios de todas las filas para poder así, clasificar de manera posterior.

```
In [ ]: seleccion_dfmt.iloc[:,7]
```

Vector de clasificacion

```
In [ ]: x_reg = seleccion_dfmt.iloc[:,7]
y_reg = seleccion_dfmt.iloc[:,7]
x_reg.shape
x
```

Mediante el comando x.shape, se busco que de que tamaño había quedado finalmente el conjunto x.

```
In [ ]: seleccion_dfmt.iloc[:,6]
```

```
In [ ]: seleccion_dfmt.iloc[:,6]
```

Como se logra ver, en el dataframe "seleccion_dfmt" en el que se encuentran todos los features de "seleccion_df" menos price_range, generamos una separacion, para así, poder dar lugar a la regresión.

Hicimos una separación entre memoria interna del celular y los otros atributos como power_battery, three_g, ram, pc, px_width y wifi.

```
In [ ]: x_reg = seleccion_dfmt.iloc[:,6]
y_reg = seleccion_dfmt.iloc[:,6]
x_reg.shape
```

Con el comando escrito a continuación, se normalizaron los valores tanto del conjunto x y el conjunto x_reg, donde en el primero se tomaban power_battery, three_g, ram, pc, px_width, wifi e int_memory. Mientras que en el segundo se tomaban power_battery, three_g, ram, pc, px_width y wifi.

```
In [ ]: ## El standard busca normalizar los valores, tomando una muestra de una
## columna x, y restandole el promedio de los valores de esa columna.
## Posteriormente, el resultado es cero, es decir, se divide por la
## desviación estándar de esa columna.
## Por otro lado, el fit calcula la media y desviación estándar.
## Para ser utilizada posteriormente para escalar.
## Y el transform, utiliza los datos obtenidos del fit, para así, poder normalizar.
## Por ultimo, el fit transform hace ambas de una vez.
x_reg = StandardScaler().fit_transform(x)
x_reg = StandardScaler().fit_transform(x_reg)
```

PARTE 2

```
In [ ]: X_entrenamiento, X_prueba, Y_entrenamiento, Y_prueba = train_test_split(x,y, test_size=0.20)
```

Mediante el comando escrito anteriormente, se definen 4 variables, donde X_entrenamiento e Y_entrenamiento son los que se usaran para entrenar (en los cuales tendremos aproximadamente el 80% de cada grupo de datos), mientras que X_prueba e Y_prueba se usaran para testear (en los cuales tendremos aproximadamente el 20% de cada grupo de datos).

```
In [ ]: Gaussian_NB = GaussianNB()
```

Se establece como "Gaussian_NB" al metodo de predicción de Naive Bayes o GaussianNB.

Se utilizará el clasificador Naive Bayes, en específico el Gaussian, por 3 motivos principalmente, en primera instancia este clasificador funciona de gran manera para trabajar con características independientes entre si, como lo son las características de los celulares. Por otro lado, tenemos el tiempo, Naive Bayes Gaussian nos permite tener una gran eficacia para tratar con una cantidad limitada de datos, y por ultimo este es un metodo que nos permite entender de manera relativamente facil los resultados.

```
In [ ]: Gaussian_NB.fit(X_entrenamiento, Y_entrenamiento)
```

Se ajustan el GaussianNB a los datos de entrenamiento X_entrenamiento e Y_entrenamiento.

```
In [ ]: y_pred = Gaussian_NB.predict(X_prueba)
```

A traves del ajuste de GaussianNB, se prueba e intenta predecir los rangos de precios que tendran celulares segun sus caracteristicas.

```
In [ ]: print("Accuracy : %.3f" % (accuracy_score(Y_prueba, y_pred, normalize=True)*100), "%")
```

Se compara Y_prueba con la predicción realizada a traves del entrenamiento, a traves de esta, se logra aproximadamente 80%.

```
In [ ]: from sklearn.metrics import confusion_matrix
seleccion_dfmt = seleccion_dfmt.drop('price_range', axis=1)
print("Matrix Confusion :")
print(Matrix Confusion)
```

