

Experiment 1

Aim: Introduction to Data science and Data preparation using Pandas steps.

Theory:

Data science is the study of data that helps us derive useful insight for business decision making. Data Science is all about using tools, techniques, and creativity to uncover insights hidden within data. It combines math, computer science, and domain expertise to tackle real-world challenges in a variety of fields.

Data science involves these key steps:

- **Data Collection:** Gathering raw data from various sources, such as databases, sensors, or user interactions.
- **Data Cleaning:** Ensuring the data is accurate, complete, and ready for analysis.
- **Data Analysis:** Applying statistical and computational methods to identify patterns, trends, or relationships.
- **Data Visualization:** Creating charts, graphs, and dashboards to present findings clearly.
- **Decision-Making:** Using insights to inform strategies, create solutions, or predict outcomes.

Dataset Overview:

The dataset consists of air pollution readings across various cities in India over the last five years. Below are the key attributes:

- **City:** The city where pollution data was recorded.
- **Date:** The timestamp of the measurement.
- **PM2.5 & PM10:** Particulate matter concentration. (The numeric figure represents the diameter in micro-meter)
- **NO, NO₂, NO_x, NH₃:** nitrogen-based pollutants.
- **CO, SO₂, O₃:** Harmful environmental pollutants.
- **Benzene, Toluene, Xylene:** Hazardous air pollutants, usually generated by industries and power plants.
- **AQI:** Air Quality Index representing overall pollution level.

- **AQI_Bucket:** Categorized pollution levels (Good, Moderate, Poor, etc.).
-

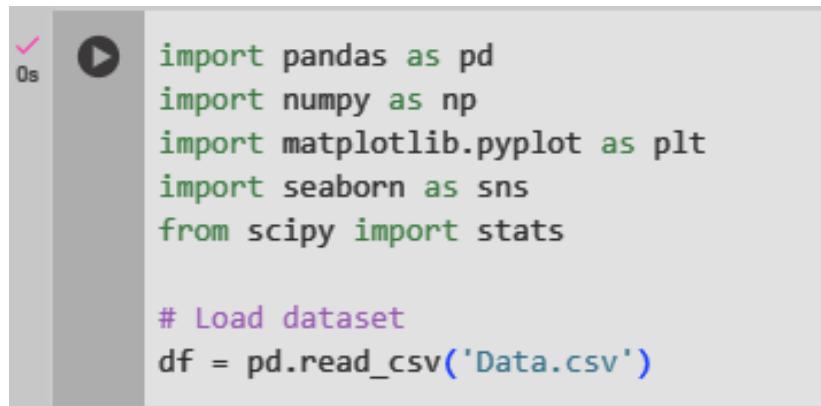
Problem Statement:

The objective is to analyze air pollution trends across Indian cities and identify key pollutants affecting air quality. Since the dataset provides information over cities of the past 5 years, we can use this information to predict air quality of a particular region in the future.

- Understanding variations in AQI across cities and time periods.
- Identifying major pollutants contributing to poor air quality.
- Visualizing trends, drawing meaningful conclusions and attempting future analysis from the dataset.

Code:

Loading the Dataset



```
✓ 0s  play
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

# Load dataset
df = pd.read_csv('Data.csv')
```

Basic Dataset Information

df.shape(): This returns a tuple indicating the number of rows and columns in the DataFrame.

df.info(): This prints "dataset info" and displays the DataFrame's structure including data types and non-null counts.

df.describe(): This prints "dataset description" and shows summary statistics like mean, standard deviation, and percentiles for numerical columns

```
✓ 0s   print("Dataset Shape:", df.shape)
      print("\nDataset Info:")
      df.info()
      print("\nDataset Description:")
      print(df.describe())

→ Dataset Shape: (29531, 16)
```

Removing Duplicate Entries

```
✓ 0s   df = df.drop_duplicates()
```

Creating Dummy Variables (One-Hot Encoding) for AQI Bucket:

This creates dummy data out of “AQI Bucket” for the various severity levels of pollution. This helps to convert categorical data to numerical data and helps in analysis in the algorithm

We use **df.head()** to verify this.

```
✓ 0s [11] df = pd.get_dummies(df, columns=['AQI_Bucket'], drop_first=True)

✓ 0s   print(df.head(10))
```

	City	Date	PM2.5	PM10	NO	NO2	NOx	NH3	CO	\
2123	Amaravati	25-11-2017	81.40	124.50	1.44	20.50	12.08	10.72	0.12	
2124	Amaravati	26-11-2017	78.32	129.06	1.26	26.00	14.85	10.28	0.14	
2125	Amaravati	27-11-2017	88.76	135.32	6.60	30.85	21.77	12.91	0.11	
2126	Amaravati	28-11-2017	64.18	104.09	2.56	28.07	17.01	11.42	0.09	
2127	Amaravati	29-11-2017	72.47	114.84	5.23	23.20	16.59	12.25	0.16	
2128	Amaravati	30-11-2017	69.80	114.86	4.69	20.17	14.54	10.95	0.12	
2129	Amaravati	01-12-2017	73.96	113.56	4.58	19.29	13.97	10.95	0.10	
2130	Amaravati	02-12-2017	89.90	140.20	7.71	26.19	19.87	13.12	0.10	
2131	Amaravati	03-12-2017	87.14	130.52	0.97	21.31	12.12	14.36	0.15	
2132	Amaravati	04-12-2017	84.64	125.00	4.02	26.98	17.58	14.41	0.18	
	S02	O3	Benzene	Toluene	Xylene	AQI	AQI_Bucket_Moderate			\
2123	15.24	127.09	0.20	6.50	0.06	184.0				True
2124	26.96	117.44	0.22	7.95	0.08	197.0				True
2125	33.59	111.81	0.29	7.63	0.12	198.0				True
2126	19.00	138.18	0.17	5.02	0.07	188.0				True
2127	10.55	109.74	0.21	4.71	0.08	173.0				True
2128	14.07	118.09	0.16	3.52	0.06	165.0				True
2129	13.90	123.80	0.17	2.85	0.04	191.0				True
2130	19.37	128.73	0.25	2.79	0.07	191.0				True
2131	11.41	114.80	0.23	3.82	0.04	227.0				False
2132	9.84	112.41	0.31	3.53	0.09	168.0				True
	AQI_Bucket_Poor							AQI_Bucket_Satisfactory	AQI_Bucket_Severe	\
2123		False						False	False	
2124		False						False	False	
2125		False						False	False	
2126		False						False	False	
2127		False						False	False	
2128		False						False	False	
2129		False						False	False	
2130		False						False	False	
2131		True						False	False	
2132		False						False	False	
	AQI_Bucket_Very_Poor									\
2123			False							
2124			False							
2125			False							
2126			False							
2127			False							

Identifying Outliers manually using the Standardization Approach (Z-Score Method)

To identify outliers manually we use the standardization approach (z score method). We find mean and standard deviation of the vehicle weight and calculate its z score; if it's less than -3 or greater than 3 means it's an outlier.

```

#By Z-score method
mean_aqi = df['AQI'].mean()
std_aqi = df['AQI'].std()

print (f"Mean of AQI: {mean_aqi}")
print (f"Standard Deviation of AQI: {std_aqi}")

df['Z_Score'] = (df['AQI'] - mean_aqi) / std_aqi
print(df[['AQI', 'Z_Score']])

# Identify outliers based on the Z-score
outliers = df[df['Z_Score'].abs() > 3]
print (outliers)

[5969 rows x 2 columns]

AQI      Z_Score
2123    184.0   0.496612
2124    197.0   0.638464
2125    198.0   0.649376
2126    188.0   0.540259
2127    173.0   0.376584
...
29523   86.0    -0.572733
29524   77.0    -0.670938
29525   47.0    -0.998288
29526   41.0    -1.063758
29527   70.0    -0.747319
[5969 rows x 2 columns]

AQI_Bucket_Moderate  AQI_Bucket_Poor  AQI_Bucket_Satisfactory \
3308                 False          False           False
4265                 False          False           False
10229                False          False           False
10230                False          False           False
10521                False          False           False
...
14880                ...           ...
14881                False          False           False
14994                False          False           False
14995                False          False           False
25531                False          False           False

AQI_Bucket_Severe  AQI_Bucket_Very_Poor  Z_Score
3308                  True          False  3.540971
4265                  True          False  3.704647
10229                 True          False  3.639176
10230                 True          False  3.442766
10521                 True          False  3.159062
...
14880                 ...           ...
14881                 True          False  3.802852
14994                 True          False  3.966527
14995                 True          False  3.311826
25531                 True          False  3.388208

```

Normalizing AQI using Min-Max Scaling

We normalize the data across the AQI on a scale of 0 to 1.

```
[25]: min_aqi = df['AQI'].min()
        max_aqi = df['AQI'].max()
        df['AQI_normalized'] = (df['AQI'] - min_aqi) / (max_aqi - min_aqi)

[26]: print(df[['AQI', 'AQI_normalized']])

      AQI  AQI_normalized
2123  184.0       0.246177
2124  197.0       0.266055
2125  198.0       0.267584
2126  188.0       0.252294
2127  173.0       0.229358
...
29523   86.0       0.096330
29524   77.0       0.082569
29525   47.0       0.036697
29526   41.0       0.027523
29527   70.0       0.071865

[5969 rows x 2 columns]
```

Conclusion

This experiment focused on preparing and analyzing air pollution data in India by addressing common data quality issues. Missing values were managed through replacement and removal techniques, while duplicate entries were eliminated to ensure data integrity. Outliers in AQI values were identified using the Z-score method, and Min-Max scaling was applied to normalize the data for better comparison. These preprocessing steps helped create a more structured and reliable dataset, which will allow for meaningful analysis of pollution trends.

Experiment 2

Aim: Data Visualization / Exploratory Data Analysis using Matplotlib and Seaborn

Introduction

Exploratory Data Analysis (EDA) is a crucial step in the data analysis pipeline. It involves visually and statistically exploring the data to gain insights and understand its underlying patterns, distributions, and relationships. In this section, we will use Matplotlib and Seaborn, two popular Python libraries, to create visualizations that help in uncovering these insights.

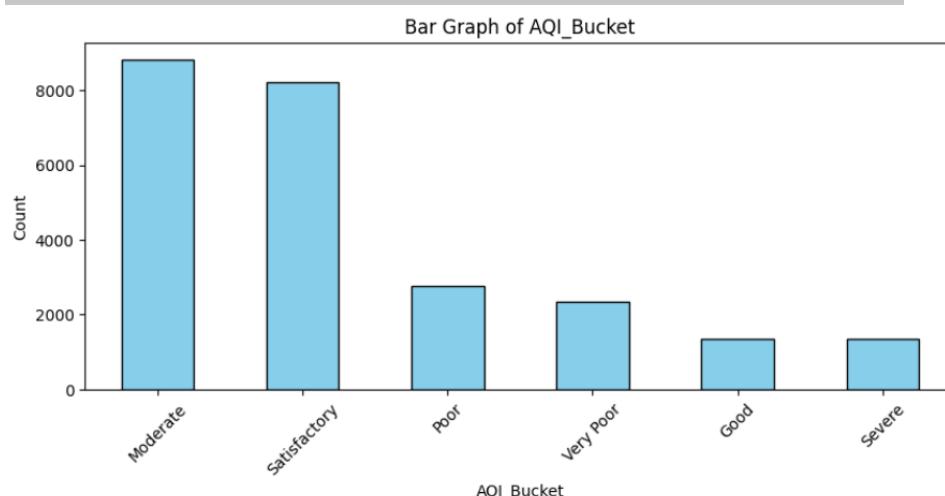
Matplotlib: A comprehensive library for creating static, animated, and interactive visualizations in Python.

Seaborn: Built on top of Matplotlib, Seaborn provides a high-level interface for drawing attractive and informative statistical graphics.

The primary objective of this analysis is to explore and visualize key patterns in the vehicle dataset, focusing on understanding relationships between different attributes and identifying significant trends.

Bar Graph and Contingency Table

```
# Bar Graph
plt.figure(figsize=(10, 4))
df[feature1].value_counts().plot(kind='bar', color='skyblue', edgecolor='black')
plt.xlabel(feature1)
plt.ylabel("Count")
plt.title(f"Bar Graph of {feature1}")
plt.xticks(rotation=45)
plt.show()
```



Observation: The majority of air quality readings fall under "Moderate" and "Satisfactory" categories, indicating generally acceptable pollution levels, but the presence of "Poor" and "Severe" levels suggests occasional hazardous conditions.

Contingency Table

A contingency table helps analyze the relationship between PM2.5 (particulate matter) and AQI_Bucket.

```
✓ 0s  feature1 = "PM2.5"
    feature2 = "AQI_Bucket"

    contingency_table = pd.crosstab(df[feature1], df[feature2])
    print("\nContingency Table:\n", contingency_table)

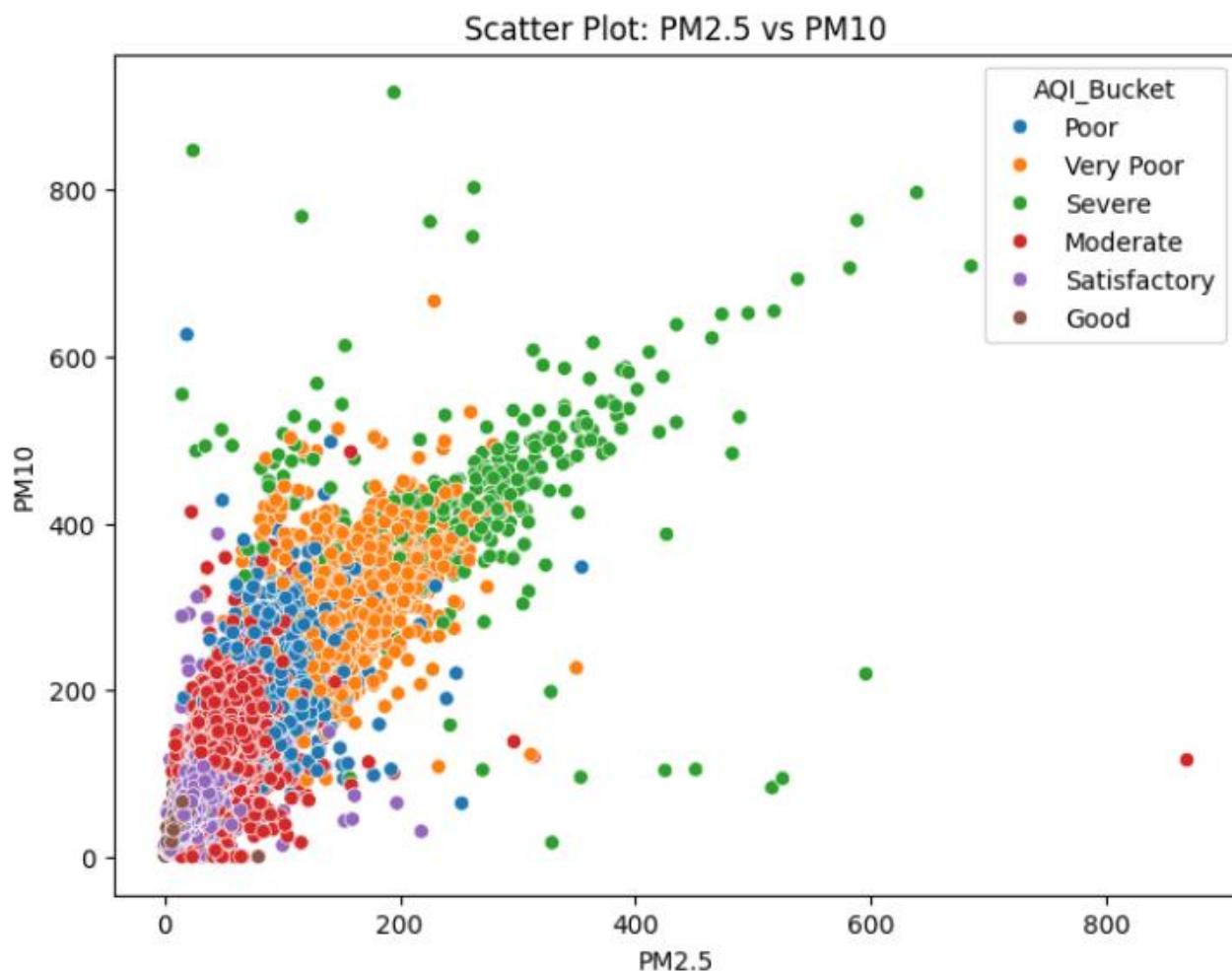
→
Contingency Table:
   AQI_Bucket  Good  Moderate  Poor  Satisfactory  Severe  Very Poor
PM2.5
0.04          0      0       0           1       0       0
0.16          1      0       0           0       0       0
0.24          1      0       0           0       0       0
0.28          1      0       0           0       0       0
0.98          1      0       0           0       0       0
```

Observation: The contingency table shows the distribution of PM2.5 values across different AQI categories, indicating a relationship between particulate matter concentration and air quality levels. If dust matter is finer ($0.16 > 0.04$), the AQI reduces.

Scatter Plot, Box Plot, and Heatmap

A scatter plot helps in visualizing the relationship between particulate matter of varying micro-meters and AQI.

```
# Scatter Plot
plt.figure(figsize=(8, 6))
sns.scatterplot(x=df["PM2.5"], y=df["PM10"], hue=df["AQI_Bucket"])
plt.title("Scatter Plot: PM2.5 vs PM10")
plt.xlabel("PM2.5")
plt.ylabel("PM10")
plt.show()
```

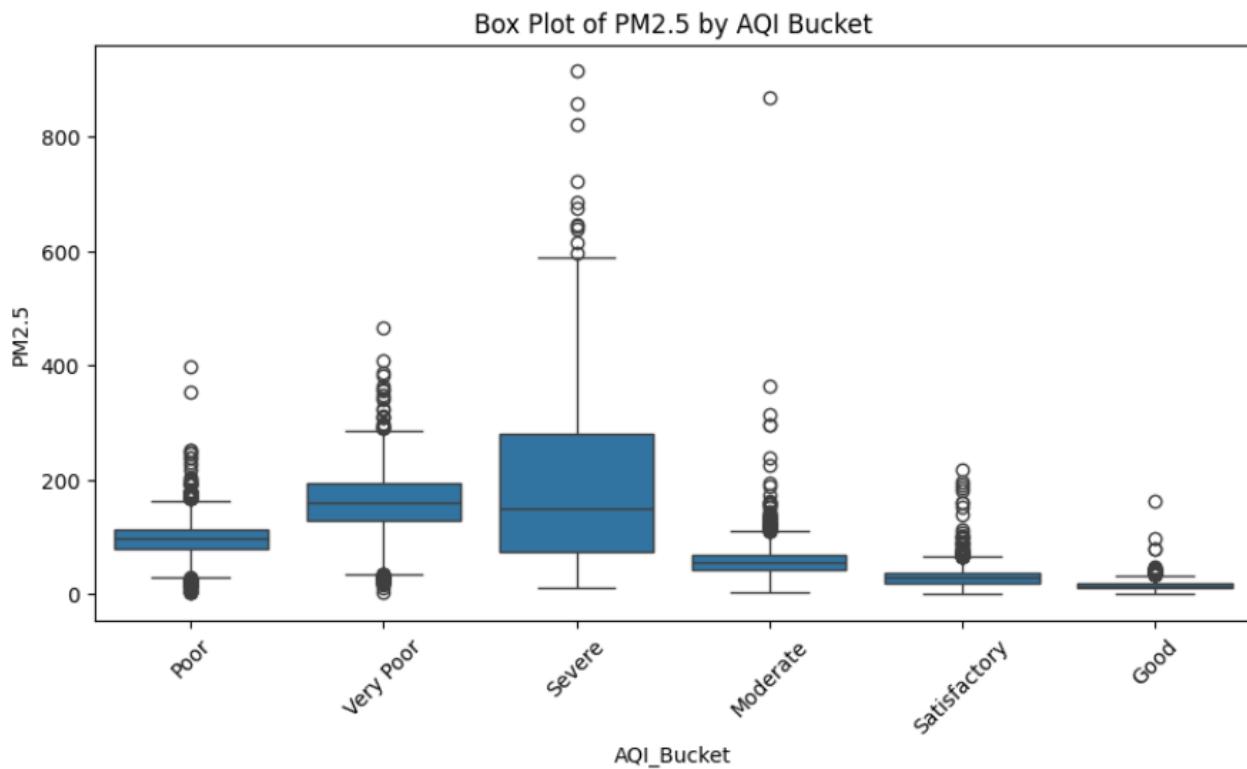


Observation: The scatter plot shows a positive correlation between PM2.5 and PM10, with higher pollutant levels corresponding to worse AQI categories, confirming that fewer particulates result in better air quality. Dots at the edge of the plot can be considered as outliers.

Box Plot

A box plot is used to summarize the distribution of the Data_Value column, helping to identify outliers, the median, and the interquartile range.

```
# Box Plot
plt.figure(figsize=(10, 5))
sns.boxplot(x=df["AQI_Bucket"], y=df["PM2.5"])
plt.title("Box Plot of PM2.5 by AQI Bucket")
plt.xticks(rotation=45)
plt.show()
```



The lines on every plot represent the upper bound, Q3 (75%ile), Q2 (50%ile), Q1 (25%ile) and the lower bound. Circles represent outliers.

After removing outliers using IQR (Inter Quartile Range Method)

```

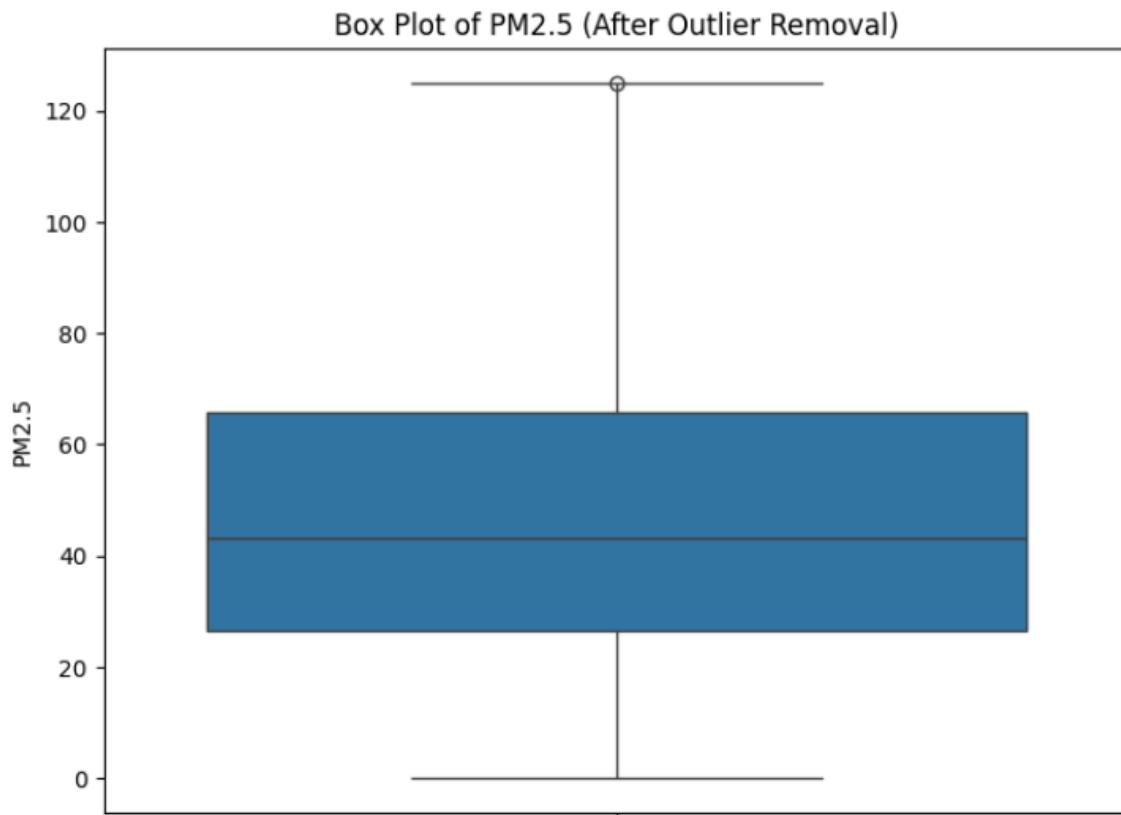
Q1 = df["PM2.5"].quantile(0.25)
Q3 = df["PM2.5"].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

outliers = df_cleaned_5[(df["PM2.5"] < lower_bound) | (df["PM2.5"] > upper_bound)]
#print("\nOutliers in PM2.5:\n", outliers)

df_cleaned = df[(df["PM2.5"] >= lower_bound) & (df["PM2.5"] <= upper_bound)]

plt.figure(figsize=(8, 6))
sns.boxplot(y=df_cleaned["PM2.5"])
plt.title("Box Plot of PM2.5 (After Outlier Removal)")
plt.show()

```



Observation: The box plot of PM2.5 (after outlier removal) shows that most values are concentrated within a reasonable range, but a few high values still exist as minor outliers, indicating occasional spikes in pollution levels.

Heatmap for Correlation Analysis

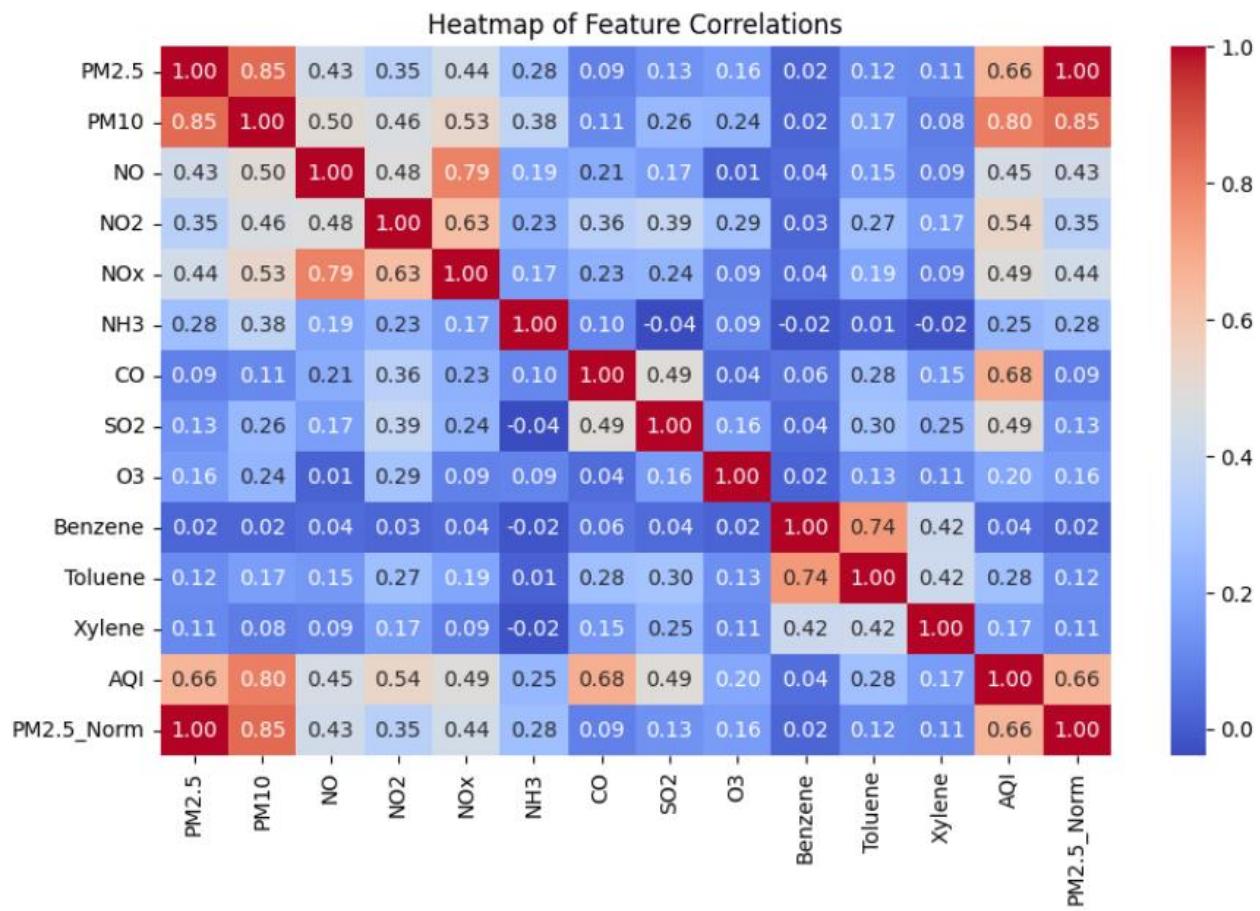
A heatmap is used to visualize the contingency table, which represents the relationship between two categorical variables.

```
# Heatmap
plt.figure(figsize=(10, 6))

numeric_df = df.select_dtypes(include=['number'])

sns.heatmap(numeric_df.corr(), annot=True, cmap="coolwarm", fmt=".2f")

plt.title("Heatmap of Feature Correlations")
plt.show()
```



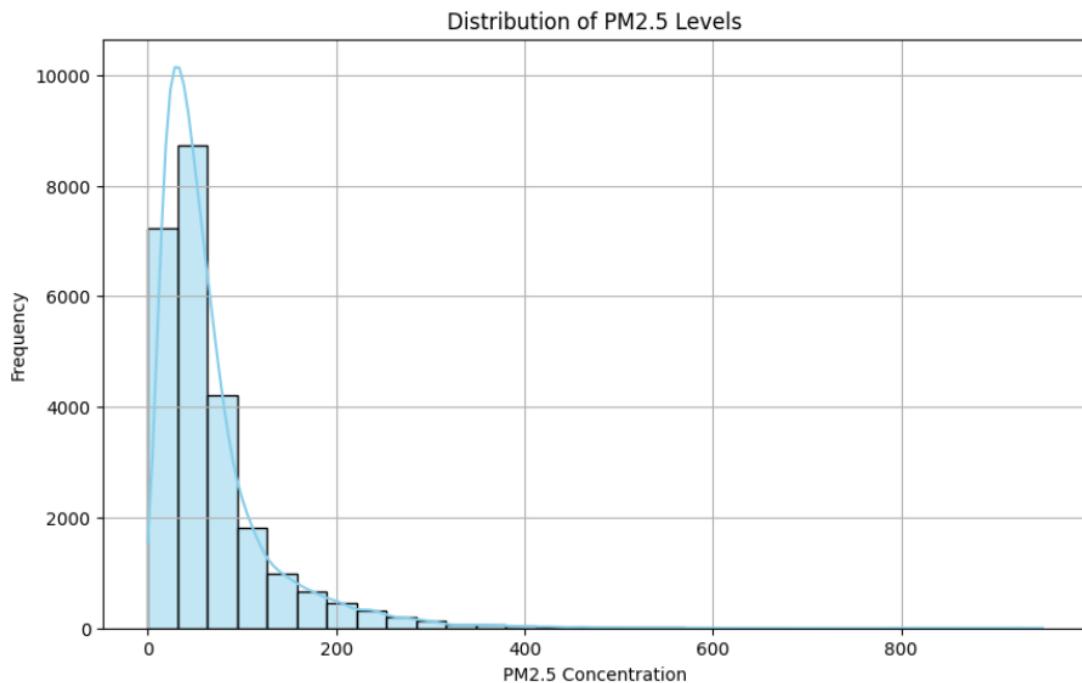
Observation: The heatmap of feature correlations shows that PM2.5 has a strong positive correlation with PM10 (0.85) and AQI (0.66), indicating that higher particle levels are associated with poorer air quality.

Histogram and Normalized Histogram

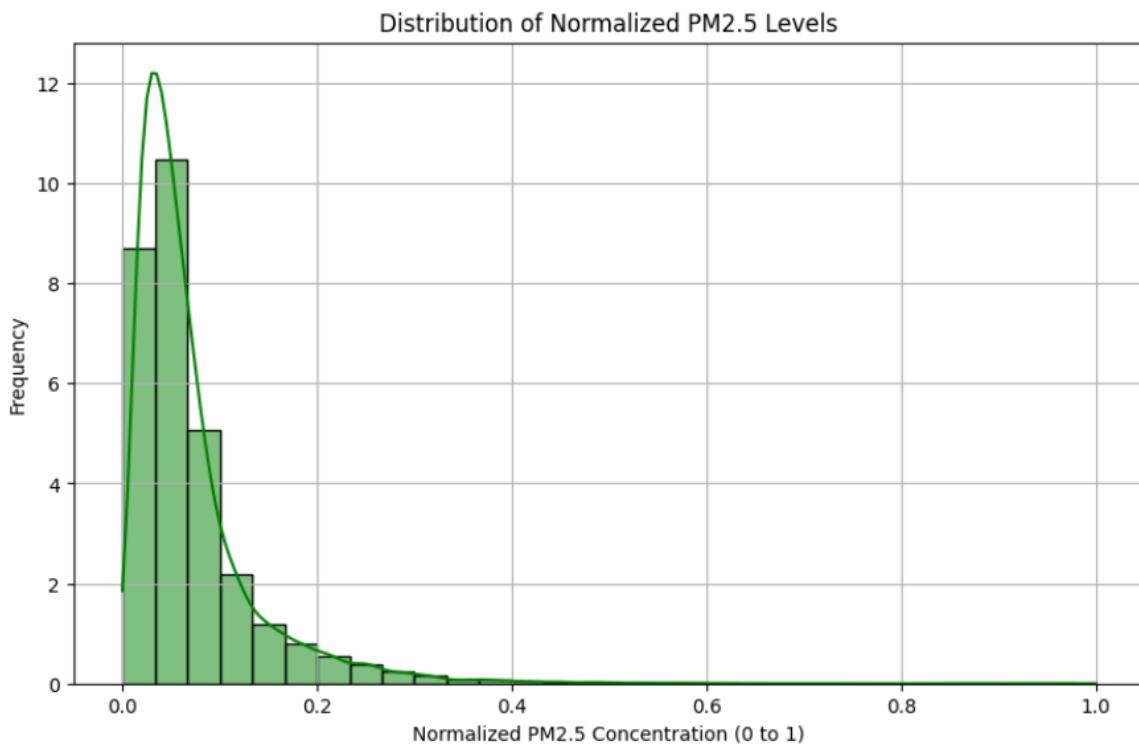
A histogram is used to visualize the distribution of Data_Value, showing how frequently different values occur.

```
# Plot histogram for PM2.5
plt.figure(figsize=(10, 6))
sns.histplot(df["PM2.5"], bins=30, kde=True, color="skyblue")

plt.xlabel("PM2.5 Concentration")
plt.ylabel("Frequency")
plt.title("Distribution of PM2.5 Levels")
plt.grid(True)
plt.show()
```



A normalized histogram represents data in terms of density instead of raw frequency, ensuring that the total area under the bars equals 1. The `stat='density'` parameter in `sns.histplot` normalizes the counts. The bin width is calculated, and the total area of the histogram is verified to confirm normalization.



Conclusion

This experiment provided valuable insights into the dataset through various visualizations. The contingency table highlighted the relationship between PM2.5 levels and AQI categories, demonstrating that lower particulate concentrations correspond to better air quality. The scatter plot revealed a strong correlation between PM2.5 and PM10, reinforcing the idea that both pollutants contribute significantly to air quality degradation. The heatmap further confirmed these correlations, showing strong associations between PM2.5, PM10, and AQI. The box plot helped detect and visualize outliers, emphasizing the need for proper handling to ensure accurate statistical analysis. Identifying and managing these extreme values is crucial for maintaining data integrity. Overall, this study reinforced the importance of exploratory data analysis (EDA) in understanding air pollution trends, detecting patterns, and ensuring data-driven decision-making for effective environmental monitoring.

Experiment 3

Aim: To perform data modeling.

Theory:

Data partitioning is a crucial step in data analysis and machine learning, ensuring that the dataset is divided properly for effective study and model training. Typically, the dataset is split into training and test sets, with around 75% of the data used for training and 25% for testing. This prevents overfitting and allows for unbiased evaluation.

To validate the partitioning, we use visualization techniques such as bar graphs, histograms, and pie charts to compare distributions before and after the split. Counting the records ensures that the dataset is divided correctly.

A statistical validation step, such as a two-sample Z-test, is performed to compare the means of numerical features in both subsets. If the p-value is greater than 0.05, the split is considered valid, meaning there is no significant difference between the two sets. If the p-value is below 0.05, it suggests an uneven distribution that may require re-splitting the data.

Ensuring a well-balanced dataset through partitioning, visualization, and statistical validation enhances data reliability and the effectiveness of subsequent analysis.

Steps:

Partitioning the dataset using train_test_split:

The process of dividing a dataset into two subsets: a training set and a test set. Typically, 75% of the data is used for training, where the model learns patterns, and 25% is used for testing to evaluate performance on unseen data. This division ensures that the model generalizes well and does not simply memorize the training examples. Proper partitioning helps in reducing overfitting and provides a fair evaluation of the model's effectiveness. The `train_test_split` function from `sklearn.model_selection` is commonly used to achieve this.

```
 0s  ⏎ import pandas as pd
 0s  import numpy as np
 0s  import matplotlib.pyplot as plt
 0s  import seaborn as sns
 0s  from sklearn.model_selection import train_test_split
 0s  from statsmodels.stats.weightstats import ztest

 0s  df = pd.read_csv('Data.csv')

 0s [33] train_df, test_df = train_test_split(df, test_size=0.25, random_state=42)

 0s  print("Training Set Shape:", train_df.shape)
 0s  print("Testing Set Shape:", test_df.shape)
```

Visualizing the distribution of training and test sets

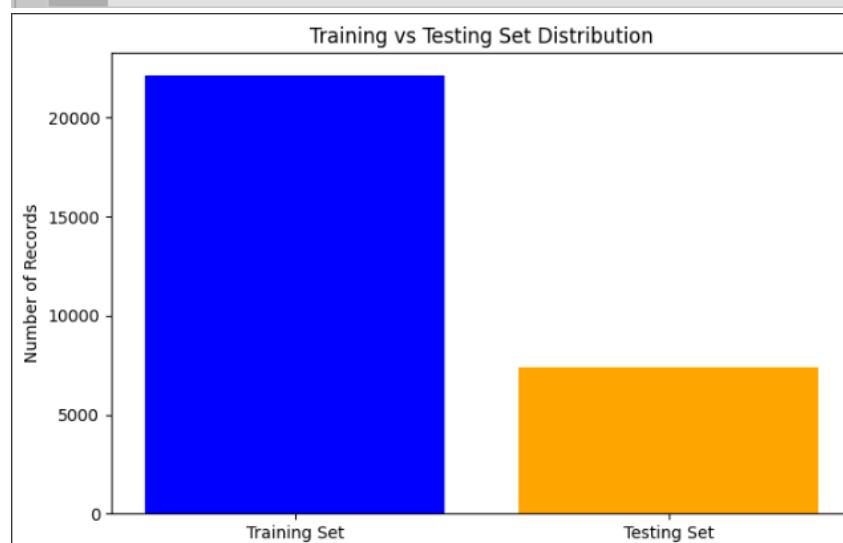
This ensures that the split maintains the original dataset's characteristics.

Bar graphs can be used to compare the number of records in both sets, while histograms and pie charts help check whether numerical and categorical feature distributions remain balanced.

If a class or feature is disproportionately represented in either subset, the split may need adjustment.

The `matplotlib.pyplot` library in Python helps create such visualizations to confirm a proper split.

```
 0s  ⏎ plt.figure(figsize=(8,5))
 0s  plt.bar(['Training Set', 'Testing Set'], [len(train_df), len(test_df)], color=['blue', 'orange'])
 0s  plt.title("Training vs Testing Set Distribution")
 0s  plt.ylabel("Number of Records")
 0s  plt.show()
```



Counting records

Counting the number of records in both training and test sets ensures that the split has been performed correctly. The expected number of samples in each set is calculated using simple percentage formulas, such as **Training Size = Total Data × 0.75** and **Testing Size = Total Data × 0.25**. By printing the lengths of the training and test sets after splitting, we can verify if the proportions match the intended split. This step helps in detecting potential errors in dataset partitioning.

```
✓ 0s   ➔ print("Training Set Records:", len(train_df))
      print("Testing Set Records:", len(test_df))

      ➔ Training Set Records: 22148
          Testing Set Records: 7383
```

Performing a two-sample Z-test to compare AQI values in both sets

It is used to statistically verify whether the training and test sets come from the same distribution. It compares the means of numerical features in both subsets and checks for significant differences.

If the p-value from the Z-test is greater than 0.05, the split is valid, meaning there is no significant difference between the two sets. However, if the p-value is below 0.05, the dataset may not be evenly distributed, requiring a reassessment of the split.

The `scipy.stats.ztest` function in Python is commonly used to perform this validation.

```
✓ 0s   ➔ # Perform Z-test for 'AQI' between training and test datasets
      train_aqi = train_df['AQI']
      test_aqi = test_df['AQI']

      # Z-test for independent samples
      z_stat, p_val = ztest(train_aqi, test_aqi)

      # Display Z-test results
      print(f"Z-statistic: {z_stat}")
      print(f"P-value: {p_val}")

      # Interpret the result
      if p_val < 0.05:
          print("There is a significant difference between the training and test data for AQI.")
      else:
          print("There is no significant difference between the training and test data for AQI.")
```

→ Z-statistic: -2.425627418979989

P-value: 0.015281950139779434

There is no significant difference between the training and test data for AQI.

Conclusion

Proper dataset partitioning, visualization, and statistical validation ensure a balanced and unbiased split, leading to reliable model performance. This approach minimizes overfitting and improves the model's generalization to real-world data.

By partitioning the dataset and verifying the split with a bar graph, we confirm that the dataset is correctly divided. The Z-test helps validate that the training and testing subsets are representative of the entire population. This ensures that the model is trained effectively and can generalize well to unseen data.

AIDS Exp 04

Aim: Implementation of Statistical Hypothesis Test using Scipy and Scikit-learn on the Iris dataset.

1. Introduction

Air quality datasets are crucial for understanding environmental pollution and its impact on public health. This dataset consists of various air quality parameters measured across different dates and cities, including pollutants like PM2.5, PM10, NO, NO2, CO, SO2, O3, Benzene, and Toluene, along with the Air Quality Index (AQI). The dataset helps analyze the correlation between these pollutants and their effect on air quality. This experiment aims to assess the relationship between different pollutants using statistical tests, including Pearson's, Spearman's, and Kendall's correlation coefficients. Additionally, the Chi-Squared test will be performed to evaluate the dependency between AQI categories and pollutant concentration levels.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	City	Date	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	Toluene	AQI	AQI_Bucke	AQI_Bucke	AQI_Bucke	AQI_Bucke	AQI_Buck
2	Ahmedaba	15-05-2019	37.55	122.41	15.08	85.12	58.72	25.24913	15.08	163.01	48.23	16.44	85.54	281	FALSE	TRUE	FALSE	FALSE	FALSE
3	Ahmedaba	16-05-2019	33.97	116.32	14.67	79.71	55.61	25.24913	14.67	91.26	51.86	15.55	83.89	330	FALSE	FALSE	FALSE	FALSE	TRUE
4	Ahmedaba	17-05-2019	35.48	130.07	18.02	77.61	58.41	25.24913	18.02	98.35	38.99	15.88	83.83	356	FALSE	FALSE	FALSE	FALSE	TRUE
5	Ahmedaba	18-05-2019	34.11	138.31	13.27	75.23	51.83	25.24913	13.27	88.66	42.22	15.93	82.73	359	FALSE	FALSE	FALSE	FALSE	TRUE
6	Ahmedaba	19-05-2019	33.69	111.73	34.56	68.9	69.77	25.24913	34.56	80.9	36.95	15.53	84.17	547	FALSE	FALSE	FALSE	TRUE	FALSE
7	Ahmedaba	20-05-2019	42.31	118.65	17.47	81.84	59.84	25.24913	17.47	89.57	46.68	15.98	83.87	813	FALSE	FALSE	FALSE	TRUE	FALSE
8	Ahmedaba	21-05-2019	24.6	103.88	11.03	81.24	52.21	25.24913	11.03	80.74	46.65	15.31	82.95	321	FALSE	FALSE	FALSE	FALSE	TRUE
9	Ahmedaba	22-05-2019	27.93	103.3	11.44	76.75	50.49	25.24913	11.44	86.48	54.34	15.6	84.17	270	FALSE	TRUE	FALSE	FALSE	FALSE
10	Ahmedaba	23-05-2019	41.39	135.65	14.29	89.1	59.76	25.24913	14.29	105.96	49.7	16.33	83.95	323	FALSE	FALSE	FALSE	FALSE	TRUE
11	Ahmedaba	24-05-2019	46.79	148	14.31	93.27	61.82	25.24913	14.31	131.04	56.31	15.21	82.4	344	FALSE	FALSE	FALSE	FALSE	TRUE
12	Ahmedaba	25-05-2019	51.63	156.97	17.96	89.18	63.98	25.24913	17.96	134.22	49.62	15.72	83.2	404	FALSE	FALSE	FALSE	TRUE	FALSE
13	Ahmedaba	26-05-2019	63.15	177.87	28.03	100.08	80.78	25.24913	28.03	122.43	50.32	17.05	84.59	558	FALSE	FALSE	TRUE	FALSE	FALSE
14	Ahmedaba	27-05-2019	57.47	163.36	21.39	112.68	79.36	25.24913	21.39	143.3	52.19	16.47	84.11	435	FALSE	FALSE	TRUE	FALSE	FALSE
15	Ahmedaba	28-05-2019	50.27	156.63	21.02	103.05	74.24	25.24913	21.02	118.55	49.86	16.92	85.28	440	FALSE	FALSE	TRUE	FALSE	FALSE
16	Ahmedaba	29-05-2019	42.02	140.66	16.43	71.51	53.58	25.24913	16.43	82.75	52.44	16.99	85.25	374	FALSE	FALSE	FALSE	TRUE	FALSE
17	Ahmedaba	30-05-2019	48.74	153.69	19.89	91.02	67.08	25.24913	19.89	138.12	52.86	17.04	83.72	515	FALSE	FALSE	TRUE	FALSE	FALSE
18	Ahmedaba	31-05-2019	46.51	136.34	16.75	85.37	60.74	25.24913	16.75	145.55	44.53	15.68	84.55	360	FALSE	FALSE	FALSE	TRUE	TRUE
19	Ahmedaba	01-06-2019	48.1	142.06	23.83	70.71	61.67	25.24913	23.83	142.52	40.22	15.5	84.08	467	FALSE	FALSE	TRUE	FALSE	FALSE
20	Ahmedaba	02-06-2019	41.38	119.91	21.8	70.35	59.18	25.24913	21.8	134.76	41.5	15.9	83.78	402	FALSE	FALSE	TRUE	FALSE	FALSE
21	Ahmedaba	03-06-2019	46.46	134.2	22.33	74.31	61.71	25.24913	22.33	137.83	47.63	14.5	81.45	419	FALSE	FALSE	TRUE	FALSE	FALSE

2. Theoretical Background

2.1 Pearson's Correlation Coefficient (r)

Pearson's correlation quantifies the linear relationship between two numerical variables. It ranges from -1 to 1, where:

Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r = correlation coefficient

x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable

- $r > 0 \rightarrow$ Positive relationship
- $r < 0 \rightarrow$ Negative relationship
- $r = 0 \rightarrow$ No correlation

Importance:

- Useful for identifying linear dependencies.
- Requires normally distributed data.

2.2 Spearman's Rank Correlation (ρ)

Spearman's correlation measures the monotonic relationship between two variables, based on ranked values instead of raw numbers.

Formula

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

ρ = Spearman's rank correlation coefficient

d_i = difference between the two ranks of each observation

n = number of observations

- Works for non-linear relationships.
- Less affected by outliers compared to Pearson's correlation.

Importance:

- Ideal for datasets that do not follow a normal distribution.

- Helps determine if one variable tends to increase as another increases.

2.3 Kendall's Rank Correlation (τ)

Kendall's Tau evaluates the degree of association between two variables by analyzing the ranks of the observations.

Formula:

$$\tau = \frac{C - D}{C + D}$$

Where:

- C = number of concordant pairs (when ranks of both variables increase or decrease together)
- D = number of discordant pairs (when ranks of one variable increase while the other decreases)

Interpretation:

- $\tau > 0 \rightarrow$ Positive association
- $\tau < 0 \rightarrow$ Negative association
- $\tau = 0 \rightarrow$ No association

Importance:

- Measures consistency in ranking.
- More effective for smaller datasets.

2.4 Chi-Squared Test (χ^2)

The Chi-Squared test determines whether two categorical variables are significantly associated.

Formula

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

χ^2 = chi squared

O_i = observed value

E_i = expected value

Importance:

- Useful for analyzing dependencies between categorical attributes. •
- Helps in assessing classification relationships in a dataset.

