

Exploratory Data Analysis

▼ Problem Statement:

We have used Cars dataset from kaggle with features including make, model, year, engine, and other properties of the car used to predict its price.

▼ 1. Importing the necessary 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 scipy import stats
import warnings
warnings.filterwarnings("ignore")
```

▼ 2. Download the dataset and load into dataframe

5 points

Please download the dataset from [here](#) and extract the csv file. Load the csv file as pandas dataframe.

```
## load the csv file
df = pd.read_csv("/content/data.csv")
```

Now we observe the each features present in the dataset.

Make: The Make feature is the company name of the Car.

Model: The Model feature is the model or different version of Car models.

Year: The year describes the model has been launched.

Engine Fuel Type: It defines the Fuel type of the car model.

Engine HP: It's say the Horsepower that refers to the power an engine produces.

Engine Cylinders: It define the nos of cylinders in present in the engine.

Transmission Type: It is the type of feature that describe about the car transmission type i.e Manual or automatic.

Driven_Wheels: The type of wheel drive.

No of doors: It defined nos of doors present in the car.

Market Category: This features tells about the type of car or which category the car belongs.

Vehicle Size: It's say about the about car size.

Vehicle Style: The feature is all about the style that belongs to car.

highway MPG: The average a car will get while driving on an open stretch of road without stopping or starting, typically at a higher speed.

city mpg: City MPG refers to driving with occasional stopping and braking.

Popularity: It can refered to rating of that car or popularity of car.

MSRP: The price of that car.

```
## print the head of the dataframe
df.head()
```

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	Number of Doors
0	BMW	Series M	2011	premium unleaded (required)	335.0	6.0	MANUAL	rear wheel drive	2.0
1	BMW	Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0
2	BMW	Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0
3	BMW	Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0
4	BMW	Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0

3. Check the datatypes

2 points

```
# Get the datatypes of each columns number of records in each column.
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11914 entries, 0 to 11913
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Make                  11914 non-null  object
1   Model                 11914 non-null  object
2   Year                  11914 non-null  int64
3   Engine Fuel Type      11911 non-null  object
4   Engine HP             11845 non-null  float64
5   Engine Cylinders      11884 non-null  float64
6   Transmission Type     11914 non-null  object
7   Driven_Wheels         11914 non-null  object
8   Number of Doors       11908 non-null  float64
9   Market Category       8172 non-null   object
10  Vehicle Size          11914 non-null  object
11  Vehicle Style         11914 non-null  object
12  highway MPG           11914 non-null  int64
13  city mpg              11914 non-null  int64
14  Popularity            11914 non-null  int64
15  MSRP                  11914 non-null  int64
dtypes: float64(3), int64(5), object(8)
memory usage: 1.5+ MB
```

▼ 4. Dropping irrelevant columns

Reference:-<https://www.youtube.com/watch?v=cRurWEfmxC0>

5 points

If we consider all columns present in the dataset then unnecessary columns will impact on the model's accuracy.

Not all the columns are important to us in the given dataframe, and hence we would drop the columns that are irrelevant to us. It would reflect our model's accuracy so we need to drop them. Otherwise it will affect our model.

The list `cols_to_drop` contains the names of the cols that are irrelevant, drop all these cols from the dataframe.

```
cols_to_drop = ["Engine Fuel Type", "Market Category", "Vehicle Style", "Popularity", "Number of Doors", "Vehicle Size"]
```

These features are not necessary to obtain the model's accuracy. It does not contain any relevant

```
# initialise cols_to_drop
cols_to_drop = ["Engine Fuel Type", "Market Category", "Vehicle Style", "Popularity", "Number"]

# drop the irrelevant cols and print the head of the dataframe
df = df.drop(cols_to_drop,axis=1)
df.head()
```

	Make	Model	Year	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	highway MPG	city mpg	M
0	BMW	Series 1 M	2011	335.0	6.0	MANUAL	rear wheel drive	26	19	46
1	BMW	Series 1	2011	300.0	6.0	MANUAL	rear wheel drive	28	19	40

5. Renaming the columns

5 points

Now, It's time for renaming the feature to useful feature name. It will help to use them in model training purpose.

We have already dropped the unnecessary columns, and now we are left with useful columns. One extra thing that we would do is to rename the columns such that the name clearly represents the essence of the column.

The given dict represents (in key value pair) the previous name, and the new name for the dataframe columns

```
rename_cols = {"Engine HP": "HP", "Engine Cylinders": "Cylinders", "Transmission Type": "Transmission Type",
               "Driven_Wheels": "Drive Mode", "highway MPG": "MPG_H", "city mpg": "MPG-C", "MSRP": "MSRP"}

# use a pandas function to rename the current columns -
df = df.rename(columns = rename_cols)

# Print the head of the dataframe
df.head()
```

	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	300.0	6.0	MANUAL	rear wheel drive	28	19	40650
2	BMW	1 Series	2011	300.0	6.0	MANUAL	rear wheel drive	28	20	36350

▼ 6. Dropping the duplicate rows

Reference: <https://www.youtube.com/watch?v=bFVMR1qfzXo>

10 points

There are many rows in the dataframe which are duplicate, and hence they are just repeating the information. Its better if we remove these rows as they don't add any value to the dataframe.

For given data, we would like to see how many rows were duplicates. For this, we will count the number of rows, remove the duplicated rows, and again count the number of rows.

```
# number of rows before removing duplicated rows
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
```

```
# drop the duplicated rows
df = df.drop_duplicates()
df.head(5)
```

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

Count Number of rows after deleting duplicated rows

```
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

```

▼ 7. Dropping the null or missing values

10 points

Missing values are usually represented in the form of Nan or null or None in the dataset.

Finding whether we have null values in the data is by using the `isnull()` function.

There are many values which are missing, in pandas dataframe these values are referred to as `np.nan`. We want to deal with these values because we can't use nan values to train models. Either we can remove them to apply some strategy to replace them with other values.

To keep things simple we will be dropping nan values

```

# check for nan values in each columns
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

```

As we can see that the HP and Cylinders have null values of 69 and 30. As these null values will impact on models' accuracy. So to avoid the impact we will drop the these values. As these values are small comparing with dataset that will not impact any major affect on model accuracy so we will drop the values.