Se calcula la matriz de confusion, la cual es un modelo de clasificación que evalua la precisión lograda a traves de GNB, esto se hará mediante la comparacion de Y_prueba (conjunto de rango de precios que serán utilizados para testeo) e "Y_pred" (Conjunto de rango de precios que fueron calculados mediante el Gaussian NB con el conjunto de características de celulares "X_prueba").

```
In [ ]: class_names = ['0.0', '1.0', '2.0', '3.0']
plot_confusion_matrix(Gaussian_NB, X_prueba, Y_prueba, display_labels=class_names, cmap=plt.cm.Blues, normalize=True)
```

Para entender esta matriz de confusión, iremos desglosando por filas, en primer lugar, en la primera fila, tenemos un 92% de casos en los que se predijo que las características de celulares correspondían a un rango de precio 0.0 y realmente eran de un 8.4% de los casos, se predijo que eran de un rango de precio 1.0 y realmente eran de un rango de precio 0.0.

En segundo lugar, tenemos que hubo alrededor de un 1% de casos en los que se predijo que serían del rango de precio 0.0 y eran realmente del rango de precio 1.0. Luego, tenemos que hubo un 70% de casos en los que se predijo que las características de celular correspondían al rango de precio 1.0 y realmente eran del rango de precio 1.0. Por ultimo, en un 15% de los casos, se predijo que las características de celulares correspondían a un rango de precio 2.0 y realmente eran de un rango de precio 1.0.

En tercer lugar, tenemos que hubo alrededor de un 23% de casos en los que se predijo que ciertas características de celulares correspondían a un rango de precio 1.0 cuando realmente eran de un rango de precio 2.0. Luego, cada vez se disminuye más llegando a una baja cada vez que aumentamos la memoria. Luego tenemos que la cantidad de memoria interna en los celulares con rangos de precio 1.0 se mantiene entre los 16 y 32GB, para posteriormente tenga un descenso. Llegando así al precio de rango 2.0 en donde se pueden apreciar dos picos importantes en primer lugar en los de 16GB y luego llegando a los 60 GB siendo este ultimo pico, mucho menor que el anterior caso donde el promedio era de 32GB. Por otro lado, tenemos que el fit transform, nos ayuda a ajustar y transformar los datos de precio 2.0 y realmente eran de un rango de precio 1.0.

En cuarto lugar, tenemos un 8.3% de casos donde se predijo que ciertas características de celulares eran acorde a un rango de precio 2.0 cuando realmente eran de un rango de precio 3.0, por ultimo, tenemos un 92% de casos donde se predijo que características de celulares iban acorde a un rango de precio 3.0 y lo eran en realidad.

Para concluir, se extrae de esta matriz de confusion que los rangos de precios en los que mejor se desempeña el clasificador es en los rangos de precios 0.0 y 3.0.

```
In [ ]: X_entrenamiento, X_prueba, Y_entrenamiento, Y_prueba = train_test_split(x_reg, y_reg, test_size=0.30)
print(X_entrenamiento.shape, X_prueba.shape, Y_entrenamiento.shape, Y_prueba.shape)
```

Se definen los parametros X_entrenamiento, X_prueba, Y_entrenamiento e Y_prueba que se utilizaran para hacer la regresión a traves del Decision Tree Regressor. Los parametros X_entrenamiento e Y_entrenamiento corresponden a un 70% de los datos del data frame mientras que los parametros X_prueba e Y_prueba corresponden a un 30% de los datos del dataframe.

```
In [ ]: regr_1 = DecisionTreeRegressor(max_depth=70)
regr_2 = DecisionTreeRegressor(max_depth=50)
regr_3 = DecisionTreeRegressor(max_depth=5)
```

Creamos tres arboles de regresion con distintas profundidades, las cuales fueron de un 70%, 90% y 95%.

```
In [ ]: regr_1.fit(X_entrenamiento, Y_entrenamiento)
regr_2.fit(X_entrenamiento, Y_entrenamiento)
regr_3.fit(X_entrenamiento, Y_entrenamiento)
```

