

Exploratory Data Analysis

In [16]:

```
1 # Importing Libraries :
2
3 # Import pandas library and make it as pd:
4 import pandas as pd
5
6 # Import Numerical Python Library and make it as np :
7 import numpy as np
8
9 # Import pyplot from Matplotlib Library and make it as plt:
10 import matplotlib.pyplot as plt
11
12
13 # Import Seaborn Library and Make it as SNS:
14 import seaborn as sns
15
16 # Inorder to Supress the Warnings Import FilterWarnings:
17 from warnings import filterwarnings
18 filterwarnings('ignore')
```

In []:

```
1
```

What – is – Data?

Data is units of information in structured or unstructured format.

- Collection of relevant tweets
- Records of yield in a farm over a period of time
- Records of stock price every minute
- Records of performance of a sports person

CLASSIFICATION – OF – DATA :

1. Time Series Data

2. Cross Sectional Data

3. Pooled Data

4. Panel Data

```
In [ ]: ⏷ 1
```



```
In [ ]: ⏷ 1
```

1. Time Series Data

- Time series data is the set of observations on a variable at different time points
- The data may be collected daily, weekly, monthly, or annually

Time-series data	
Year	Numbers of cars rented by alamo
2000	10000
2001	8000
2002	11000
2003	5000
2004	12000

	A	B	C
1	Year	Quarter	Sales
2	2012	1	\$165,000.00
3		2	\$253,000.00
4		3	\$316,000.00
5		4	\$287,000.00
6	2013	1	\$257,000.00
7		2	\$308,000.00
8		3	\$376,000.00
9		4	\$351,000.00

2. Cross Sectional Data

Cross-sectional data are set of observations on two or more variables at the same time point

City	Date	MaxTemperature	Humidity	Wind
NYC	1/1/2015	55	45%	4 mph
NYC	1/1/2014	30	39%	16 mph
NYC	1/1/2013	47	65%	21 mph
SFO	1/1/2015	70	35%	21 mph
SFO	1/1/2014	75	23%	2 mph
SFO	1/1/2013	71	39%	13 mph
Boston	1/1/2015	34	39%	16 mph
Boston	1/1/2014	26	17%	27 mph
Boston	1/1/2013	45	46%	18 mph

3. Pooled Data

Pooled data is combination of both time series data and cross sectional data.

Test questions	Original n	Original M	Original SE	Pooled M	Pooled SE	χ^2	Asymptotic significant ^a	Exact significant (two-tailed) ^b
Anchor Item Q1	290	0.272	0.026	0.297	0.027	2.828	.985	1.000
Anchor Item Q2	285	0.179	0.023	0.204	0.028	7.142	.712	1.000
Anchor Item Q3	252	0.202	0.025	0.215	0.032	9.938	.446	1.000
Anchor Item Q4	250	0.304	0.029	0.340	0.037	11.377	.329	1.000
Anchor Item Q5	271	0.472	0.030	0.489	0.030	2.488	.991	1.000
Anchor Item Q6	259	0.467	0.031	0.473	0.032	3.517	.967	1.000
Anchor Item Q7	282	0.298	0.027	0.332	0.031	5.493	.856	1.000
Anchor Item Q8	282	0.294	0.027	0.316	0.030	4.749	.907	1.000
Anchor Item Q9	290	0.172	0.022	0.207	0.030	9.886	.451	1.000

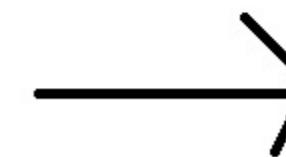
^aKruskal-Wallis test, $df = 10$.

^bBinomial distribution used for McNemar.

4. Panel Data

- Panel data is type of pooled data
- The same cross-sectional unit is surveyed over time
- Also known as longitudinal or micro-panel data

Person ID	Age	Sex	Year	<i>Income</i>
1	27	1	2015	1600
1	28	1	2016	1500
2	42	2	2015	1900
2	43	2	2016	2000
2	44	2	2017	2100
3	34	1	2015	3300



	A	B	C	D
1		FIRM A	FIRM B	FIRM C
2	3/01/1995	0.191439	0.780436	0.329861
3	4/01/1995	-0.04008	0.332088	0.306904
4	5/01/1995	-0.14042	0.681645	0.820281
5	6/01/1995	-5.65741	0.272911	0.637198
6	9/01/1995	5.898009	0.212024	0.173159
7	10/01/1995	-3.1429	0.719781	0.078806
8	11/01/1995	-6.70905	0.230741	0.090504
9	12/01/1995	6.139131	0.201295	0.409763
10	13/01/1995	-3.19934	0.635427	0.561843
11	16/01/1995	3.282443	0.636448	0.43751
12	17/01/1995	-0.05193	0.655553	0.791657
13	18/01/1995	15.7144	0.013025	0.057594
14	19/01/1995	11.65196	0.535783	0.657062
15	20/01/1995	0.220995	0.65229	0.872535
16	23/01/1995	-14.755	0.092412	0.171942
17	24/01/1995	2.882044	0.68886	0.360246
18	25/01/1995	0.186261	0.544323	0.677453
19				

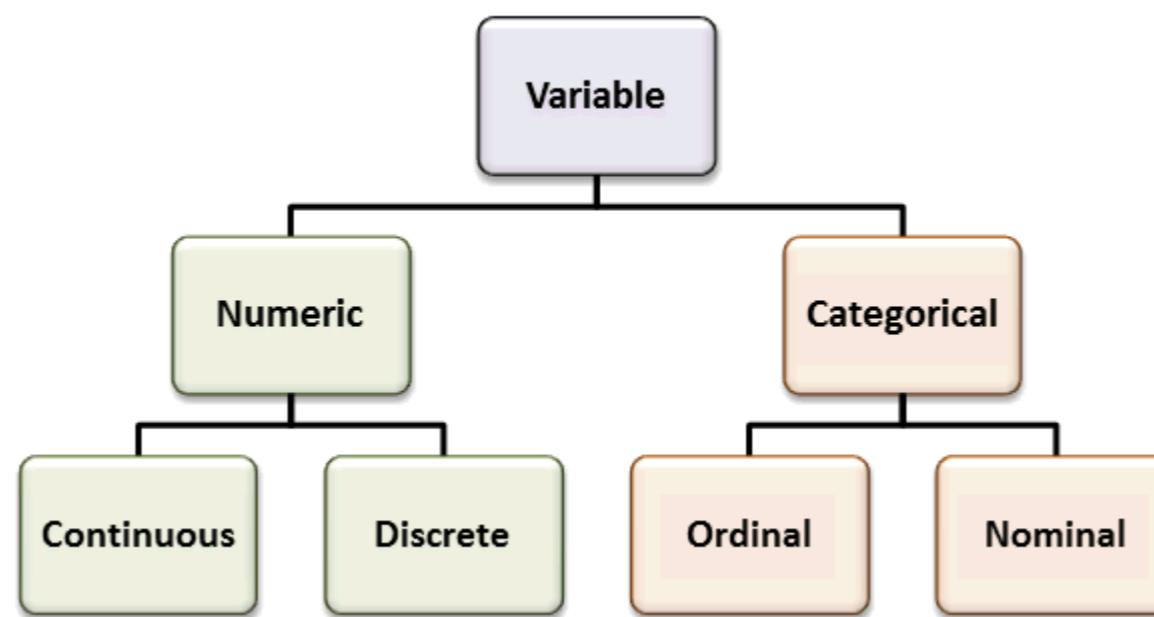
original

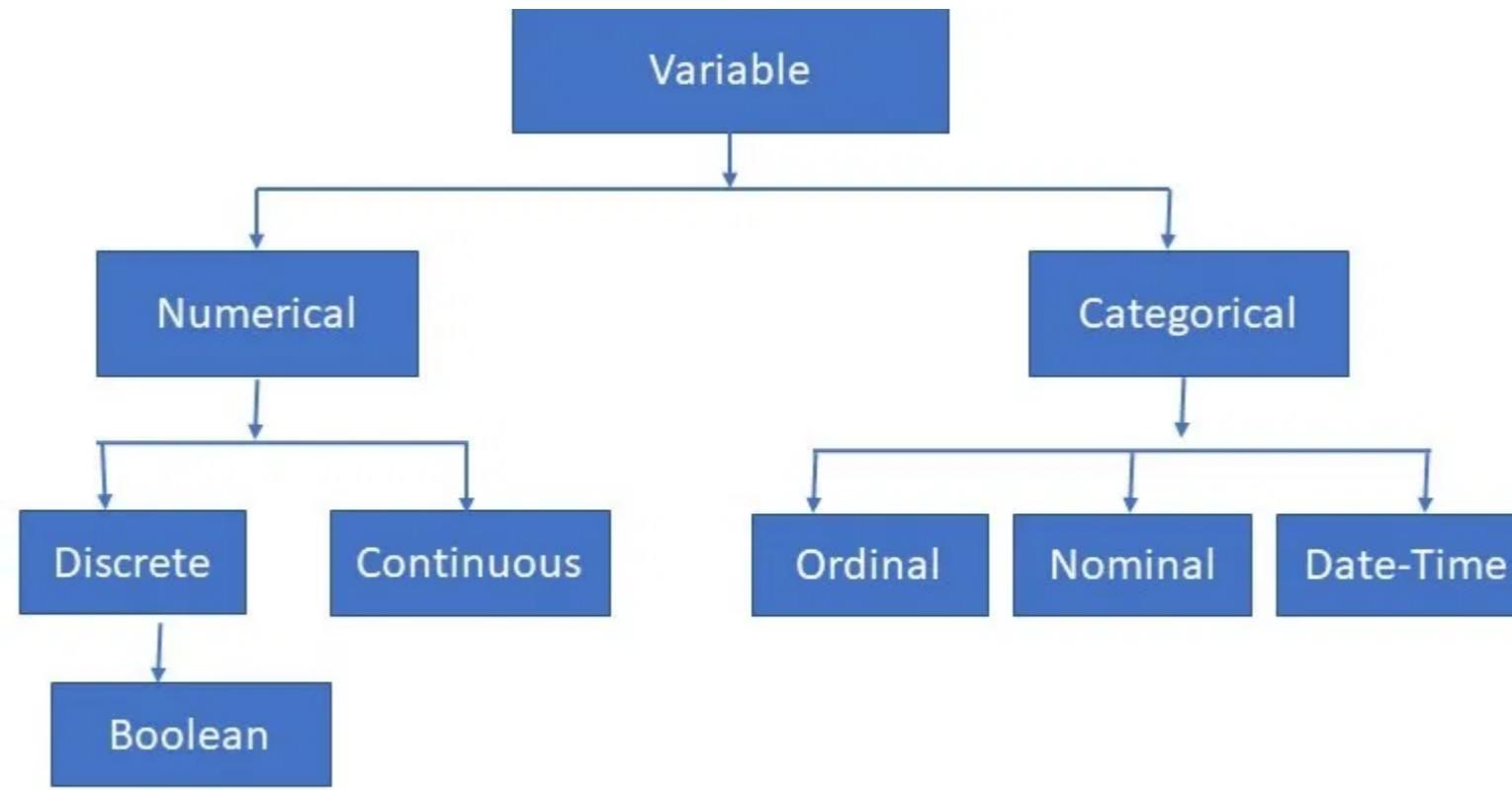
3/01/1995	FIRM A	0.821848
4/01/1995	FIRM A	-0.04008
5/01/1995	FIRM A	-0.14042
6/01/1995	FIRM A	-5.65741
9/01/1995	FIRM A	5.898009
10/01/1995	FIRM A	-3.1429
11/01/1995	FIRM A	-6.70905
3/01/1995	FIRM B	6.139131
4/01/1995	FIRM B	-3.19934
5/01/1995	FIRM B	3.282443
6/01/1995	FIRM B	-0.05193
9/01/1995	FIRM B	15.7144
10/01/1995	FIRM C	11.65196
11/01/1995	FIRM C	0.220995
23/01/1995	FIRM C	-14.755
24/01/1995	FIRM C	2.882044
25/01/1995	FIRM C	0.186261

panel form I want

Understand the variable:

To understand the descriptive statistics, we need to first understand the type of each variable or column





Categorical Variable

Nominal Data:

- Nominal data has no order and has two or more than two categories
- It represents discrete units and are used to label variables
- Example: Housing type - Apartment, Bungalow, Penthouse

Binary Data:

- Binary data has no order and has strictly two categories
- Also known as dichotomous variable
- Example: Presence or absence of a disease

Categorical Variable

Ordinal Data:

- Ordinal data is ordered nominal data
- It represents discrete and ordered
- Example: Education - Primary, Secondary, High School, College and University which are ordered

IMPORTANT:

- At times data may appear to be numeric but is actually categorical
- Consider the adjoining table. The 'Gender' column has value 1 and 0
- The numeric values used to represent a categorical variable cannot be used for numerical computation

```
In [12]: 1 import pandas as pd
```

```
In [14]: 1 _dict_ = {'Name' : ['Riya', 'Arya', 'Sam', 'Joe', 'Harry', 'Sarah'],
2           'Gender' : [0, 0, 1, 1, 1, 0],
3           'Age' : [45, 23, 34, 54, 12, 43]}
4
5 DataFrame = pd.DataFrame(_dict_)
6 DataFrame
```

Out[14]:

	Name	Gender	Age
0	Riya	0	45
1	Arya	0	23
2	Sam	1	34
3	Joe	1	54
4	Harry	1	12
5	Sarah	0	43

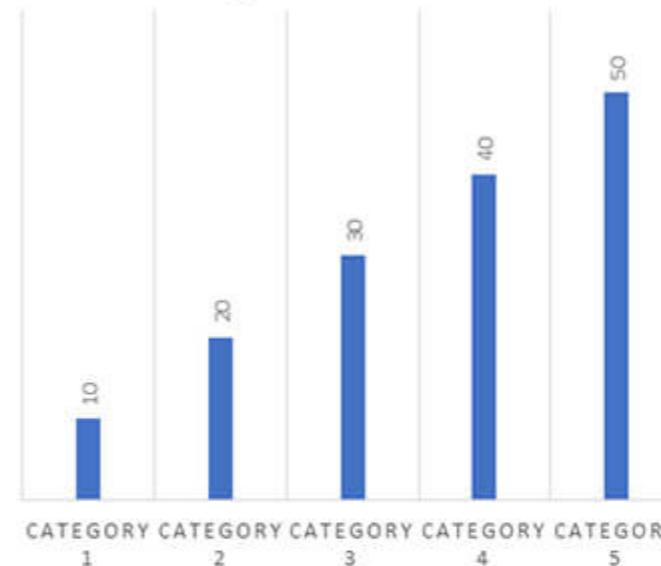
Data Analysis:

- Descriptive
- Inferential

Descriptive:

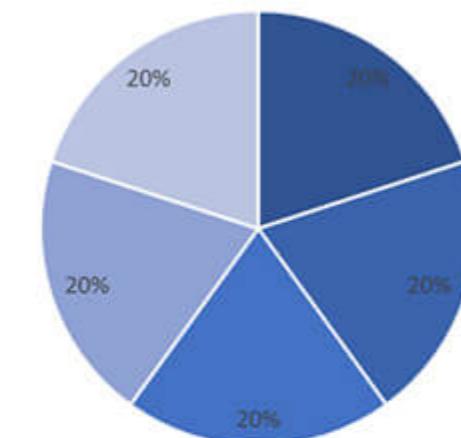
- Descriptive statistics is the term given to the analysis of the data that helps describe, visualize or summarize the data
- Helps exhibit the patterns in the data
- However, descriptive statistics only help in understanding the data and not make conclusions regarding any hypothesis
- We can not generalise the facts, they are confined only to the data at hand (that may be the sample)

Descriptive Statistics



- Describes the characteristics of data
- **Example:** Population, Frequency of the variables

Inferential Statistics



- 1. Category ■ 2. Category ■ 3. Category ■ 4. Category ■ 5. Category

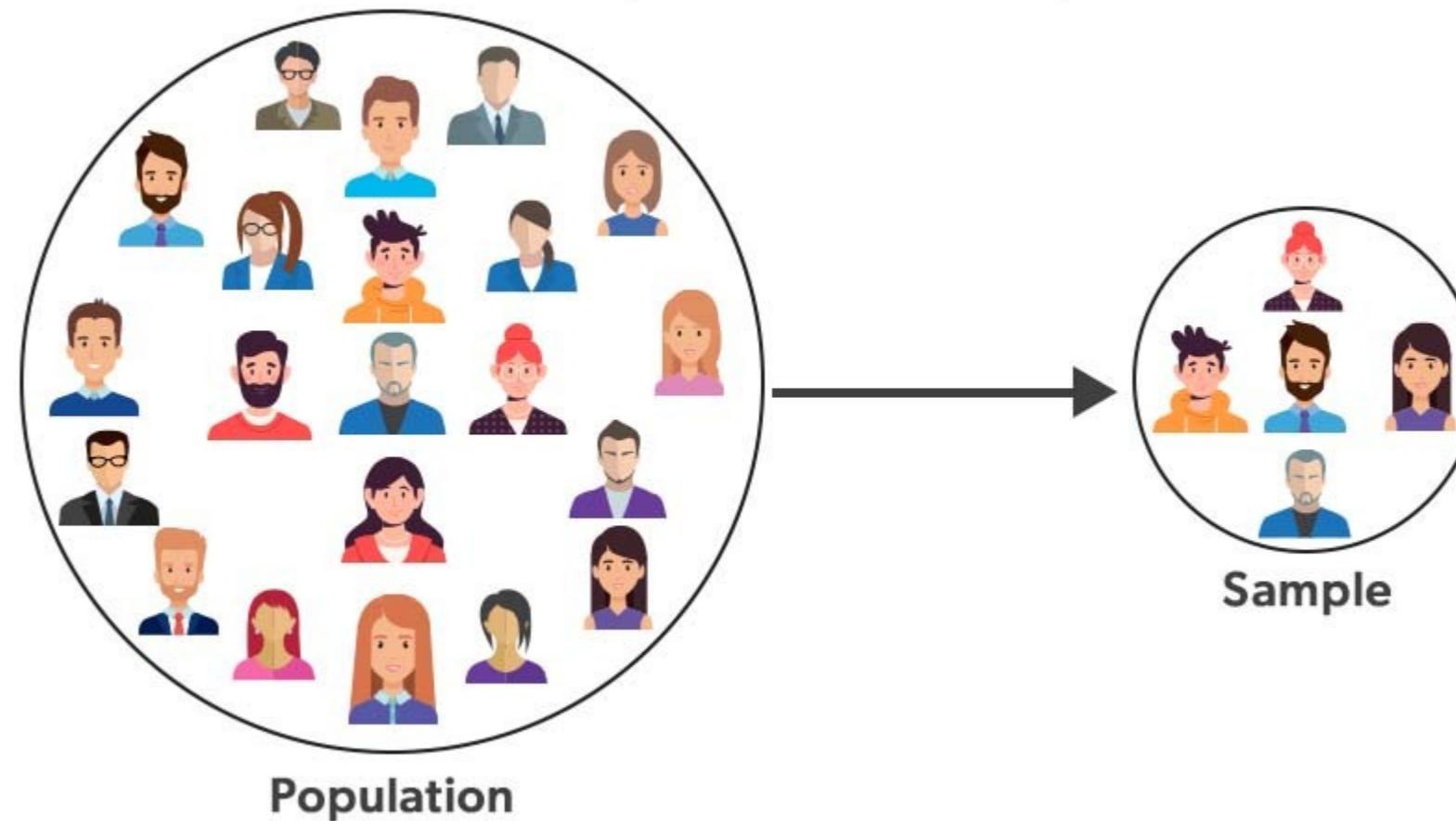
- Studies the sample of the same data
- **Example:** Grade, Percentile



Population & Sample:

- Population is a collection of all the individuals or objects
- Sample is a subset of the population which is a representative of the population

Population and Sample



Example:

Suppose a company producing electric bulbs wants to know the average life of a bulb. If the details on average life of all bulbs are available, then this information is regarded as the population. It would be challenging for the company to test each and every bulb produced. In such a scenario, the company would draw a sample from the produced bulbs to test.

Descriptive statistics

The major characteristics we observe in each variable of the data set are:

1. Measures of central tendency
2. Distribution of the data
3. Measures of dispersion

Measures of Central Tendency

- A measure of central tendency is a single value that identifies the central position of the data
- It includes mean, median, mode and partition values

Mean

- Mean is a central tendency of the data
- It is defined as the sum of all observations divided by the number of observations: where N is the total number of observations
- The whole data is spread out around the mean

Mean

Merits:

- Based on all observation

Demerits:

- Affected by extreme observations (outliers)
- Arithmetically it is not possible to obtain the mean of a data with missing values
- Mean can not be calculated for categorical data

Outlier

- An outlier is a value that behaves differently than other observations
- It does not follow the usual pattern

Obtain the mean

Each column of a DataFrame is a pandas Series.

Let us create a series and find its mean.

In []: 1

Trimmed mean

- Arithmetic mean obtained by ignoring lowest and highest $\alpha\%$ observations is called as the $\alpha\%$ - trimmed mean
- Using a trimmed mean, it is possible to eliminate the influence of outliers
- Note the observations are arranged either in increasing or decreasing order
- Also called as truncated mean

In [2]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt

In [76]: 1 prices_missing = pd.Series([1, 1, 4, 56, 7, 8, 90, 34, 67, 4, 76, 32, 8, 45, 35, 87,
2 np.NaN, 89, 65, 90, 74, 88, 34, 99, 88, 67])
3
4 print('Mean : ', prices_missing.mean())
5 print('Median : ', prices_missing.median())

Mean : 49.96
Median : 56.0

In [14]: 1 # Obtain a 20% of trimmed Mean:
2
3 import scipy
4 from scipy import stats
5
6 scipy.stats.trim_mean(prices_missing, proportiontocut = 0.20)

Out[14]: 54.0

In []: 1

Median

- Median is the middle most observation in the data when it is arranged either in ascending or descending order of their values
- It divides the data into two equal parts. Thus, it is a positional average
- Median will be the middle term, if the number of observations is odd
- Median will be average of middle two terms, if number of observations is even

Obtain the median

Using mean and median

- For a numeric variable, the missing value is replaced by the mean if there are no outliers present

- In case if outliers are present the missing values are replaced by the median
- For symmetric data median = mean, so either of the value can be used

In []: 1

Mode:

- Mode of the data is the value which has the highest frequency
- In other words, it is most repeated observation
- A set of observations having two modes is said to be bimodal
- A set of observations may have more than two modes. Such data is said to be multimodal

Mode Merits:

- It is applicable to both numeric and categorical data
- It is not affected by extreme observations

Demerits:

- It is not based on all observations

In []: 1

OBTAİN MODE:

In []: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt

In [61]: 1 # Create a Mode:
2
3 prices_pens = pd.Series([0, 35, 67, 34, 10, 56, 100, 67,
4 99, 45, 100, 45, 36, 99, 56])
5
6 prices_pens.value_counts()

Out[61]: 67 2
56 2
100 2
99 2
45 2
0 1
35 1
34 1
10 1
36 1
dtype: int64

```
In [62]: 1 # print the Mode :  
2  
3 prices_pens.mode()
```

```
Out[62]: 0    45  
1    56  
2    67  
3    99  
4   100  
dtype: int64
```

Measure of Dispersion:

- The measure of Dispersion refers to the Variability with the Data.
- Variability is the Measure of how Close or Far the Data from the Central Value.
- Reliability of a Measure of Central Tendency is more dispersion in the Data is less.
- It includes Range, variance, Standard Deviation, Coefficient of Variation and Interquartile Range.

Range:

- Range is the Difference B/w the Largest and Smallest Observation.
- Range = $X_n - X_1$
- Where X_n is the Largest Value and X_1 is the Smallest value.

Merits:

- It is a basic way to measure the Variability

Demerits:

- If Either of the Smallest or the Largest Values turns out to be an Outlier then the Range Value is Unreliable.
- The Other in-Between Observation do not affect the range Data.
- It Doesn't give the Clear picture on the Data.

Obtain the Range:

```
In [63]: 1 year_establishment = pd.Series((1955, 1990, 1930, 2000))  
2  
3 Range = year_establishment.max() - year_establishment.min()  
4 Range
```

```
Out[63]: 70
```

Variance:

Variance is the arithmetic mean of squares of deviations taken from the mean.

Interpretation: The variance looks at how spread out the observations are from the mean. The higher the variance more the data is spread out.

Obtain the variance:

```
In [66]: 1 | prices_pens.var()
```

```
Out[66]: 1044.6857142857143
```

The Unit of measurement of variance

- The unit of measurement of the variance is not the same as that of the original data
- Consider an industry producing screws. The diameter of each screw is noted in millimeters
- The variance in the diameter is measured in squared Millimeters
- Thus a usual practice is to take square root of the variance which would have the same unit as the original data
- The value thus obtained is called the standard deviation of the data

Interquartile range

- The interquartile range is defined as the difference between the third and the first quartile
- It gives the range in which the middle 50% of the data lies
- The value of IQR is used to remove outliers in the data

```
In [78]: 1 | # Obtain the First Quartiles:  
2 | Q1 = prices_missing.quantile(0.25)  
3 |  
4 | # Obtain the Quantile:  
5 | Q3 = prices_missing.quantile(0.75)  
6 |  
7 | IQR = Q3 - Q1  
8 |  
9 | # Print the IQR:  
10 | IQR
```

```
Out[78]: 79.0
```

Interpretation : The InterQuartile range of prices_missing is 79.0

Z-score

- Every value in the dataset has the z-score
- To find the z-score of a value, subtract the mean from it and divide it by the standard
- The z-scores are extensively used in statistics especially in hypothesis testing
- It is also used to normalize the data, ie is it scales the data to get the mean of the scaled data as 0 and variance as 1

The z-score tells how much is the value away from the mean

```
1 0 ----> The corresponding value is same as the mean
2 < 0 ----> The corresponding value is less than the mean
3 > 0 ----> The corresponding value is greater than the mean
4
5 < -3 or > 3 ----> The corresponding value is potential outlier
```

Covariance

- Covariance is a measure of how much two random variables vary together
- Covariance is similar to variance, but variance explains how a single variable varies, and covariance explains how two variables vary together
- It cannot tell whether the value indicates a strong relationship or weak relationship, since the covariance can take any value

```
In [16]: # Create Imports Series:
1 imports_raw_material = pd.Series([10, 11, 14, 15, 20, 22, 16, 12, 15, 13])
2
3 # Create Export Series:
4 export_finished_products = pd.Series([11, 12, 14, 18, 21, 26, 44, 34, 12])
5
6 # Calculate the Covariance:
7 imports_raw_material.cov(export_finished_products)
```

Out[16]: 14.0

Correlation

- Correlation is the extent of linear relationship among numeric variables
- It indicates the extent to which two variables increase or decrease in parallel
- The value of a correlation coefficient ranges between -1 and 1
- A positive correlation exists if variable increase/decrease simultaneously
- A negative correlation exists if one variable increases, while the other variable decreases

```
In [17]: imports_raw_material.corr(export_finished_products)
```

Out[17]: 0.30909991855425206

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

Let us now do hands on practise with dataset

- We will work on mpg dataset
- mpg stands for miles per gallon
- The data is generally used to predict the miles per gallon of a car

```
In [7]: 1 import pandas as pd  
2 import numpy as np  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns
```

```
In [9]: 1 # Import the Mpg Data set:  
2 df_mpg = pd.read_csv('auto-mpg.csv')
```

Problem Statement

The data is generally used to predict the miles per gallon of a car

```
In [10]: 1 # Display head of the dataset:  
2 # Head() " Displays the First Five rows:  
3 df_mpg.head()
```

Out[10]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino

```
In [11]: 1 # Tail () " Displays the Last Five rows:  
2 df_mpg.tail()
```

Out[11]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
393	27.0	4	140.0	86	2790	15.6	82	1	ford mustang gl
394	44.0	4	97.0	52	2130	24.6	82	2	vw pickup
395	32.0	4	135.0	84	2295	11.6	82	1	dodge rampage
396	28.0	4	120.0	79	2625	18.6	82	1	ford ranger
397	31.0	4	119.0	82	2720	19.4	82	1	chevy s-10

```
In [12]: 1 # Entire Dataset:  
2  
3 df_mpg
```

```
Out[12]:   mpg cylinders displacement horsepower weight acceleration model year origin car name  
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu  
1 15.0 8 350.0 165 3693 11.5 70 1 buick skylark 320  
2 18.0 8 318.0 150 3436 11.0 70 1 plymouth satellite  
3 16.0 8 304.0 150 3433 12.0 70 1 amc rebel sst  
4 17.0 8 302.0 140 3449 10.5 70 1 ford torino  
... ... ... ... ... ... ... ... ...  
393 27.0 4 140.0 86 2790 15.6 82 1 ford mustang gl  
394 44.0 4 97.0 52 2130 24.6 82 2 vw pickup  
395 32.0 4 135.0 84 2295 11.6 82 1 dodge rampage  
396 28.0 4 120.0 79 2625 18.6 82 1 ford ranger  
397 31.0 4 119.0 82 2720 19.4 82 1 chevy s-10
```

398 rows × 9 columns

Interpretation : There are (398 Rows and 9 Columns) has been presented in this DataSet.

Check the Size and Variable Types:

```
In [27]: 1 # Check the Size of the Data:  
2 df_mpg.shape
```

```
Out[27]: (398, 9)
```

```
In [28]: 1 # Check the Variable Type :  
2 df_mpg.dtypes
```

```
Out[28]: mpg          float64  
cylinders      int64  
displacement    float64  
horsepower      object  
weight          int64  
acceleration    float64  
model year      int64  
origin          int64  
car name        object  
dtype: object
```

```
In [29]: 1 # Information about the Dataset:  
2 df_mpg.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 398 entries, 0 to 397  
Data columns (total 9 columns):  
 #   Column      Non-Null Count  Dtype     
---  --          -----          ----  
 0   mpg         398 non-null    float64  
 1   cylinders   398 non-null    int64  
 2   displacement 398 non-null    float64  
 3   horsepower   398 non-null    object  
 4   weight       398 non-null    int64  
 5   acceleration 398 non-null    float64  
 6   model year   398 non-null    int64  
 7   origin       398 non-null    int64  
 8   car name     398 non-null    object  
dtypes: float64(3), int64(4), object(2)  
memory usage: 28.1+ KB
```

```
In [13]: 1 # Check for Duplicate Values in the DataSet:  
2  
3 df_mpg.duplicated().sum().sum()
```

Out[13]: 0

```
In [41]: 1 # Check for Null Values in our DataSet:  
2 df_mpg.isnull().sum().sum()
```

Out[41]: 0

```
In [20]: 1 # Descriptive Statistics:  
2  
3 df_mpg.describe()
```

Out[20]:

	mpg	cylinders	displacement	weight	acceleration	model year	origin
count	398.000000	398.000000	398.000000	398.000000	398.000000	398.000000	398.000000
mean	23.514573	5.454774	193.425879	2970.424623	15.568090	76.010050	1.572864
std	7.815984	1.701004	104.269838	846.841774	2.757689	3.697627	0.802055
min	9.000000	3.000000	68.000000	1613.000000	8.000000	70.000000	1.000000
25%	17.500000	4.000000	104.250000	2223.750000	13.825000	73.000000	1.000000
50%	23.000000	4.000000	148.500000	2803.500000	15.500000	76.000000	1.000000
75%	29.000000	8.000000	262.000000	3608.000000	17.175000	79.000000	2.000000
max	46.600000	8.000000	455.000000	5140.000000	24.800000	82.000000	3.000000

```
In [22]: 1 # Descriptive Stats for Cat Col:  
2  
3 df_mpg.describe(include = 'object')
```

```
Out[22]:
```

	horsepower	car name
count	398	398
unique	94	305
top	150	ford pinto
freq	22	6

```
In [29]: 1 # Skewness for my Entire DataSet:  
2  
3 from warnings import filterwarnings  
4 filterwarnings('ignore')  
5  
6 df_mpg.skew()
```

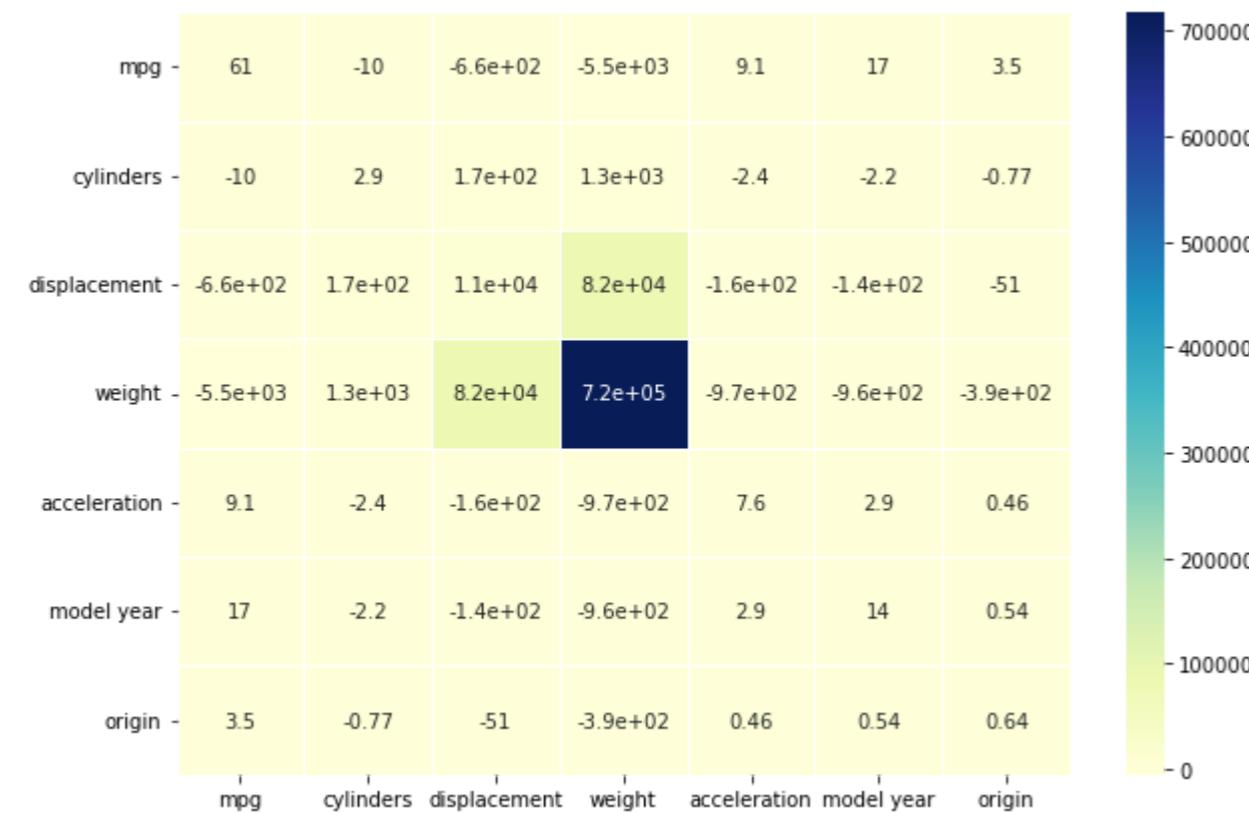
```
Out[29]: mpg          0.457066  
cylinders      0.526922  
displacement   0.719645  
weight         0.531063  
acceleration   0.278777  
model year     0.011535  
origin          0.923776  
dtype: float64
```

```
In [30]: 1 # Obtain kurt():  
2  
3 df_mpg.kurt()
```

```
Out[30]: mpg          -0.510781  
cylinders      -1.376662  
displacement   -0.746597  
weight         -0.785529  
acceleration   0.419497  
model year     -1.181232  
origin          -0.817597  
dtype: float64
```

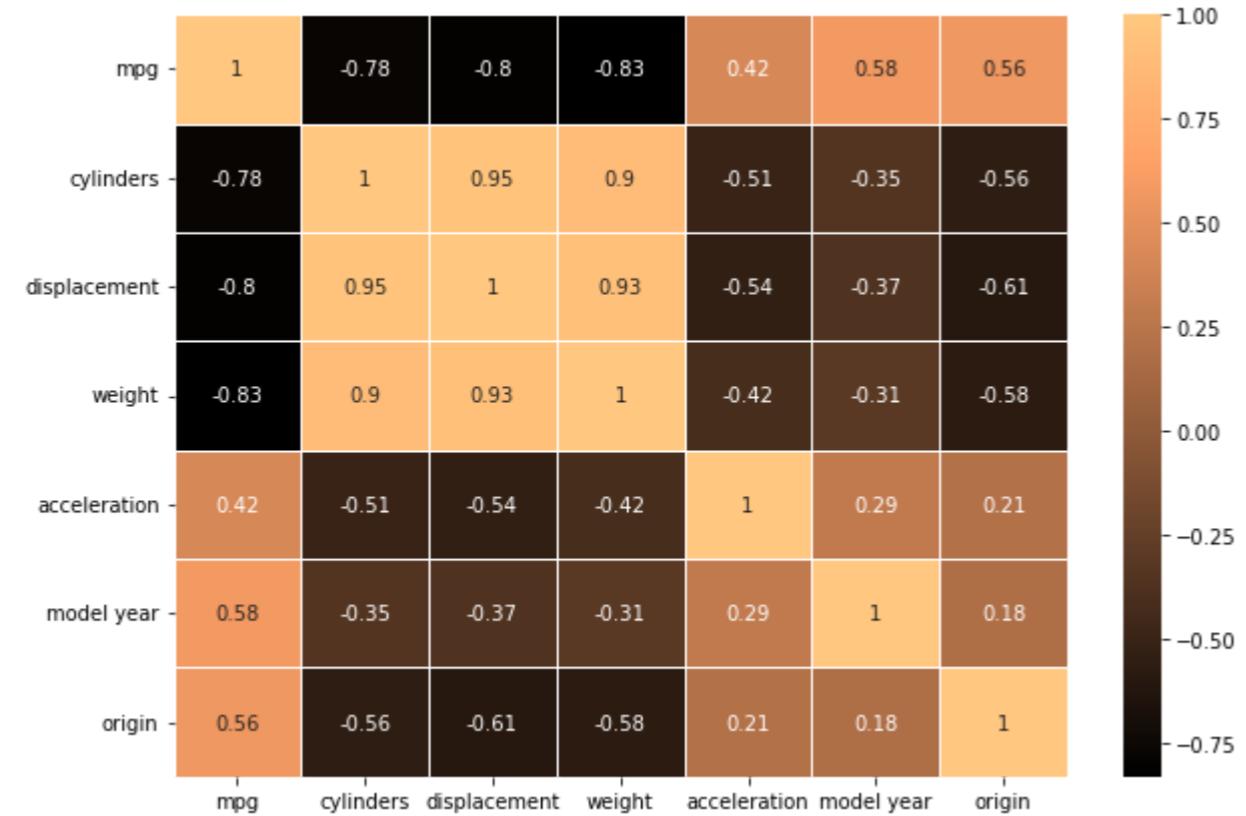
In [33]:

```
1 # Covariance:  
2  
3 cov = df_mpg.cov()  
4 cov  
5  
6 # Set the plot Size:  
7 fig, ax = plt.subplots(figsize = (10, 7))  
8  
9 # Plot a Heatmap for the Correlation Matrix  
10 # annot : print Values in each Cell  
11 # linewidths : specify width of the Line and Specifying the Plot  
12 # vmin : Minimum Value of the Variable  
13 # vmax : Maximum Value of the Variable  
14 # Cmap : colour Code for the Plot  
15 # fmt : set the Decimal place of the Annot  
16 sns.heatmap(cov, annot = True, linewidths = 0.95,  
17             cmap = 'YlGnBu', fmt = '.2g')  
18  
19 # Display the Plot:  
20 plt.show()
```



In [35]:

```
1 # Correlation:  
2  
3 corr = df_mpg.corr()  
4 corr  
5  
6 # Set the plot Size:  
7 fig, ax = plt.subplots(figsize = (10, 7))  
8  
9 # Plot a Heatmap for the Correlation Matrix  
10 # annot : print Values in each Cell  
11 # linewidths : specify width of the Line and Specifying the Plot  
12 # vmin : Minimum Value of the Variable  
13 # vmax : Maximum Value of the Variable  
14 # Cmap : colour Code for the Plot  
15 # fmt : set the Decimal place of the Annot  
16 sns.heatmap(corr, annot = True, linewidths = 0.95,  
17               cmap = 'copper', fmt = '.2g')  
18  
19 # Display the Plot:  
20 plt.show()
```



UNIVARIATE ANALYSIS

- Univariate Analysis is the Simplest Form of Statistical Analysis
- It is Study of a single (uni) variable
- The Key Purpose is to Describe the Data using different numerical and Visualization Techniques

1. Histogram
2. Density Plot
3. Box Plot
4. Violin plot
5. Cumulative Distribution

Univariate --Analysis:

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
In [3]: 1 df_car = pd.read_csv('auto-mpg.csv')
2 df_car.head()
```

Out[3]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino

```
In [6]: 1 # Univariate Analysis:
```

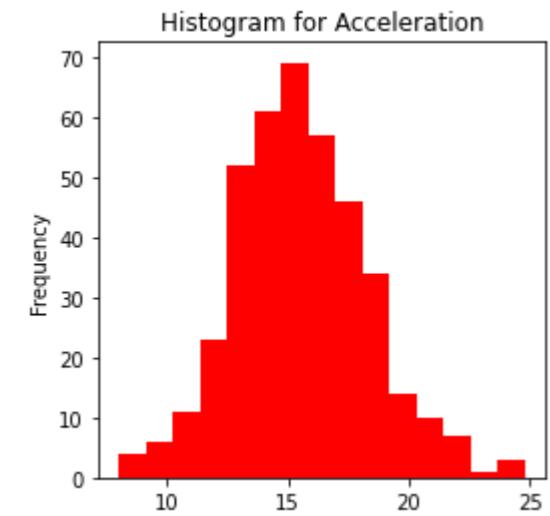
- Considering the Variable 'Acceleration'
- Check the Number of Observations in the Data using the len()

```
In [7]: 1 acceleration = df_car['acceleration']
```

```
In [8]: 1 len(acceleration)
```

```
Out[8]: 398
```

```
In [9]: 1 # Bins = 15, Creates 15 class Intervals
2 plt.hist(acceleration, bins = 15, color = 'r')
3
4 # Set the Title:
5 plt.title('Histogram for Acceleration')
6
7 # Set Label for y - axis :
8 plt.ylabel('Frequency')
9
10 # Display the Plot:
11 plt.show()
```

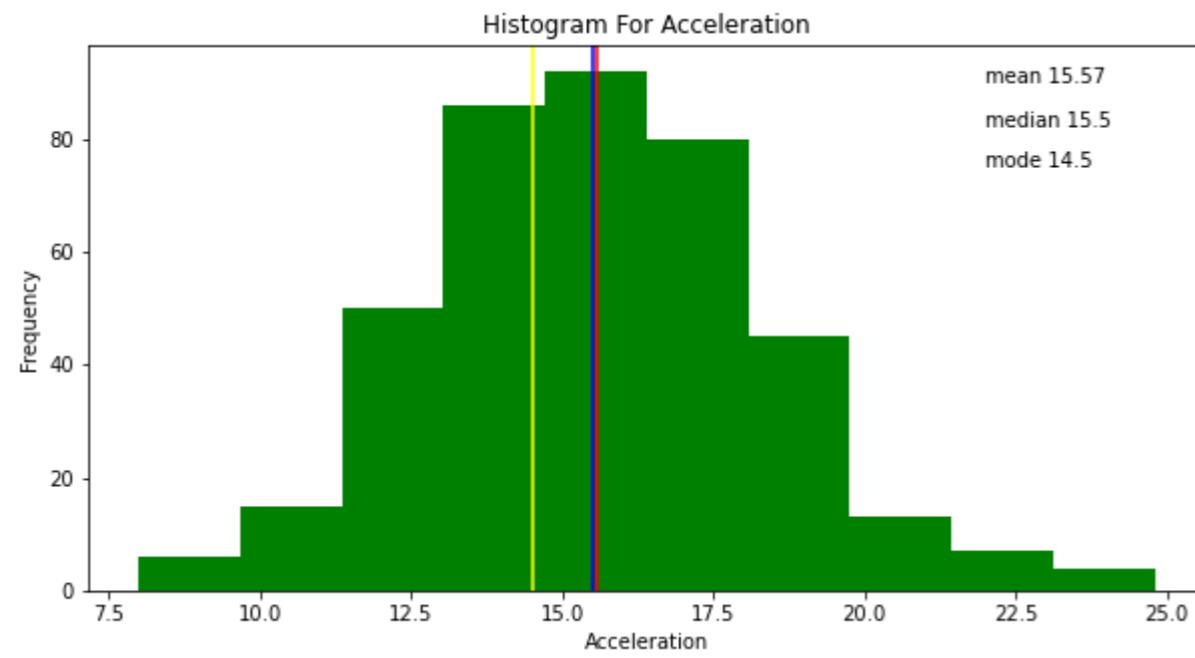


Interpretation : The Histogram divides the Data into Bins, counts the Data points in each bin

It shows the Bins on the x - axis and the Counts on the y - axis

The Variable "acceleration" has Symmetrical Distribution

```
In [10]: 1 # Set the Figure size:
2 plt.figure(figsize = (10, 5))
3
4 # Bins - 10, Creates the class Interval:
5 plt.hist(acceleration, bins = 10, color = 'g')
6
7 # Plot the Lines of the Mean, Median and Mode on my Histogram PLOT:
8 # Specify Different colors for each Line along by using the 'Color' Parameter
9 plt.axvline(acceleration.mean(), color = 'red')
10 plt.axvline(acceleration.median(), color = 'blue')
11 plt.axvline(acceleration.mode()[0], color = 'Yellow')
12
13
14 # Add the Values in the Plot:
15 plt.text(22, 90, 'mean'+ ' '+ str(round(acceleration.mean(),2)))
16 plt.text(22, 82, 'median'+ ' '+ str(round(acceleration.median(),2)))
17 plt.text(22, 75, 'mode'+ ' '+ str(round(acceleration.mode()[0],2)))
18
19 # Set title:
20 plt.title('Histogram For Acceleration')
21
22 # Set the Label of x- axis :
23 plt.xlabel('Acceleration')
24
25 # Set the Label for y - axis :
26 plt.ylabel('Frequency')
27
28 # Display the Plot:
29 plt.show()
```



```
In [2]: 1 # Import the Library:
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

What is Density Plot:

- A density Plot is a Smoothened Version of Histogram Estimated from the Data
- The 'kde'(Kernel Density Estimator) is drawn at every individual data point and all of these curves are the added together to create a single smooth density estimation.

In [24]:

```
1 # Set the Figure Size:  
2 plt.rcParams['figure.figsize'] =[10, 10]  
3  
4 # Distplot() - a plot for kernel Density Estimator:  
5 sns.distplot(acceleration)  
6  
7 # Set the Title:  
8 plt.title('Density plot for Acceleration')  
9  
10  
11 # Display the PLOT:  
12 plt.show()  
13
```

```
NameError                                 Traceback (most recent call last)  
~\AppData\Local\Temp\ipykernel_47640\296992581.py in <module>  
      3  
      4 # Distplot() - a plot for kernel Density Estimator:  
----> 5 sns.distplot(acceleration)  
      6  
      7 # Set the Title:
```

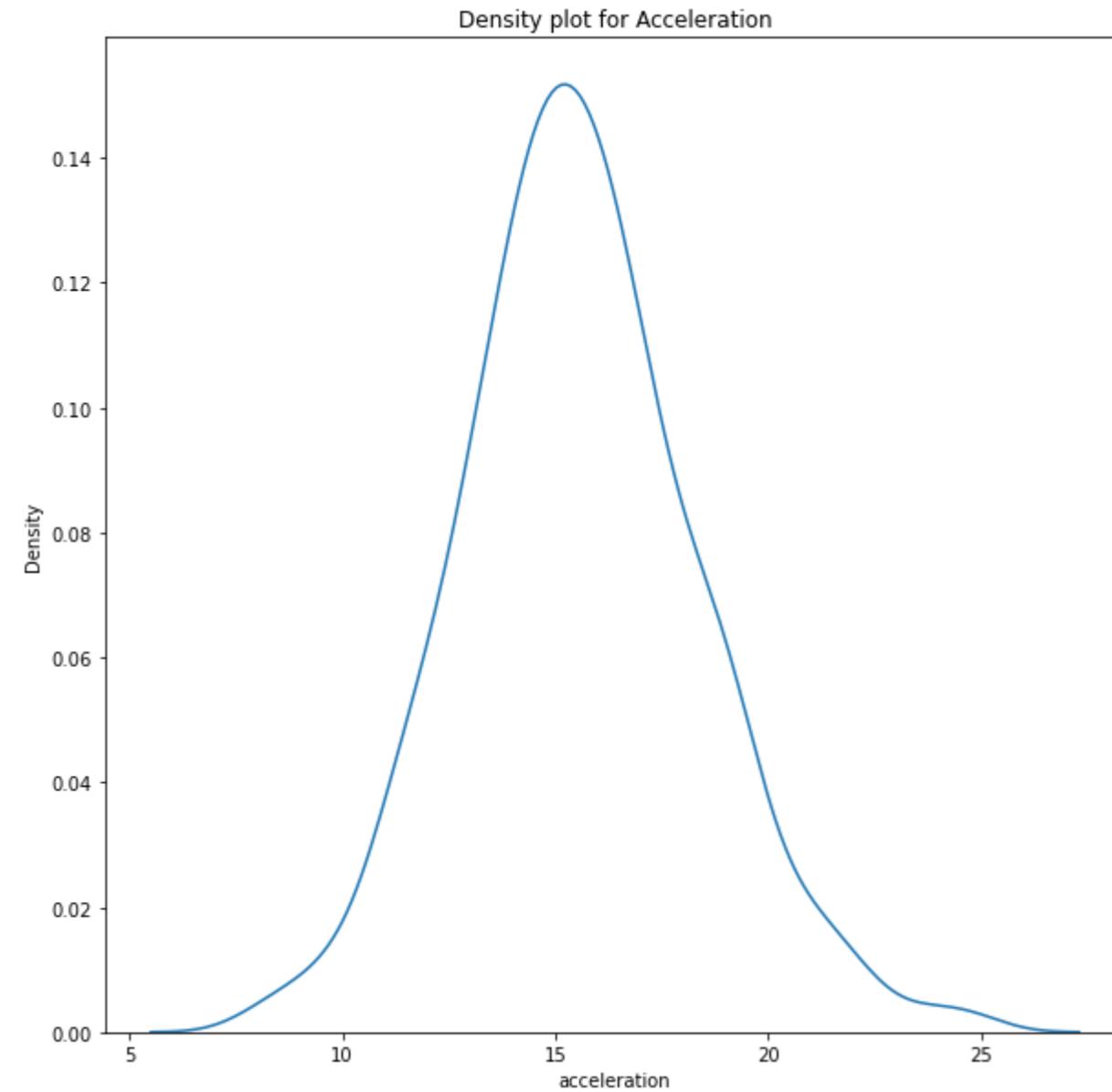
NameError: name 'acceleration' is not defined

Interpretation :The Dist of a acc is near to normal

In [14]:

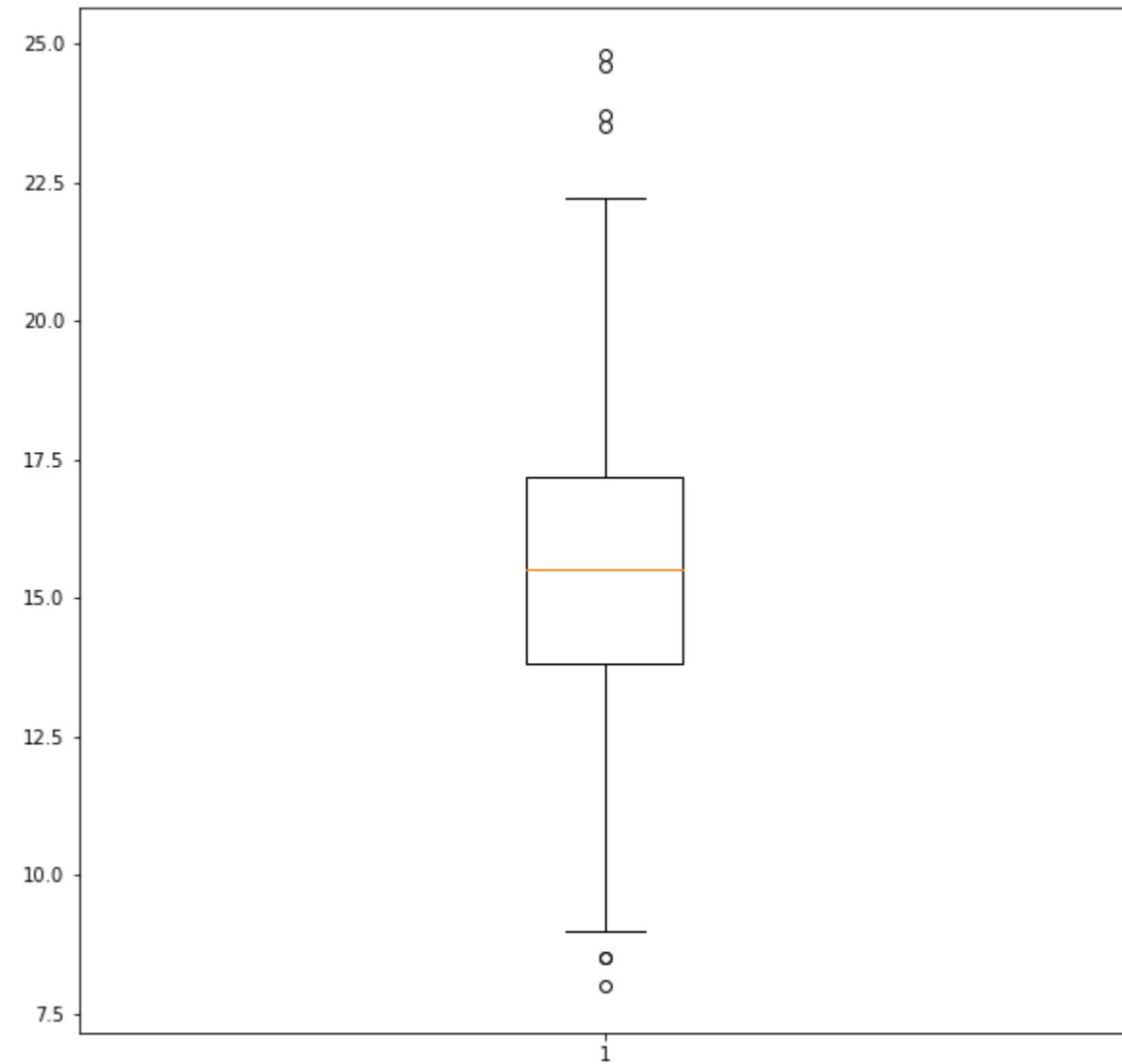
```
1 # Distplot without Histogram:  
2  
3 # Distplot() - a plot for kernel Density Estimator:  
4 sns.distplot(acceleration, hist = False)  
5  
6 # Set the Title:  
7 plt.title('Density plot for Acceleration')  
8  
9  
10 # Display the Plot:  
11 plt.show()  
12
```

C:\Users\bhuva\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)



In [16]:

```
1 # Box plot:  
2  
3 plt.boxplot(acceleration)  
4  
5 # Display the plot:  
6 plt.show()
```

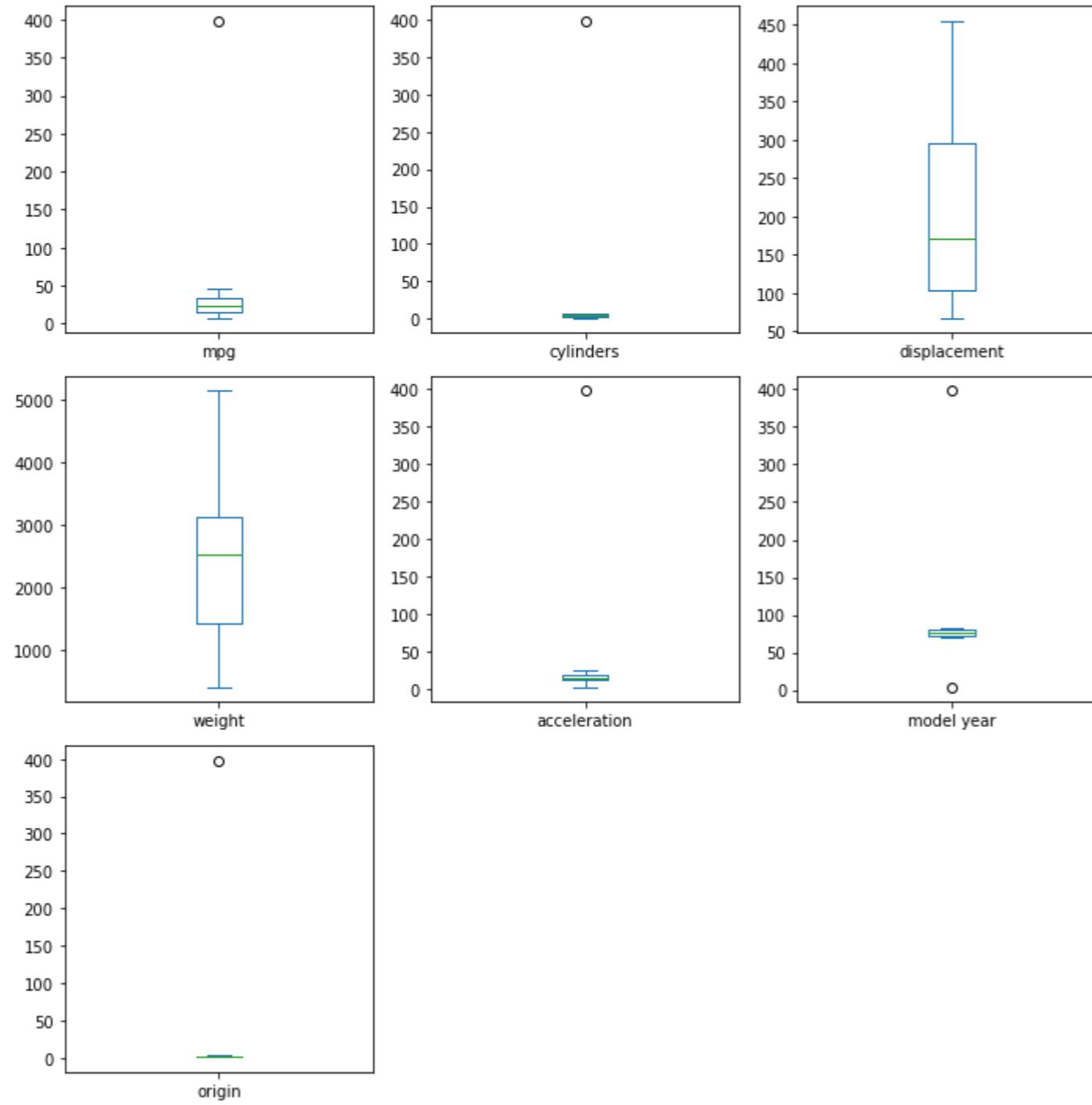


Interpretation : The Boxplot of column Acceleration has Outliers in it, thus we can able to remove by using IQR method.

In [20]:

```
1 # Only Numerical Colms from my Dataset:  
2  
3  
4 df_numerical = df_car.describe(include = [np.number])
```

```
In [21]: 1 df_numerical.plot(kind = 'box', subplots = True, layout = (3, 3))
2
3
4
5 # To give the Specified padding from the subplots:
6 plt.tight_layout()
7
8 # Display the Plot:
9 plt.show()
```



```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

Distribution of the Plot by Using Distplot:

Hint : We can also the plot the Distribution plot by using without Histogram

```
In [ ]: 1
```

```
In [ ]: 1
```

MULTIVARIATE ANALYSIS

- It is the analysis of two or more variables
- Used to study the relationships between the Variables
- Useful in determining the effect of one variable on other variables

1. Scatter plot
2. Correlation Matrix
3. Pair Plot
4. Grouped Box Plot

Handling Non-numeric Data :

Variable Encoding

- The dataset may contain numerical and / or Categorical Variables. Most of the algorithms are designed to work on numeric data.
- There are Several Methods (label Encoding , Dummy Encoding) to convert the Categorical data into numerical Data.

Types of Encoding

1. N-1 Dummy Encoding
2. One-Hot Encoding
3. Label encoding
4. Ordinal Encoding
5. Frequency Encoding
6. Target Encoding

```
In [11]: 1 df_car.head(5)
```

```
Out[11]:   mpg cylinders displacement horsepower weight acceleration model year origin          car name
0 18.0           8        307.0       130     3504         12.0        70      1  chevrolet chevelle malibu
1 15.0           8        350.0       165    3693          11.5        70      1  buick skylark 320
2 18.0           8        318.0       150    3436          11.0        70      1  plymouth satellite
3 16.0           8        304.0       150    3433          12.0        70      1  amc rebel sst
4 17.0           8        302.0       140    3449          10.5        70      1  ford torino
```

```
In [12]: 1 df_car['origin'].value_counts()
```

```
Out[12]: 1    249
3     79
2     70
Name: origin, dtype: int64
```

```
In [13]: 1 df_car['car name'].value_counts()
```

```
Out[13]:  ford pinto      6
toyota corolla      5
amc matador        5
ford maverick      5
chevrolet chevette  4
..
chevrolet monza 2+2  1
ford mustang ii     1
pontiac astro        1
amc pacer            1
chevy s-10            1
Name: car name, Length: 305, dtype: int64
```

In [14]:

```

1 # Create Dummy Variables for 'Car name' :
2 # Drop_first = 'True' creates (n - 1) Dummy Variables from Categories
3 pd.get_dummies(df_car, columns = ['car name'], drop_first = True)

```

Out[14]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name_amc ambassador dpl	car name_amc ambassador sst	...	car name_volvo 145e (sw)	car name_volvo 244dl	car name_volvo 245	car name_volvo 264gl	car name_volvo diesel	car name_vw dasher (diesel)	car name_vw pickup	car name_vw rabbit	na
0	18.0	8	307.0	130	3504	12.0	70	1	0	0	...	0	0	0	0	0	0	0	0	0
1	15.0	8	350.0	165	3693	11.5	70	1	0	0	...	0	0	0	0	0	0	0	0	0
2	18.0	8	318.0	150	3436	11.0	70	1	0	0	...	0	0	0	0	0	0	0	0	0
3	16.0	8	304.0	150	3433	12.0	70	1	0	0	...	0	0	0	0	0	0	0	0	0
4	17.0	8	302.0	140	3449	10.5	70	1	0	0	...	0	0	0	0	0	0	0	0	0
...
393	27.0	4	140.0	86	2790	15.6	82	1	0	0	...	0	0	0	0	0	0	0	0	0
394	44.0	4	97.0	52	2130	24.6	82	2	0	0	...	0	0	0	0	0	0	0	1	0
395	32.0	4	135.0	84	2295	11.6	82	1	0	0	...	0	0	0	0	0	0	0	0	0
396	28.0	4	120.0	79	2625	18.6	82	1	0	0	...	0	0	0	0	0	0	0	0	0
397	31.0	4	119.0	82	2720	19.4	82	1	0	0	...	0	0	0	0	0	0	0	0	0

398 rows × 312 columns



One Hot Encoding

In [15]:

```

1 # One Hot Encoding:
2
3 pd.get_dummies(df_car, columns = ['origin'])

```

Out[15]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	car name	origin_1	origin_2	origin_3
0	18.0	8	307.0	130	3504	12.0	70	chevrolet chevelle malibu	1	0	0
1	15.0	8	350.0	165	3693	11.5	70	buick skylark 320	1	0	0
2	18.0	8	318.0	150	3436	11.0	70	plymouth satellite	1	0	0
3	16.0	8	304.0	150	3433	12.0	70	amc rebel sst	1	0	0
4	17.0	8	302.0	140	3449	10.5	70	ford torino	1	0	0
...
393	27.0	4	140.0	86	2790	15.6	82	ford mustang gl	1	0	0
394	44.0	4	97.0	52	2130	24.6	82	vw pickup	0	1	0
395	32.0	4	135.0	84	2295	11.6	82	dodge rampage	1	0	0
396	28.0	4	120.0	79	2625	18.6	82	ford ranger	1	0	0
397	31.0	4	119.0	82	2720	19.4	82	chevy s-10	1	0	0

398 rows × 11 columns

In [20]:

```
1 # Sklearn Library:
2
3 # Import The OneHotEncoder:
4
5 from sklearn.preprocessing import OneHotEncoder
6
7 # Creating an Instance of One-Hot-Encoder:
8 encode = OneHotEncoder()
9
10 # fit_transform : Fit to data and Returns a Transformed Version:
11 # toarray() : Returns the Numpy array
12 # Columns : Add the Column names
13
14 df_encode = pd.DataFrame(encode.fit_transform(df_car[['origin']]).toarray(), columns = ['origin_europe', 'origin_japan', 'origin_usa'])
15
16
17 # Merge with main Data DataFrame ( df_car):
18 # Axis = 1: it stands for columns
19 df_encode = pd.concat([df_car, df_encode], axis = 1)
20
21
22 # Print the First 2 Rows of the DataSet:
23 df_encode
```

Out[20]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name	origin_europe	origin_japan	origin_usa
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu	1.0	0.0	0.0
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320	1.0	0.0	0.0
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite	1.0	0.0	0.0
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst	1.0	0.0	0.0
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino	1.0	0.0	0.0
...
393	27.0	4	140.0	86	2790	15.6	82	1	ford mustang gl	1.0	0.0	0.0
394	44.0	4	97.0	52	2130	24.6	82	2	vw pickup	0.0	1.0	0.0
395	32.0	4	135.0	84	2295	11.6	82	1	dodge rampage	1.0	0.0	0.0
396	28.0	4	120.0	79	2625	18.6	82	1	ford ranger	1.0	0.0	0.0
397	31.0	4	119.0	82	2720	19.4	82	1	chevy s-10	1.0	0.0	0.0

398 rows × 12 columns

Label Encoding

The Label Encoder Consider as level in a Categorical Variable by Alphabetical Order

In [27]:

```
1 df_car.head(2)
```

Out[27]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320

In [6]:

```
1 # Example:  
2  
3 # Use Sklearn Library:  
4 from sklearn.preprocessing import LabelEncoder  
5  
6 # Create an Instance:  
7 labelencoder = LabelEncoder()  
8  
9 # fit the encoder:  
10 df_car['encoded_performance_of_my_car'] = labelencoder.fit_transform(df_car.displacement)  
11  
12 # display the Data:  
13 df_car
```

Out[6]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name	encoded_performance_of_my_car
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu	69
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320	72
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite	70
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst	67
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino	66
...
393	27.0	4	140.0	86	2790	15.6	82	1	ford mustang gl	41
394	44.0	4	97.0	52	2130	24.6	82	2	vw pickup	17
395	32.0	4	135.0	84	2295	11.6	82	1	dodge rampage	40
396	28.0	4	120.0	79	2625	18.6	82	1	ford ranger	34
397	31.0	4	119.0	82	2720	19.4	82	1	chevy s-10	33

398 rows × 10 columns

In []:

```
1 # Interpretation  
2  
3 # Low performance ranges from 17 - 27  
4 # Average Performance ranges from 37 - 43  
5 # Super Performance ranges from 61 - 72
```

In [4]:

```
1 df_car['encoded_performance_of_my_car'].value_counts()
```

Out[4]:

17	21
19	18
72	18
70	17
61	17
..	
22	1
26	1
37	1
27	1
43	1

Name: encoded_performance_of_my_car, Length: 82, dtype: int64

Ordinal Encoding

```
In [12]: 1 # Import the Ordinal Encoder:  
2 from sklearn.preprocessing import OrdinalEncoder  
3  
4 # Create an Instance:  
5 orderencoding = OrdinalEncoder(categories = [['Average', 'Good', 'Excellent']])  
6  
7 # Fit the Encoder:  
8 df_car['Ordered_performance'] = orderencoding.fit_transform(df_car.mpg.values.reshape(-1, 1))  
9  
10 # Display the Data:  
11 df_car.head(3)
```

```
-----  
ValueError                                Traceback (most recent call last)  
~\AppData\Local\Temp\ipykernel_50112\1775378340.py in <module>  
      6  
      7 # Fit the Encoder:  
----> 8 df_car['Ordered_performance'] = orderencoding.fit_transform(df_car.mpg.values.reshape(-1, 1))  
      9  
     10 # Display the Data:  
  
~\anaconda3\lib\site-packages\sklearn\base.py in fit_transform(self, X, y, **fit_params)  
  865     if y is None:  
  866         # fit method of arity 1 (unsupervised transformation)  
--> 867         return self.fit(X, **fit_params).transform(X)  
  868     else:  
  869         # fit method of arity 2 (supervised transformation)  
  
~\anaconda3\lib\site-packages\sklearn\preprocessing\_encoders.py in fit(self, X, y)  
1292  
1293     # `__fit__` will only raise an error when `self.handle_unknown="error"`  
-> 1294     self._fit(X, handle_unknown=self.handle_unknown, force_all_finite="allow-nan")  
1295  
1296     if self.handle_unknown == "use_encoded_value":  
  
~\anaconda3\lib\site-packages\sklearn\preprocessing\_encoders.py in _fit(self, X, handle_unknown, force_all_finite, return_counts)  
104             cats = result  
105         else:  
--> 106             cats = np.array(self.categories[i], dtype=Xi.dtype)  
107             if Xi.dtype.kind not in "OUS":  
108                 sorted_cats = np.sort(cats)  
  
ValueError: could not convert string to float: 'Average'
```

Frequency Encoding

skill Encoded Feature

Python 0.44 python = 0.44 Python 0.44 Python 0.44 Python 0.44 R 0.33 R = 0.33 R 0.33 R 0.33 SQL 0.22 SQLI 0.22 SQL = 0.22 SQLI 0.22

import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns

```
In [14]: 1 df_car.head()
```

```
Out[14]:   mpg cylinders displacement horsepower weight acceleration model year origin car name
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu
1 15.0 8 350.0 165 3693 11.5 70 1 buick skylark 320
2 18.0 8 318.0 150 3436 11.0 70 1 plymouth satellite
3 16.0 8 304.0 150 3433 12.0 70 1 amc rebel sst
4 17.0 8 302.0 140 3449 10.5 70 1 ford torino
```

```
In [3]: 1 df_car.head(1)
```

```
Out[3]:   mpg cylinders displacement horsepower weight acceleration model year origin car name
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu
```

```
In [5]: 1 # Frequency Encoding:
2 # Size of the each Category:
3 encoding = df_car.groupby('displacement').size()
4
5 # Get the Frequency of Each Category:
6 encoding = encoding / len(df_car)
7
8 # encode the Column:
9 df_car['frequency_displacement'] = df_car.displacement.map(encoding)
10
11 # Display the Data:
12 df_car.head()
```

```
Out[5]:   mpg cylinders displacement horsepower weight acceleration model year origin car name frequency_displacement
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu 0.007538
1 15.0 8 350.0 165 3693 11.5 70 1 buick skylark 320 0.045226
2 18.0 8 318.0 150 3436 11.0 70 1 plymouth satellite 0.042714
3 16.0 8 304.0 150 3433 12.0 70 1 amc rebel sst 0.017588
4 17.0 8 302.0 140 3449 10.5 70 1 ford torino 0.027638
```

Target Encoding

```
In [3]: 1 df_smoker = pd.DataFrame({  
2     'Smoker' : ['yes', 'yes', 'no', 'no', 'yes', 'no', 'yes', 'yes', 'no'],  
3     'Target' : [1, 0, 1, 0, 0, 0, 1, 0, 0]  
4 })  
5 df_smoker
```

Out[3]:

	Smoker	Target
0	yes	1
1	yes	0
2	no	1
3	no	0
4	yes	0
5	no	0
6	yes	1
7	yes	0
8	no	0

```
In [4]: 1 mean = df_smoker.groupby('Smoker')['Target'].mean()
```

```
In [9]: 1 df_smoker['Smoker_encoded'] = df_smoker.Smoker.map(mean)  
2  
3  
4 df_smoker
```

Out[9]:

	Smoker	Target	Smoker_encoded
0	yes	1	0.40
1	yes	0	0.40
2	no	1	0.25
3	no	0	0.25
4	yes	0	0.40
5	no	0	0.25
6	yes	1	0.40
7	yes	0	0.40
8	no	0	0.25

Feature Scaling:

- Feature Scaling is also Called as Data Normalization.
- It is the technique used to Transform the Data into a common scale. Since the Features have various ranges, it becomes a necessary step in the data preprocessing while using machine learning algorithms.

Methods to Perform Feature Scaling:

1. Z-score Normalization (or) Standardization

- Standardization transforms the Data such that the data has mean 0 and Unit Variance.
- The Procedure involves subtracting the mean from Observation and dividing by the Standard Deviation
- $(x_{\text{new}}) = \frac{x - \mu}{\sigma}$

```
In [15]: 1 df_car.head(1)
```

```
Out[15]:   mpg cylinders displacement horsepower weight acceleration model year origin      car name
0    18.0          8        307.0       130     3504         12.0         70      1  chevrolet chevelle malibu
```

```
In [5]: 1 # Import Standard Scaler:
2 from sklearn.preprocessing import StandardScaler
3
4 # Minimum and maximum Values of 'Weight':
5 # '\n' - add Space
6
7 print('Minimum value Before Transformation : ', df_car['weight'].min(), '\n'
8      'Maximum value Before Transformation : ', df_car['weight'].max(), '\n')
9
10 # Instantiate the StandardScaler:
11 standard_scale = StandardScaler()
12
13 # Fit the StandardScaler:
14 # fit_transform() : returns a Transformed Data
15 df_car['Scaled_Weight'] = standard_scale.fit_transform(df_car[['weight']])
16
17 print('Minimum value after Transformation : ', df_car['Scaled_Weight'].min(), '\n'
18      'Maximum value after Transformation : ', df_car['Scaled_Weight'].max(), '\n')
```

```
Minimum value Before Transformation :  1613
Maximum value Before Transformation :  5140

Minimum value after Transformation :  -1.6049434405635041
Maximum value after Transformation :  2.565185359572092
```

Scaled value mean = 0

----->> IMPORTANT

Scaled value Standard Deviation = 1

```
In [6]: 1 print('Mean : ', df_car['Scaled_Weight'].mean())
2 print('\n')
3 print('Standard Deviation : ', df_car['Scaled_Weight'].std())
```

Mean : -9.902743058842477e-17

Standard Deviation : 1.0012586537392114

```
In [4]: 1 df_car
```

Out[4]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name	Scaled_Weight
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu	0.630870
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320	0.854333
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite	0.550470
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst	0.546923
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino	0.565841
...
393	27.0	4	140.0	86	2790	15.6	82	1	ford mustang gl	-0.213324
394	44.0	4	97.0	52	2130	24.6	82	2	vw pickup	-0.993671
395	32.0	4	135.0	84	2295	11.6	82	1	dodge rampage	-0.798585
396	28.0	4	120.0	79	2625	18.6	82	1	ford ranger	-0.408411
397	31.0	4	119.0	82	2720	19.4	82	1	chevy s-10	-0.296088

398 rows × 10 columns

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

2. Min-Max Normalization

- Performs linear transformation on the original data
- The min-max normalization is given as
- $[(x_{norm}) = X - X_{min} / X_{max} - X_{min}]$

```
In [11]: 1 # IMPORT LIBRARIES:
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```
In [12]: 1 df_car.head(2)
```

```
Out[12]:   mpg cylinders displacement horsepower weight acceleration model year origin car name min_max_scaled_weight
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu 0.536150
1 15.0 8 350.0 165 3693 11.5 70 1 buick skylark 320 0.589736
```

```
In [10]: 1 # Import MinMaxScaler:
2 from sklearn.preprocessing import MinMaxScaler
```

```
3
4 # Creating an Instance:
5 min_max = MinMaxScaler()
6
7 # Fit the MinMaxScaler:
8 df_car['min_max_scaled_weight'] = min_max.fit_transform(df_car[['weight']])
9
10 # Minimum and Maximum of Normalized Weight:
11 df_car['min_max_scaled_weight'].min(), df_car['min_max_scaled_weight'].max()
```

```
Out[10]: (0.0, 1.0)
```

Interpretation :

The New Data Has been Generated for the Particular Column weight b/w 0 - 1

Data Transformation:

1. Log Transformation

- Reduce the **skewness in the Distribution of the Original Data**
- Makes the Data more interpretable
- The Arithmetic mean of the log-transformed data is the geometric mean of the original Data
- It Converts the Exponential Growth of the linear Growth

```
In [13]: 1 df_car.head(1)
```

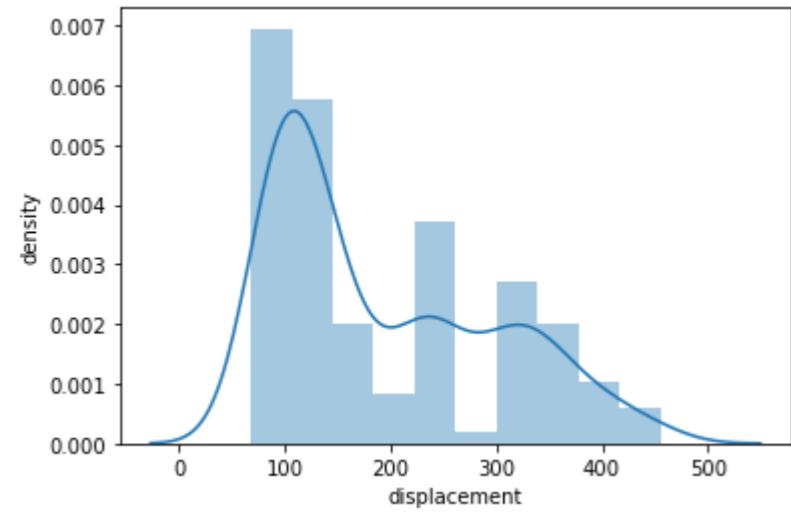
```
Out[13]:   mpg cylinders displacement horsepower weight acceleration model year origin car name min_max_scaled_weight
0 18.0 8 307.0 130 3504 12.0 70 1 chevrolet chevelle malibu 0.536150
```

In [14]:

```
1 # Set the Figure Size:  
2 plt.rcParams['figure.figsize'] = [6, 4]  
3  
4 # Distribution of the Displacement:  
5 sns.distplot(df_car['displacement'])  
6  
7 plt.ylabel('density')  
8  
9 # Coefficient of Skewness:  
10 print('Skewness : ', df_car['displacement'].skew())  
11  
12 # Display the Plot:  
13 plt.show()
```

C:\Users\bhuva\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

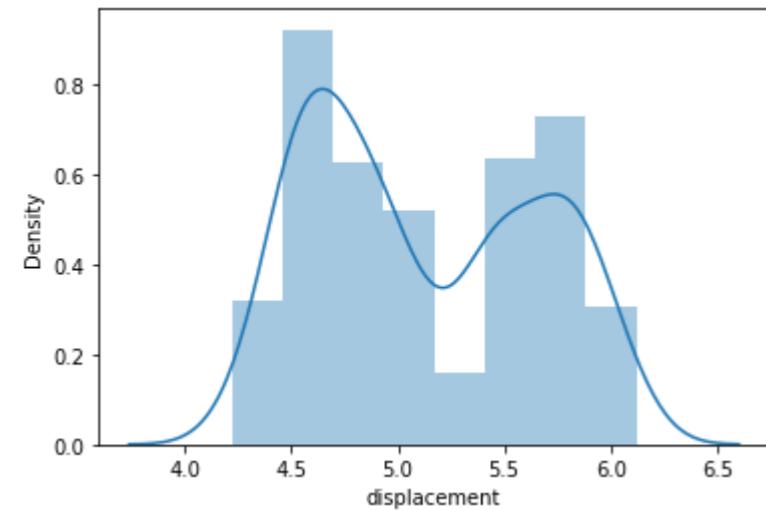
Skewness : 0.7196451643005952



```
In [15]: 1 # Apply Natural Log Transformation with (base 'e')
2 log_displacement = np.log(df_car['displacement'])
3
4
5 # Coefficient of Skewness:
6 print('Skewness : ', log_displacement.skew())
7
8 # Distribution of Log_transformed Variable:
9 sns.distplot(log_displacement)
10
11 # Set Label for y - Axis :
12 plt.ylabel('Density')
13
14 # Display the Plot :
15 plt.show()
```

Skewness : 0.22600298495225188

C:\Users\bhuva\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



```
In [26]: 1 df_car.shape
2
```

Out[26]: (398, 10)

In []: 1

In []: 1

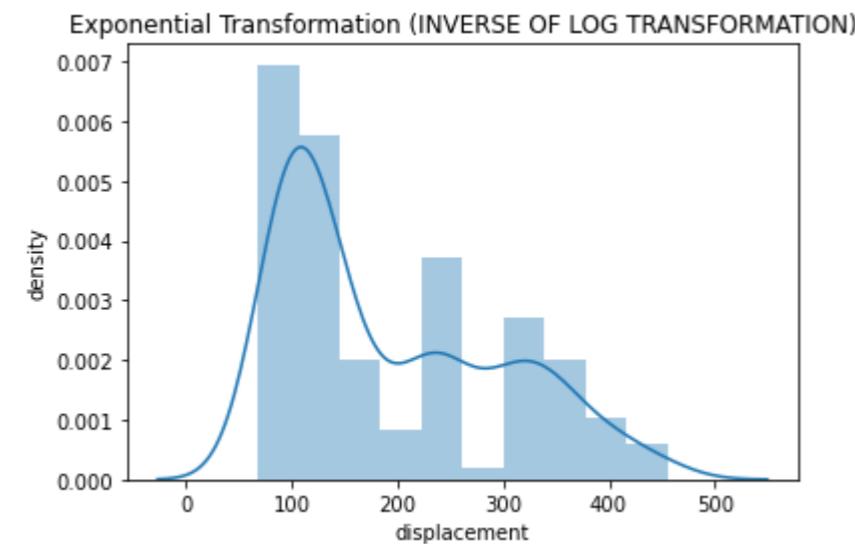
2. Exponential Transformation

- It is the **Inverse transformation** of the Log Transformation
- It is used to Convert the log Transformed values to their original Units

In [16]:

```
1 # Anti - Log of 'Displacement' :
2 displacement = np.exp(log_displacement)
3
4 # Plot The Distribution :
5 sns.distplot(displacement)
6
7 # Set the Label y - axis:
8 plt.ylabel('density')
9
10 # Set the Title:
11 plt.title('Exponential Transformation (INVERSE OF LOG TRANSFORMATION)')
12
13 # Display the Plot:
14 plt.show()
```

C:\Users\bhuva\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



3. Box-Cox Transformation

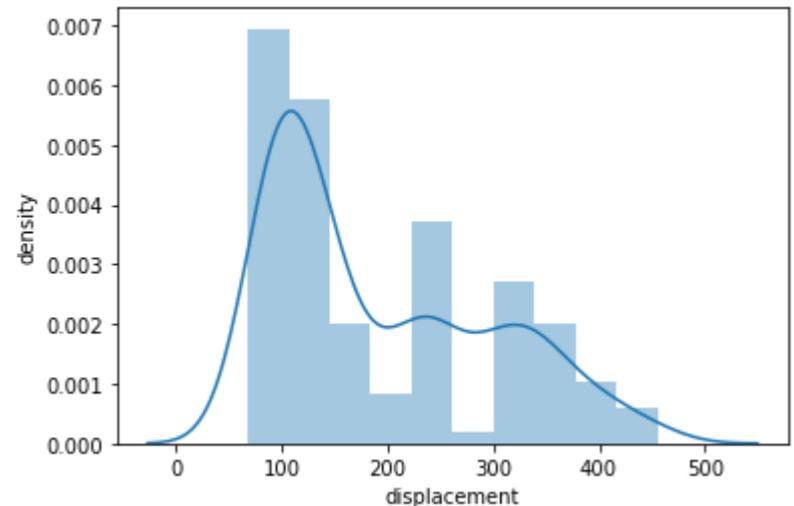
- Generalized version of Log transformation
- Reduce the SKewness of the data making it more Symmetrical
- The Box-Cox transformation for Variable x with Positive Values is defined as:

$$x(\lambda) = x^{\lambda} - 1 / \lambda \text{ Whereas, } \lambda \neq 0$$

In [37]:

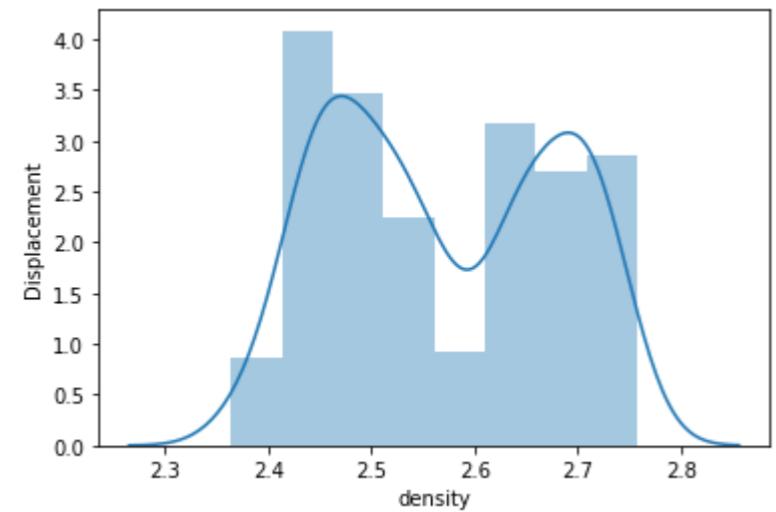
```
1 # Set the Figure Size:  
2 plt.rcParams['figure.figsize'] = [6, 4]  
3  
4 # Distribution of the Displacement:  
5 sns.distplot(df_car['displacement'])  
6  
7 # Set the Label for y - axis:  
8 plt.ylabel('density')  
9  
10 # Coefficient of the Skewness:  
11 print('Skewness : ', df_car['displacement'])  
12  
13 # Display the Plot:  
14 plt.show()
```

```
Skewness : 0      307.0  
1      350.0  
2      318.0  
3      304.0  
4      302.0  
...  
393    140.0  
394    97.0  
395    135.0  
396    120.0  
397    119.0  
Name: displacement, Length: 398, dtype: float64
```



In [41]:

```
1 # Box Cox Transformation:  
2  
3 # Import scipy Library:  
4  
5 import scipy  
6  
7 from scipy import stats  
8  
9 # Scipy.stats.boxcox() returns tuple with transformation:  
10 box_displacement = scipy.stats.boxcox(df_car['displacement'])  
11  
12 # Distribution of Transformed Variable:  
13 # [0] refers to array of the indexing particular Column in Displacement:  
14 sns.distplot(box_displacement[0])  
15  
16 # Set axes label:  
17 plt.ylabel('Displacement')  
18 plt.xlabel('density')  
19  
20 # Display the plot:  
21 plt.show()  
22  
23
```



```
In [49]: 1 df_car_price = pd.read_csv('CarPrice_Assignment.csv')
2
3 # Shape of the Data:
4 print(df_car_price.shape)
5
6
7 # Print First Five Rows:
8 df_car_price.head()
```

(205, 26)

Out[49]:

	car_ID	symboling	CarName	fuelytype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	...	enginesize	fuelsystem	boreratio	stroke	compressionratio	horsepower	peakrpm	citympg
0	1	3	alfa-romero giulia	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21
1	2	3	alfa-romero stelvio	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21
2	3	1	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154	5000	19
3	4	2	audi 100 ls	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	3.40	10.0	102	5500	24
4	5	2	audi 100ls	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	8.0	115	5500	18

5 rows × 26 columns

Read the Dataset Data Information Data Type Conversion

Missing Value Treatment

- Replace by Frequency
- Replac by Mean
- Replace by Median
- Repalce by bfill() and ffill()

Outlier Treatment

- z - Score
- Interquartile Range Methodology

Feature Engineering

Train Test Split

Supervised Learning:

Linear Regression Algorithms

Unsupervised Learning:

Clustering

Hyper parameter tuning

Decision Tree

Random Forest

MISSING VALUES :

- Standard Missing Values
- Non Standard Missing Values

In [62]:

```

1 # Import Libraries:
2
3
4 # Importing pandas library and make it as pd
5 import pandas as pd
6
7
8
9 import numpy as np
10 import matplotlib.pyplot as plt
11 import seaborn as sns

```

In [63]:

```

1 df = pd.read_csv('k-circlesales-cleaned.csv')
2 df.head()

```

Out[63]:

	Unnamed: 0	Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales	Profit
0	0	FDA15	9.30	Low Fat	0.016047	Dairy	249.8	OUT049	1999	Medium	Tier 2	Supermarket Type1	3735.1380	11.5
1	1	DRC01	5.92	Regular	0.019278	Soft Drinks	48.3	OUT018	2009	Medium	Tier 2	Supermarket Type2	443.4228	14.3
2	2	FDN15	17.50	Low Fat	0.016760	Meat	141.6	OUT049	1999	Medium	Tier 2	Supermarket Type1	2097.2700	14.5
3	3	FDX07	19.20	Regular	0.000000	Fruits and Vegetables	182.1	OUT010	1998	NaN	NaN	Grocery Store	732.3800	13.6
4	4	NCD19	8.93	Low Fat	0.000000	Household	53.9	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052	14.1

1. Data Summary:

```
In [64]: 1 # Find the Shape of the Dataset:  
2 df.shape
```

Out[64]: (8523, 14)

```
In [65]: 1 # Summary of all the Variable in the Dataset:  
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 8523 entries, 0 to 8522  
Data columns (total 14 columns):  
 #   Column           Non-Null Count  Dtype     
---  --  
 0   Unnamed: 0        8523 non-null   int64    
 1   Item_Identifier   8523 non-null   object    
 2   Item_Weight        7774 non-null   float64   
 3   Item_Fat_Content   8523 non-null   object    
 4   Item_Visibility    8523 non-null   float64   
 5   Item_Type          8523 non-null   object    
 6   Item_MRP           8523 non-null   float64   
 7   Outlet_Identifier  8523 non-null   object    
 8   Outlet_Establishment_Year  8523 non-null   int64    
 9   Outlet_Size         6113 non-null   object    
 10  Outlet_Location_Type 6473 non-null   object    
 11  Outlet_Type        8523 non-null   object    
 12  Item_Outlet_Sales  8523 non-null   float64   
 13  Profit             8523 non-null   float64  
dtypes: float64(5), int64(2), object(7)  
memory usage: 932.3+ KB
```

The Statistical Summary of the Dataset for Numerical Variables

```
In [66]: 1 # Summary of Numerical Variables:  
2  
3 df.describe()  
4  
5 # By default describe(), gives you five point summaryb only for Numerical Columns.  
6  
7  
8  
9
```

Out[66]:

	Unnamed: 0	Item_Weight	Item_Visibility	Item_MRP	Outlet_Establishment_Year	Item_Outlet_Sales	Profit
count	8523.000000	7774.000000	8523.000000	8523.000000	8523.000000	8523.000000	8523.000000
mean	4261.000000	11.676740	0.066132	140.998838	1997.831867	2181.288914	13.414514
std	2460.522505	5.776851	0.051598	62.258099	8.371760	1706.499616	1.701840
min	0.000000	0.000000	0.000000	31.300000	1985.000000	33.290000	0.100000
25%	2130.500000	7.720000	0.026989	93.800000	1987.000000	834.247400	13.150000
50%	4261.000000	11.800000	0.053931	142.700000	1999.000000	1794.331000	13.900000
75%	6391.500000	16.500000	0.094585	185.650000	2004.000000	3101.296400	14.300000
max	8522.000000	21.350000	0.328391	266.900000	2009.000000	13086.964800	24.000000

The Statistical Summary of the Dataset for Categorical Variables

```
In [67]: 1 # Get Summary for Cat columns:  
2 df.describe(include = 'object')
```

Out[67]:

	Item_Identifier	Item_Fat_Content	Item_Type	Outlet_Identifier	Outlet_Size	Outlet_Location_Type	Outlet_Type
count	8523	8523	8523	8523	6113	6473	8523
unique	1559	2	16	10	3	8	4
top	FDW13	Low Fat	Fruits and Vegetables	OUT027	Medium	Tier 2	Supermarket Type1
freq	10	5517	1232	935	2793	2793	5577

```
In [68]: 1 # Check for Null Values:  
2  
3 df.isnull().sum()
```

Out[68]:

Unnamed:	0
Item_Identifier	0
Item_Weight	749
Item_Fat_Content	0
Item_Visibility	0
Item_Type	0
Item_MRP	0
Outlet_Identifier	0
Outlet_Establishment_Year	0
Outlet_Size	2410
Outlet_Location_Type	2050
Outlet_Type	0
Item_Outlet_Sales	0
Profit	0
dtype:	int64

```
In [69]: 1 # Total Null Values in the Dataset  
2 # Sum of all the Null values:  
3  
4 df.isnull().sum().sum()
```

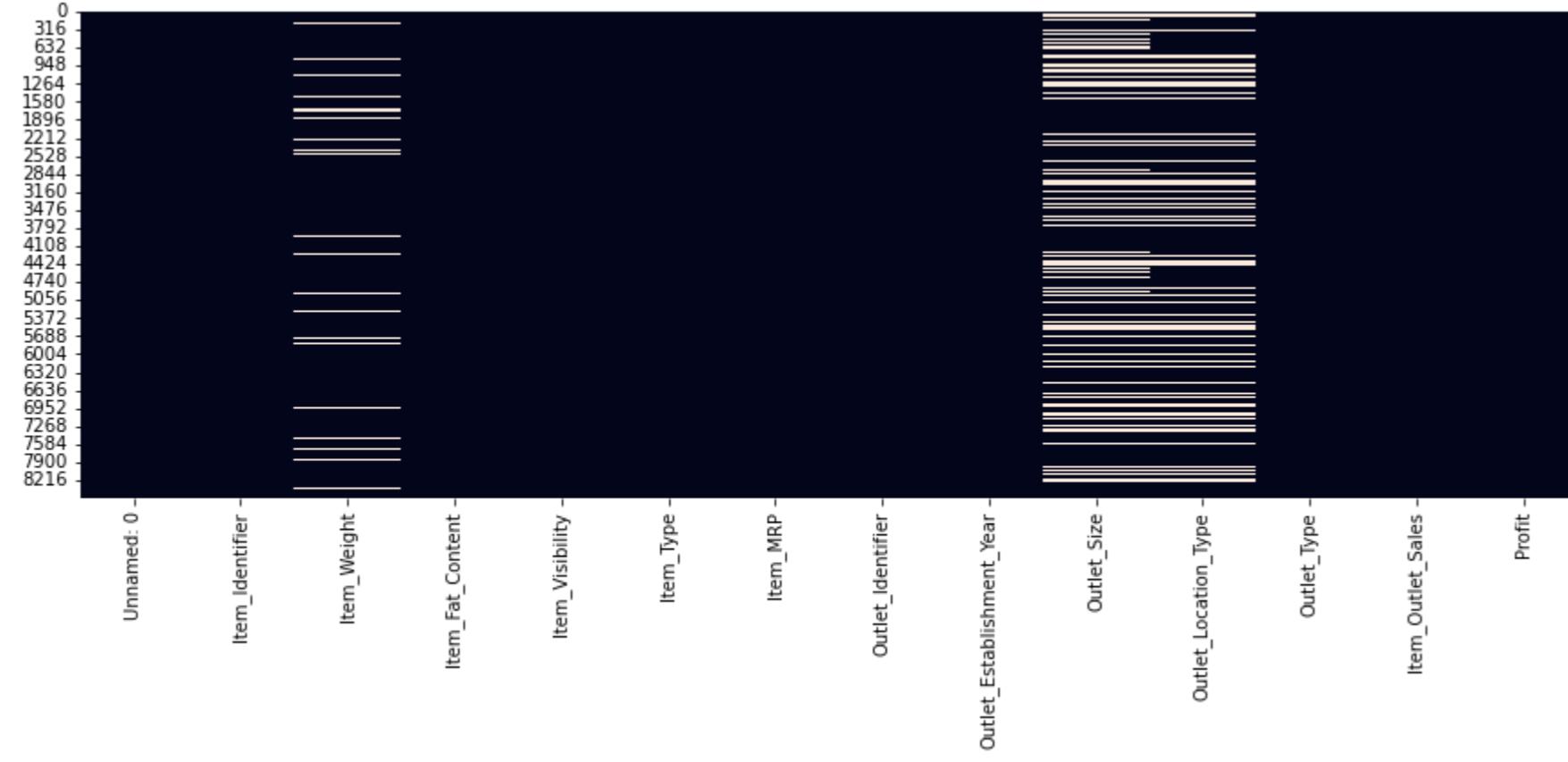
Out[69]: 5209

```
In [70]: 1 # Sanity Check whether our dataset has null values:  
2 df.isnull().values.any()
```

Out[70]: True

In [71]:

```
1 # Let us plot the Heatmap to visualize the missing values:  
2  
3  
4  
5 # Set the Figure Size:  
6 plt.rcParams['figure.figsize'] = [15, 5]  
7  
8 # Plot the heatmap:  
9 sns.heatmap(df.isnull(), cbar = False)  
10  
11 # Plot the map:  
12 plt.show()
```



```
In [72]: 1 # get the count of Missing Values:  
2 missing_values = df.isnull().sum()  
3  
4 # check for missing Values:  
5 total = df.isnull().sum().sort_values(ascending = False)  
6  
7 # Calculate percentage of Null Values:  
8 percent = ((df.isnull().sum() / df.shape[0])*100)  
9  
10 # Sort the Values in descending order  
11 percent = percent.sort_values(ascending = False)  
12  
13 # Concatenate the Total Missing Values and Percentage of Missing Values:  
14 missing_data = pd.concat([total, percent], axis = 1,  
15                         keys = ['Total', 'percentage'])  
16 # Add the Data types:  
17 missing_data['Type'] = df[missing_data.index].dtypes  
18  
19 # View the Missing Data:  
20 missing_data
```

Out[72]:

	Total	percentage	Type
Outlet_Size	2410	28.276428	object
Outlet_Location_Type	2050	24.052564	object
Item_Weight	749	8.787985	float64
Unnamed: 0	0	0.000000	int64
Item_Identifier	0	0.000000	object
Item_Fat_Content	0	0.000000	object
Item_Visibility	0	0.000000	float64
Item_Type	0	0.000000	object
Item_MRP	0	0.000000	float64
Outlet_Identifier	0	0.000000	object
Outlet_Establishment_Year	0	0.000000	int64
Outlet_Type	0	0.000000	object
Item_Outlet_Sales	0	0.000000	float64
Profit	0	0.000000	float64

```
In [73]: 1 df.shape
```

Out[73]: (8523, 14)

```
In [74]: 1 # drop the rows having missing values:  
2 df_sales_drop = df.dropna()  
3  
4 # Display the DataFrame  
5 df_sales_drop.shape
```

Out[74]: (5364, 14)

```
In [75]: 1 df.shape
```

```
Out[75]: (8523, 14)
```

```
In [76]: 1 # Make a copy of the Dropped DataFrame:  
2 df_Sales2 = df_sales_drop.copy()
```

Let us Consider each Variable Separately for missing value Treatment

1. Outlet_size

```
In [77]: 1 # check Number of Null Values:  
2 df.Outlet_Size.isnull().sum()
```

```
Out[77]: 2410
```

```
In [78]: 1 df.Outlet_Size.value_counts()
```

```
Out[78]: Medium    2793  
Small     2388  
High      932  
Name: Outlet_Size, dtype: int64
```

```
In [79]: 1 # Let's Check the Head of the Column:  
2 df.Outlet_Size
```

```
Out[79]: 0      Medium  
1      Medium  
2      Medium  
3        NaN  
4      High  
...  
8518     High  
8519      NaN  
8520     Small  
8521     Medium  
8522     Small  
Name: Outlet_Size, Length: 8523, dtype: object
```

```
In [80]: 1 # Obtain the Mode Value:  
2 df.Outlet_Size.mode()
```

```
Out[80]: 0      Medium  
dtype: object
```

```
In [81]: 1 # Replace all the Missing Values with 'Medium':  
2 df.Outlet_Size.replace(np.NaN, 'Medium', inplace = True)
```

```
In [82]: 1 # Sanity Check for the Missing Values:  
2 df.Outlet_Size.isnull().sum()
```

Out[82]: 0

```
In [83]: 1 # check for Value Counts:  
2 df.Outlet_Size.value_counts()
```

```
Out[83]: Medium    5203  
Small     2388  
High      932  
Name: Outlet_Size, dtype: int64
```

2. The Variable Outlet_Location_Type:

```
In [84]: 1 df.Outlet_Location_Type.isnull().sum()
```

Out[84]: 2050

```
In [85]: 1 # Check for the Count of Sub-categories:  
2 df.Outlet_Location_Type.value_counts()
```

```
Out[85]: Tier 2    2793  
Tier1    2388  
Tier 3    932  
?        120  
--       109  
-        67  
na       48  
NAN      16  
Name: Outlet_Location_Type, dtype: int64
```

```
In [86]: 1 # replace '?' with NaN  
2 # to_replace : value that will be replaced  
3 # Value : Value to replace values matching 'to replace' with  
4 df.Outlet_Location_Type.replace(to_replace = '?', value = np.NaN, inplace = True)  
5  
6  
7 # replace ' - ' with NaN  
8 # to_replace : value that will be replaced  
9 # Value : Value to replace values matching 'to replace' with  
10 df.Outlet_Location_Type.replace(to_replace = 'na', value = np.NaN, inplace = True)  
11  
12  
13 # replace ' - ' with NaN  
14 # to_replace : value that will be replaced  
15 # Value : Value to replace values matching 'to replace' with  
16 df.Outlet_Location_Type.replace(to_replace = 'NAN', value = np.NaN, inplace = True)
```

```
In [87]: 1 # replace '--' with NaN  
2 # to_replace : value that will be replaced  
3 # Value : Value to replace values matching 'to replace' with  
4 df.Outlet_Location_Type.replace(to_replace = ' --', value = np.NaN, inplace = True)
```

```
In [88]: 1 # replace '-' with NaN  
2 # to_replace : value that will be replaced  
3 # Value : Value to replace values matching 'to replace' with  
4 df.Outlet_Location_Type.replace(to_replace = ' -', value = np.NaN, inplace = True )
```

```
In [89]: 1 # Sanity Check:  
2 df.Outlet_Location_Type.value_counts()
```

```
Out[89]: Tier 2    2793  
Tier1     2388  
Tier 3     932  
Name: Outlet_Location_Type, dtype: int64
```

```
In [90]: 1 # Sanity Check for Missing Values:  
2 df.Outlet_Location_Type.isnull().sum()
```

```
Out[90]: 2410
```

```
In [91]: 1 df.head(1)
```

```
Out[91]: Unnamed:  
          0 Item_Identifier Item_Weight Item_Fat_Content Item_Visibility Item_Type Item_MRP Outlet_Identifier Outlet_Establishment_Year Outlet_Size Outlet_Location_Type Outlet_Type Item_Outlet_Sales Profit  
          0          0        FDA15         9.3      Low Fat       0.016047     Dairy      249.8        OUT049           1999      Medium      Tier 2 Supermarket Type1      3735.138     11.5
```

```
In [92]: 1 # Considering the Variable Item_Outlet_sales and Outlet_Location_Type  
2 # Dropna() : drops rows having missing values:  
3  
4 df_impute = df[['Item_Outlet_Sales', 'Outlet_Location_Type']].dropna()  
5  
6  
7 # Group the Data by the Loaction Type:  
8 # Mean() : return the mean of the Values:  
9 df_impute = df_impute.groupby(df_impute['Outlet_Location_Type']).mean()  
10  
11 # Display the DataFrame:  
12 df_impute
```

```
Out[92]: Item_Outlet_Sales  
          Outlet_Location_Type  
          Tier 2    2681.603542  
          Tier 3    2298.995256  
          Tier1     1912.149161
```

```
In [93]: 1 # Create a variable 'ind' for all the indexes:  
2 index = np.where(df.Outlet_Location_Type.isnull())[0]  
3  
4  
5 # Impute the Values using For Loop:  
6 for ind in index:  
7     if df.loc[ind, 'Item_Outlet_Sales'] <= 2100:  
8         df.loc[ind, 'Outlet_Location_Type'] = 'Tier1'  
9     elif df.loc[ind, 'Item_Outlet_Sales'] >= 2500:  
10        df.loc[ind, 'Outlet_Location_Type'] = 'Tier2'  
11    else:  
12        df.loc[ind, 'Outlet_Location_Type'] = 'Tier3'
```

```
In [94]: 1 df.Outlet_Location_Type.isnull().sum()
```

```
Out[94]: 0
```

```
In [95]: 1 df.Item_Weight.describe()
```

```
Out[95]: count    7774.000000  
mean      11.676740  
std       5.776851  
min       0.000000  
25%      7.720000  
50%      11.800000  
75%      16.500000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

```
In [96]: 1 # Replace the Value with NaN  
2 df.Item_Weight.replace(0, np.NaN, inplace = True)
```

```
In [97]: 1 df.Item_Weight.describe()
```

```
Out[97]: count    7060.000000  
mean      12.857645  
std       4.643456  
min       4.555000  
25%      8.773750  
50%      12.600000  
75%      16.850000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

Choice of Making Median or Mean for Numerical Variable

In [98]:

```
1 # Import Necessary Libraries
2 from matplotlib import gridspec
3
4 # Set the Plot Size:
5 plt.rcParams['figure.figsize'] = [15, 5]
6
7 # Specify the Geometry of the Grid that a Subplot is Placed in:
8 # Split the Plot into 2 Rows and 3 Columns:
9 gs = gridspec.GridSpec(2, 3, width_ratios = [.5, .5, .5], height_ratios = [2, 15])
10
11
12 # Step = 1: Specify the Plot Location by calling 'gs' initiative above:
13 # Step = 2: Write the Plot Text:
14 # Write the Text in the Plot using text()
15 # x : Location on x-axis where the Text is to be Written
16 # y : Location on y-axis where the Text is to be Written
17 # s : Text to be Written
18 # Fontsize: set the Font Size
19 # Step = 3: use.axis('off') to Hide the X and Y axes:
20
21 # Plot the 1st Row and 1st column:
22 a11 = plt.subplot(gs[0, 0])
23 a11.text(x = 0.1, y = 0.03, s = 'Original Data', fontsize = 15)
24 a11.axis('off')
25
26 # Plot the 1st Row and 2nd Column:
27 a12 = plt.subplot(gs[0, 1])
28 a12.text(x = 0.1, y = 0.03, s = 'Data Imputed with Mean', fontsize = 15)
29 a12.axis('off')
30
31 # Plot the 1st Row and 3rd Column:
32 a12 = plt.subplot(gs[0, 2])
33 a12.text(x = 0.1, y = 0.03, s = 'Data Imputed with Median', fontsize = 15)
34 a12.axis('off')
35
36 # SUMMARY STATISTICS:
37 # Step = 1: Specify the Location of the Plot Calling 'gs' Initiative Above in:
38 # Step = 2: Specify the Text along with Location:
39 # Step = 3: use.axis('off') to Hide the x and y axes:
40
41 # Plot in 2nd Row and 1st Column:
42 # Original Data:
43 a12 = plt.subplot(gs[1, 0])
44 a12.text(0.05, .3, s = str(df['Item_Weight'].describe()), fontsize = 15)
45 a12.axis('off')
46
47 # Plot in 2nd Row and 2nd Column:
48 # fill the missing Values with Mean
49 # Obtain the Mean of the Data
50 mu = df['Item_Weight'].mean()
51 a12 = plt.subplot(gs[1, 1])
52 a12.text(0.05, 0.3, s = str(df['Item_Weight'].fillna(mu).describe()), fontsize = 15)
53 a12.axis('off')
54
55
56 # Fill the Missing Value with Median:
57 me = df['Item_Weight'].median()
58 a12 = plt.subplot(gs[1, 2])
59 a12.text(0.05, 0.3, s = str(df['Item_Weight'].fillna(me).describe()), fontsize = 15)
60 a12.axis('off')
61
```

```
62 # Display the Plot:  
63 plt.show()  
64
```

Original Data

```
count    7060.000000  
mean     12.857645  
std      4.643456  
min      4.555000  
25%     8.773750  
50%     12.600000  
75%     16.850000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

Data Imputed with Mean

```
count    8523.000000  
mean     12.857645  
std      4.226124  
min      4.555000  
25%     9.310000  
50%     12.857645  
75%     16.000000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

Data Imputed with Median

```
count    8523.000000  
mean     12.813420  
std      4.227240  
min      4.555000  
25%     9.310000  
50%     12.600000  
75%     16.000000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

```
In [99]: 1 df['Item_Weight'].fillna(me, inplace = True)
```

```
In [100]: 1 df.Item_Weight.describe()
```

```
Out[100]: count    8523.000000  
mean     12.813420  
std      4.227240  
min      4.555000  
25%     9.310000  
50%     12.600000  
75%     16.000000  
max      21.350000  
Name: Item_Weight, dtype: float64
```

In [101]:

```
1 # Sanity check:  
2  
3  
4 # get the count of Missing Values:  
5 missing_values = df.isnull().sum()  
6  
7 # check for missing Values:  
8 total = df.isnull().sum().sort_values(ascending = False)  
9  
10 # Calculate percentage of Null Values:  
11 percent = ((df.isnull().sum() / df.shape[0])*100)  
12  
13 # Sort the Values in descending order  
14 percent = percent.sort_values(ascending = False)  
15  
16 # Concatenate the Total Missing Values and Percentage of Missing Values:  
17 missing_data = pd.concat([total, percent], axis = 1,  
18                         keys = ['Total', 'percentage'])  
19 # Add the Data types:  
20 missing_data['Type'] = df[missing_data.index].dtypes  
21  
22 # View the Missing Data:  
23 missing_data
```

Out[101]:

	Total	percentage	Type
Unnamed: 0	0	0.0	int64
Item_Identifier	0	0.0	object
Item_Weight	0	0.0	float64
Item_Fat_Content	0	0.0	object
Item_Visibility	0	0.0	float64
Item_Type	0	0.0	object
Item_MRP	0	0.0	float64
Outlet_Identifier	0	0.0	object
Outlet_Establishment_Year	0	0.0	int64
Outlet_Size	0	0.0	object
Outlet_Location_Type	0	0.0	object
Outlet_Type	0	0.0	object
Item_Outlet_Sales	0	0.0	float64
Profit	0	0.0	float64

OUTLIER DETECTION

- 1. Based on Boxplots

```
In [102]: 1 # Filtering the Numerical Column from the DataSet:  
2 df_num = df.select_dtypes(include = [np.number])
```

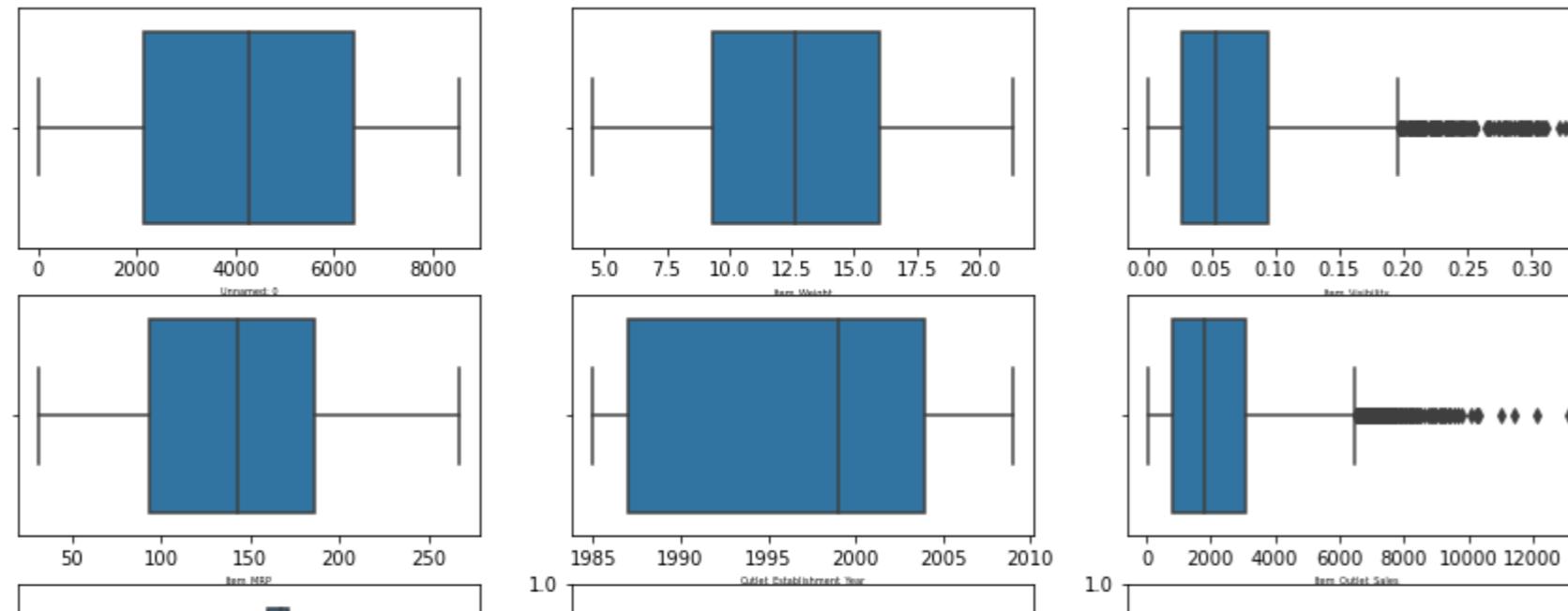
```
In [103]: 1 df_num.head(2)
```

```
Out[103]: Unnamed: 0 Item_Weight Item_Visibility Item_MRP Outlet_Establishment_Year Item_Outlet_Sales Profit  
0 0 9.30 0.016047 249.8 1999 3735.1380 11.5  
1 1 5.92 0.019278 48.3 2009 443.4228 14.3
```

```
In [104]: 1 df_num.columns
```

```
Out[104]: Index(['Unnamed: 0', 'Item_Weight', 'Item_Visibility', 'Item_MRP',  
'Outlet_Establishment_Year', 'Item_Outlet_Sales', 'Profit'],  
dtype='object')
```

```
In [105]: 1 # Plot the Boxplot for each Columns:  
2 # Subplots() : plot Subplots:  
3 # Figsize() : set the Figure Size:  
4 fig, ax = plt.subplots(3, 3, figsize = (15, 8))  
5  
6  
7 # Plot the Boxplot using boxplot() from Seaborn Library:  
8 # Z : let the Variable z defines the Boxplot:  
9 # X : data from which the Boxplot is to be Predicted:  
10 # orient : 'h' specifies the Horizontal Boxplot  
11 # Whis : Proportion of the IQR, It past the Low and High Quartiles to extend the plot  
12 # ax : specifies the axes object to draw the plot  
13 # Set_xlabel() : set the x - axis Label  
14 # FontSize () : set the Font size of the x - axis  
15 for variable, subplot in zip(df_num.columns, ax.flatten()):  
16     z = sns.boxplot(x = df_num[variable], orient = 'h', whis = 1.5, ax = subplot)  
17     z.set_xlabel(variable, fontsize = 5)
```



- 1. Based on IQR Method

```
In [106]: 1 # Obtain the First Quartiles:  
2 Q1 = df_num.quantile(0.25)  
3  
4 # Obtain the Quartiles:  
5 Q3 = df_num.quantile(0.75)  
6  
7 # Obtain the IQR:  
8 IQR = Q3 - Q1  
9  
10 # Print the IQR:  
11 IQR
```

```
Out[106]: Unnamed: 0          4261.000000  
Item_Weight           6.690000  
Item_Visibility      0.067596  
Item_MRP             91.850000  
Outlet_Establishment_Year 17.000000  
Item_Outlet_Sales    2267.049000  
Profit               1.150000  
dtype: float64
```

```
In [107]: 1 # Filter out the Outlier Values:  
2 # ~ : Select all the rows which do not satisfy the Condition  
3 # any() : returns whether the element is True over the Column  
4 # axis = 1: Indicates should Select the alternate columns('0' for index Positions)  
5  
6 from warnings import filterwarnings  
7 filterwarnings('ignore')  
8  
9  
10 df_sales_iqr = df[(df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))].any(axis = 1)
```

```
In [108]: 1 df_sales_iqr.shape
```

```
Out[108]: (7515, 14)
```

```
In [109]: 1 df.shape
```

```
Out[109]: (8523, 14)
```

Based on Z - Score

z - score of a value is the Difference Between that value and the mean, divided by the Standard Deviation. If the Z - score greater than +3 or less than -3 it indicates an Outlier Value in our Dataset.

```
In [110]: 1 df.head(1)
```

```
Out[110]: Unnamed:  
          0 Item_Identifier Item_Weight Item_Fat_Content Item_Visibility Item_Type Item_MRP Outlet_Identifier Outlet_Establishment_Year Outlet_Size Outlet_Location_Type Outlet_Type Item_Outlet_Sales Profit  
          0      0        FDA15       9.3      Low Fat     0.016047    Dairy      249.8      OUT049           1999      Medium      Tier 2 Supermarket Type1      3735.138   11.5
```

```
In [111]: 1 # Import the Library:  
2 import scipy  
3  
4 # from scipy import Stats module:  
5 from scipy import stats  
6  
7 # z - score are defined for each observations in a Variable  
8 # Compute the z - scores using the method zscore from the Scipy Library  
9 z_scores_profit = scipy.stats.zscore(df_num['Profit'])  
10  
11 # display the z - scores:  
12 z_scores_profit
```

```
Out[111]: 0      -1.125033  
1       0.520342  
2       0.637869  
3       0.108998  
4       0.402815  
...  
8518    0.402815  
8519    0.461578  
8520   -2.300301  
8521    0.461578  
8522    0.696632  
Name: Profit, Length: 8523, dtype: float64
```

```
In [112]: 1 # Print the Rows where the z - score is less than -3  
2 row_index_less = np.where(z_scores_profit < - 3)  
3  
4 print(row_index_less)
```

```
(array([ 41, 50, 144, 217, 320, 324, 406, 425, 432, 457, 546,  
       602, 607, 611, 641, 716, 761, 816, 892, 921, 926, 957,  
       960, 986, 1066, 1133, 1178, 1212, 1237, 1304, 1409, 1459, 1533,  
       1551, 1561, 1582, 1623, 1645, 1649, 1659, 1669, 1715, 1784, 1788,  
       1805, 1816, 1835, 1867, 1893, 1900, 1916, 1956, 1974, 1984, 1989,  
       2000, 2019, 2060, 2062, 2130, 2201, 2224, 2246, 2289, 2332, 2338,  
       2425, 2449, 2477, 2510, 2515, 2526, 2537, 2543, 2572, 2636, 2792,  
       2836, 2919, 2925, 2969, 2987, 2992, 2996, 3025, 3035, 3098, 3105,  
       3223, 3254, 3296, 3356, 3415, 3416, 3476, 3513, 3619, 3668, 3675,  
       3759, 3785, 3849, 3852, 3904, 3917, 3932, 3936, 3990, 4073, 4091,  
       4146, 4197, 4213, 4407, 4487, 4493, 4556, 4571, 4705, 4721, 4724,  
       4764, 4800, 5049, 5127, 5231, 5274, 5286, 5313, 5336, 5378, 5516,  
       5540, 5657, 5658, 5728, 5730, 5750, 5788, 5883, 5891, 5906, 5971,  
       5987, 6093, 6137, 6211, 6247, 6271, 6372, 6390, 6419, 6507, 6534,  
       6541, 6558, 6606, 6610, 6696, 6759, 6760, 6776, 6778, 6897, 6931,  
       6972, 7042, 7180, 7242, 7267, 7304, 7351, 7384, 7432, 7629, 7684,  
       7751, 7779, 7797, 7822, 7851, 7893, 7919, 7948, 8074, 8099, 8101,  
       8237, 8273, 8366, 8376, 8395, 8429, 8433, 8456], dtype=int64),)
```

Interpretation : The Rows Corresponding to the above Displayed index are the Outliers for the Data.

```
In [113]: 1 # Print the rows where z - score is more than 3:  
2 row_index_more = np.where(z_scores_profit > 3)  
3  
4 row_index_more
```

```
Out[113]: (array([3026, 4386, 5089, 8369], dtype=int64),)
```

```
In [114]: 1 # Count of Outliers in the Variable representing Profit:  
2  
3  
4 len(row_index_less[0]) + len(row_index_more[0])  
5
```

```
Out[114]: 199
```

There are 199 Outliers in the Data by using z - score identification Methodology

```
In [115]: 1 # Filter out the Outlier Values:  
2 # ~ : select all the Rows which do not Satisfy the Condition  
3  
4 df_Zscore_profit = df['Profit'][~((z_scores_profit < -3)|(z_scores_profit > 3))]  
5
```

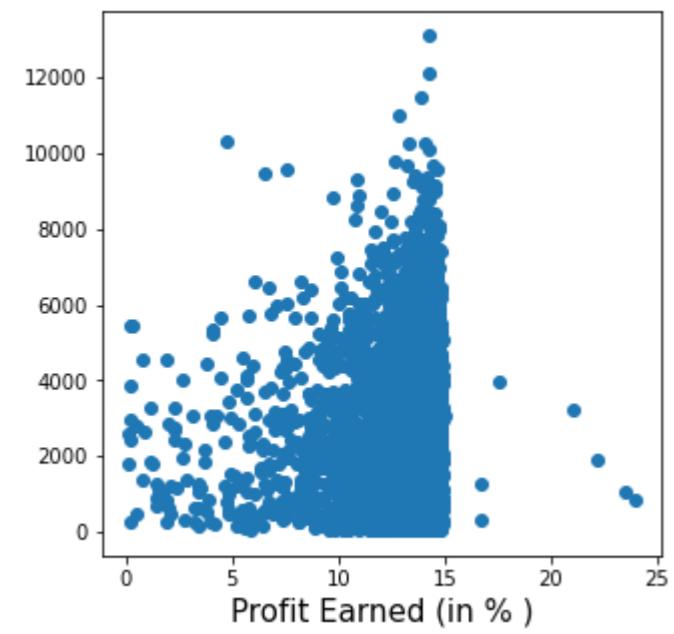
```
In [116]: 1 # Check for the Shape:  
2 df_Zscore_profit.shape
```

```
Out[116]: (8324,)
```

```
In [117]: 1 # Original Dataset Shape:  
2 df.shape
```

Out[117]: (8523, 14)

```
In [118]: 1 # Initialize the Figure  
2 # FItsize: set the Plot size:  
3  
4 fig, ax = plt.subplots(figsize = (5, 5))  
5  
6 # Plot the Scatter Plot:  
7 ax.scatter(df['Profit'], df['Item_Outlet_Sales'])  
8  
9  
10 # Set the x - axis label:  
11 # Fontsize : set the font size:  
12 ax.set_xlabel('Profit Earned (in % )', fontsize = 15)  
13  
14 # Set the y - axis  
15 ax.set_ylabel('')  
16 plt.show()
```



```
In [2]: 1 import pandas as pd  
2 import numpy as np  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns
```

```
In [124]: 1 df.head(5)
```

	Unnamed: 0	Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales	Profit
0	0	FDA15	9.30	Low Fat	0.016047	Dairy	249.8	OUT049	1999	Medium	Tier 2	Supermarket Type1	3735.1380	11.5
1	1	DRC01	5.92	Regular	0.019278	Soft Drinks	48.3	OUT018	2009	Medium	Tier 2	Supermarket Type2	443.4228	14.3
2	2	FDN15	17.50	Low Fat	0.016760	Meat	141.6	OUT049	1999	Medium	Tier 2	Supermarket Type1	2097.2700	14.5
3	3	FDX07	19.20	Regular	0.000000	Fruits and Vegetables	182.1	OUT010	1998	Medium	Tier1	Grocery Store	732.3800	13.6
4	4	NCD19	8.93	Low Fat	0.000000	Household	53.9	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052	14.1

```
In [123]: 1 # import the sklearn Library:
```

```
2 import sklearn
3
4 # import the train_test_split module from sklearn
5 from sklearn.model_selection import train_test_split
```

```
In [125]: 1 # Select the target column:
```

```
2 Y = df['Item_Outlet_Sales']
3
4 # select the Independent Column:
5 # By dropping the Target Column
6 X = df.drop(['Item_Outlet_Sales'], axis = 1)
```

```
In [130]: 1 # Let us now split the Dataset into Test and Train
```

```
2 # Test size : The Proportion of the Data to be included in the Testing set.
```

```
3
4 x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size = 0.25,
5                                                 random_state = 100)
6 print('X_train : ', x_train.shape)
7 print('X_test : ', x_test.shape)
8 print('Y_train : ', y_train.shape)
9 print('Y_test : ', y_test.shape)
```

```
X_train : (6392, 13)
X_test : (2131, 13)
Y_train : (6392,)
Y_test : (2131,)
```

Linear Regression in Machine Learning

- Linear regression is one of the easiest and most popular Machine Learning algorithms.
- It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables such as [sales](#), [salary](#), [age](#), [product price](#), etc.
- Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression.
- Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

Types of Linear Regression

Linear regression can be further divided into two types of the algorithm:

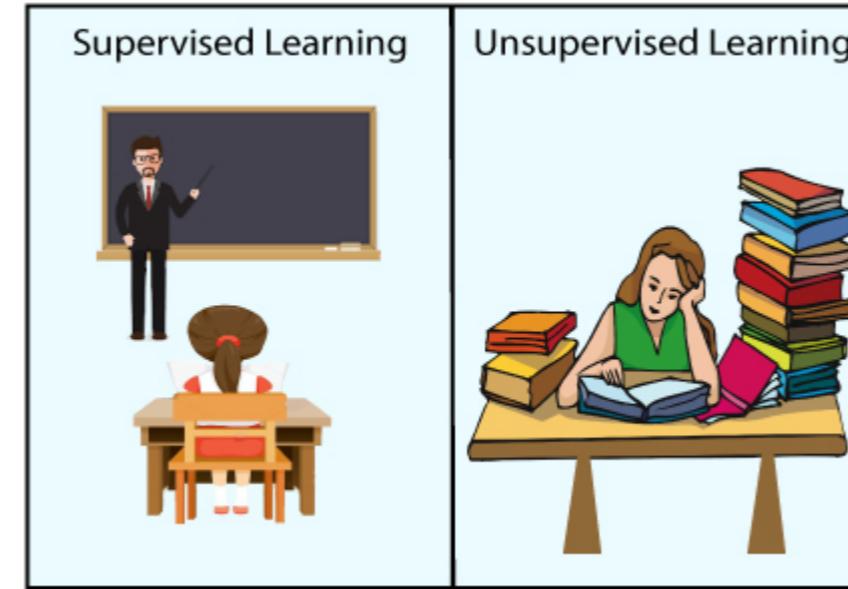
Simple Linear Regression:

If a single independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Simple Linear Regression.

Multiple Linear regression:

If more than one independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Multiple Linear Regression.

Supervised and Unsupervised Learning



Supervised Machine Learning:

Supervised learning is a machine learning method in which models are trained using labeled data. In supervised learning, models need to find the mapping function to map the input variable (X) with the output variable (Y).

Supervised Machine learning Supervised learning needs supervision to train the model, which is similar to as a student learns things in the presence of a teacher. Supervised learning can be used for two types of problems: Classification and Regression.

Example: Suppose we have an image of different types of fruits. The task of our supervised learning model is to identify the fruits and classify them accordingly. So to identify the image in supervised learning, we will give the input data as well as output for that, which means we will train the model by the shape, size, color, and taste of each fruit. Once the training is completed, we will test the model by giving the new set of fruit. The model will identify the fruit and predict the output using a suitable algorithm

Supervised learning is the types of machine learning in which machines are trained using well "labelled" training data, and on basis of that data, machines predict the output. The labelled data means some input data is already tagged with the correct output.

In supervised learning, the training data provided to the machines work as the supervisor that teaches the machines to predict the output correctly. It applies the same concept as a student learns in the supervision of the teacher.

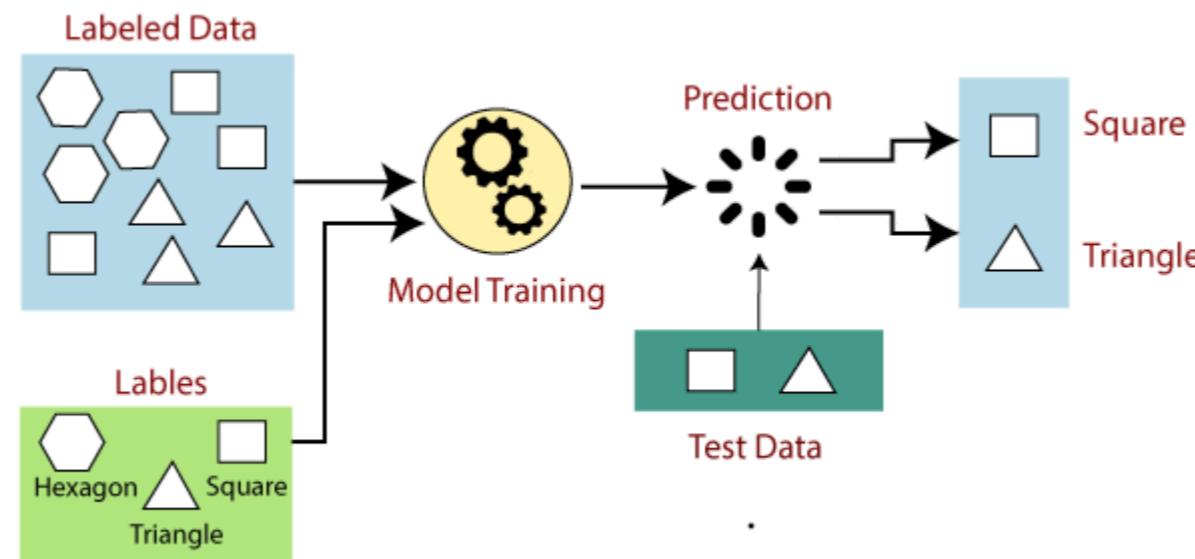
Supervised learning is a process of providing input data as well as correct output data to the machine learning model. The aim of a supervised learning algorithm is to find a mapping function to map the input variable(x) with the output variable(y).

In the real-world, supervised learning can be used for Risk Assessment, Image classification, Fraud Detection, spam filtering, etc

How Supervised Learning Works?

In supervised learning, models are trained using labelled dataset, where the model learns about each type of data. Once the training process is completed, the model is tested on the basis of test data (a subset of the training set), and then it predicts the output.

The working of Supervised learning can be easily understood by the below example and diagram



Suppose we have a dataset of different types of shapes which includes square, rectangle, triangle, and Polygon. Now the first step is that we need to train the model for each shape.

If the given shape has four sides, and all the sides are equal, then it will be labelled as a Square.

If the given shape has three sides, then it will be labelled as a triangle.

If the given shape has six equal sides then it will be labelled as hexagon.

Now, after training, we test our model using the test set, and the task of the model is to identify the shape.

The machine is already trained on all types of shapes, and when it finds a new shape, it classifies the shape on the bases of a number of sides, and predicts the output.

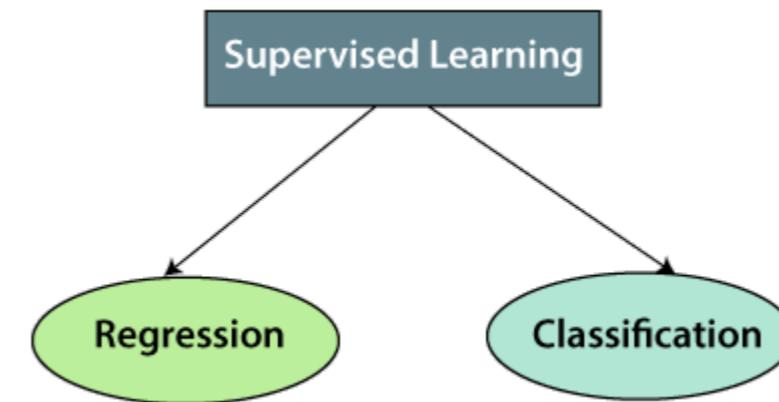
Steps Involved in Supervised Learning:

- First Determine the type of training dataset
- Collect/Gather the labelled training data.
- Split the training dataset into training dataset, test dataset, and validation dataset.
- Determine the input features of the training dataset, which should have enough knowledge so that the model can accurately predict the output.
- Determine the suitable algorithm for the model, such as support vector machine, decision tree, etc.
- Execute the algorithm on the training dataset. Sometimes we need validation sets as the control parameters, which are the subset of training datasets.

- Evaluate the accuracy of the model by providing the test set. If the model predicts the correct output, which means our model is accurate.

Types of supervised Machine learning Algorithms:

Supervised learning can be further divided into two types of problems:



1. Regression

Regression algorithms are used if there is a relationship between the input variable and the output variable. It is used for the prediction of continuous variables, such as Weather forecasting, Market Trends, etc. Below are some popular Regression algorithms which come under supervised learning:

- Linear Regression
- Regression Trees
- Non-Linear Regression
- Bayesian Linear Regression
- Polynomial Regression

2. Classification

Classification algorithms are used when the output variable is categorical, which means there are two classes such as Yes-No, Male-Female, True-false, etc.

Spam Filtering,

- Random Forest
- Decision Trees
- Logistic Regression
- Support vector Machines

Advantages of Supervised learning:

- With the help of supervised learning, the model can predict the output on the basis of prior experiences.

- In supervised learning, we can have an exact idea about the classes of objects.
- Supervised learning model helps us to solve various real-world problems such as fraud detection, spam filtering, etc.

Disadvantages of supervised learning:

- Supervised learning models are not suitable for handling the complex tasks.
- Supervised learning cannot predict the correct output if the test data is different from the training dataset.
- Training required lots of computation times.
- In supervised learning, we need enough knowledge about the classes of object.

Unsupervised Machine Learning

In the previous topic, we learned supervised machine learning in which models are trained using labeled data under the supervision of training data. But there may be many cases in which we do not have labeled data and need to find the hidden patterns from the given dataset. So, to solve such types of cases in machine learning, we need unsupervised learning techniques.

As the name suggests, unsupervised learning is a machine learning technique in which models are not supervised using training dataset. Instead, models itself find the hidden patterns and insights from the given data. It can be compared to learning which takes place in the human brain while learning new things. It can be defined as:

Unsupervised learning cannot be directly applied to a regression or classification problem because unlike supervised learning, we have the input data but no corresponding output data. **The goal of unsupervised learning is to find the underlying structure of dataset, group that data according to similarities, and represent that dataset in a compressed format.**

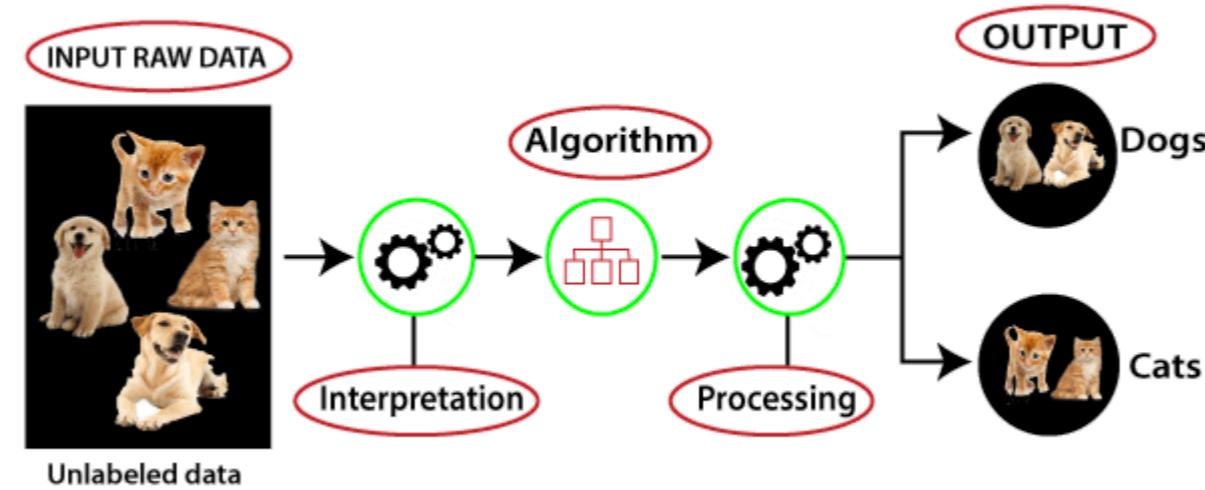
Example: Suppose the unsupervised learning algorithm is given an input dataset containing images of different types of cats and dogs. The algorithm is never trained upon the given dataset, which means it does not have any idea about the features of the dataset. The task of the unsupervised learning algorithm is to identify the image features on their own. Unsupervised learning algorithm will perform this task by clustering the image dataset into the groups according to similarities between images.



Why use Unsupervised Learning?

Below are some main reasons which describe the importance of Unsupervised Learning:

- Unsupervised learning is helpful for finding useful insights from the data.
- Unsupervised learning is much similar as a human learns to think by their own experiences, which makes it closer to the real AI.
- Unsupervised learning works on unlabeled and uncategorized data which make unsupervised learning more important.
- In real-world, we do not always have input data with the corresponding output so to solve such cases, we need unsupervised learning.