```
# drop missing values
df = df.dropna()

# Make sure that missing values are removed
# check number of nan values in each col again
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
```

```
#Describe statistics of df
df.describe()
```

	Year	HP	Cylinders	MPG_H	MPG-C	Price
count	10827.000000	10827.000000	10827.000000	10827.000000	10827.000000	1.082700e+04
mean	2010.896370	254.553062	5.691604	26.308119	19.327607	4.249325e+04
std	7.029534	109.841537	1.768551	7.504652	6.643567	6.229451e+04
min	1990.000000	55.000000	0.000000	12.000000	7.000000	2.000000e+03
25%	2007.000000	173.000000	4.000000	22.000000	16.000000	2.197250e+04
50%	2015.000000	240.000000	6.000000	25.000000	18.000000	3.084500e+04
75%	2016.000000	303.000000	6.000000	30.000000	22.000000	4.330000e+04
max	2017.000000	1001.000000	16.000000	354.000000	137.000000	2.065902e+06

▼ 8. Removing outliers

Reference: <https://www.youtube.com/watch?v=yxTRB0boTVg>

Sometimes a dataset can contain extreme values that are outside the range of what is expected and unlike the other data. These are called outliers and often machine learning modeling and model skill in general can be improved by understanding and even removing these outlier values.

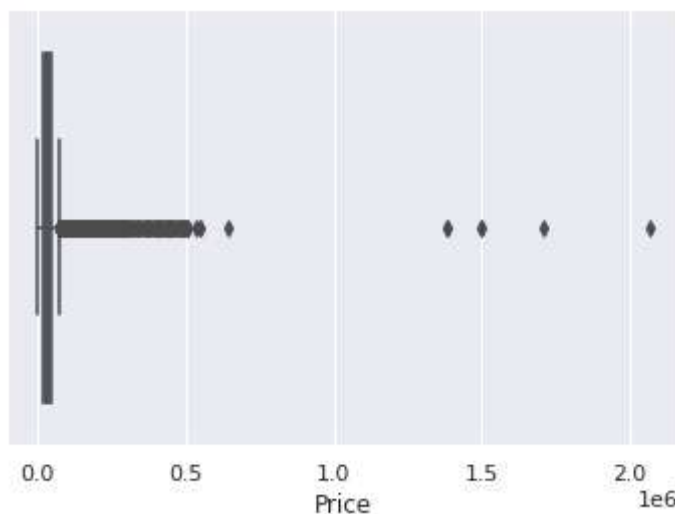
▼ Detecting outliers

There many techniques to detect outliers. Let us first see the simplest form of visualizing outliers.

Box plots are a graphical depiction of numerical data through their quantiles. It is a very simple but effective way to visualize outliers. Think about the lower and upper whiskers as the boundaries of the data distribution. Any data points that show above or below the whiskers, can be considered outliers or anomalous.

15 points

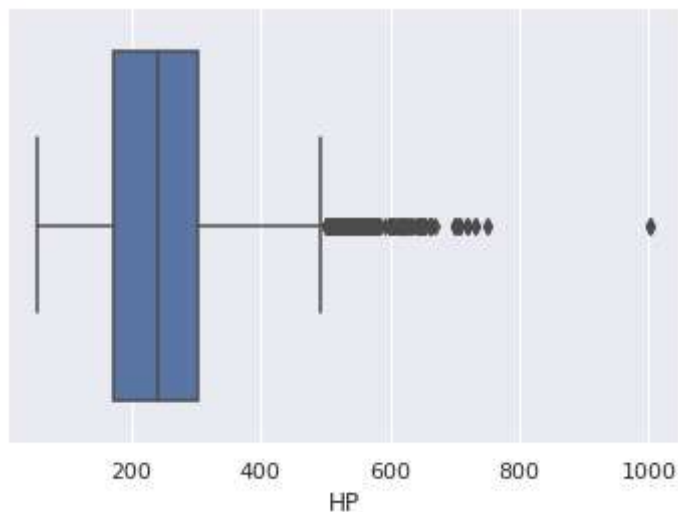
```
## Plot a boxplot for 'Price' column in dataset.  
sns.boxplot( x=df['Price'] )  
plt.show()
```



Observation:

Here as you see that we got some values near to 1.5 and 2.0 . So these values are called outliers. Because there are away from the normal values. Now we have detect the outliers of the feature of Price. Similarly we will checking of anothers features.

```
## PLoT a boxplot for 'HP' columns in dataset  
sns.boxplot( x=df['HP'] )  
plt.show()
```

Observation:

Here boxplots show the proper distribution of 25 percentile and 75 percentile of the feature of HP.

print all the columns which are of int or float datatype in df.

Hint: Use loc with condition

```
df.loc[:, df.dtypes != object]
```

	Year	HP	Cylinders	MPG_H	MPG-C	Price
0	2011	335.0	6.0	26	19	46135
1	2011	300.0	6.0	28	19	40650
2	2011	300.0	6.0	28	20	36350
3	2011	230.0	6.0	28	18	29450
4	2011	230.0	6.0	28	18	34500
...
11909	2012	300.0	6.0	23	16	46120
11910	2012	300.0	6.0	23	16	56670
11911	2012	300.0	6.0	23	16	50620
11912	2013	300.0	6.0	23	16	50920
11913	2006	221.0	6.0	26	17	28995

10827 rows × 6 columns

Save the column names of the above output in variable list named 'l'

```
l=list(df.loc[:, df.dtypes != object].columns)
```

▼ Outliers removal techniques

Reference: <https://www.youtube.com/watch?v=A3gClkblXK8>

25 points

1. Using IQR Technique

Here comes cool Fact for you!

IQR is the first quartile subtracted from the third quartile; these quartiles can be clearly seen on a box plot on the data.

The anatomy of boxplot is given below.



- Calculate IQR and give a suitable threshold to remove the outliers and save this new dataframe into df2.

Let us help you to decide threshold: Outliers in this case are defined as the observations that are below $(Q1 - 1.5 \times IQR)$ or above $(Q3 + 1.5 \times IQR)$

```
## Your code here
Q1 = df.quantile(0.25)
Q3 = df.quantile(0.75)
IQR = Q3 - Q1 #interquantile range
df2 = df[~((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))).any(axis=1)]
```

2. Outlier removal using Z-score function

Reference:- <https://www.youtube.com/watch?v=rfmfXa0kdrQ>

- The intuition behind Z-score is to describe any data point by finding their relationship with the Standard Deviation and Mean of the group of data points.

We will use Z-score function defined in scipy library to detect the outliers in dataframe df having columns which are in variable 'l'

```
z = np.abs(stats.zscore(df[1]))
print(z)
```

	Year	HP	Cylinders	MPG_H	MPG-C	Price
0	0.014743	0.732425	0.174386	0.041059	0.049314	0.058463
1	0.014743	0.413769	0.174386	0.225455	0.049314	0.029591
2	0.014743	0.413769	0.174386	0.225455	0.101214	0.098621
3	0.014743	0.223542	0.174386	0.225455	0.199843	0.209390
4	0.014743	0.223542	0.174386	0.225455	0.199843	0.128320
...
11909	0.157006	0.413769	0.174386	0.440829	0.500900	0.058222
11910	0.157006	0.413769	0.174386	0.440829	0.500900	0.227587
11911	0.157006	0.413769	0.174386	0.440829	0.500900	0.130463
11912	0.299270	0.413769	0.174386	0.440829	0.500900	0.135279
11913	0.696575	0.305482	0.174386	0.041059	0.350371	0.216695

```
[10827 rows x 6 columns]
```

Hey buddy! do you understand the above output? Difficult right? let's try and define a threshold to identify an outlier so that we get a clear picture of whats going on.

We will not spare you without a good fact! ;)

In most of the cases a threshold of 3 or -3 is used i.e if the Z-score value is greater than or less than 3 or -3 respectively, that data point will be identified as outliers.

```
# print the values in dataframe which are less than the threshold and save this dataframe as
threshold = 3
df3 = df[(z < threshold).all(axis=1)]
df3
```

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

print the shape difference of df df2 and df3.

```

1      BMW      1      2011      335.0      6.0      MANUAL      rear wheel      26      19      46135
print(df.shape)
print(df2.shape)
print(df3.shape)

(10827, 10)
(9191, 10)
(10338, 10)

```

drive

Interesting right? Bam! you have removed 489 rows from the dataframe which was detected as outlier by Z-score technique. and removed 1636 rows from the dataframe which was detected as outlier by IQR technique.

By the way there are many other techniques by which you can remove outliers. You can explore on more interesting techniques available.

We know you must be having many questions in you mind like:

- Which technique we should use and why?
- Is it necessary that whatever detected as outlier are really outliers?

Don't worry these dilemma is faced by many data analyst. We provide you with good references below for you to explore further on this

- <https://www.theanalysisfactor.com/outliers-to-drop-or-not-to-drop/>
- <https://www.researchgate.net/post/Which-is-the-best-method-for-removing-outliers-in-a-data-set>

Lets find unique values and there counts in each column in df using value counts function.

Reference: <https://www.youtube.com/watch?v=8wRfskrpTlk>

```

for i in df.columns:
    print ("----- %s -----" % i)
    print(df[i].value_counts())

19      793
20      742
14      603
22      571

```

```
21      551
13      537
23      425
25      392
24      372
12      282
27      243
26      207
11      187
28      160
30      127
31      116
29      98
10      76
9       33
32      21
34      20
36      20
40      19
44      18
42      17
41      17
35      15
33      13
53      13
43      13
54      10
8        9
37        8
39        6
51        6
50        6
128       6
49        4
137       3
85        3
55        3
47        2
58        2
129       1
7         1
38        1
```

Name: MPG-C, dtype: int64

```
----- Price -----
2000      599
29995      18
25995      16
20995      15
27995      15
...
66347      1
62860      1
48936      1
68996      1
```

Visualising Univariate Distributions

Reference: -<https://www.youtube.com/watch?v=ll5-7rX3xPY>

We will use seaborn library to visualize eye catchy univariate plots.

Do you know? you have just now already explored one univariate plot. guess which one? Yeah its box plot.

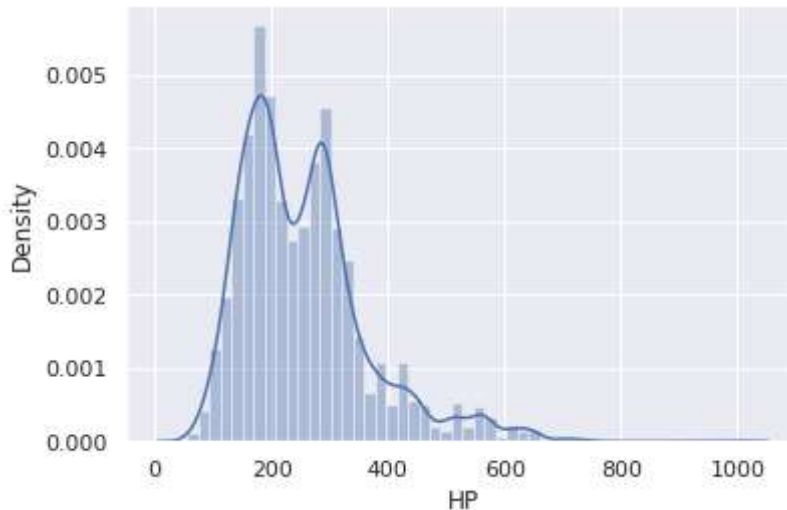
1 . Histogram & Density Plots

15 points

Histograms and density plots show the frequency of a numeric variable along the y-axis, and the value along the x-axis. The `sns.distplot()` function plots a density curve. Notice that this is aesthetically better than vanilla `matplotlib`.

```
#plotting distplot for variable HP  
sns.distplot(df['HP'])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7ffa9131d410>



Observation: We plot the Histogram of feature HP with help of distplot in seaborn.

In this graph we can see that there is max values near at 200. similarly we have also the 2nd highest value near 400 and so on.

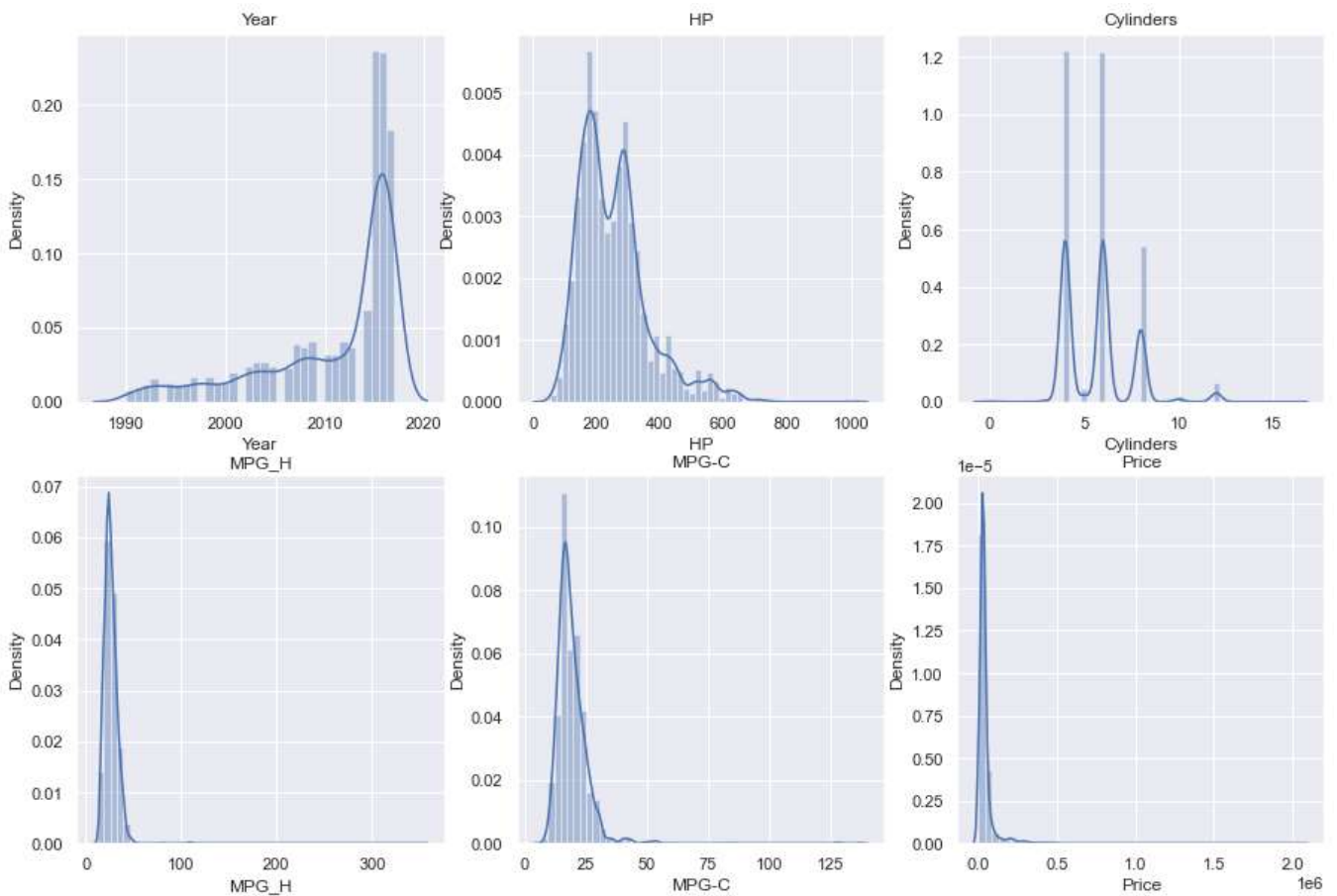
It represents the overall distribution of continuous data variables.

Since seaborn uses matplotlib behind the scenes, the usual matplotlib functions work well with seaborn. For example, you can use subplots to plot multiple univariate distributions.

- Hint: use matplotlib subplot function

Reference:- <https://www.youtube.com/watch?v=Tf-dgRR1PMA>

```
# plot all the columns present in list l together using subplot of dimention (2,3).
c=0
plt.figure(figsize=(15,10))
for i in l:
    c=c+1
    plt.subplot(2, 3, c)
    plt.title(i)
    sns.distplot(df[i])
plt.show()
```

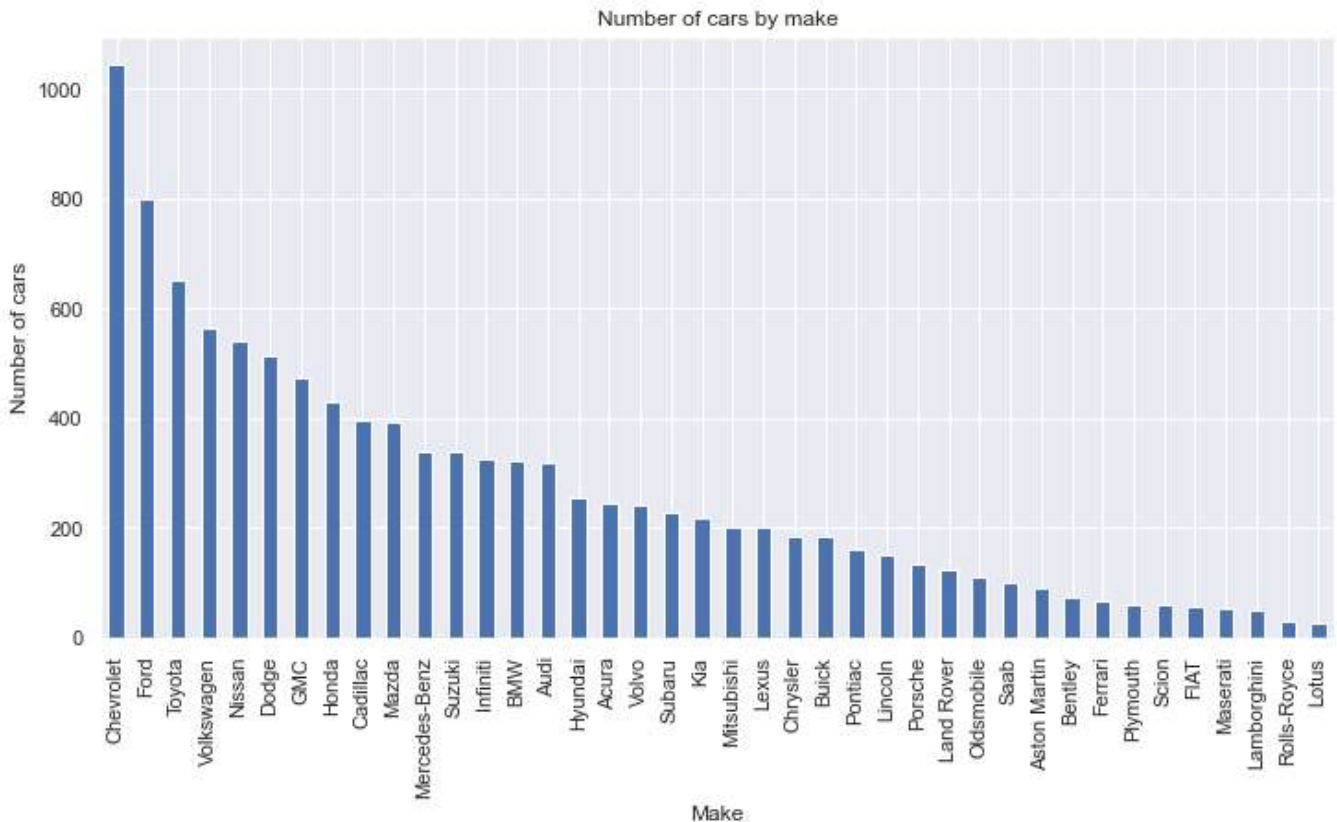


2. Bar plots

10 points

Plot a histogram depicting the make in X axis and number of cars in y axis.

```
plt.figure(figsize = (12,8))
df.Make.value_counts().nlargest(40).plot(kind='bar', figsize=(12,6))
plt.title("Number of cars by make")
plt.ylabel('Number of cars')
plt.xlabel('Make')
```



Observation: In this plot we can see that we have plot the bar plot with the cars model and nos. of cars.

3. Count Plot

Reference: - <https://www.youtube.com/watch?v=8U5h3EJuu8M>

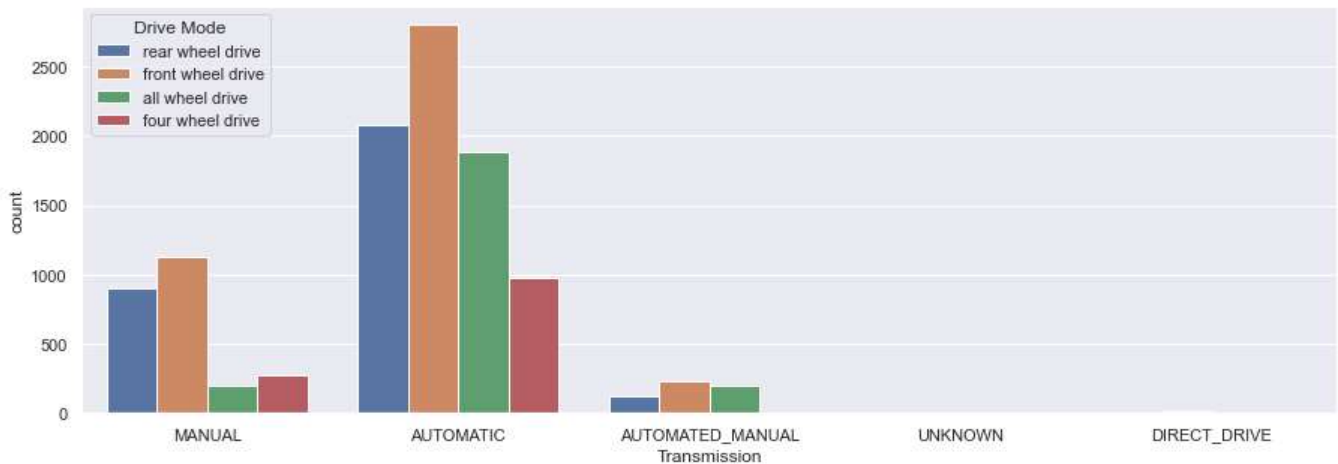
10 points

A count plot can be thought of as a histogram across a categorical, instead of quantitative, variable.

Plot a countplot for a variable Transmission vertically with hue as Drive mode

```
plt.figure(figsize=(15,5))

sns.countplot(x='Transmission', hue="Drive Mode", data=df)
plt.show()
# 'Cylinders', y='Price'
```



Observation: In this count plot, We have plot the feature of Transmission with help of hue. We can see that the the nos of count and the transmission type and automated manual is plotted. Drive mode as been given with help of hue.

▼ Visualising Bivariate Distributions

Bivariate distributions are simply two univariate distributions plotted on x and y axes respectively. They help you observe the relationship between the two variables.

▼ 1. Scatterplots

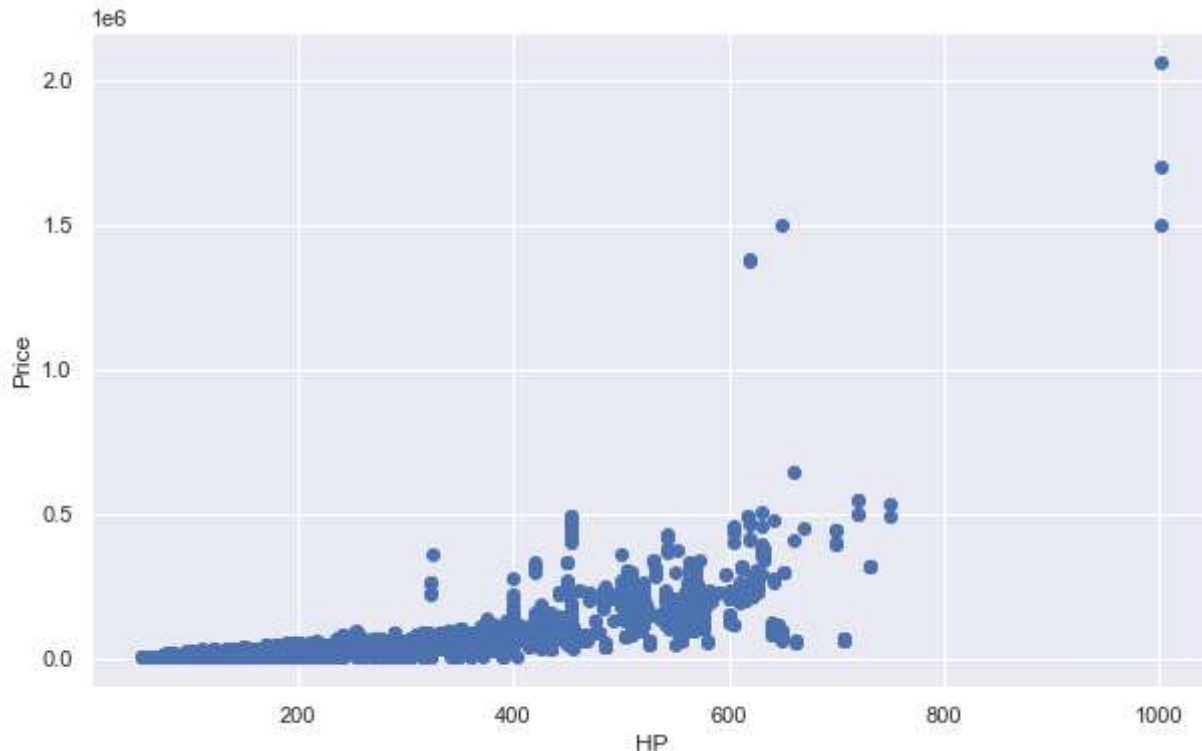
Reference:- <https://www.youtube.com/watch?v=ls9bc-WA-c8>

Scatterplots are used to find the correlation between two continuos variables.

10 points

Using scatterplot find the correlation between 'HP' and 'Price' column of the data.

```
## Your code here -
fig, ax = plt.subplots(figsize=(10,6))
ax.scatter(df['HP'], df['Price'])
ax.set_xlabel('HP')
ax.set_ylabel('Price')
plt.show()
```



Observation:

It is a type of plot or mathematical diagram using Cartesian coordinates to display values for typically two variables for a set of data.

We have plot the scatter plot with x axis as HP and y axis as Price.

The data points between the features should be same either wise it give errors.