Ajustamos los datos de cada conjunto de entrenamientos segun las profundidades que mencionamos anteriormente.

```
In [ ]: y_1=regr_1.predict(X_prueba)
y_2=regr_2.predict(X_prueba)
y_3=regr_3.predict(X_prueba)
```

A partir de los datos aprendidos a traves de los conjuntos de entrenamiento, se intentan predecir los valores de memorias internas segun distintos parametros como pc, battery_power, etc y se guardan en las variables y_1, y_2 e y_3.

```
In [ ]: print("max_depth = 70 score=", regr_1.score(X_prueba, y_1))
print("max_depth = 50 score=", regr_2.score(X_prueba, y_2))
print("max_depth = 5 score=", regr_3.score(X_prueba, y_3))
```

Se comparan los valores de memoria interna predichos a traves de los valores de X_prueba y se comparan con los valores reales. Se puede ver que la capacidad de predicción es nula o hasta peor que nula.

```
In [ ]: regr_1.feature_importances_
```

Se calcula la importancia de las características del celular al momento de predecir la memoria interna, sin embargo, no sabemos que tan fiable es esta medida despues de que la capacidad de predicción fue pauperima.

```
In [ ]: from sklearn.manifold import TSNE
Visualizacion_X = TSNE(n_components=1).fit_transform(X_prueba)
Visualizacion_X.shape
```

El TSNE ayuda a visualizar datos con una gran cantidad de

1.¿Como resolvieron el tema de los valores nulos? Justifique.(0.3)

R// Para solucionar el problema de los valores nulos, lo que hicimos fue eliminarlos a traves del comando df=df.dropna(), pues, consideramos que estos generaban un problema al momento de recopilar informacion precisa sobre las características de los celulares y sus respectivos rangos de precios, dando lugar así, a sesgos.

2.¿Que normalizaste, filas o columnas?¿Por que?¿Para que sirve normalizar los datos?¿Que tipo de normalizacion usaste y por que?Justifique. (0.3)

R//En este caso, normalizamos columnas, pues eran los datos que se relacionaban entre sí, no hubiera tenido sentido sacar el promedio de diferentes tipos de características con diferentes tipos de medidas. Normalizar los datos sirve para poder llevar todos los datos a una misma medida y facilitar al clasificador el manejo de ellos. Para esto, nosotros utilizamos StandardScaler().fit_transform(x) , pues, creemos que era fácil de entender y bastante eficiente para el uso que necesitabamos.

3.¿Que graficos utilizaste para caracterizar los datos?¿Por que?¿Que observaste de tus datos, encontraste alguna característica particular? Justifique.(0.3)

R// Utilizamos graficos de barras, líneas, puntos y caja, debido a que eran faciles de entender en su mayoría y permitian obtener conclusiones más facilmente. En general se pudo observar que los datos eran consistentes con la tendencia de los rangos de precios,es decir, a medida que aumentaban los rangos de precios, tambien tendian a aumentar las medidas de cada característica. Nos parecio curioso que en todos los rangos de precios habian celulares con hasta solo 400 miliamperios, como podria servir eso??

4.¿Por que se separan los datos en set de entrenamiento y set de pruebas?¿Que proporcion de los datos utilizaste para cada uno y por que?Justifique.(0.3)

R// Estos se separan debido a que, si usas los mismos datos para entrenar y luego predices esos mismos datos, no tendria sentido alguno, pues la ia solo estaria memorizando, mas no aprendiendo. Nosotros utilizamos un 80% de entrenamiento y un 20% de prueba en el clasificador dado que fue el que mejor nos dio resultados a la hora de predecir. Por otro lado, utilizamos un 70% de entrenamiento y un 30% de prueba en el clasificador, pues, ya que los datos no parecian tener gran relacion, una menor cantidad de datos para entrenar, daba lugar a un error mas leve, o, puesto de otra manera, una mejor precision.

In []: