

An
Industrial Training Report
On
Linear Regression with Machine Learning
Submitted in Partial Fulfilment
For the award of the
Degree of Bachelor of Technology
In Department of Computer Science Engineering



Session 2023-2024

Submitted To:
Dr. Subash Chandra
HOD, CSE

Submitted By:
Ashita
20ERWCS011

Department of Computer Science Engineering
Rajasthan College of Engineering for Women

DECLARATION

I hereby declare that the discussion entitled “Linear Regression with Machine Learning” being submitted by me toward the partial fulfilment of the Degree of Bachelor, In Department of Computer Science Engineering is a project work carried out by me. Under the supervision of Mr. Vinod Todwal, and haven’t been submitted anywhere else. I will be solely responsible if any kind of plagiarisms is found.

Ashita

20ERWCS011

Rajasthan College of Engineering for Women

Counter Signed by

Mr. Vinod Todwal

Certificate of Training

Ashita Rustagi

from Rajasthan College Of Engineering For Women has successfully completed a 6-week online training on **Machine Learning**. The training consisted of Introduction to Machine Learning, Data, Introduction to Python, Data Exploration and Pre-processing, Linear Regression, Introduction to Dimensionality Reduction, Logistic Regression, Decision Tree, Ensemble Models, and Clustering (Unsupervised Learning) modules.

In the final assessment, Ashita scored 68% marks.

We wish Ashita all the best for future endeavours.

A handwritten signature in black ink, appearing to read 'Javed'.

ACKNOWLEDGEMENT

I like to share my sincere gratitude toward all those who helps me in the completion of the project. During the project I have faced many challenges because of lack of experience but this training process helps me to get over from all the challenges and in the final compilation of my idea to shaped sculpture.

Minute aspects of the project work.

In the last I would like to thanks the management of Rajasthan College of Engineering for Women for Providing me such an opportunity to learn from these experiences.

I am also thankful to my friends and parents who have inspired me to face all the challenges and win over them in life.

TABLE OF CONTENTS

DECLARATION	ii
CERTIFICATE	iii
ACKNOWLEDGEMENT	iv
ABSTRACT	1
Chapter 1. INTRODUCTION	2
1.1 Background of Project	2
1.2 Introduction of Machine Learning	2
1.3 Introduction of Data	3
1.4 Introduction of Python	5
Chapter 2. LITERATURE REVIEW	6
2.1 Machine Learning	6
2.2 Objective of Machine learning	6
2.3 Scope of Machine Learning	6
2.4 Kinds of Machine Learning	7
2.4.1 Supervised Machine Learning	7
2.4.2 Unsupervised Machine Learning	9
2.4.3 Semi Supervised Machine learning	11
2.4.4 Reinforcement Machine learning	12
2.4.5 Deep Learning	12

Chapter 3 Design of Project	13
3.1 Importing of Python Libraries	13
3.2 Importing of Data	14
3.3 Operation on Training Data	17
3.4 Manipulation of Data	20
3.5 Representation of Data	23
3.6 Prediction from Training Data	30
3.7 Introduction to Dimensionality Reduction	32
3.8 Decision Tree	33
3.9 Naïve Bayes Classification Theorem	34
3.10 Conclusion	35
3.11 References	35

TABLE OF FIGURES

Types of Machine Learning	3
Difference Between Qualitative and Quantitative Data	4
Linear Relationship between Dependent and Independent Variables	8
Logistics Regression with Multiple Variables	8
Regression vs Classification	9
K Means Clustering	10
Hierarchical Clustering	11
Importing of Python Libraries	13
Importing of Data	15
Output of Head Operation	17
Output of Tail Operation	18
Output of Shape Operation	18
Output of Describe and Info Operation	19
Dealing and Deleting Missing Values	20
Updated Info Operation	21
Value Count Operation	22
Relation Between Various Components of Data	23
Box Plot Representation	26
Correlation Between Variables	27
Multicollinearity	28
Dimensional Reduction	31
Decision Tree	32

ABSTRACT

The linear Regression is one of the simplest and most widely used Machine Learning algorithm. This algorithm is a kind of Supervised Learning which deals with the labeled data to be analyzed and make the prediction.

Sir Francis Galton is the one who first proposed the concept of Linear Regression in 1894.

Linear Regression works with two kinds of variables x which is dependent variable and y is independent Variable. Linear Regression provides the linear relationship between them. Linear Regression used the Least Square method to determine the best fit for every specific part of the data.

The Linear Regression has two types simple and multiple regression. This type of regression is applicable or applied only with the cleansed data by means of data mining theory.

Machine Learning is a collection of algorithms which not only manage the data but also in finding new patterns to find out the new and best opportunities which is beneficial for the organizations or the company. Machine Learning is a branch of AI that are mostly concerns with the current or the most recent updates of the data of any particular sector.

INTRODUCTION**1.1 BACKGROUND OF PROJECT**

Data Management on a large scale or we can say management of big data is very complex by simple techniques thus Machine Learning is used for the evaluation, analyzing and making prediction with the help of advanced algorithms and technologies which is physically impossible.

Machine Learning is simply all about data and its analyzing to make predictions and organizing the data in a proper format. Every change or the information which is related to any specific entity or topic is data for it. So, a million of byte of data is inserting, deleting and updating in every second all over the world as all the data are connected to each other through internet.

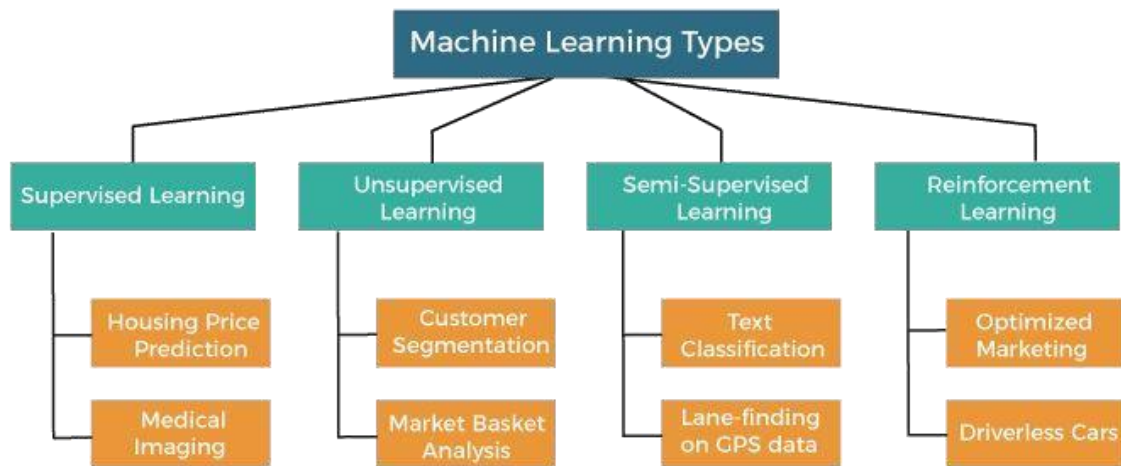
1.2 INTRODUCTION OF MACHINE LEARNING

Machine Learning is a branch of AI that enables computer to self-learn from the training data and improve over time without being explicitly programmed. Machine Learning is a programming computer to optimize a performance criterion using example data and past experiences. Machine Learning system builds the prediction model for analyzing based on algorithm applied.

The main purpose of Machine Learning is not only to train the model on previous data but also make the prediction value based on it for the best outcome. Machine Learning also helps in discovering the new patterns in the data. Machine Learning uses different kinds of algorithm to parse the data and it also helps in decision making. Machine Learning used in recognizing the practical benefits to the individuals in the real world.

There are mainly four types of Learning algorithms in Machine Learning

1. Supervised Machine Learning Algorithm
2. Unsupervised Machine Learning Algorithm
3. Semi Supervised Machine Learning Algorithm
4. Reinforcement Machine Learning Algorithm



1.3 INTRODUCTION OF DATA

Data is a raw collection of facts and figures. Data is basically any collection of different kind of facts in the forms such as numeric, alphabetical, relations, diagram, text, images, flowcharts and many more. Data is the main factor of Machine Learning on which almost all the operations and algorithms are being applied.

The unprocessed or unorganized information related to any specific is known as the data for that specific. As Machine Learning is used for prediction. So, as data is important for prediction, decision making, Problem solving, improving process to get the best outcome for the future.

Data is generally represented into two forms, which are:

1. Qualitative Data

Qualitative Data is descriptive which can be observed but not measured.

2. Quantitative Data

Quantitative Data is any data that can be counted or measured

WHAT'S THE DIFFERENCE BETWEEN QUANTITATIVE AND QUALITATIVE DATA?

Quantitative Data

- Countable or measurable, relating to numbers.
- Tells us how many, how much, or how often.
- Fixed and universal, "factual."
- Gathered by measuring and counting things.
- Analyzed using statistical analysis.

Qualitative Data

- Descriptive, relating to words and language.
- Describes certain attributes, and helps us to understand the "why" or "how" behind certain behaviors.
- Dynamic and subjective, open to interpretation.
- Gathered through observations and interviews.
- Analyzed by grouping the data into meaningful themes or categories.

1.4 INTRODUCTION OF PYTHON

Python was developed by Guido Van Rossum in 1991.

Python is a general-purpose language.

Python is a high-level interpreted language.

The most recent version of Python is Python3.

Python is object-oriented language.

Python is simple as English in terms of reading and understanding.

There are different python libraries that are used during the project such as numpy, pandas, Matplotlib.pyplot, seaborn and many more.

LITERATURE REVIEW

2.1 MACHINE LEARNING

Machine Learning is simply a branch of AI which uses different kinds of algorithm to maintain and organize a very vast amount of data in a proper format. Machine Learning builds different models according to algorithm used which help the user in decision making and predictions. Machine Learning allows the data to self-learn.

2.2 OBJECTIVES OF MACHINE LEARNING

The main purpose of Machine Learning is to train the data for the model to predict the value of some quantity which gives the best outcome. The main purpose of Machine Learning also includes is to discover the new pattern in your data which leads you toward decision making.

2.3 SCOPE OF MACHINE LEARNING

Machine Learning algorithms are used in almost all the fields. Some of the application of the Machine Learning are:

1. Spam Filtering
2. Fraud Detection
3. Smart Healthcare System
4. Speech Recognition
5. Computer Vision
6. Smart Transportation

2.4 KINDS OF MACHINE LEARNING

Machine Learning is a collection of different algorithms which are classified into different types of Machine Learning.

2.4.1 SUPERVISED MACHINE LEARNING

Supervised Machine Learning algorithm is designed to train the data for the model to get the desired output. The algorithm measures its accuracy through the loss function and manipulates until the error has been sufficiently minimized.

There are mainly two types of Machine Learning algorithms are:

A. REGRESSION

Regression is a kind of supervised algorithm which deals with the numerical type data. Regression is a technique which is used to predict the continuous data. Regression algorithm is also of two kinds. They are:

A.1 LINEAR REGRESSION

Linear Regression is a model used to determine the regression analysis. Regression analysis tells us the relationship between dependent and independent variables. Linear Regression only support the numerical type data. The regression is statistical in nature. It is used for predictive analysis. Linear Regression can be calculated with the equation:

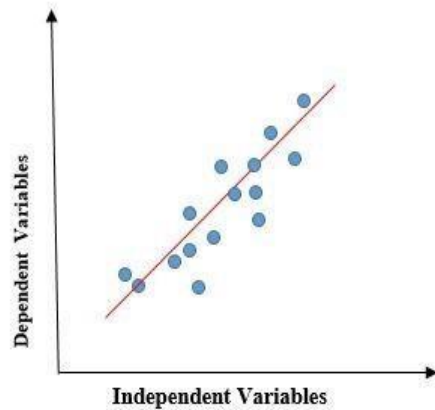
$$y = mx + c$$

y=dependent variable

x= independent variable

c= constant and

m = slope



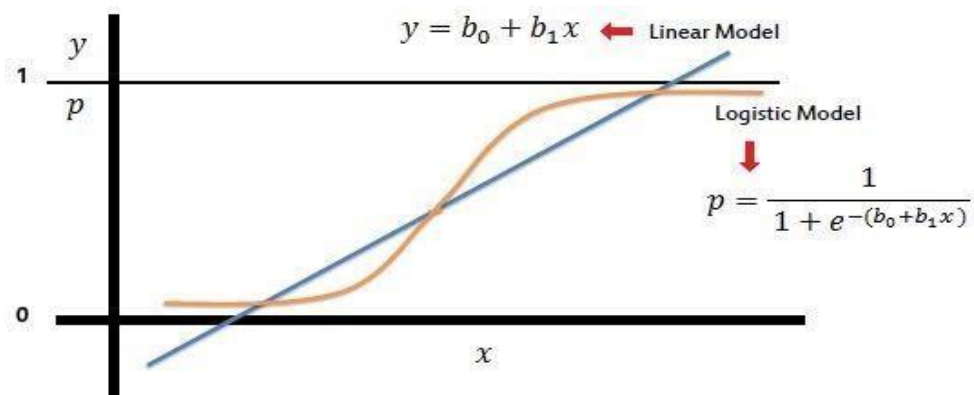
A.2 LOSS FUNCTION FOR LINEAR REGRESSION

The Loss Function is the difference between the predicted value and the actual value. It is the mean squared error value between them. This function is also known as cost function.

A.3 LOGISTIC REGRESSION

Logistic Regression is used for categorical variables. The output of Logistic Regression is either 0 or 1. This algorithm is used to solve the classification problem. The Logistic Regression equation is:

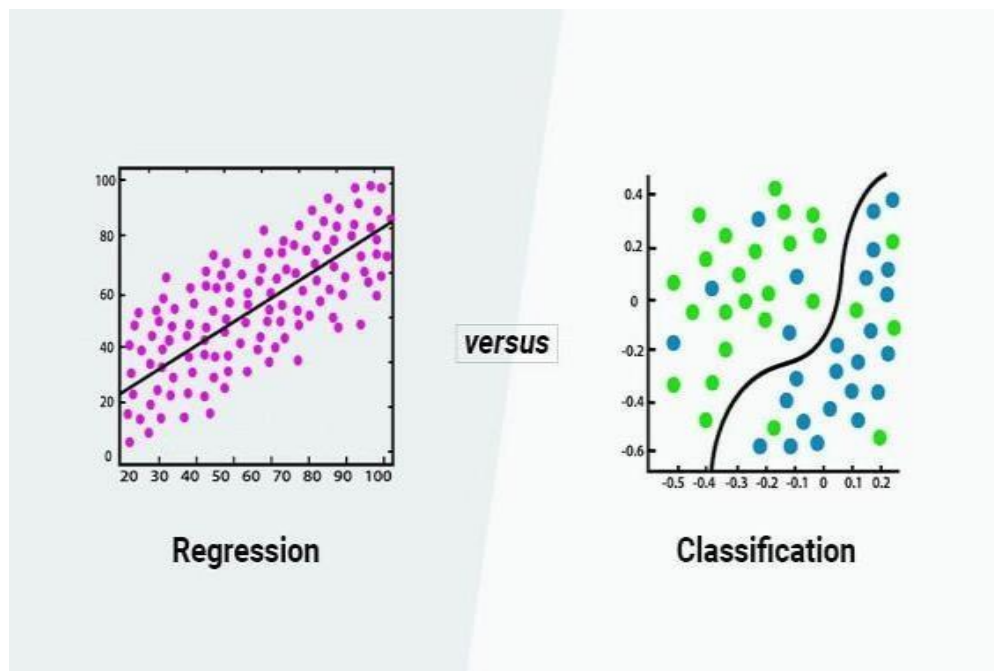
$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$



B. CLASSIFICATION

In Classification, the model is built on training data but is also applied on test data before the use in general. Classification is a technique in which the data having same characteristics are placed in one group and the prediction is done accordingly on unseen data. The main algorithm that comes under the classification are:

Naïve Bayes, K Nearest Neighbor, Decision Tree, Support Vector Machine, Random Forest and so on.



2.4.2 UNSUPERVISED MACHINE LEARNING

Unsupervised learning is a type of Machine Learning that learns pattern from unlabeled data. It is used to find hidden pattern. There are mainly two types of Unsupervised Learning:

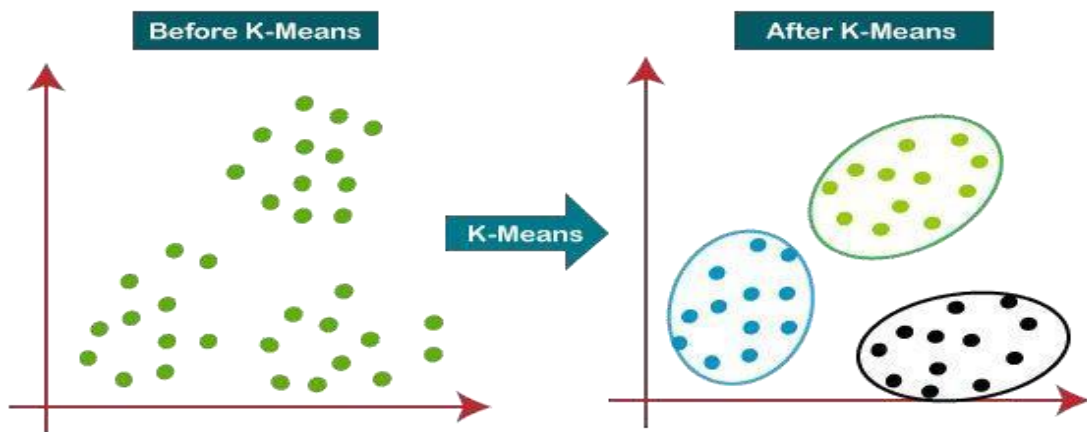
A. CLUSTERING

Clustering is a process in which the data having the same characteristics are placed into the same group with the help of different algorithms. They are:

A.1 K MEANS CLUSTERING ALGORITHM

K Means Clustering algorithm is used to group the unlabeled data into different clusters. In this algorithm K defines the number of predefined values of clusters that are needed to be create in the process. This algorithm is also called Centre Base or Centroid algorithm. In this algorithm we have to find the Euclidean Distance for every new data from the Centroid to decide in which cluster the data is placed. The formula for Euclidean Distance is:

$$K=((x-x1)^2 + (y-y1)^2)^{1/2}$$

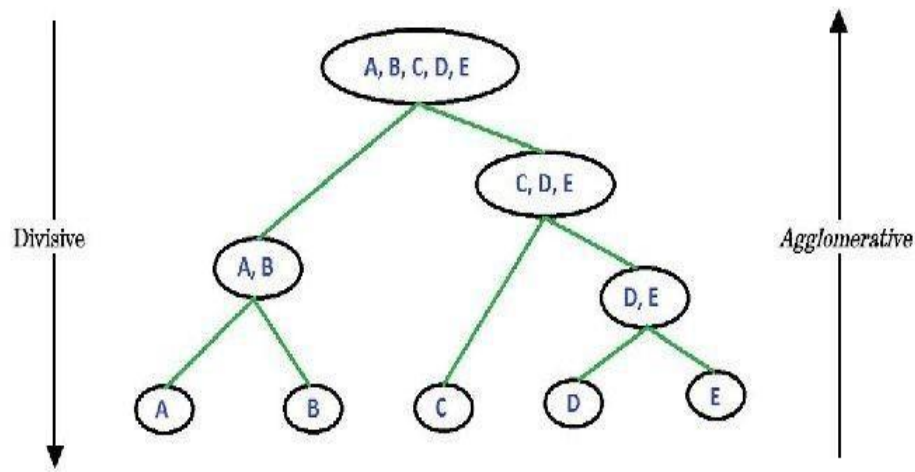


A.2 HIERARCHICAL CLUSTERING

Hierarchical Clustering is a type of clustering in which the algorithm is followed by the data as inherit of the characteristics of the data. There are mainly two types of Hierarchical clustering. They are:

Agglomerative

Divisive



B. ASSOCIATIVE RULE MINING

Associative Rule Mining is a kind of Unsupervised Learning algorithm which is used to check the dependency of one data element on another data element. The two main factor of this algorithm is:

B.1 SUPPORT

Support is a term which is defined as the percentage of the dataset in which a particular set of values is considered at a time.

B.2 CONFIDENCE

The percentage of transactions that contain a particular item or a set of items.

2.4.3 SEMI SUPERVISED MACHINE LEARNING

Semi Supervised learning is a collection of supervised and unsupervised learning. This type of learning is applicable for labeled and unlabeled data to train the model.

2.4.4 REINFORCEMENT MACHINE LEARNING

Reinforcement learning is a type of Machine Learning which allows machine and software agents to automatically determine the ideal behavior within a specific context in order to maximize its performance.

2.4.5 DEEP LEARNING

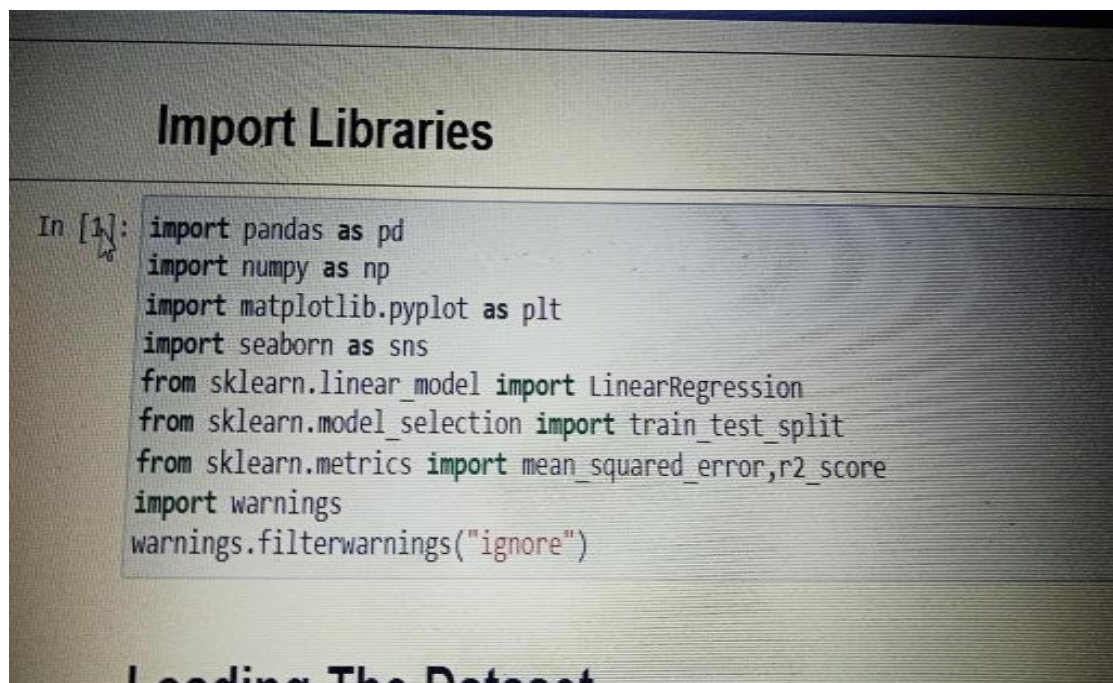
Deep Learning is a type of Machine Learning which guides the computer to work or process the data as a human brain does. Deep Learning have advance algorithm and features in it which help this learning to recognize the complex pictures, audio, video, text, graph and any other data source that lead to the better result and better prediction for future. Some application of Deep learning is as follow:

1. Digital Assistants
2. Fraud Detection
3. Automatic Face Recognition
4. Voice Activated Controller

DESIGN OF PROJECT

The project in the training is of Linear Regression in Machine Learning. This project is Completed with stepwise process which include many steps. The project of Linear Regression is about how the data is managed to make the right decision as well as Prediction for the best.

3.1 IMPORTING PYTHON LIBRARIES



PANDAS

Pandas is a python library which is used in data analysis, manipulation and its cleansing. Pandas supports operations like sorting, re-indexing and concatenation.

NUMPY

Numpy is a python library which is used for numerical data. This library supports large Metrics and multi-dimensional data. This library consists of in-built mathematical function For easy computations.

MATPLOTLIB.PYPILOT

This library is used for analysis of data. Mat plot library helps to plot the data in the graph representation and helps us to determine the relationship between different elements.

SEABORN

Seaborn library is a type of python library which provide a large interface for drawing attractive and informative statistical graphics.

SKLEARN

Sklearn is a python library used for the implementation of Machine Learning Model and Statistical model.

3.2 IMPORTING OF DATA

The data is imported by its address. The data is collected from different sources to analyze the real status of the environment and provides the appropriate suggestion.

WPS Home QualityAir

Home Insert Page Layout Formulas Data Review View Tools Smart Toolbox

Format Painter Bold Italic Underline Text Color Fill Color Orientation Merge and Center General Rows and Columns Conditional Formatting AutoSum AutoFilter

Last chance! Black Friday ending soon—Buy one year of WPS Pro, get one year FREE! Don't miss out! Get Now

J4 fx

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	Place	PM2.5	PM10	SO2	NO2	O3	CO	AQI	Quality														
2	Haryana	92.1	59.8	11.3	39.8	6	3.8	92	Moderate														
3	Delhi	27.9	18.1	3.3	53.9	7.7	3	53	Moderate														
4	Uttar Prad	119.2	100.2	11.5	52.3	8.1	4.9	119	Unhealthy														
5	Karnataka	25.3	20.8	3.8	48.1	10.7	1.9	48	Good														
6	Kerala	146.9	116.2	11	53.4	6	4.8	146	Unhealthy														
7	Madhya Pr	159.2	123.1	12.1	52.3	5	4.9	159	Unhealthy														
8	Maharashtra	37.4	25.8	18.3	28.3	4.7	2.5	37	Good														
9	Andhra Pradesh	96.3	71.3	7.4	51	14.5	2.7	96	Moderate														
10	Arunachal	82.8	81.4	4.4	33.9	6.3	2.8	82	Moderate														
11	Assam	48.1	24	2.4	23.7	2.8	2.8	48	Good														
12	Bihar	20.7	20	3.8	38.2	4.7	2.4	38	Good														
13	Chhattisgarh	67	50.3	5.1	46.3	5.6	2.1	67	Moderate														
14	Goa	100	76.2	4.8	51.1	11	2.5	100	Moderate														
15	Gujarat	129.6	103.1	9.6	52.1	5.4	5.3	129	Unhealthy														
16	Himanchal	72	44.8	3.8	38.3	6	1.4	72	Moderate														
17	Jammu and	69.1	55.9	7.4	34.2	1.8	1.5	69	Moderate														
18	Jharkhand	34	30.5	3.8	31.4	4.6	2.6	34	Good														
19	Manipur	37	31.1	20.2	33.1	4.2	3.4	37	Good														
20	Meghalaya	41.1	38.1	11	26.7	3.8	1.5	41	Good														
21	Mizoram	124.4	98.8	12.6	52.3	5.9	5	124	Unhealthy														
22	Nagaland	134.7	110	10.4	54.2	10.4	4.2	134	Unhealthy														
23	Odisha	72.9	43.3	7	33.4	5.5	2.3	72	Moderate														
24	Punjab	25.3	20.8	3.8	48.1	10.7	1.9	48	Good														
25	Rajasthan	69.1	34	8.5	15	7	2.2	69	Moderate														
26	Sikkim	31.6	45.4	7.4	38.4	11.5	2.2	45	Good														
27	Tamil Nadu	46.1	27.9	5.8	7.3	14.2	5.3	46	Good														
28	Telangan	179.3	103.7	9.6	53.4	6.5	5.3	179	Unhealthy														

Sheet1

23°C Haze 2:20 PM 12/5/2023

WPS Qualityxloz

Home Insert Page Layout Formulas Data Review View Tools Smart Toolbox

Calibri 11 A A

General Rows and Columns Fill Sort Conditional Formatting AutoSum AutoFilter

Last chance! Black Friday ending soon—Buy one year of WPS Pro, get one year FREE! Don't miss out! Get Now

JS4 X ✓ fx

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
8	Maharashtra	37.4	25.8	18.3	28.3	4.7	2.5		37	Good													
9	Andhra Pradesh	96.3	71.3	7.4	51	14.5	2.7		96	Moderate													
10	Arunachal Pradesh	82.8	81.4	4.4	33.9	6.3	2.8		82	Moderate													
11	Assam	48.1	24	2.4	23.7	2.8	2.8		48	Good													
12	Bihar	20.7	20	3.8	38.2	4.7	2.4		38	Good													
13	Chhattisgarh	67	50.3	5.1	46.3	5.6	2.1		67	Moderate													
14	Goa	100	76.2	4.8	51.1	11	2.5		100	Moderate													
15	Gujarat	129.6	103.1	9.6	52.1	5.4	5.3		129	Unhealthy													
16	Himachal Pradesh	72	44.8	3.8	38.3	6	1.4		72	Moderate													
17	Jammu and Kashmir	69.1	55.9	7.4	34.2	1.8	1.5		69	Moderate													
18	Jharkhand	34	30.5	3.8	31.4	4.6	2.6		34	Good													
19	Manipur	37	31.1	20.2	33.1	4.2	3.4		37	Good													
20	Meghalaya	41.1	38.1	11	26.7	3.8	1.5		41	Good													
21	Mizoram	124.4	98.8	12.6	52.3	5.9	5		124	Unhealthy													
22	Nagaland	134.7	110	10.4	54.2	10.4	4.2		134	Unhealthy													
23	Odisha	72.9	43.3	7	33.4	5.5	2.3		72	Moderate													
24	Punjab	25.3	20.8	3.8	48.1	10.7	1.9		48	Good													
25	Rajasthan	69.1	34	8.5	15	7	2.2		69	Moderate													
26	Sikkim	31.6	45.4	7.4	38.4	11.5	2.2		45	Good													
27	Tamil Nadu	46.1	27.9	5.8	7.3	14.2	5.3		46	Good													
28	Telangana	129.3	103.7	9.6	53.4	6.5	5.3		129	Unhealthy													
29	Tripura	126.4	70.4	13.1	50.7	0.9	4.6		126	Unhealthy													
30	Uttarakhand	19.9	11.5	0.4	3.1	30	1		30	Good													
31	West Bengal	24.2	13	0.3	2.6	17.6	1.4		24	Good													
32																							
33																							
34																							
35																							

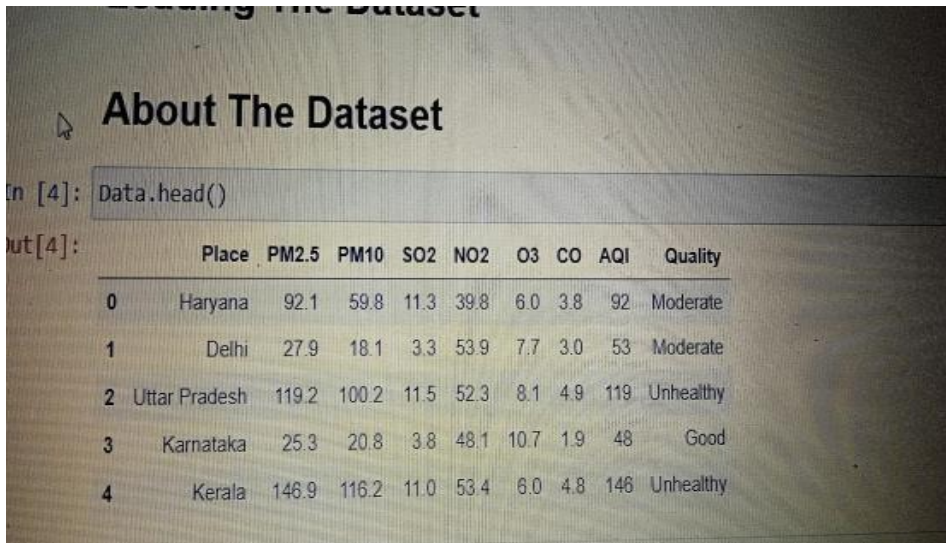
Sheet1

23°C Haze 2:20 PM 12/5/2023

3.3 OPERATION ON TRAINING DATA

3.3.1 HEAD OPERATION

Head Operation is used to show the topmost rows of the data according to the value entered by the user. The default value of head is 5.



About The Dataset

```
In [4]: Data.head()
```

Out[4]:

	Place	PM2.5	PM10	SO2	NO2	O3	CO	AQI	Quality
0	Haryana	92.1	59.8	11.3	39.8	6.0	3.8	92	Moderate
1	Delhi	27.9	18.1	3.3	53.9	7.7	3.0	53	Moderate
2	Uttar Pradesh	119.2	100.2	11.5	52.3	8.1	4.9	119	Unhealthy
3	Karnataka	25.3	20.8	3.8	48.1	10.7	1.9	48	Good
4	Kerala	146.9	116.2	11.0	53.4	6.0	4.8	146	Unhealthy

3.3.2 TAIL OPERATION

Tail Operation is used to display the bottommost rows of the data according to the value entered by the user. The default value for tail is 5.

[5]: Data.tail()

	Place	PM2.5	PM10	SO2	NO2	O3	CO	AQI	Quality
25	TamilNadu	46.1	27.9	5.8	7.3	14.2	5.3	46	Good
26	Telangana	129.3	103.7	9.6	53.4	6.5	5.3	129	Unhealthy
27	Tripura	126.4	70.4	13.1	50.7	0.9	4.6	126	Unhealthy
28	Uttarakhand	19.9	11.5	0.4	3.1	30.0	1.0	30	Good
29	West Bengal	24.2	13.0	0.3	2.6	17.6	1.4	24	Good

3.3.3 SHAPE OPERATION

Shape Operation is used to find out the numbers of rows and column in the data.

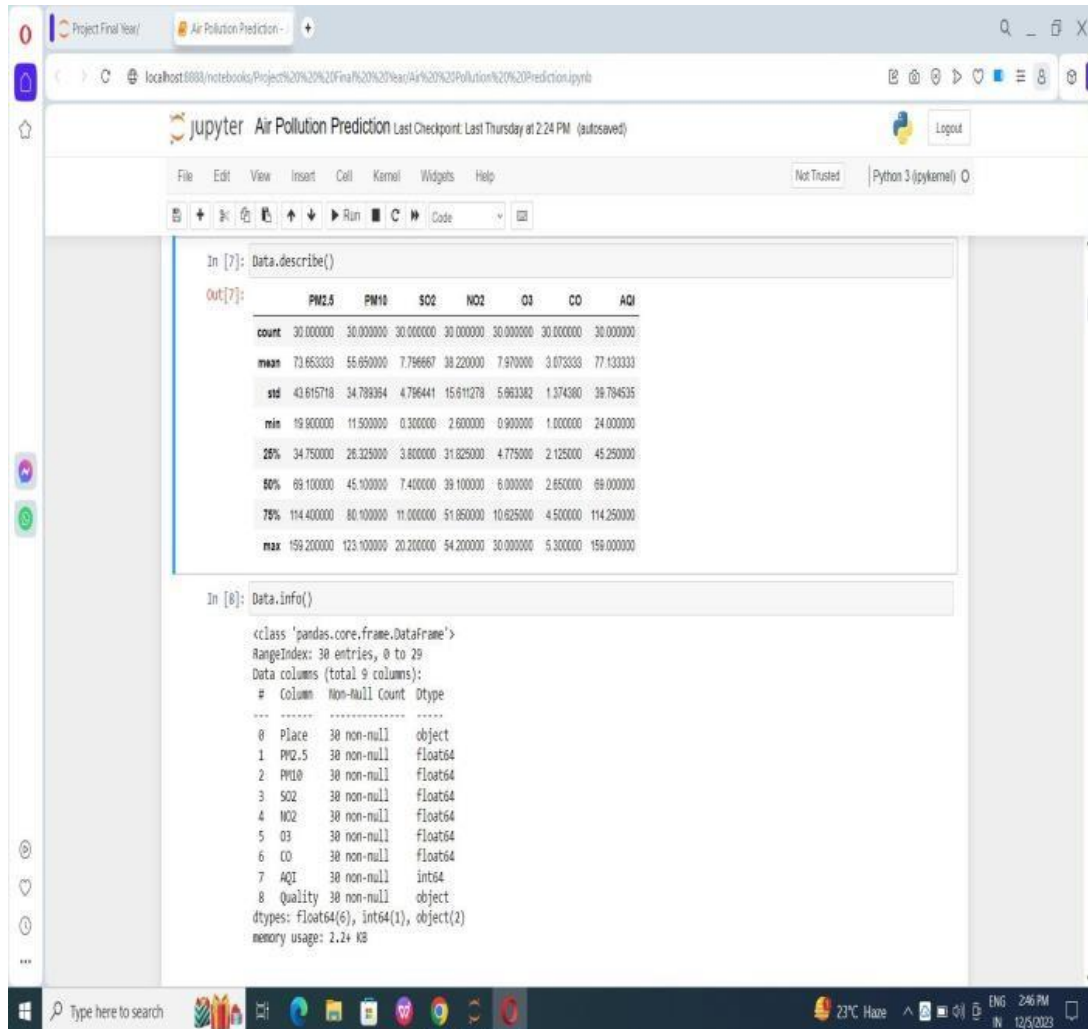
```
In [6]: Data.shape
Out[6]: (30, 9)
```

3.3.4 DESCRIBE OPERATION

Describe Operation helps in finding the standard value for each of the element of operation like mean, mode, minimum, maximum, standard deviation and many more.

3.3.5 INFO OPERATION

Info Operation helps the user to find out all the data types with the number of not null values for each and every row and column along with the size of datatypes.



The screenshot displays a Jupyter Notebook interface titled "Air Pollution Prediction". The notebook is running on a local host (localhost:8888) and shows two code cells. The first cell, labeled "In [7]:", contains the command `Data.describe()`. The output, labeled "Out[7]:", is a summary statistics table for the dataset. The second cell, labeled "In [8]:", contains the command `Data.info()`. The output shows the data structure, including the number of entries (30), the number of columns (9), and the data types for each column.

	PM2.5	PM10	SO2	NO2	O3	CO	AQI
count	30.000000	30.000000	30.000000	30.000000	30.000000	30.000000	30.000000
mean	73.653333	55.650000	7.796667	30.220000	7.970000	3.073333	77.133333
std	43.615718	34.789364	4.796441	15.611278	5.061382	1.074380	38.784535
min	19.900000	11.500000	0.300000	2.600000	0.900000	1.000000	24.000000
25%	34.750000	26.325000	3.800000	31.825000	4.775000	2.125000	45.250000
50%	69.100000	45.100000	7.400000	39.100000	6.000000	2.850000	69.000000
75%	114.400000	80.100000	11.000000	51.850000	10.625000	4.500000	114.250000
max	159.200000	123.100000	20.200000	54.200000	30.000000	5.300000	159.000000

The second code cell, labeled "In [8]:", contains the command `Data.info()`. The output shows the data structure, including the number of entries (30), the number of columns (9), and the data types for each column.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 9 columns):
#   Column  Non-null Count  Dtype
---  -
0   Place   30 non-null      object
1   PM2.5   30 non-null      float64
2   PM10    30 non-null      float64
3   SO2     30 non-null      float64
4   NO2     30 non-null      float64
5   O3      30 non-null      float64
6   CO      30 non-null      float64
7   AQI     30 non-null      int64
8   Quality 30 non-null      object
dtypes: float64(6), int64(1), object(2)
memory usage: 2.2+ KB
```

3.4 MANIPULATION OF DATA

3.4.1 IS NULL FUNCTION

The is null method return all the element having their values are replaced with True or False.

3.4.2 DROPNA FUNCTION

Dropna Function removes the cells which contains the null values.

The screenshot shows a web browser window displaying a Jupyter Notebook titled "Air Pollution Prediction". The notebook's file path is visible as `localhost:8888/notebooks/Project%20Final%20Notebook(Air%20Pollution%20Prediction).ipynb`. The interface includes standard Jupyter controls like "File", "Edit", "View", "Insert", "Cell", "Kernel", "Widgets", and "Help" menus. A toolbar contains icons for running code, saving, and other functions. Below the menu bar, there are tabs for "Not Trusted" and "Python 3 (ipykernel)". The main area displays a table with 26 rows representing Indian states. Each row has a state name followed by several columns of numerical data, most of which are "NaN".

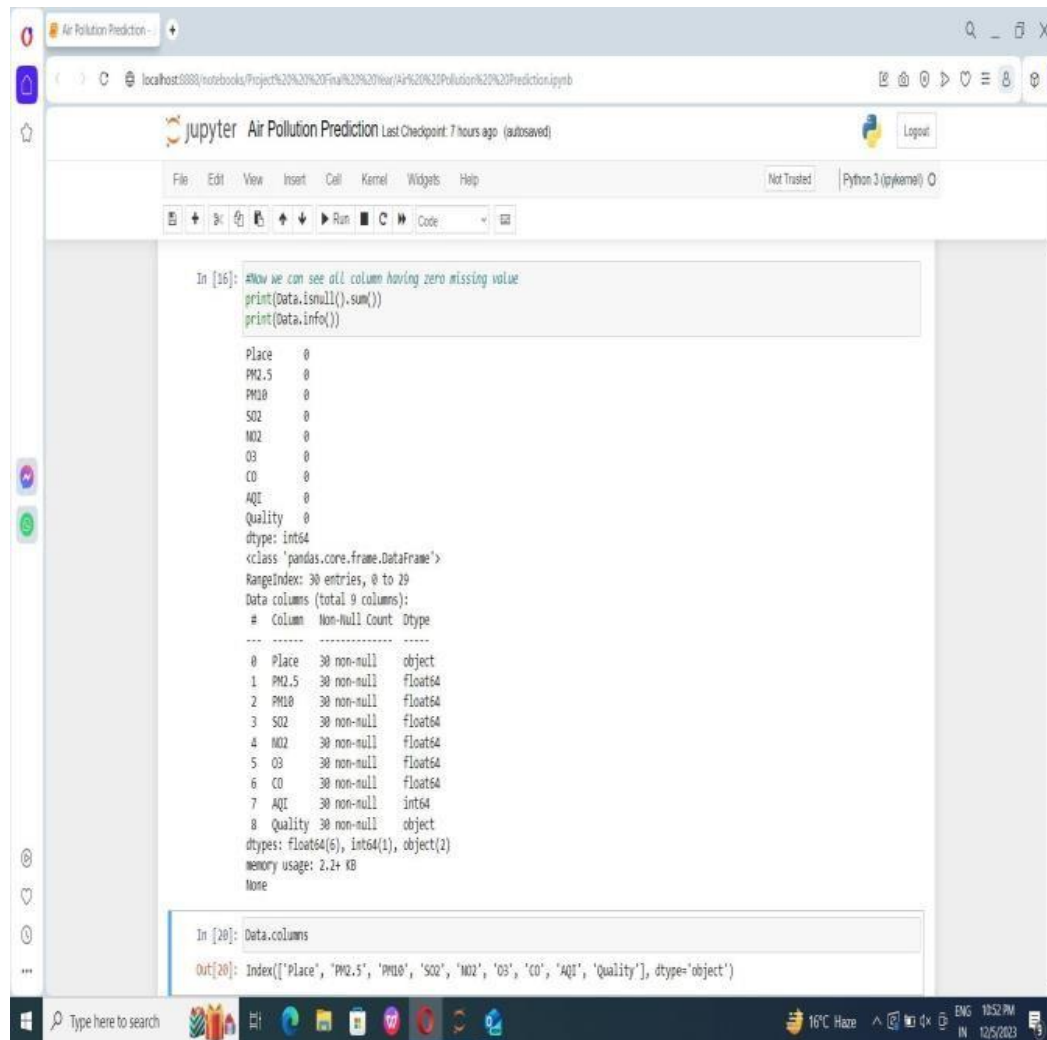
1	Arunachal Pradesh	66.1	NaN	NaN	NaN	NaN	NaN	NaN
2	Assam	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	Bihar	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	Chattisgarh	NaN	NaN	NaN	NaN	NaN	NaN	NaN
5	Delhi	NaN	NaN	NaN	NaN	NaN	NaN	NaN
6	Goa	NaN	NaN	NaN	NaN	NaN	NaN	NaN
7	Gujarat	NaN	NaN	NaN	NaN	NaN	NaN	NaN
8	Haryana	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9	Himachal Pradesh	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10	Jammu and Kashmir	NaN	NaN	NaN	NaN	NaN	NaN	NaN
11	Jharkhand	NaN	NaN	NaN	NaN	NaN	NaN	NaN
12	Karnataka	NaN	NaN	NaN	NaN	NaN	NaN	NaN
13	Kerala	NaN	NaN	NaN	NaN	NaN	NaN	NaN
14	Madhya Pradesh	NaN	NaN	NaN	NaN	NaN	NaN	NaN
15	Maharashtra	NaN	NaN	NaN	NaN	NaN	NaN	NaN
16	Manipur	NaN	NaN	NaN	NaN	NaN	NaN	NaN
17	Meghalaya	NaN	NaN	NaN	NaN	NaN	NaN	NaN
18	Mizoram	NaN	NaN	NaN	NaN	NaN	NaN	NaN
19	Nagaland	NaN	NaN	NaN	NaN	NaN	NaN	NaN
20	Odisha	NaN	NaN	NaN	NaN	NaN	NaN	NaN
21	Punjab	NaN	NaN	NaN	NaN	NaN	NaN	NaN
22	Rajasthan	NaN	NaN	NaN	NaN	NaN	NaN	NaN
23	Sikkim	NaN	NaN	NaN	NaN	NaN	NaN	NaN
24	Tamil Nadu	NaN	NaN	NaN	NaN	NaN	NaN	NaN
25	Telangana	NaN	NaN	NaN	NaN	NaN	NaN	NaN

3.4.3 INFO FUNCTION

Info Operation is used to determine the type of data and its size of each and every element of data.

3.4.4 COLUMN FUNCTION

Column function is used to find the number and names of the data.



The screenshot shows a Jupyter Notebook titled "Air Pollution Prediction" running on a local host. The notebook contains two code cells. The first cell, labeled "In [16]:", executes the following code:

```
#How we can see all column having zero missing value
print(Data.isnull().sum())
print(Data.info())
```

The output of the first cell is as follows:

```
Place      0
PM2.5      0
PM10       0
SO2        0
NO2        0
O3         0
CO         0
AQI        0
Quality    0
dtype: int64
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 9 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0   Place   30 non-null     object
 1   PM2.5   30 non-null     float64
 2   PM10    30 non-null     float64
 3   SO2     30 non-null     float64
 4   NO2     30 non-null     float64
 5   O3      30 non-null     float64
 6   CO      30 non-null     float64
 7   AQI     30 non-null     int64
 8   Quality 30 non-null     object
dtypes: float64(6), int64(1), object(2)
memory usage: 2.2+ KB
None
```

The second cell, labeled "In [20]:", executes the following code:

```
Data.columns
```

The output of the second cell is:

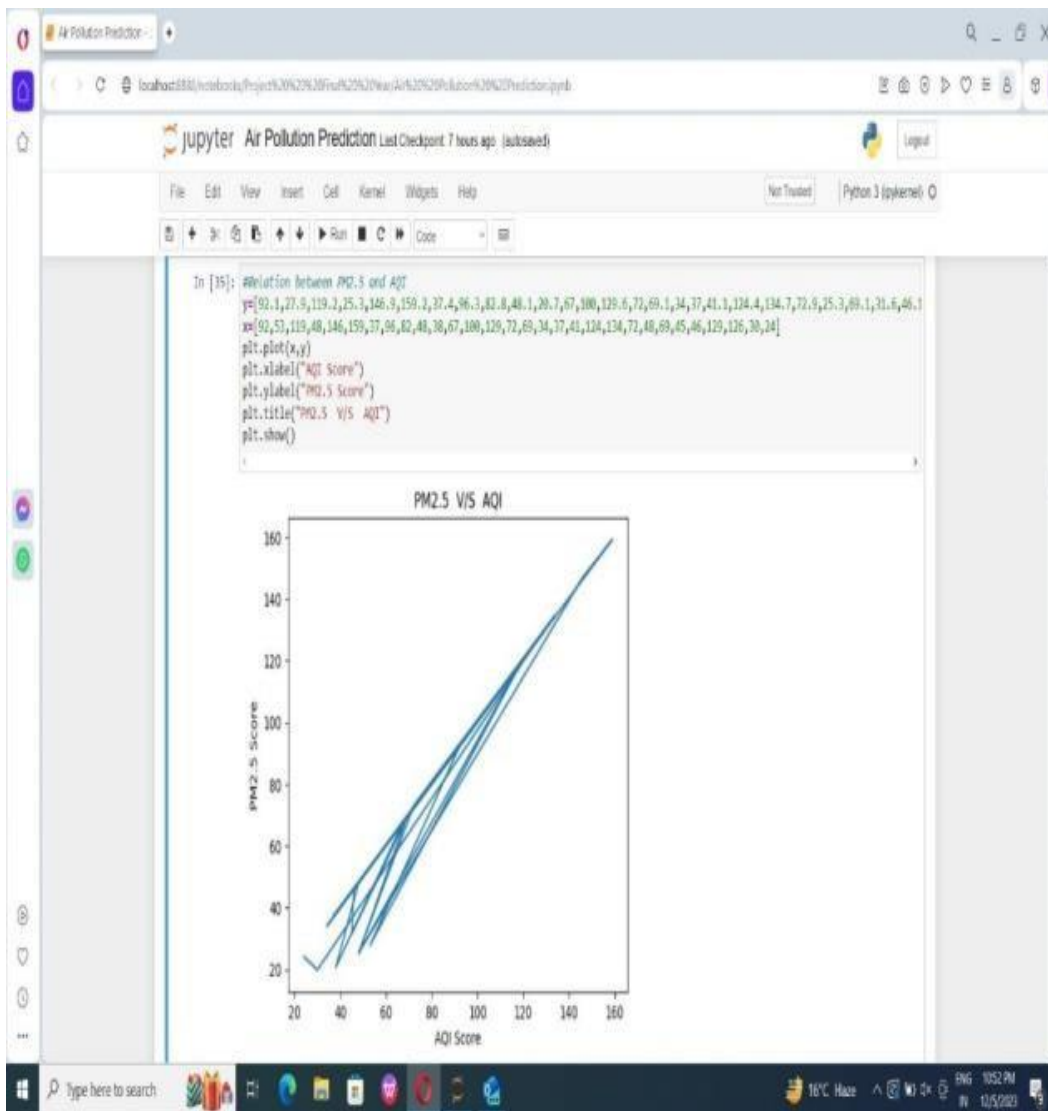
```
Out[20]: Index(['Place', 'PM2.5', 'PM10', 'SO2', 'NO2', 'O3', 'CO', 'AQI', 'Quality'], dtype='object')
```

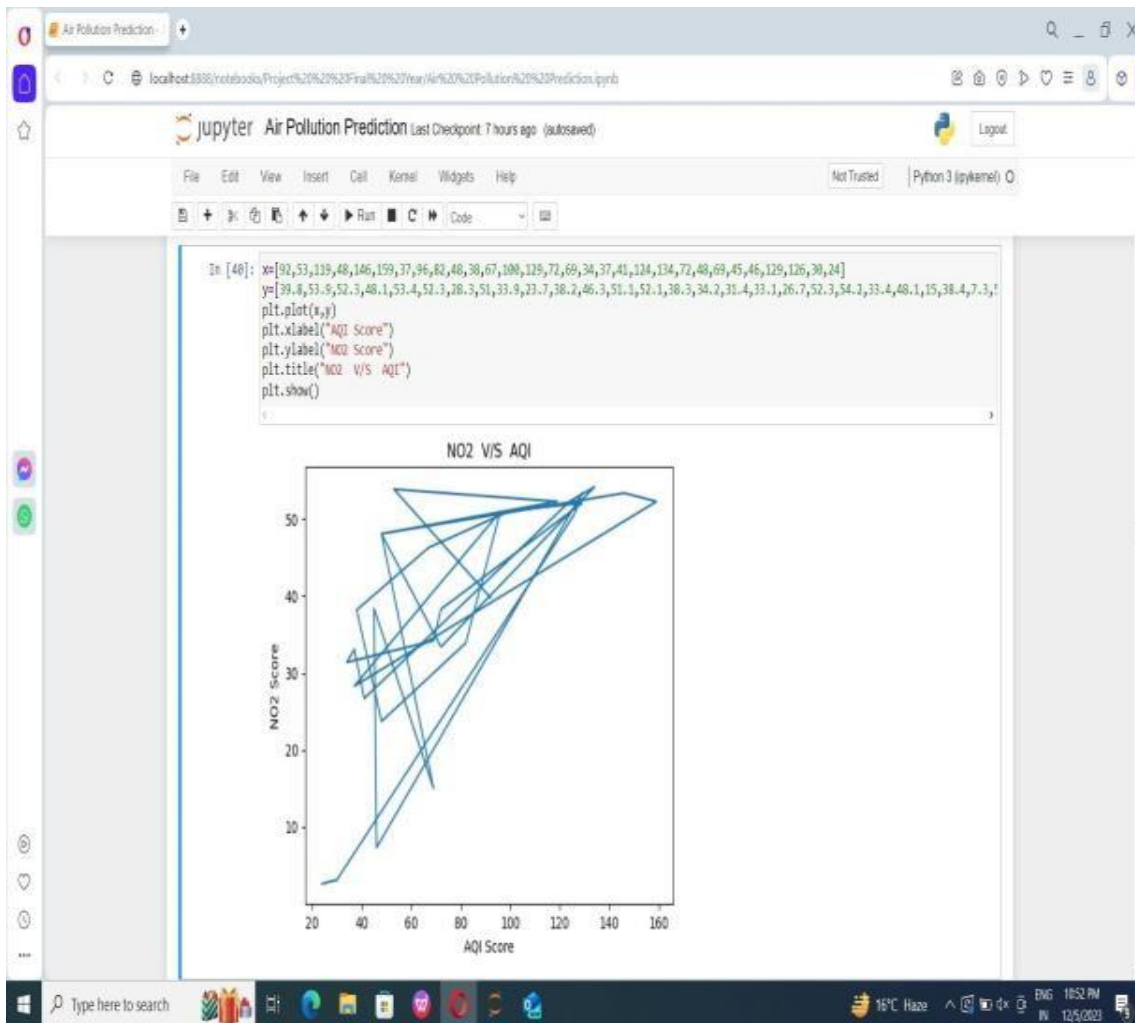
3.4.5 VALUE COUNT FUNCTION

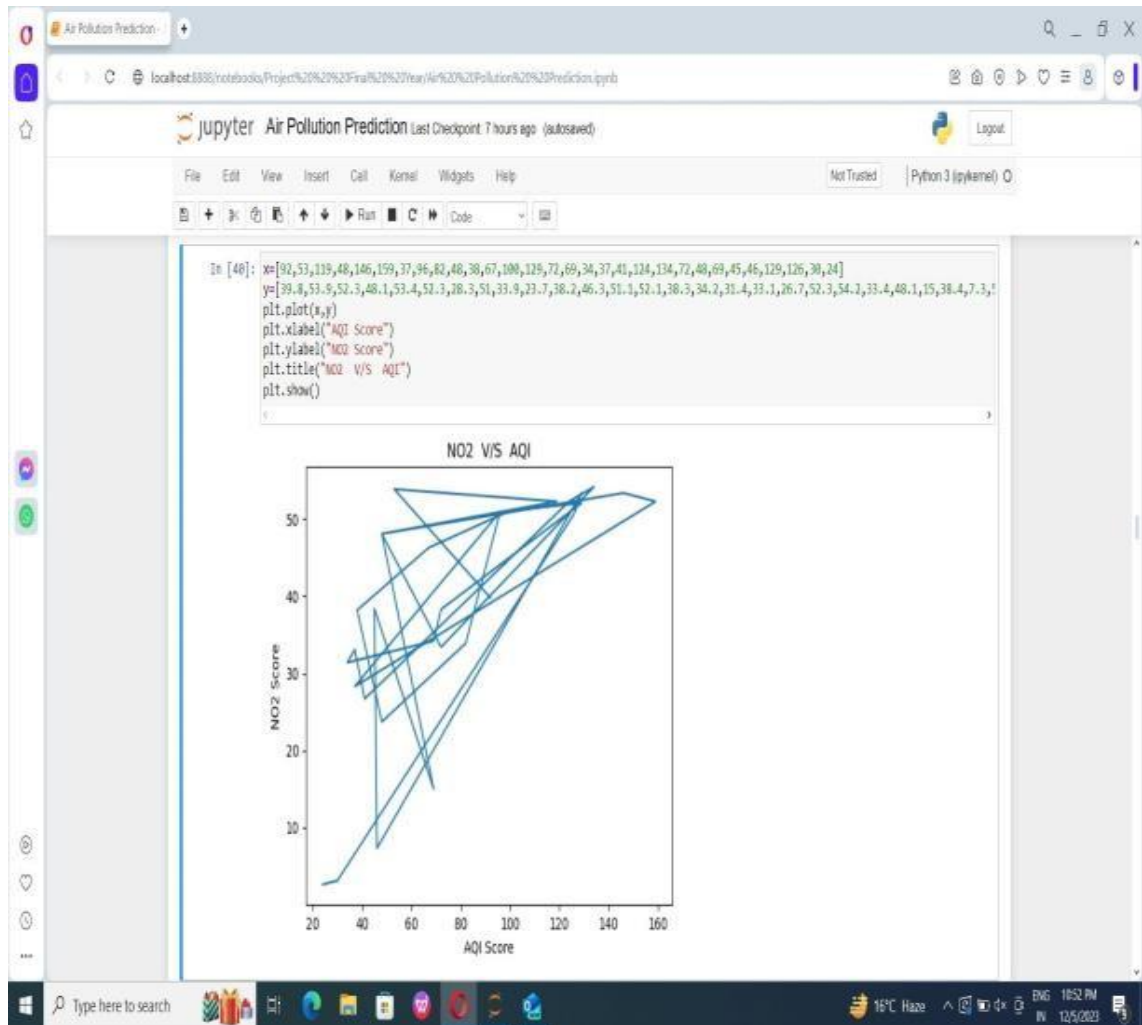
Value count function return the series containing the numbers and value of the unique data of the training dataset.

[illegible]

3.5 REPRESENTATION OF DATA

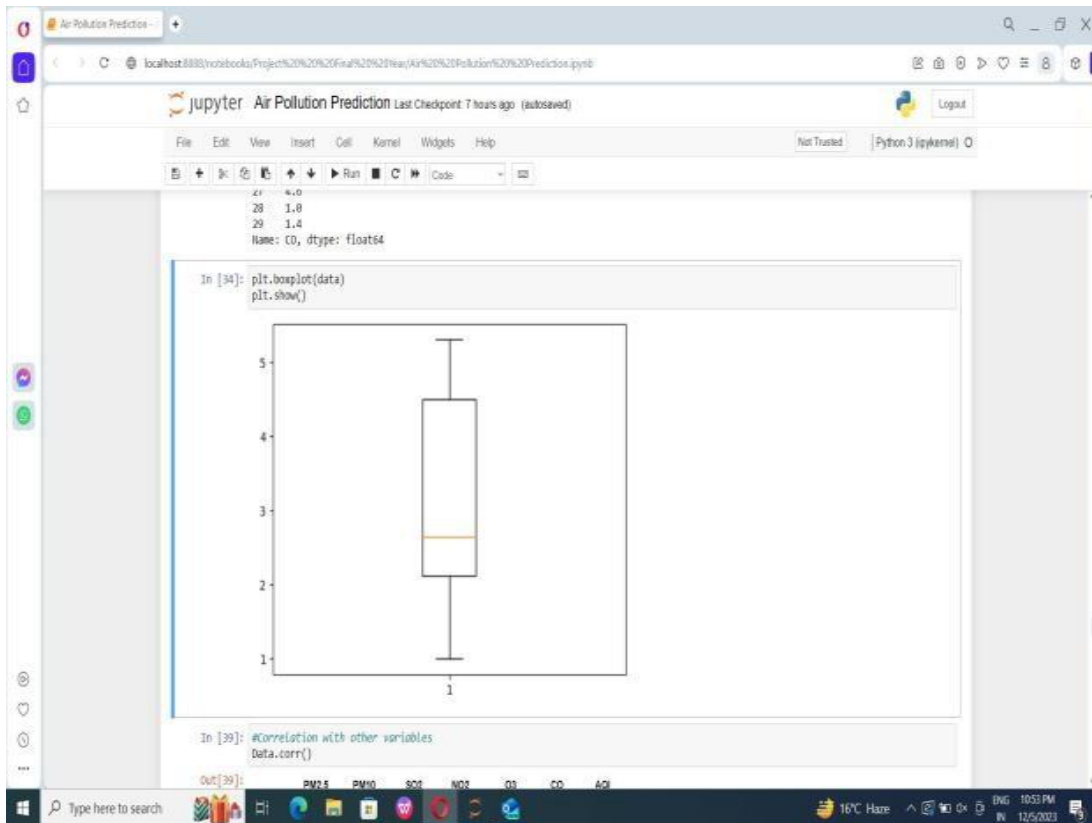






3.5.1 BOX PLOT REPRESENTATION

Box plot is used to represent the relation between the data and the differences between the expected and actual value.



3.5.2 CORRELATION

Correlation gives the direction of the linear line between the quantitative variables

Air Pollution Prediction

localhost:8888/notebooks/Project%20%20Final%20%20Year/Air%20%20Pollution%20%20Prediction.ipynb

jupyter Air Pollution Prediction (last checkpoint: 7 hours ago) (auto saved)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) 0

1

In [39]:

```
data.corr()
```

Out[39]:

	PM2.5	PM10	SO2	NO2	O3	CO	AQI
PM2.5	1.000000	0.969099	0.432085	0.607077	-0.266008	0.734967	0.969529
PM10	0.969099	1.000000	0.424189	0.665596	-0.265129	0.717082	0.965280
SO2	0.432085	0.424189	1.000000	0.309036	-0.471621	0.503989	0.382221
NO2	0.607077	0.665596	0.309036	1.000000	-0.429421	0.482017	0.691967
O3	-0.266008	-0.265129	-0.471621	-0.429421	1.000000	-0.273270	-0.273438
CO	0.734967	0.717082	0.503989	0.482017	-0.273270	1.000000	0.740099
AQI	0.969529	0.965280	0.382221	0.691967	-0.273438	0.740099	1.000000

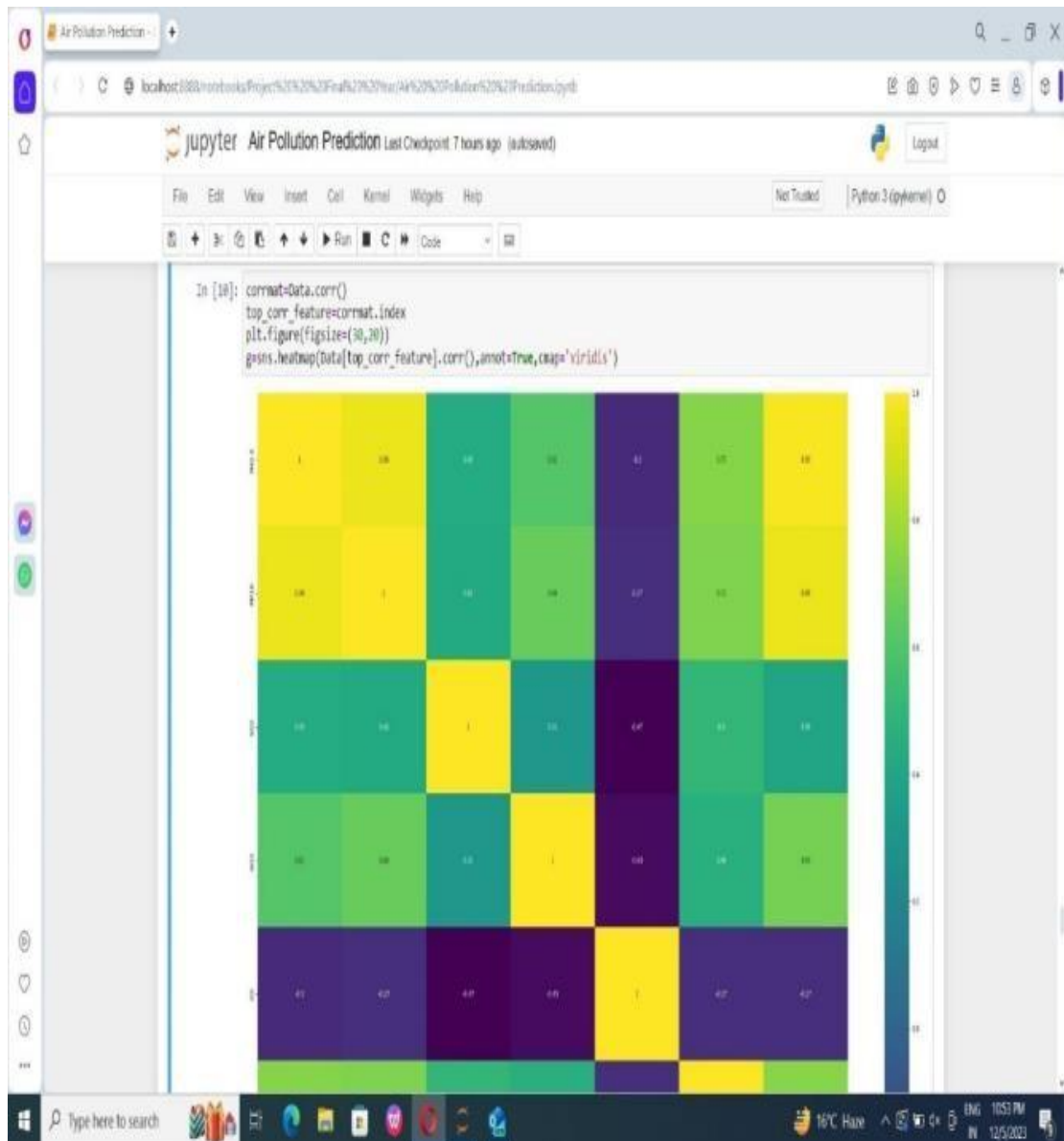
In [23]:

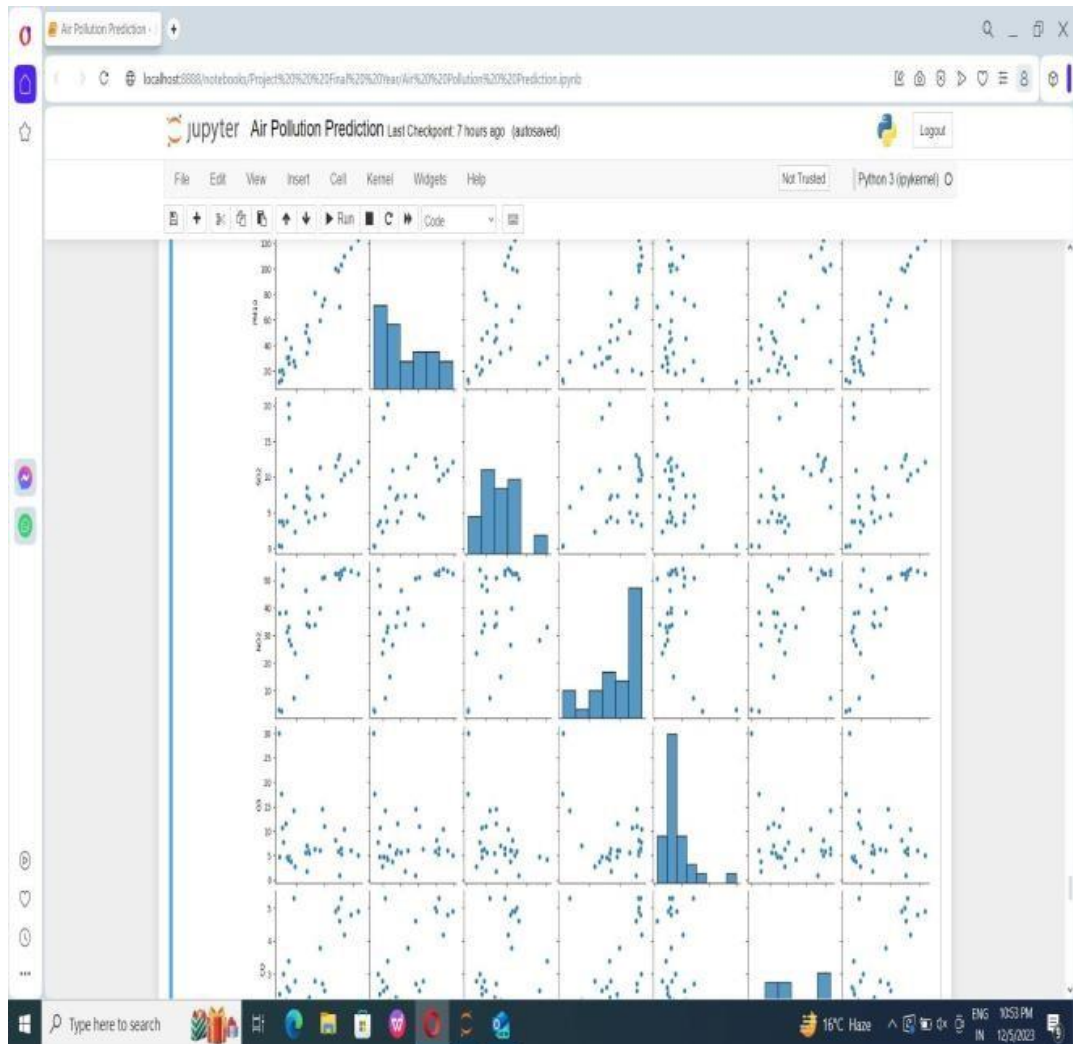
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression as lr
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score
import warnings
warnings.filterwarnings("ignore")
```

In [6]:

```
Data=pd.read_excel("C:\\Users\\Ashita\\Desktop\\Quality.xlsx")
```

16°C, Haze 10:53 PM 12/5/2023





3.6 PREDICTION FROM TRAINING DATA

3.6.1 LINEAR REGRESSION

Linear Regression is used to show the relationship between the dependent and the Independent variables.

3.6.2 TRAIN VALUES AND TEST VALUES

Train values and test values are very important evaluation of Linear Regression models by evaluating the training data with the concerns toward the requirement or the best fit outcomes.

3.6.3 FEATURES IN LINEAR REGRESSION

A. PARAMETERS

Linear Regression has two main parameters slope and intercepts. The slope represents the change in the dependent variable for a unit of change in independent variable. The intercept is the value of the variables when the independent variable is zero.

B. ASSUMPTIONS

Linear Regression is a parametric technique that relies on parameters learned from the data. The data must fulfil certain assumptions to obtain reliable results.

C. GRADIENT DESCENT

A Linear Regression model can be trained using the optimization algorithm gradient descent. The algorithm iteratively modifies the model parameters to reduce the mean square error of the model on a training set.

D. ORDINARY LEAST SQUARES

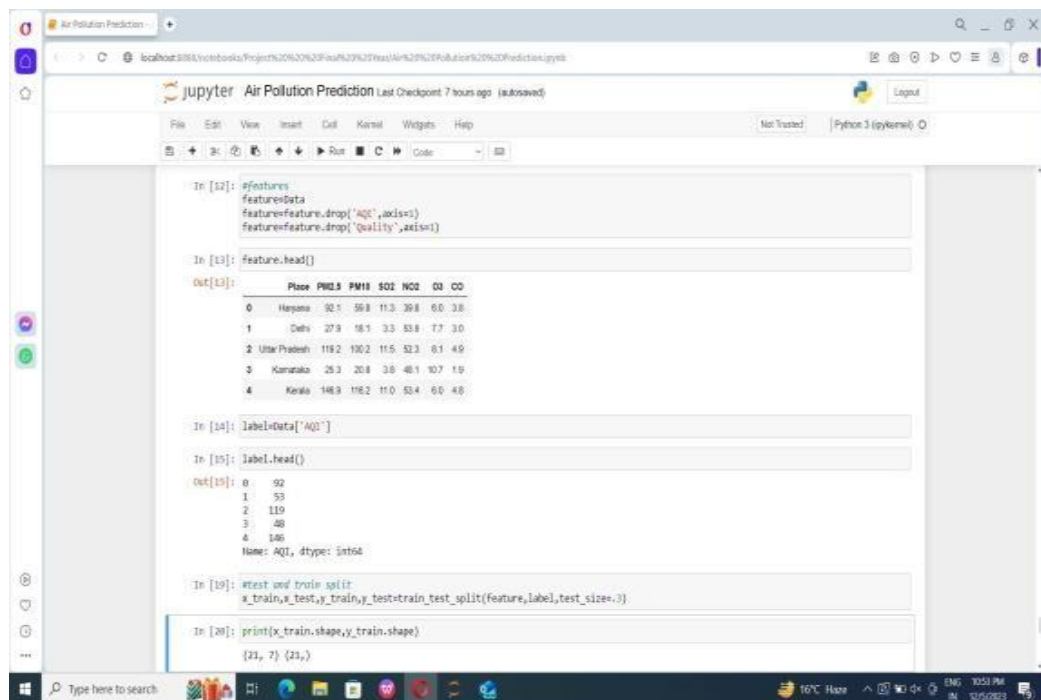
The ordinary least square procedure seeks to minimize the sum of the square residual.

E. COST FUNCTION

The cost function helps to figure out the best possible value which provide the best fit line for the training data.

F. MULTICOLLINEARITY

Multicollinearity in regression analysis occur when two or more predictors or independent Variable is highly correlated.



```
In [12]: #features
features=Data
features=features.drop('AQI',axis=1)
features=features.drop('Quality',axis=1)

In [13]: features.head()

Out[13]:
```

	Place	PM2.5	PM10	SO2	NO2	O3	CO
0	Haryana	92.1	59.8	11.3	39.8	6.0	3.8
1	Delhi	27.9	18.1	3.3	53.8	7.7	3.0
2	Uttar Pradesh	119.2	130.2	11.5	52.3	6.1	4.9
3	Karnataka	25.3	20.8	3.8	48.1	10.7	1.9
4	Kerala	148.9	116.2	11.0	52.4	6.0	4.8

```
In [14]: label=Data['AQI']

In [15]: label.head()

Out[15]:
```

0	92
1	53
2	119
3	48
4	146

```

Name: AQI, dtype: int64

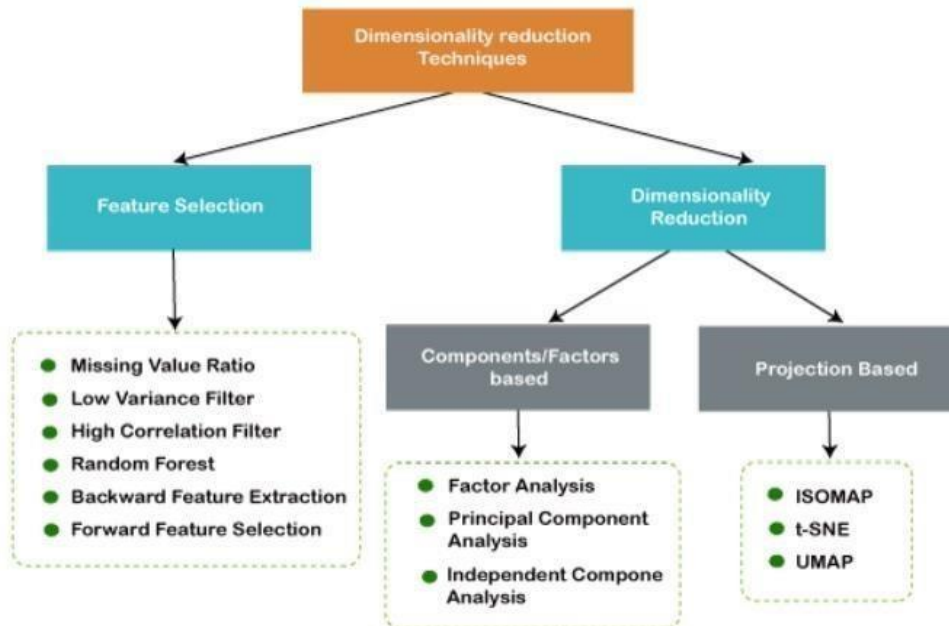
In [19]: #test and train split
x_train,x_test,y_train,y_test=train_test_split(features,label,test_size=.3)

In [20]: print(x_train.shape,y_train.shape)

(22, 7) (22,)
```

3.7 INTRODUCTION TO DIMENSIONALITY REDUCTION

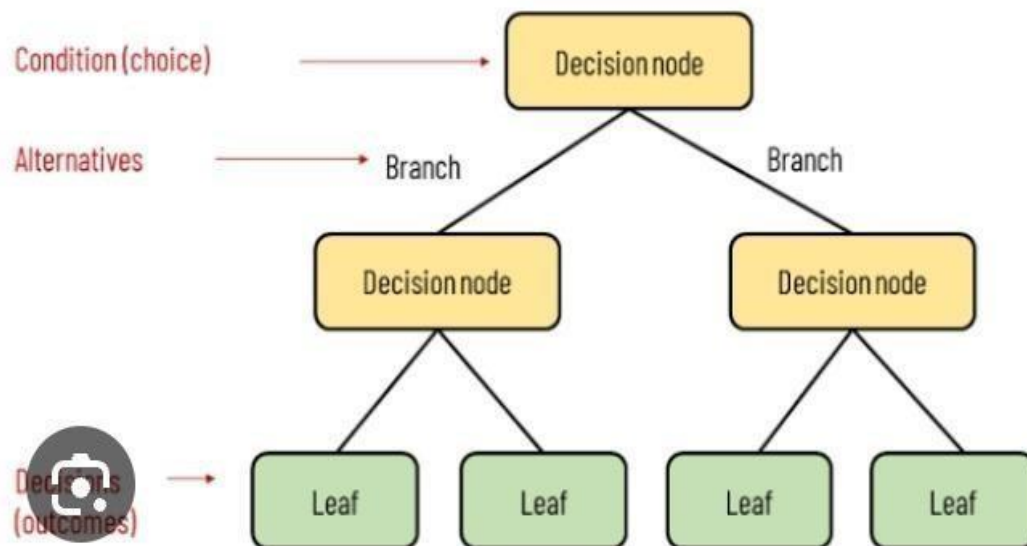
Dimensionality Reduction is mainly the collection of the technique used to reduce the number of features, reduce dimension in a dataset. It improves the performance of any model by reducing its complexity and helps in better visualization.



3.8 DECISION TREE

Decision Tree can be used for both regression and classification problem but mainly preferred for classification problems. In a decision tree, decision nodes are used to make any decision and have multiple branches. Leaf nodes are the output of those decision. In order to build a tree we use CART algorithm.

Elements of a decision tree



3.9 NAÏVE BAYES CLASSIFICATION THOREM

It is a probabilistic classifier algorithm which means its prediction is done on the basis of Probability of an object. Bayes theorem is also known as Bayes rule which is used to Determine the probability of a hypothesis. It depends on the conditional probability.

$$P(A/B) = (P(B/A) * P(A)) / P(B)$$

$P(A/B)$ =Posterior Probability

$P(A)$ =Prior Probability

$P(B)$ = Marginal Probability

$P(B/A)$ =Likelihood Probability

WORKING PROCESS OF NAÏVE BAYES ARE:

1. Convert the given dataset into frequency table.
2. Generate likelihood table by finding the probability of the feature.
3. Now use the Bayes formula to calculate posterior probability

3.10 CONCLUSION

Linear Regression is able to predict values for both dependent and independent variables and determines the linear relationship between them along with the discussion of null or alternative hypothesis.

The cost function provides the best feasible value which leads to the best fit.

But this Linear Regression is not applicable for complex computation.

It is very difficult to find the linear relationship between the variables.

3.11 REFERENCES

1. Types of Machine Learning from Java Point website
2. Quantitative vs Qualitative Data: What is the difference? By Emily Stevens
3. Writing Null Hypothesis in Research and Statistics by Joseph Quinones and Jennifer Mueller.
4. Logistics Regression from the Saedsayad.com website.
5. Regression vs Classification by Dhanush V
6. K Means Clustering from Java Point website.