4. joint distributions

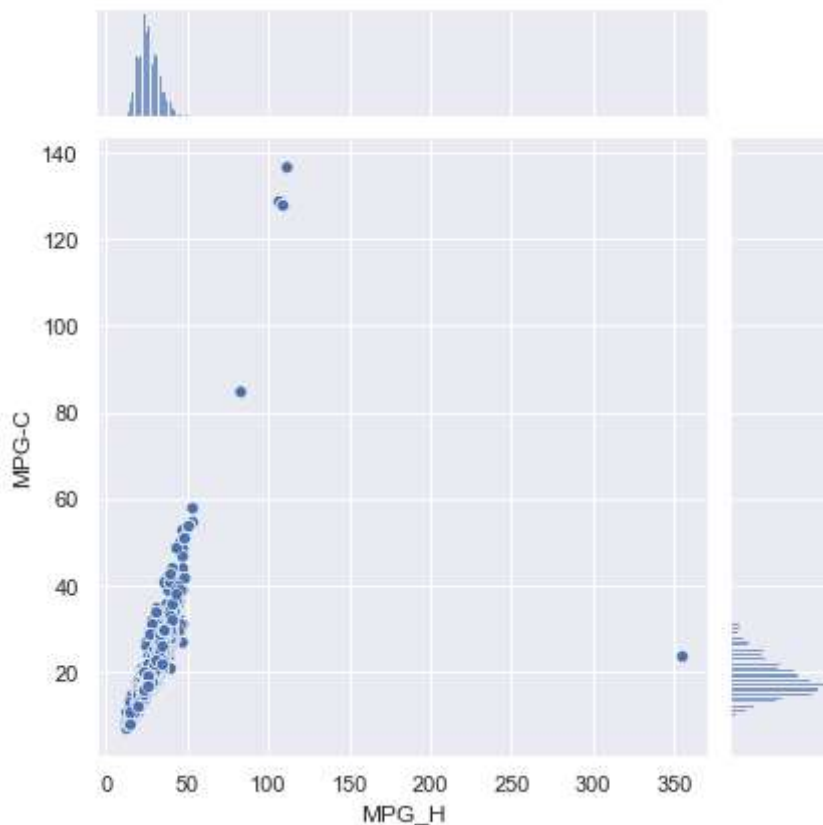
Reference:-<https://www.youtube.com/watch?v=LrSLBbe6pGY>

30 points

Seaborn's jointplot displays a relationship between 2 variables (bivariate) as well as 1D profiles (univariate) in the margins. This plot is a convenience class that wraps JointGrid

```
# joint plots of MPG_H and MPG-C
sns.jointplot('MPG_H', 'MPG-C', df)
```

```
plt.show()
```



Observations:

Jointplot is library specific and can be used to quickly visualize and analyze the relationship between two variables and describe their individual distributions on the same plot.

In this plot we can see the relationship of MPG-C abd MPG_H.

You can adjust the arguments of the jointplot() to make the plot more readable.

5. Plotting Aggregated Values across Categories

Reference:-<https://www.youtube.com/watch?v=yWYWmeuH7no>

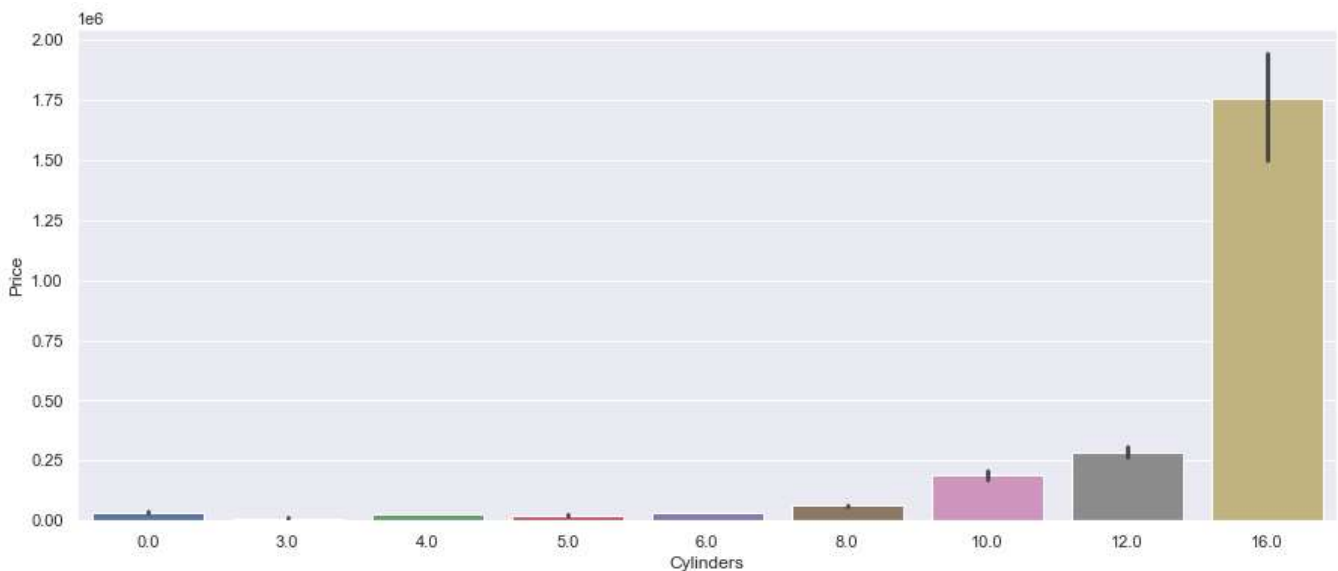
Bar Plots - Mean, Median and Count Plots

30 points

Bar plots are used to **display aggregated values** of a variable, rather than entire distributions. This is especially useful when you have a lot of data which is difficult to visualise in a single figure.

For example, say you want to visualise and *compare the Price across Cylinders*. The `sns.barplot()` function can be used to do that.

```
# bar plot with default statistic=mean between Cylinder and Price
plt.figure(figsize=(15,6))
sns.barplot(x='Cylinders', y='Price', data=df)
plt.show()
```



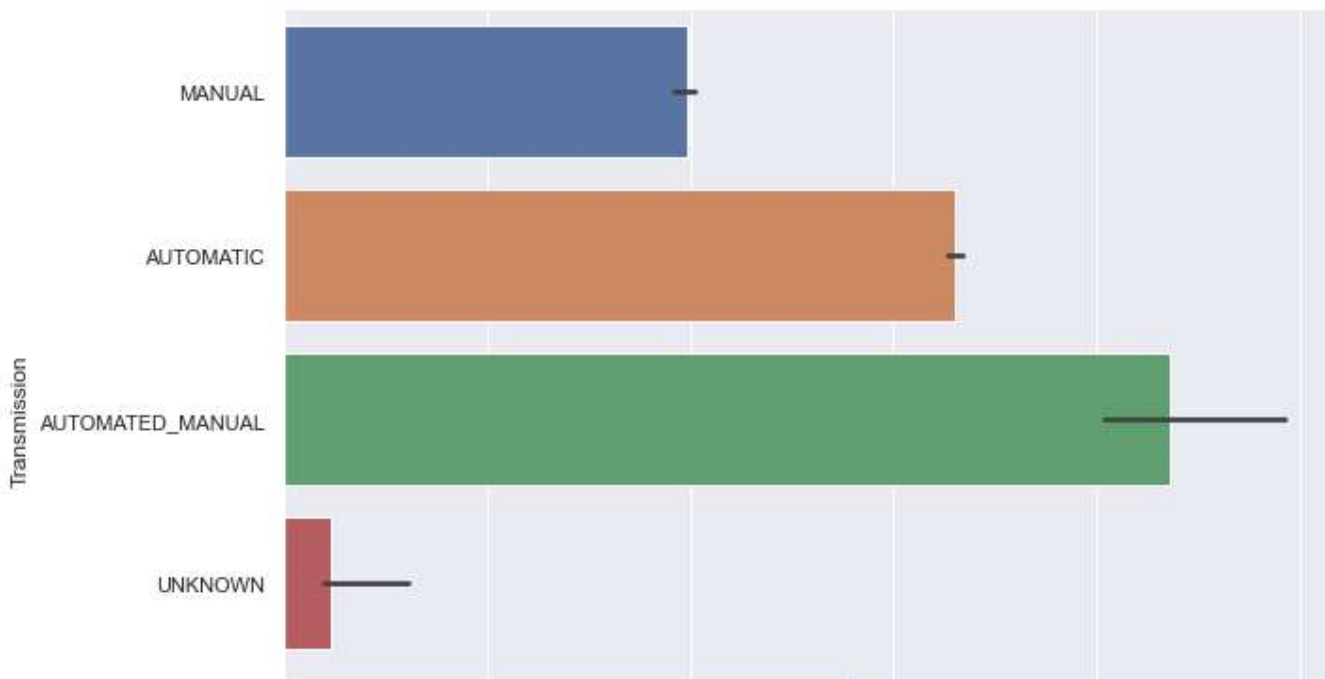
Observation:

By default, seaborn plots the mean value across categories, though you can plot the count, median, sum etc.

Also, barplot computes and shows the confidence interval of the mean as well.

When you want to visualise having a large number of categories, it is helpful to plot the categories across the y-axis. Let's now *drill down into Transmission sub categories*.

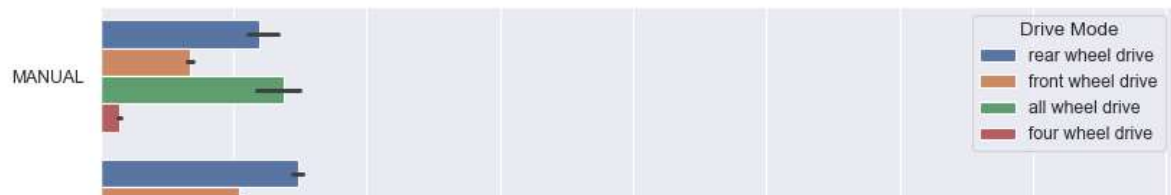
```
# Plotting categorical variable Transmission across the y-axis
plt.figure(figsize=(10, 8))
sns.barplot(x='Price', y="Transmission", data=df, estimator=np.median)
plt.show()
```



Plot bar plot for Price and Transmission with hue="Drive Mode"

```
plt.figure(num=None, figsize=(12, 8), dpi=80, facecolor='w', edgecolor='k')
```

```
# specify hue="Drive Mode"
sns.barplot(x='Price', y='Transmission', hue="Drive Mode", data=df, estimator=np.median)
plt.show()
```



These plots look beautiful, isn't it? In Data Analyst life, such charts are their unavoidable friend.:)

☐

Multivariate Plots



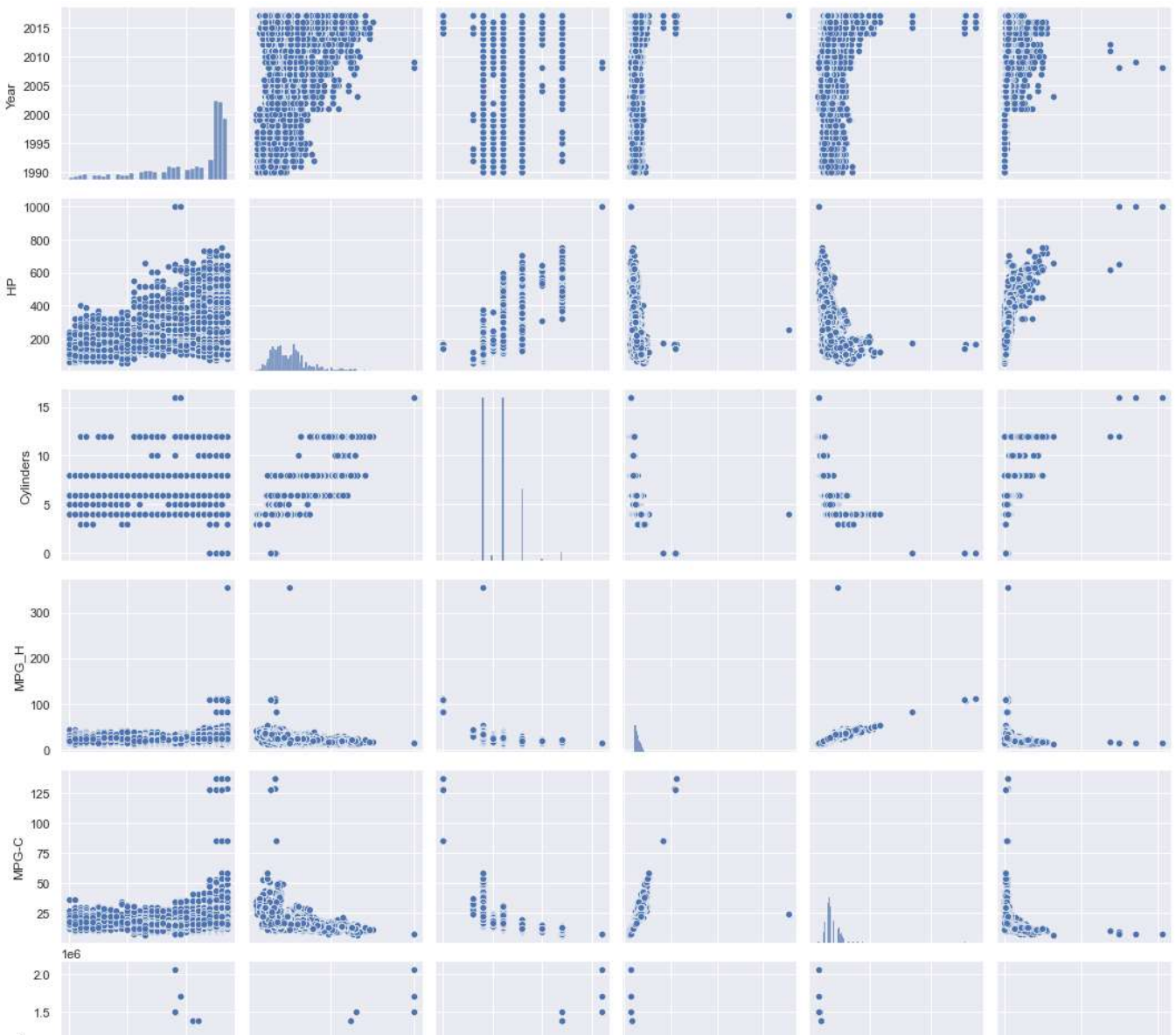
▼ 1. Pairplot

Reference:- <https://www.youtube.com/watch?v=TexdD7t0IKU>

10 points

Plot a pairplot for the dataframe df.

```
sns.pairplot(df)
plt.show()
```



Observation:

To plot multiple pairwise bivariate distributions in a dataset, you can use the `pairplot()` function. This shows the relationship for $(n, 2)$ combination of variable in a DataFrame as a matrix of plots and the diagonal plots are the univariate plots.