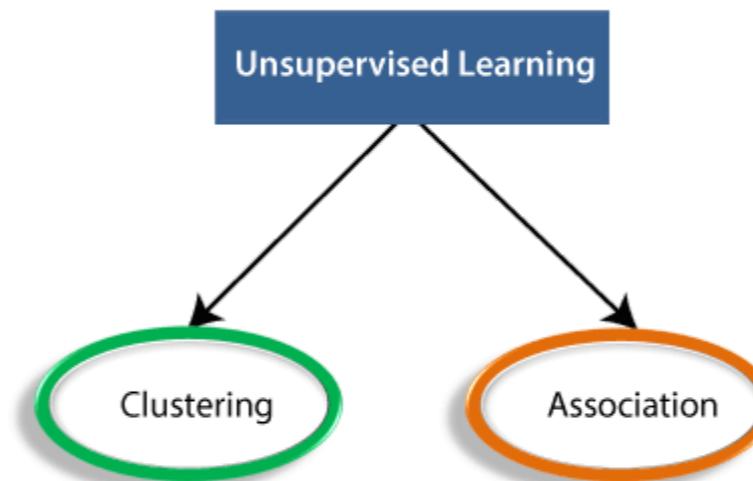


Here, we have taken an unlabeled input data, which means it is not categorized and corresponding outputs are also not given. Now, this unlabeled input data is fed to the machine learning model in order to train it. Firstly, it will interpret the raw data to find the hidden patterns from the data and then will apply suitable algorithms such as k-means clustering, Decision tree, etc.

Once it applies the suitable algorithm, the algorithm divides the data objects into groups according to the similarities and difference between the objects.

Types of Unsupervised Learning Algorithm:

The unsupervised learning algorithm can be further categorized into two types of p



Clustering: Clustering is a method of grouping the objects into clusters such that objects with most similarities remain into a group and has less or no similarities with the objects of another group. Cluster analysis finds the commonalities between the data objects and categorizes them as per the presence and absence of those commonalities.

Association: An association rule is an unsupervised learning method which is used for finding the relationships between variables in the large database. It determines the set of items that occurs together in the dataset. Association rule makes marketing strategy more effective. Such as people who buy X item (suppose a bread) are also tend to purchase Y (Butter/Jam) item. A typical example of Association rule is Market Basket Analysis.

Advantages of Unsupervised Learning

- Unsupervised learning is used for more complex tasks as compared to supervised learning because, in unsupervised learning, we don't have labeled input data.
- Unsupervised learning is preferable as it is easy to get unlabeled data in comparison to labeled data.

Disadvantages of Unsupervised Learning

- Unsupervised learning is intrinsically more difficult than supervised learning as it does not have corresponding output.
- The result of the unsupervised learning algorithm might be less accurate as input data is not labeled, and algorithms do not know the exact output in advance.

```
In [4]: # Imputation of SLR and OLS Summary:  
1  
2  
3 data = pd.read_excel('Weightexcel.xlsx')  
4 data
```

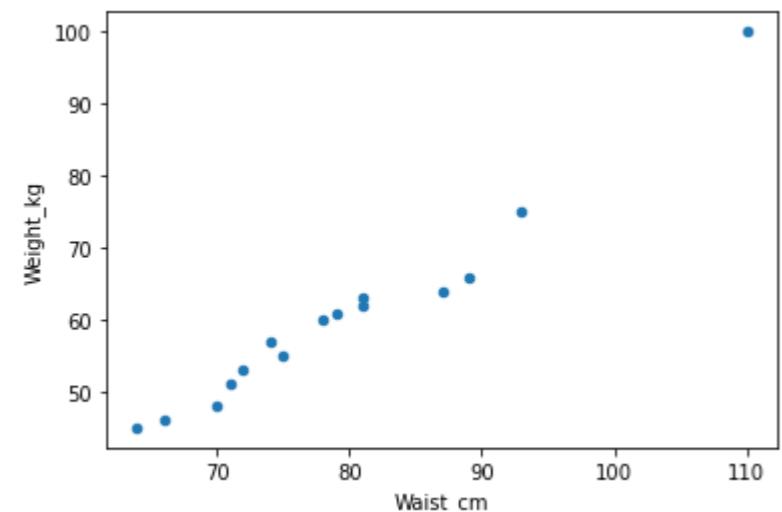
Out[4]:

	Waist_cm	Weight_kg
0	71	51
1	89	66
2	64	45
3	74	57
4	87	64
5	93	75
6	79	61
7	81	62
8	75	55
9	72	53
10	70	48
11	66	46
12	81	63
13	78	60
14	110	100

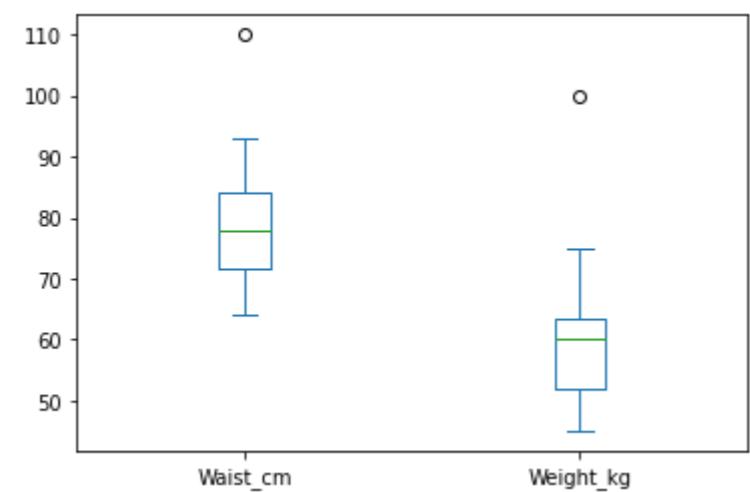
```
In [5]: data.shape
```

Out[5]: (15, 2)

```
In [6]: 1 data.plot(kind = 'scatter', x = 'Waist_cm', y = 'Weight_kg')
2 plt.show()
```



```
In [5]: 1 data.plot(kind = 'box')
2 plt.show()
```



```
In [6]: 1 data.corr()
```

Out[6]:

	Waist_cm	Weight_kg
Waist_cm	1.000000	0.981216
Weight_kg	0.981216	1.000000

```
In [13]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
In [14]: 1 # Change to the DataFrame variable:  
2 waist = pd.DataFrame(data['Waist_cm'])  
3 weight = pd.DataFrame(data['Weight_kg'])
```

```
In [15]: 1 waist
```

```
Out[15]: Waist_cm
```

	Waist_cm
0	71
1	89
2	64
3	74
4	87
5	93
6	79
7	81
8	75
9	72
10	70
11	66
12	81
13	78
14	110

```
In [16]: 1 import sklearn  
2  
3 from sklearn.linear_model import LinearRegression
```

```
In [17]: 1 # Making Instances:  
2 lm = LinearRegression()  
3 model = lm.fit(waist, weight)
```

```
In [18]: 1 model.coef_
```

```
Out[18]: array([[1.13470708]])
```

```
In [19]: 1 model.intercept_
```

```
Out[19]: array([-29.62009537])
```

```
In [20]: 1 model.score(waist, weight)
```

```
Out[20]: 0.9627843958606758
```

```
In [21]: 1 import statsmodels  
2 import statsmodels.api as sm
```

```
In [22]: 1 # Predict the new value of Weight:  
2 from warnings import filterwarnings  
3 filterwarnings('ignore')  
4  
5  
6 waist_new = np.array([97])  
7 waist_new = waist_new.reshape(-1, 1)  
8 weight_predict = model.predict(waist_new)  
9 weight_predict
```

Out[22]: array([[80.44649183]])

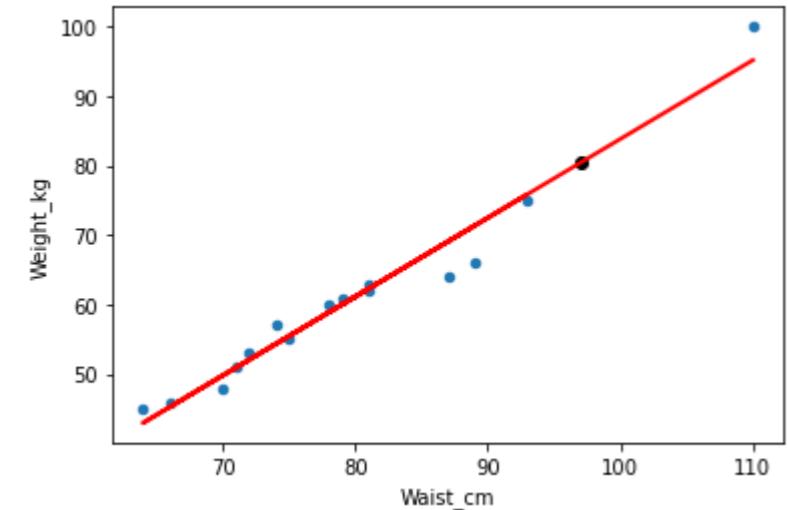
```
In [21]: 1 # Predict more Values:  
2 X = ([67, 78, 94])  
3 X = pd.DataFrame(X)  
4 Y = model.predict(X)  
5 Y = pd.DataFrame(Y)  
6  
7 df = pd.concat([X, Y], axis = 1, keys = ['Waist_new', 'Weight_predicted'])  
8 df
```

Out[21]:

	Waist_new	Weight_predicted
0	0	0
0	67	46.405279
1	78	58.887057
2	94	77.042371

In [27]:

```
1 # Visual the Result:  
2 data.plot(kind = 'scatter', x = 'Waist_cm', y = 'Weight_kg')  
3  
4  
5 # Plot the REgression Line:  
6 plt.plot(waist, model.predict(waist), color = 'r', linewidth = 2)  
7  
8  
9 # Plotting the Predicted Values:  
10 plt.scatter(waist_new, weight_predict, color = 'black')  
11  
12 # Display to Show:  
13 plt.show()
```



In [28]:

```
1 model = sm.OLS(X, Y)
```

In [29]:

```
1 fit = model.fit()
```

In [30]:

```
1 fit.pvalues
```

Out[30]:

```
0    0.002165  
dtype: float64
```

```
In [31]: 1 fit.summary()
```

```
Out[31]: OLS Regression Results
```

```
Dep. Variable: 0 R-squared (uncentered): 0.996
Model: OLS Adj. R-squared (uncentered): 0.994
Method: Least Squares F-statistic: 460.3
Date: Sat, 11 Feb 2023 Prob (F-statistic): 0.00217
Time: 21:34:05 Log-Likelihood: -9.2545
No. Observations: 3 AIC: 20.51
Df Residuals: 2 BIC: 19.61
Df Model: 1
Covariance Type: nonrobust

coef std err t P>|t| [0.025 0.975]
0 1.2931 0.060 21.455 0.002 1.034 1.552

Omnibus: nan Durbin-Watson: 0.981
Prob(Omnibus): nan Jarque-Bera (JB): 0.306
Skew: -0.222 Prob(JB): 0.858
Kurtosis: 1.500 Cond. No. 1.00
```

Notes:

- [1] R^2 is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Logistic Regression

```
In [24]: 1 df = pd.read_csv('Breast_cancer_data.csv')
2 df
```

Out[24]:

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness	diagnosis
0	17.99	10.38	122.80	1001.0	0.11840	0
1	20.57	17.77	132.90	1326.0	0.08474	0
2	19.69	21.25	130.00	1203.0	0.10960	0
3	11.42	20.38	77.58	386.1	0.14250	0
4	20.29	14.34	135.10	1297.0	0.10030	0
...
564	21.56	22.39	142.00	1479.0	0.11100	0
565	20.13	28.25	131.20	1261.0	0.09780	0
566	16.60	28.08	108.30	858.1	0.08455	0
567	20.60	29.33	140.10	1265.0	0.11780	0
568	7.76	24.54	47.92	181.0	0.05263	1

569 rows × 6 columns

```
In [25]: 1 df.diagnosis.value_counts()
```

Out[25]: 1 357
0 212
Name: diagnosis, dtype: int64

```
In [13]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 from warnings import filterwarnings
7 filterwarnings('ignore')
```

```
In [14]: 1 from sklearn.datasets import load_breast_cancer  
2  
3 dataset = load_breast_cancer(as_frame = True)  
4 dataset
```

```
Out[14]: {'data':   mean radius  mean texture  mean perimeter  mean area  mean smoothness  \  
0      17.99       10.38       122.80     1001.0      0.11840  
1      20.57       17.77       132.90     1326.0      0.08474  
2      19.69       21.25       130.00     1203.0      0.10960  
3      11.42       20.38       77.58      386.1      0.14250  
4      20.29       14.34       135.10     1297.0      0.10030  
..      ...        ...        ...        ...        ...  
564    21.56       22.39       142.00     1479.0      0.11100  
565    20.13       28.25       131.20     1261.0      0.09780  
566    16.60       28.08       108.30     858.1       0.08455  
567    20.60       29.33       140.10     1265.0      0.11780  
568    7.76        24.54       47.92      181.0      0.05263  
  
      mean compactness  mean concavity  mean concave points  mean symmetry  \  
0      0.27760       0.30010       0.14710      0.2419  
1      0.07864       0.08690       0.07017      0.1812  
2      0.15990       0.19740       0.12790      0.2069  
3      0.28390       0.24140       0.10520      0.2597  
4      0.13280       0.19800       0.10430      0.1809
```

```
In [15]: 1 dataset.data.head()
```

```
Out[15]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	worst area	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	25.38	17.33	184.60	2019.0	0.1622	0.6656	0.7119	0.2654	0.4601	0.1189
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	24.99	23.41	158.80	1956.0	0.1238	0.1866	0.2416	0.1860	0.2750	0.0890
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	23.57	25.53	152.50	1709.0	0.1444	0.4245	0.4504	0.2430	0.3613	0.0875
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	14.91	26.50	98.87	567.7	0.2098	0.8663	0.6869	0.2575	0.6638	0.1730
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	22.54	16.67	152.20	1575.0	0.1374	0.2050	0.4000	0.1625	0.2364	0.0767

5 rows × 30 columns

```
In [16]: 1 dataset.target.value_counts()
```

```
Out[16]: 1    357  
0    212  
Name: target, dtype: int64
```

```
In [17]: 1 X = dataset['data']  
2 Y = dataset['target']
```

```
In [18]: 1 # Load Library of Test Train and Split  
2 from sklearn.model_selection import train_test_split  
3  
4 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 0.25,  
5 random_state = 100)  
6  
7 print('X_train : ', X_train.shape)  
8 print('X_test : ', X_test.shape)  
9 print('Y_train : ', y_train.shape)  
10 print('Y_test : ', y_test.shape)
```

```
X_train : (426, 30)  
X_test : (143, 30)  
Y_train : (426,)  
Y_test : (143,)
```

```
In [19]: 1 from sklearn.preprocessing import StandardScaler  
2  
3 ss_train = StandardScaler()  
4 X_train = ss_train.fit_transform(X_train)  
5  
6 ss_test = StandardScaler()  
7 X_test = ss_test.fit_transform(X_test)
```

```
In [20]: 1 # Import Logistic Regression from the Sklearn Library:  
2 from sklearn.linear_model import LogisticRegression
```

```
In [21]: 1 # Create an Instance:  
2 lm = LogisticRegression()  
3 model = lm.fit(X, Y)
```

```
In [22]: 1 model
```

```
Out[22]: LogisticRegression()
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [23]: 1 predictions = model.predict(X_test)
```

In [24]:

```
1 from sklearn.metrics import confusion_matrix
2
3 cm = confusion_matrix(y_test, predictions)
4
5 TN, FP, FN, TP = confusion_matrix(y_test, predictions).ravel()
6
7 print('True Positives(TP) = ', TP)
8 print('False Positives(FP) = ', FP)
9 print('True Negative(TN) = ', TN)
10 print('False Negative(FN) = ', FN)
```

```
True Positives(TP) =  53
False Positives(FP) =  33
True Negative(TN) =  23
False Negative(FN) =  34
```

In [28]:

```
1 accuracy = (TP + TN) / (TP + FP + TN + FN)
2
3 print('Accuracy of the Binary Classifier : {:.3f}'.format(accuracy))
```

```
Accuracy of the Binary Classifier : 0.531
```

In [33]:

```
1 models = {}
2
3 # Logistic Regression:
4 from sklearn.linear_model import LogisticRegression
5 models['Logistic Regression'] = LogisticRegression()
6
7
8 # Support Vector Machines:
9 from sklearn.svm import LinearSVC
10 models['Support Vector Machines'] = LinearSVC()
11
12
13 # Decision Tree:
14 from sklearn.tree import DecisionTreeClassifier
15 models['Decision Tree'] = DecisionTreeClassifier()
16
17 # Random Forest:
18 from sklearn.ensemble import RandomForestClassifier
19 models['Random Forest'] = RandomForestClassifier()
20
21 # Naive Bayes:
22 from sklearn.naive_bayes import GaussianNB
23 models['Naive Bayes'] = GaussianNB()
24
25
26 # k - Nearestv Neighbors:
27 from sklearn.neighbors import KNeighborsClassifier
28 models['K- NN Algorithm'] = KNeighborsClassifier()
```

```
In [37]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score
2
3 # creating an instance:
4 accuracy, precision, recall = {}, {}, {}
5
6 for key in models.keys():
7
8     # Fit the Classifier:
9     models[key].fit(X_train, y_train)
10
11    # Make Predictions:
12    predictions = models[key].predict(X_test)
13
14    # Calculate metrics:
15    accuracy[key] = accuracy_score(predictions, y_test)
16    precision[key] = precision_score(predictions, y_test)
17    recall[key] = recall_score(predictions, y_test)
```

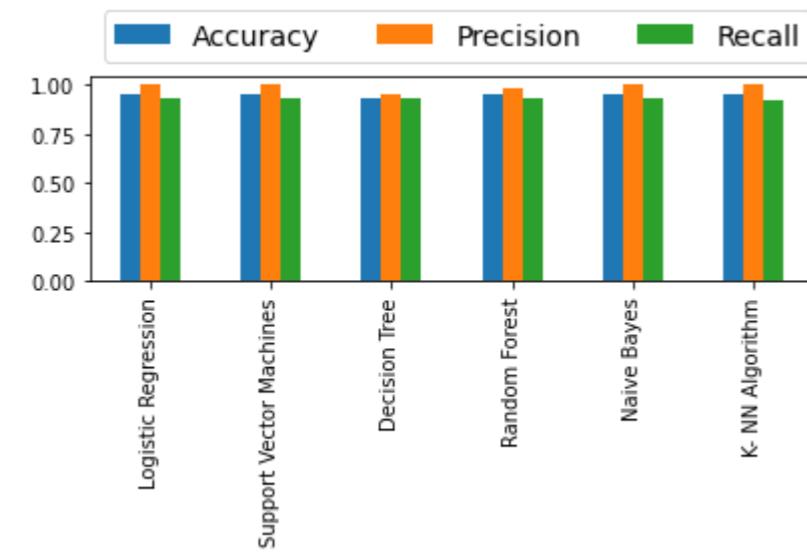
```
In [39]: 1 import pandas as pd
2
3 df_model = pd.DataFrame(index = models.keys(),
4                         columns = ['Accuracy', 'Precision', 'Recall'])
5 df_model['Accuracy'] = accuracy.values()
6 df_model['Precision'] = precision.values()
7 df_model['Recall'] = recall.values()
8
9 df_model
```

Out[39]:

	Accuracy	Precision	Recall
Logistic Regression	0.958042	1.000000	0.935484
Support Vector Machines	0.958042	1.000000	0.935484
Decision Tree	0.930070	0.954023	0.932584
Random Forest	0.951049	0.988506	0.934783
Naive Bayes	0.958042	1.000000	0.935484
K- NN Algorithm	0.951049	1.000000	0.925532

In [45]:

```
1 # Visualization of Scores:  
2  
3  
4 ax = df_model.plot.bar()  
5 ax.legend(  
6     ncol = len(models.keys()),  
7     bbox_to_anchor = (0, 1),  
8     loc = 'lower left',  
9     prop = {'size' : 14}  
10 )  
11  
12 plt.tight_layout()
```

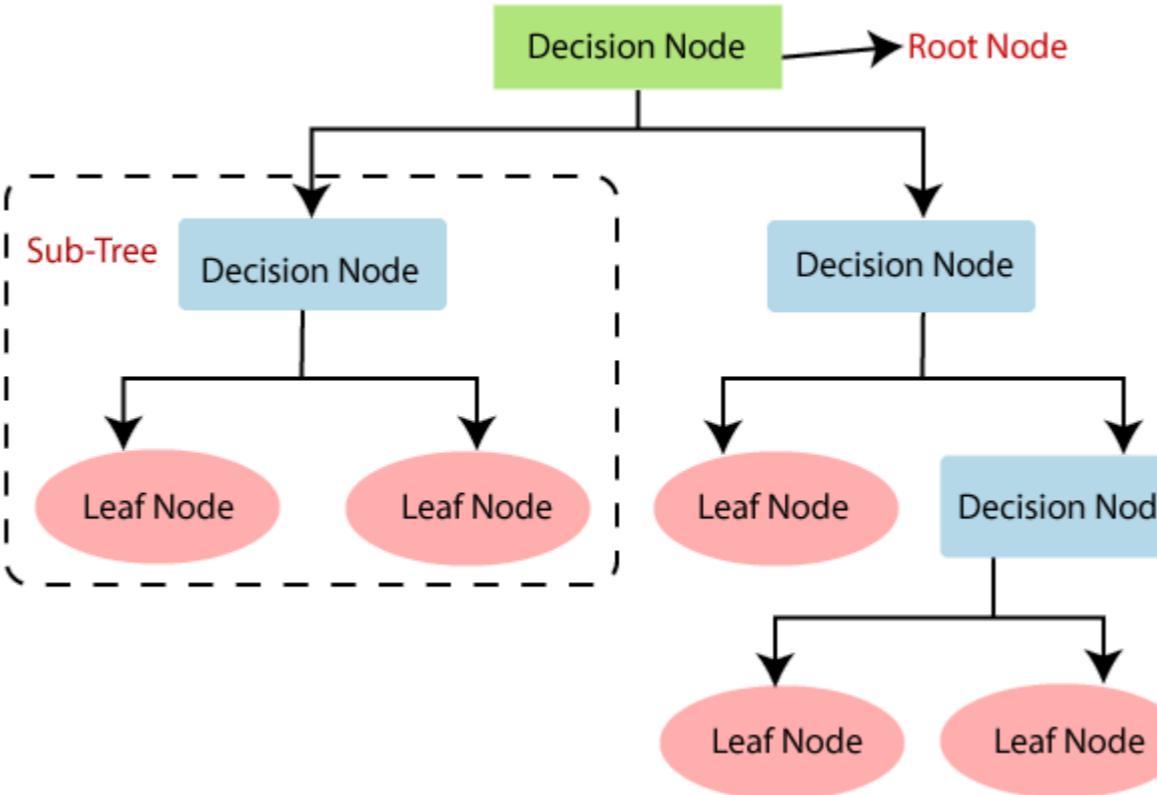


Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.

The logic behind the decision tree can be easily understood because it shows a tree-like structure.



Decision Tree Terminologies

Root Node: Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.

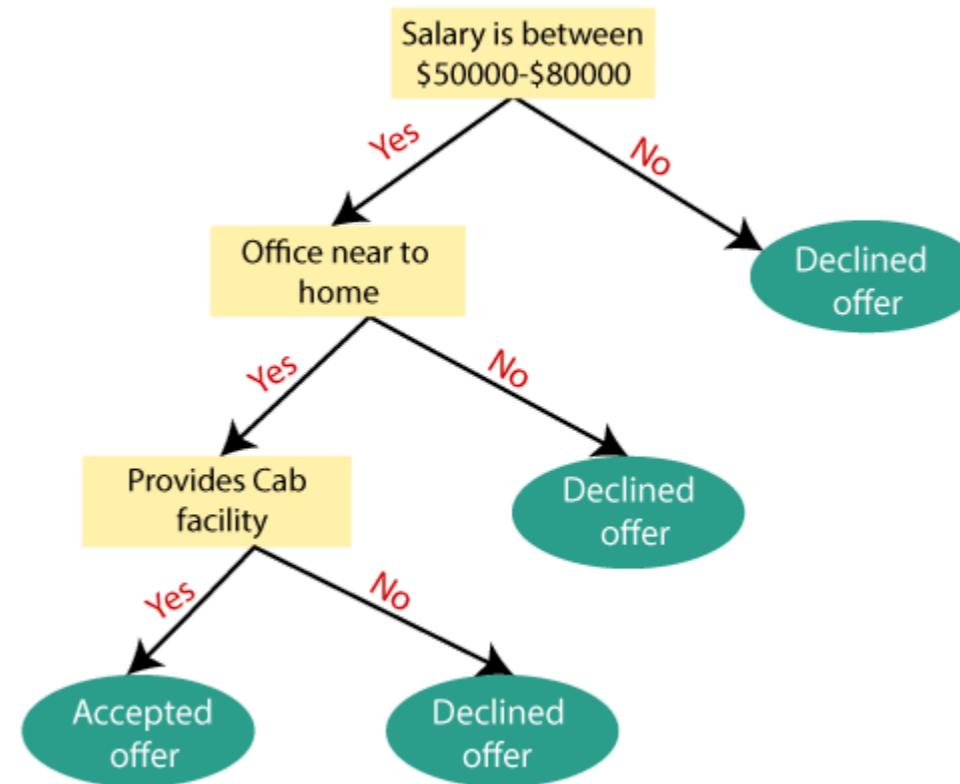
Leaf Node: Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.

Splitting: Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.

Branch/Sub Tree: A tree formed by splitting the tree.

Pruning: Pruning is the process of removing the unwanted branches from the tree. **Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.

Example: Suppose there is a candidate who has a job offer and wants to decide whether he should accept the offer or Not. So, to solve this problem, the decision tree starts with the root node (Salary attribute by ASM). The root node splits further into the next decision node (distance from the office) and one leaf node based on the corresponding labels. The next decision node further gets split into one decision node (Cab facility) and one leaf node. Finally, the decision node splits into two leaf nodes (Accepted offers and Declined offer). Consider the below diagram:



Advantages of the Decision Tree

- It is simple to understand as it follows the same process which a human follow while making any decision in real-life.
- It can be very useful for solving decision-related problems.
- It helps to think about all the possible outcomes for a problem.
- There is less requirement of data cleaning compared to other algorithms.

Disadvantages of the Decision Tree

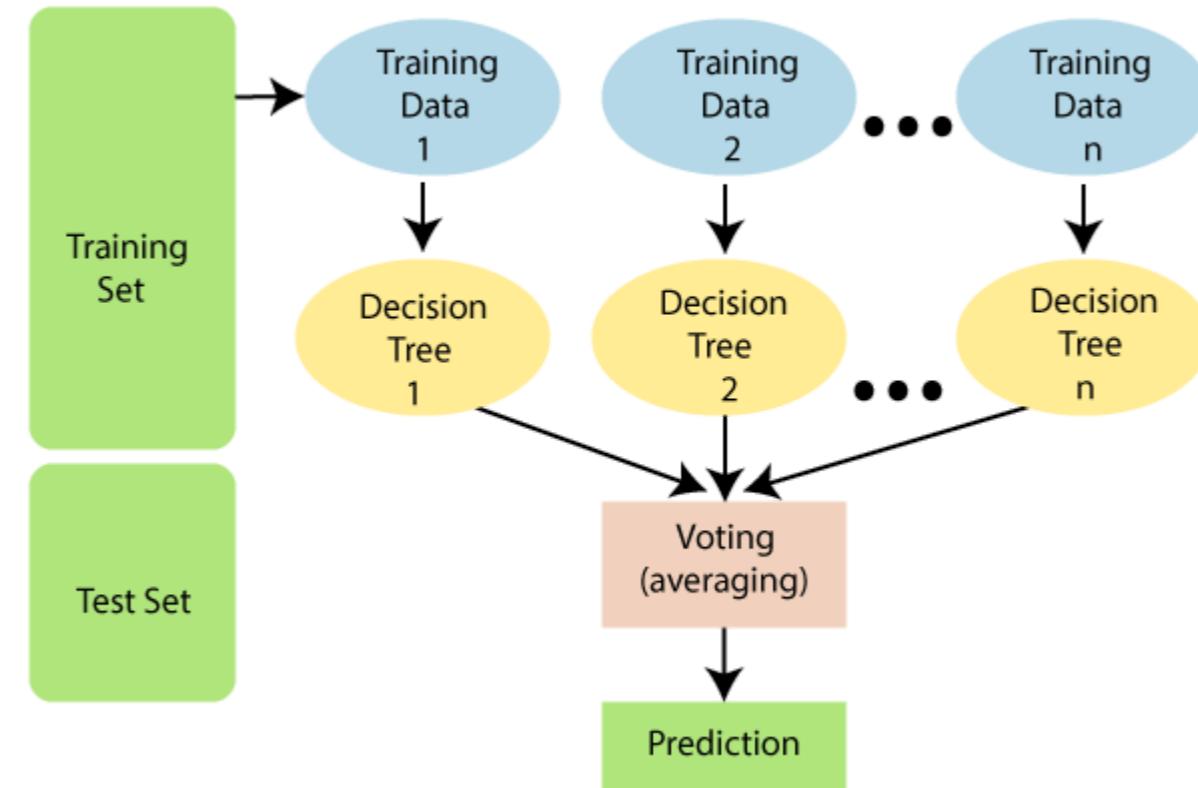
- The decision tree contains lots of layers, which makes it complex.
- It may have an overfitting issue, which can be resolved using the Random Forest algorithm.
- For more class labels, the computational complexity of the decision tree may increase.

Random Forest Algorithm

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.



Applications of Random Forest

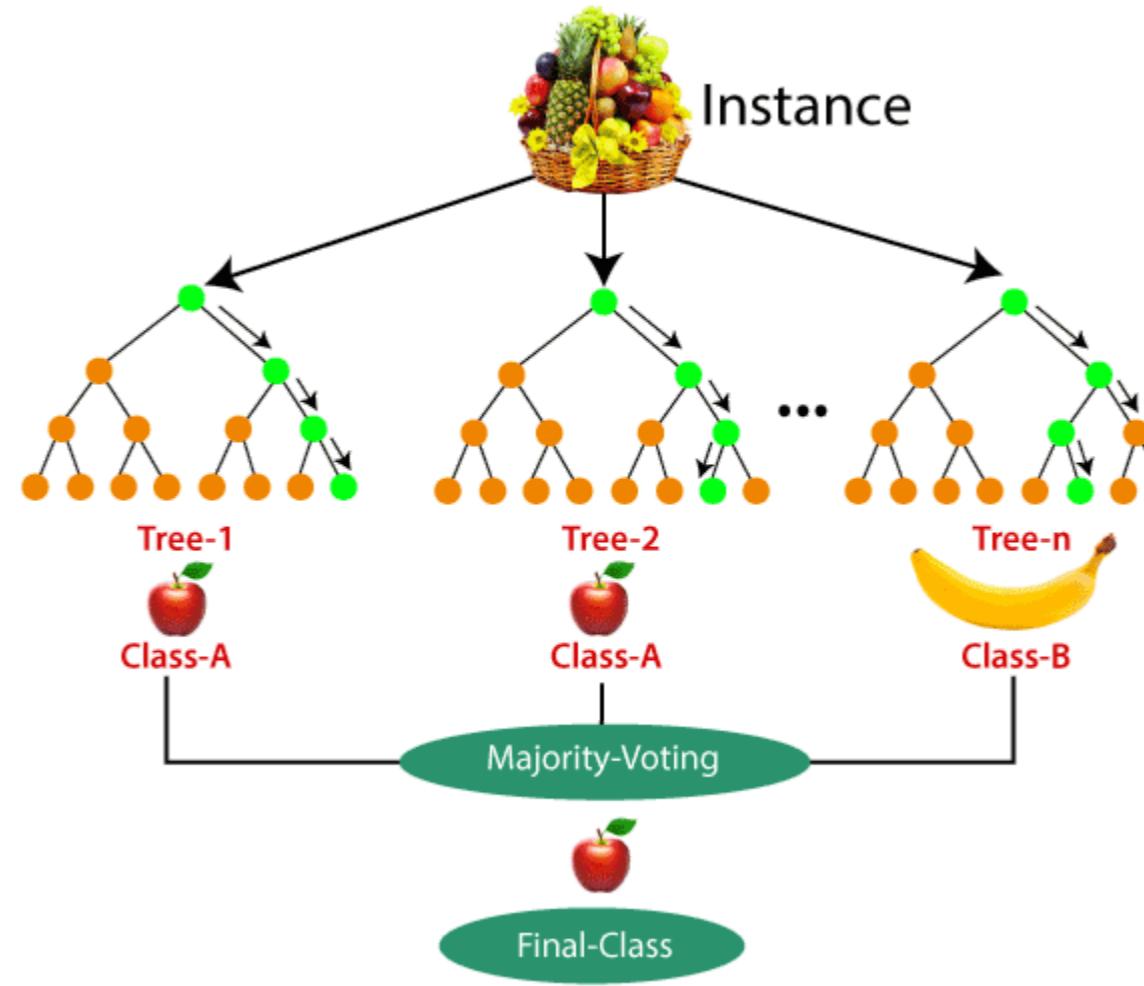
There are mainly four sectors where Random forest mostly used:

Banking: Banking sector mostly uses this algorithm for the identification of loan risk.

Medicine: With the help of this algorithm, disease trends and risks of the disease can be identified.

Land Use: We can identify the areas of similar land use by this algorithm.

Marketing: Marketing trends can be identified using this algorithm.

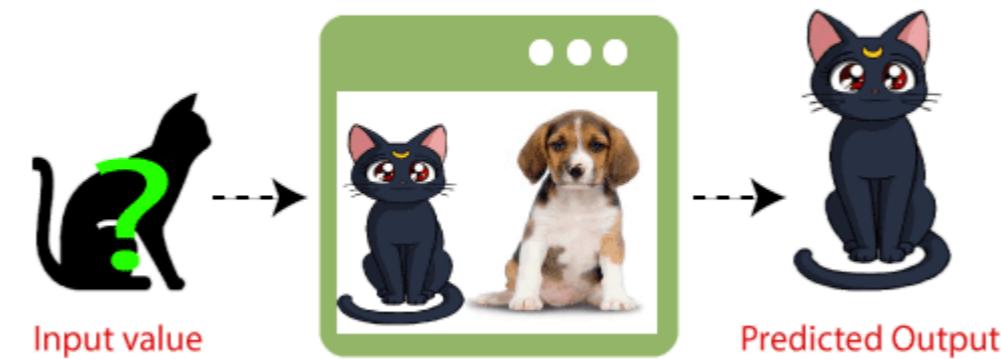


K-Nearest Neighbor(KNN) Algorithm for Machine Learning

- K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
- K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.
- K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suited category by using K- NN algorithm.
- K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.
- K-NN is a non-parametric algorithm, which means it does not make any assumption on underlying data.
- It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

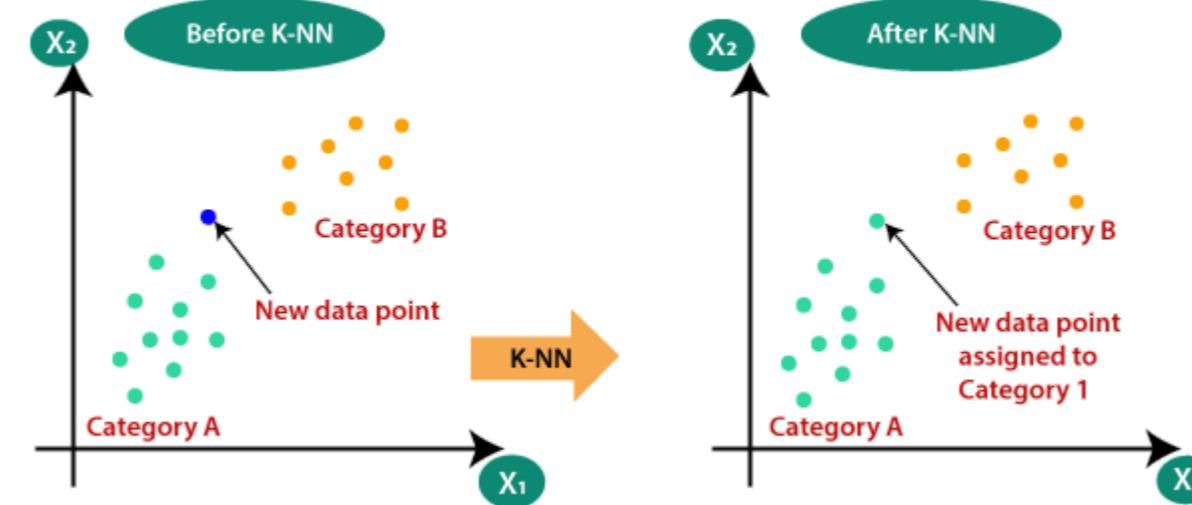
Example: Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.

KNN Classifier



Why do we need a K-NN Algorithm?

Suppose there are two categories, i.e., Category A and Category B, and we have a new data point x_1 , so this data point will lie in which of these categories. To solve this type of problem, we need a K-NN algorithm. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:



How does K-NN work? The K-NN working can be explained on the basis of the below algorithm:

Step-1: Select the number K of the neighbors

Step-2: Calculate the Euclidean distance of K number of neighbors

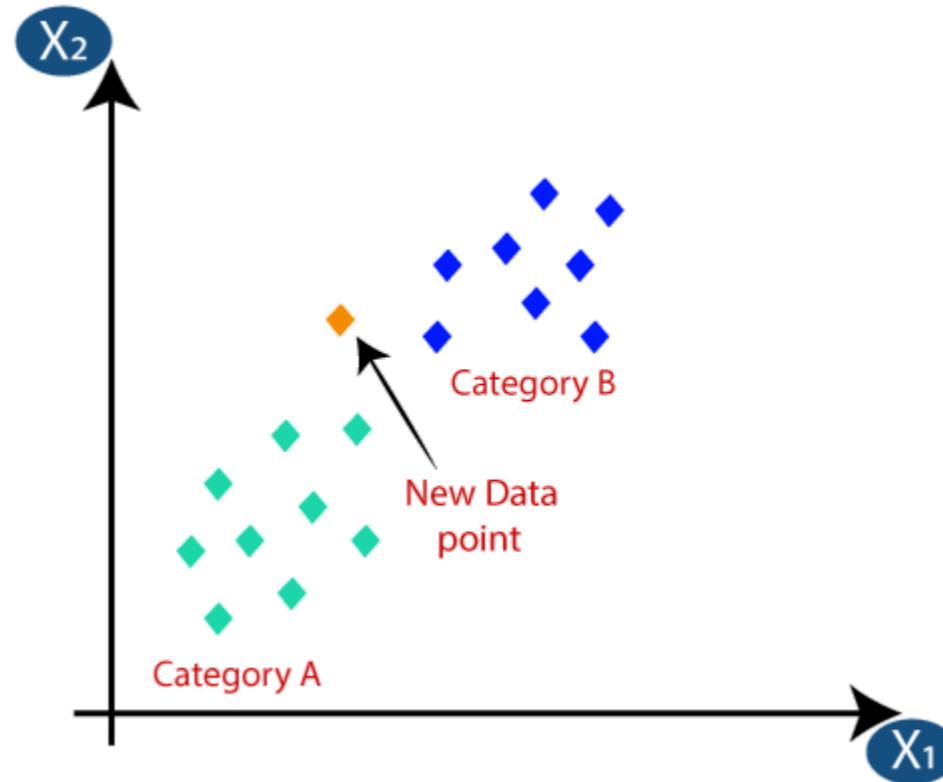
Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.

Step-4: Among these k neighbors, count the number of the data points in each category.

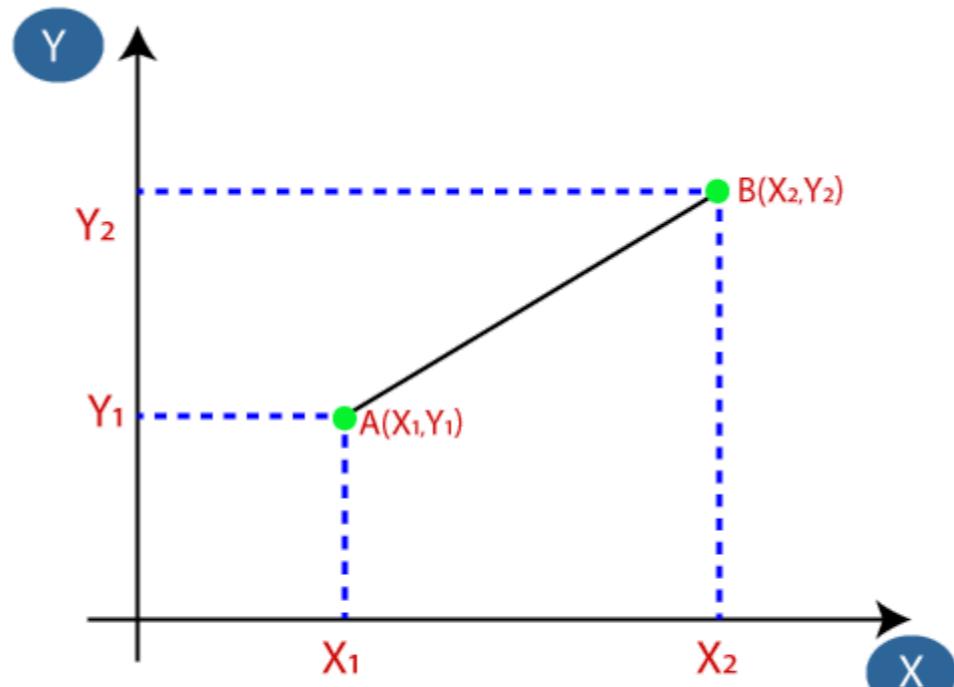
Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.

Step-6: Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:

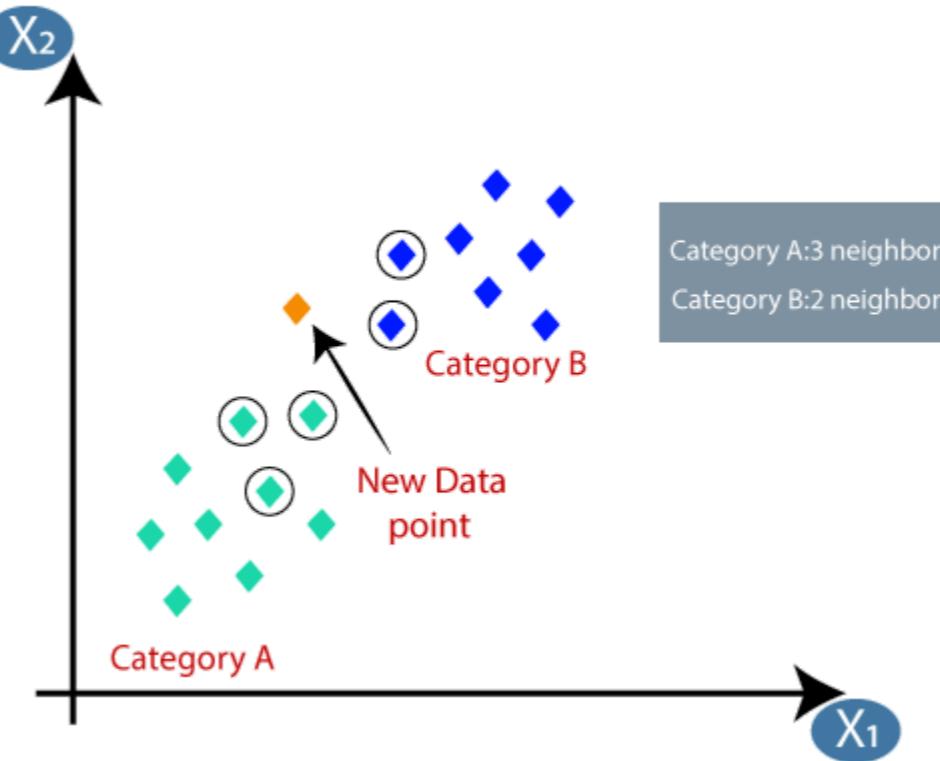


- Firstly, we will choose the number of neighbors, so we will choose the k=5.
- Next, we will calculate the Euclidean distance between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:



$$\text{Euclidean Distance between } A_1 \text{ and } B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:



How to select the value of K in the K-NN Algorithm?

- There is no particular way to determine the best value for "K", so we need to try some values to find the best out of them. The most preferred value for K is 5.
- A very low value for K such as K=1 or K=2, can be noisy and lead to the effects of outliers in the model.
- Large values for K are good, but it may find some difficulties.

Advantages of KNN Algorithm:

- It is simple to implement.
- It is robust to the noisy training data
- It can be more effective if the training data is large.

Disadvantages of KNN Algorithm:

- Always needs to determine the value of K which may be complex some time.
- The computation cost is high because of calculating the distance between the data points for all the training samples.

Naive Bayes Classifier Algorithm

- Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.
- It is mainly used in text classification that includes a high-dimensional training dataset. Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.
- It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

- Some popular examples of Naïve Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

Why is it called Naïve Bayes?

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

Naïve: It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

Bayes: It is called Bayes because it depends on the principle of Bayes' Theorem.

Bayes' Theorem:

Bayes' theorem is also known as Bayes' Rule or Bayes' law, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability. The formula for Bayes' theorem is given as:

Naïve Bayes Classifier Algorithm Where,

P(A|B) is Posterior probability: Probability of hypothesis A on the observed event B.

P(B|A) is Likelihood probability: Probability of the evidence given that the probability of a hypothesis is true.

P(A) is Prior Probability: Probability of hypothesis before observing the evidence.

P(B) is Marginal Probability: Probability of Evidence.

Working of Naïve Bayes' Classifier:

Working of Naïve Bayes' Classifier can be understood with the help of the below example:

Suppose we have a dataset of weather conditions and corresponding target variable "Play". So using this dataset we need to decide that whether we should play or not on a particular day according to the weather conditions. So to solve this problem, we need to follow the below steps:

- Convert the given dataset into frequency tables.
- Generate Likelihood table by finding the probabilities of given features.
- Now, use Bayes theorem to calculate the posterior probability.

Problem: If the weather is sunny, then the Player should play or not?

Solution: To solve this, first consider the below dataset:

Advantages of Naïve Bayes Classifier:

- Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.
- It can be used for Binary as well as Multi-class Classifications.
- It performs well in Multi-class predictions as compared to the other Algorithms.
- It is the most popular choice for text classification problems.

Disadvantages of Naïve Bayes Classifier:

- Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

Applications of Naïve Bayes Classifier:

- It is used for Credit Scoring.
- It is used in medical data classification.
- It can be used in real-time predictions because Naïve Bayes Classifier is an eager learner.
- It is used in Text classification such as Spam filtering and Sentiment analysis.

In []: ► 1