

ASSIGNMENT 4

PRATHAMESH CHIKANE

```
# Importing required libraries.
import pandas as pd
import numpy as np
import seaborn as sns #visualisation
import matplotlib.pyplot as plt #visualisation
%matplotlib inline
sns.set(color_codes=True)
from google.colab import files
```

```
uploaded_file= files.upload()
```

Choose files

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving data.csv to data.csv

```
df = pd.read_csv('data.csv')
df.head()
```

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	Number of Doors
0	BMW	1 Series M	2011	premium unleaded (required)	335.0	6.0	MANUAL	rear wheel drive	4
1	BMW	1 Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	4
2	BMW	1 Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	4
3	BMW	1 Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	4
4	BMW	1 Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	4

```
# To display the bottom 5 rows
```

```
df.tail(5)
```

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_
11909	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wt
11910	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wt
11911	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wt
11912	Acura	ZDX	2013	premium unleaded (recommended)	300.0	6.0	AUTOMATIC	all wt
11913	Lincoln	Zephyr	2006	regular unleaded	221.0	6.0	AUTOMATIC	front wt

Here we check for the datatypes because sometimes the MSRP or the price of the car would be stored as a string or object, if in that case, we have to convert that string to the integer data only then we can plot the data via a graph. Here, in this case, the data is already in integer format so nothing to worry.

```
# Checking the data type
df.dtypes
```

```
Make                object
Model              object
Year               int64
Engine Fuel Type    object
Engine HP          float64
Engine Cylinders    float64
Transmission Type   object
Driven_Wheels       object
Number of Doors     float64
Market Category     object
Vehicle Size        object
Vehicle Style       object
highway MPG         int64
city mpg            int64
Popularity          int64
MSRP               int64
dtype: object
```

This step is certainly needed in every EDA because sometimes there would be many columns that we never use in such cases dropping is the only solution. In this case, the columns such as Engine Fuel Type, Market Category, Vehicle style, Popularity, Number of doors, Vehicle Size doesn't make any sense to me so I just dropped for this instance.

```
# Dropping irrelevant columns
df = df.drop(['Engine Fuel Type', 'Market Category', 'Vehicle Style', 'Popularity', 'Number of doors', 'Vehicle Size'])
df.head(5)
```

	Make	Model	Year	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	highway MPG	city MPG
0	BMW	1 Series M	2011	335.0	6.0	MANUAL	rear wheel drive	26	23
1	BMW	1 Series	2011	300.0	6.0	MANUAL	rear wheel drive	28	24
2	BMW	1 Series	2011	300.0	6.0	MANUAL	rear wheel drive	28	24
3	BMW	1 Series	2011	230.0	6.0	MANUAL	rear wheel drive	28	24
4	BMW	1 Series	2011	230.0	6.0	MANUAL	rear wheel drive	28	24

In this instance, most of the column names are very confusing to read, so I just tweaked their column names. This is a good approach it improves the readability of the data set.

```
# Renaming the column names
df = df.rename(columns={'Engine HP': 'HP', 'Engine Cylinders': 'Cylinders', 'Transmission Type': 'Transmission'})
df.head(5)
```

	Make	Model	Year	HP	Cylinders	Transmission	Drive Mode	MPG-H	MPG-C	Price
0	BMW	1 Series M	2011	335.0	6.0	MANUAL	rear wheel drive	26	19	46135
							rear			

This is often a handy thing to do because a huge data set as in this case contains more than 10,000 rows often have some duplicate data which might be disturbing, so here I remove all the duplicate value from the data-set. For example prior to removing I had 11914 rows of data but after removing the duplicates 10925 data meaning that I had 989 of duplicate data.

Now let us remove the duplicate data because it's ok to remove them.

```
df.count()

Make          11914
Model         11914
Year          11914
HP            11845
Cylinders     11884
Transmission  11914
Drive Mode    11914
MPG-H         11914
MPG-C         11914
Price         11914
dtype: int64
```

So seen above there are 11914 rows and we are removing 989 rows of duplicate data.

```
# Dropping the duplicates
df = df.drop_duplicates()
df.head(5)
```

	Make	Model	Year	HP	Cylinders	Transmission	Drive Mode	MPG-H	MPG-C	Price
0	BMW	1 Series M	2011	335.0	6.0	MANUAL	rear wheel drive	26	19	46135
1	BMW	1 Series	2011	335.0	6.0	MANUAL	rear wheel drive	26	19	46135

```
# Counting the number of rows after removing duplicates.
df.count()
```

```
Make          10925
Model         10925
Year          10925
HP            10856
Cylinders     10895
Transmission  10925
Drive Mode    10925
MPG-H         10925
MPG-C         10925
Price         10925
dtype: int64
```

```
# Finding the null values.
print(df.isnull().sum())
```

```
Make          0
Model         0
Year          0
HP            69
Cylinders     30
Transmission  0
Drive Mode    0
MPG-H         0
MPG-C         0
Price         0
dtype: int64
```

This is the reason in the above step while counting both Cylinders and Horsepower (HP) had 10856 and 10895 over 10925 rows.

```
# Dropping the missing values.
df = df.dropna()
df.count()
```

```
Make          10827
Model         10827
Year          10827
HP            10827
Cylinders     10827
Transmission  10827
Drive Mode    10827
```

```
MPG-H          10827
MPG-C          10827
Price          10827
dtype: int64
```

Now we have removed all the rows which contain the Null or N/A values (Cylinders and Horsepower (HP)).

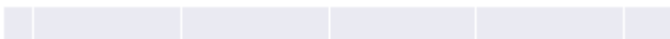
```
# After dropping the values
print(df.isnull().sum())
```

```
Make          0
Model         0
Year          0
HP            0
Cylinders     0
Transmission  0
Drive Mode    0
MPG-H         0
MPG-C         0
Price         0
dtype: int64
```

Detecting Outliers An outlier is a point or set of points that are different from other points. Sometimes they can be very high or very low. It's often a good idea to detect and remove the outliers. Because outliers are one of the primary reasons for resulting in a less accurate model. Hence it's a good idea to remove them. The outlier detection and removing that I am going to perform is called IQR score technique. Often outliers can be seen with visualizations using a box plot. Shown below are the box plot of MSRP, Cylinders, Horsepower and EngineSize. Herein all the plots, you can find some points are outside the box they are none other than outliers. The technique of finding and removing outlier that I am performing in this assignment is taken help of a tutorial from towards data science.

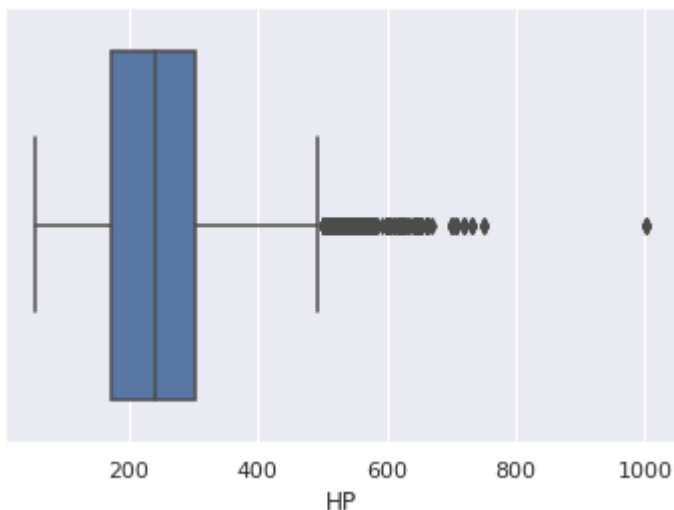
```
sns.boxplot(x=df['Price'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f45910b7950>
```



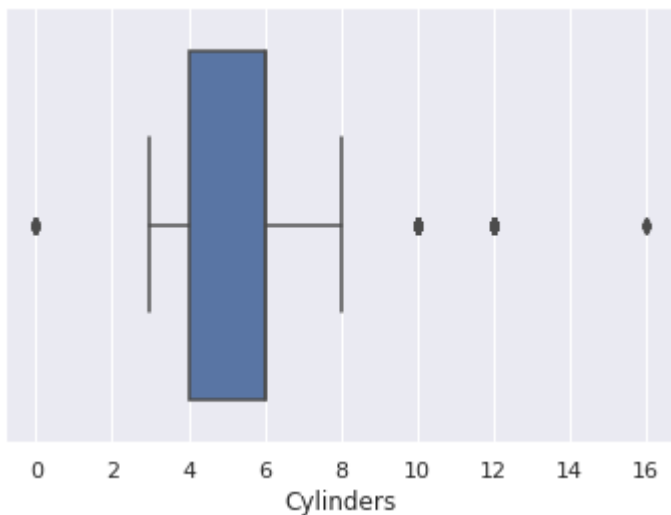
```
sns.boxplot(x=df['HP'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f45900c1490>
```



```
sns.boxplot(x=df['Cylinders'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f45
```



```
Q1 = df.quantile(0.25)
```

```
Q3 = df.quantile(0.75)
```

```
IQR = Q3 - Q1
```

```
print(IQR)
```

Year	9.0
HP	130.0
Cylinders	2.0
MPG-H	8.0
MPG-C	6.0

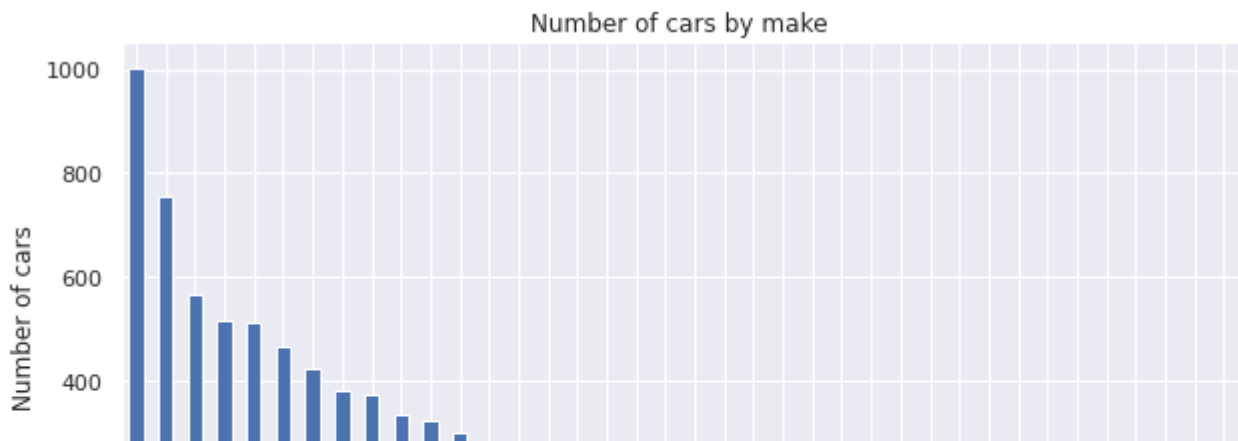
```
Price      21327.5
dtype: float64
```

```
df = df[~((df < (Q1-1.5 * IQR)) |(df > (Q3 + 1.5 * IQR))).any(axis=1)]
df.shape
```

```
(9191, 10)
```

As seen above there were around 1600 rows were outliers. But you cannot completely remove the outliers because even after you use the above technique there maybe 1–2 outlier unremoved but that ok because there were more than 100 outliers. Something is better than nothing.

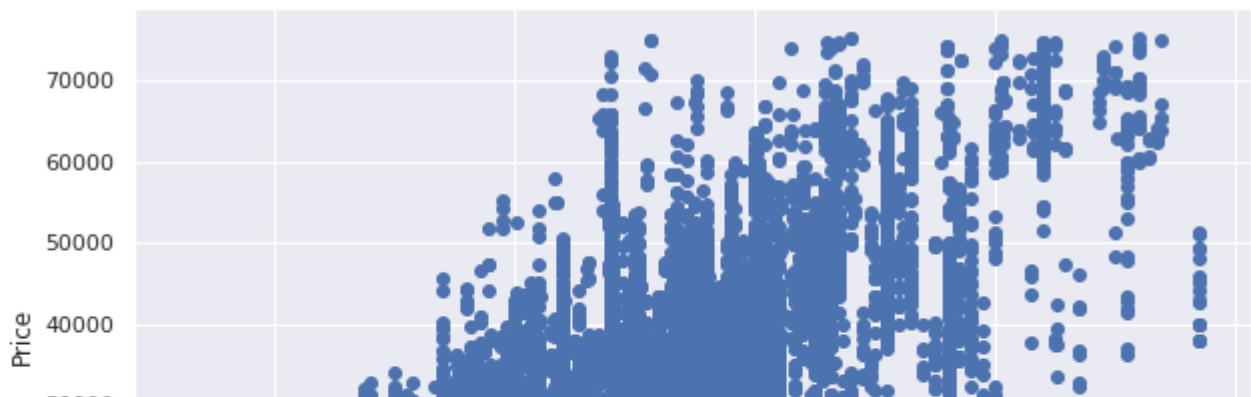
```
# Plotting a Histogram
df.Make.value_counts().nlargest(40).plot(kind='bar', figsize=(10,5))
plt.title("Number of cars by make")
plt.ylabel('Number of cars')
plt.xlabel('Make');
```



```
# Finding the relations between the variables.
plt.figure(figsize=(20,10))
c= df.corr()
sns.heatmap(c,cmap="BrBG",annot=True)
c
```


	Year	HP	Cylinders	MPG-H	MPG-C	Price
Year	1.000000	0.326726	-0.133920	0.378479	0.338145	0.592983
HP	0.326726	1.000000	0.715237	-0.443807	-0.544551	0.739042
Cylinders	-0.133920	0.715237	1.000000	-0.703856	-0.755540	0.354013
MPG-H	0.378479	-0.443807	-0.703856	1.000000	0.939141	-0.106320

```
# Plotting a scatter plot
fig, ax = plt.subplots(figsize=(10,6))
ax.scatter(df['HP'], df['Price'])
ax.set_xlabel('HP')
ax.set_ylabel('Price')
plt.show()
```



Hence the above are some of the steps involved in Exploratory data analysis, these are some general steps that you must follow in order to perform EDA. There are many more yet to come but for now, this is more than enough idea as to how to perform a good EDA given any data sets. Stay tuned for more updates. If you have some doubts then the comment section is all yours. I'll try my level best to answer your questions. Thank you.