2. Heatmaps

Reference:- <https://www.youtube.com/watch?v=ZSwXRn50lnA>

A heat map is a two-dimensional representation of information with the help of colors. Heat maps can help the user visualize simple or complex information

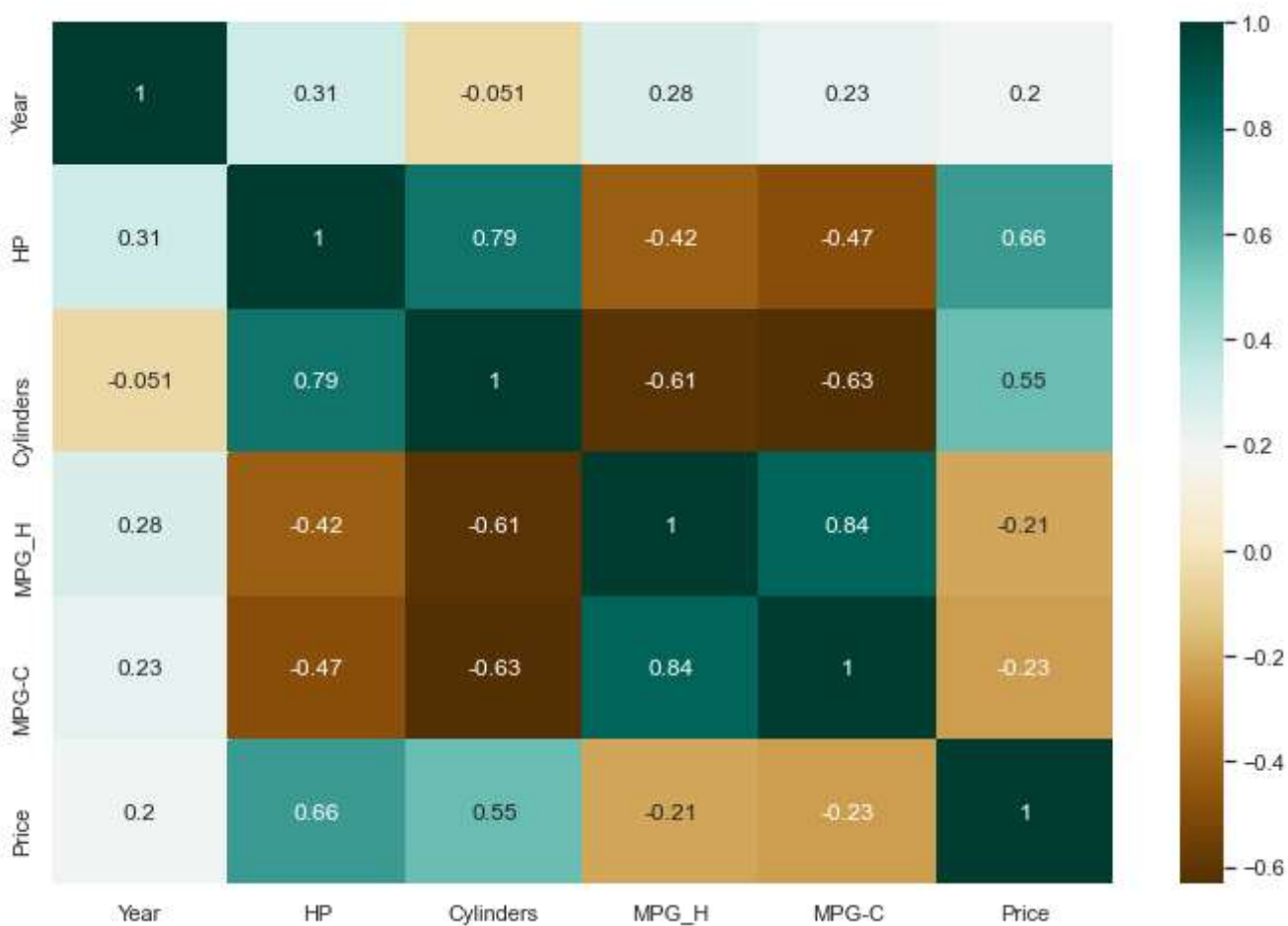
20 points

Using heatmaps plot the correlation between the features present in the dataset.

```
#find the correlation of features of the data
corr = df.corr()
corr
```

	Year	HP	Cylinders	MPG_H	MPG-C	Price
Year	1.000000	0.314971	-0.050598	0.284237	0.234135	0.196789
HP	0.314971	1.000000	0.788007	-0.420281	-0.473551	0.659835
Cylinders	-0.050598	0.788007	1.000000	-0.611576	-0.632407	0.554740
MPG_H	0.284237	-0.420281	-0.611576	1.000000	0.841229	-0.209150
MPG-C	0.234135	-0.473551	-0.632407	0.841229	1.000000	-0.234050
Price	0.196789	0.659835	0.554740	-0.209150	-0.234050	1.000000

```
# Using the correlated df, plot the heatmap
# set cmap = 'BrBG', annot = True - to get the same graph as shown below
# set size of graph = (12,8)
plt.figure(figsize=(12,8))
sns.heatmap(corr,cmap='BrBG', annot= True)
plt.show()
```



Observation:

A heatmap contains values representing various shades of the same colour for each value to be plotted. Usually the darker shades of the chart represent higher values than the lighter shade. For a very different value a completely different colour can also be used.

The above heatmap plot shows correlation between various variables in the colored scale of -1 to 1.

Amazing work done ! you have really made eye catchy visualization plots so far. Did you felt its complicate to understand the above plot?. Hey smarty don't worry, in near assignments you will have enough practise to analyse and prepare insights from such plots that you will become pro in this field.

Then soon you will be like below meme image.png

Have a sweet cookie:) Congratulations! you have completed the 6th milestone challenge too.

FeedBack

We hope you've enjoyed this course so far. We're committed to helping you use AlforAll course to its full potential so you can grow with us. And that's why we need your help in form of a feedback here

We appreciate your time for your thoughtful comment.

Comment

Rate the following in the scale of 1 to 10

Assignment Difficulty Level

Rate

Question Clarity Level

#Rate

YouTube video content quality.

#Rate