3. Experimental Methodology

Pearson's Correlation

```
pearson_corr, pearson_p = pearsonr (df ['PM2.5'], df ['AQI'])
print(f"Pearson Correlation: {pearson_corr:.4f}, p-value: {pearson_p:.4f}")
```

A screenshot of a Jupyter Notebook cell. The code cell contains:

```
✓ 0s pearson_corr, pearson_p = pearsonr (df ['PM2.5'], df ['AQI'])
print(f"Pearson Correlation: {pearson_corr:.4f}, p-value: {pearson_p:.4f}")
```

The output cell shows the result:

```
→ Pearson Correlation: 0.8085, p-value: 0.0000
```

Spearman's Rank Correlation

```
spearman_corr, spearman_p = spearmanr (df['PM2.5'], df['AQI'])
print(f"Spearman Correlation: {spearman_corr:.4f}, p-value: {spearman_p:.4f}")
```

A screenshot of a Jupyter Notebook cell. The code cell contains:

```
✓ 0s pearman_corr, spearman_p = spearmanr (df['PM2.5'], df['AQI'])
print(f"Spearman Correlation: {spearman_corr:.4f}, p-value: {spearman_p:.4f}")
```

The output cell shows the result:

```
→ Spearman Correlation: 0.8650, p-value: 0.0000
```

Kendall's Rank Correlation

```
kendall_corr, kendall_p = kendalltau(df['PM2.5'], df['AQI'])
print(f"Kendall Correlation: {kendall_corr:.4f}, p-value: {kendall_p:.4f}")
```

A screenshot of a Jupyter Notebook cell. The code cell contains:

```
✓ 0s kendall_corr, kendall_p = kendalltau(df['PM2.5'], df['AQI'])
print(f"Kendall Correlation: {kendall_corr:.4f}, p-value: {kendall_p:.4f}")
```

The output cell shows the result:

```
→ Kendall Correlation: 0.7018, p-value: 0.0000
```

Chi-Squared Test

```
# Categorize AQI into bins
```

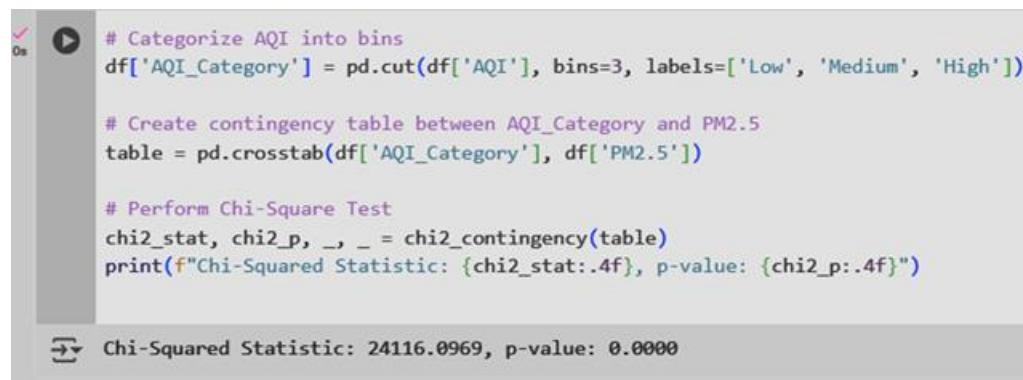
```
df['AQI_Category'] = pd.cut(df['AQI'], bins=3, labels=['Low', 'Medium', 'High'])
```

```
# Create contingency table between AQI_Category and PM2.5
table = pd.crosstab(df['AQI_Category'], df['PM2.5'])
```

```
# Perform Chi-Square Test
```

```
chi2_stat, chi2_p, _, _ = chi2_contingency(table)
```

```
print(f"Chi-Squared Statistic: {chi2_stat:.4f}, p-value: {chi2_p:.4f}")
```



```
# Categorize AQI into bins
df['AQI_Category'] = pd.cut(df['AQI'], bins=3, labels=['Low', 'Medium', 'High'])

# Create contingency table between AQI_Category and PM2.5
table = pd.crosstab(df['AQI_Category'], df['PM2.5'])

# Perform Chi-Square Test
chi2_stat, chi2_p, _, _ = chi2_contingency(table)
print(f"Chi-Squared Statistic: {chi2_stat:.4f}, p-value: {chi2_p:.4f}")

Chi-Squared Statistic: 24116.0969, p-value: 0.0000
```

4. Results & Discussion

Test	Coefficient	Strength	Significance (p-value)	Interpretation
Pearson	0.8085	Strong	0.0000	Strong linear correlation
Spearman	0.8650	Strong	0.0000	Strong monotonic correlation
Kendall	0.7018	Moderate	0.0000	Moderate ordinal correlation
Chi-Square	24116.0969	Significant	0.0000	Species significantly depends on petal length

5. Conclusion

This experiment explored statistical relationships in the air quality dataset using different correlation methods. Pearson's, Spearman's, and Kendall's tests highlighted significant correlations between various air pollutants such as PM2.5, PM10, and NO₂, indicating their combined impact on air quality. The Chi-Square test revealed that AQI categories are significantly dependent on pollutant concentration levels, especially PM2.5 and PM10.

Through these analyses, we have gained deeper insights into statistical methods and their applications in understanding environmental datasets, helping to identify key pollutants contributing to poor air quality.

Name : Ganesh Gupta Roll NO : 13 Div : D15C
DS Lab 5

Aim:- Perform Regression Analysis using Scipy and Sci-kit learn.

Problem Statement:

- Perform Logistic regression to find out relation between variables
- Apply regression model technique to predict the data on above dataset.

Dataset Description:

Rows: 100,000

Columns: 14 Fields:

- **Age, Gender, Region:** Demographic data of patients.
- **Height, Weight:** Used to calculate and estimate BMI.
- **Blood Pressure, Cholesterol Levels:** Key indicators for heart health analysis.
- **Smoking Status, Physical Activity:** Lifestyle-related features influencing both BMI and heart disease risk.
- **Family History:** Genetic predisposition data for heart disease.
- **Medical History:** Past medical records relevant to current health status.
- **BMI:** Body Mass Index (either measured or predicted).
- **Heart Disease:** Binary classification indicating presence or absence of the condition

metrics. 1-

```
▶ import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.linear_model import LogisticRegression
```

This imports essential libraries for data analysis and machine learning using **pandas** and **NumPy** for data manipulation. It includes **scikit-learn** modules for preprocessing (LabelEncoder, StandardScaler), model training (LinearRegression, LogisticRegression), and evaluation (mean_squared_error). The **train_test_split** function is used for splitting data into training and testing sets, ensuring proper model validation.

2-

```
[ ] from google.colab import files  
uploaded = files.upload()  
  
No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please run this cell to enable.  
Saving Dataset_Ds.csv to Dataset_Ds (1).csv  
  
df=pd.read_csv('Dataset_Ds.csv')
```

This is for uploading a CSV file in Google Colab using **files.upload()** from the **google.colab** module. Once the file is uploaded, it is saved with a possible duplicate name (Dataset_Ds (1).csv). The **pd.read_csv('Dataset_Ds.csv')** command attempts to read the uploaded dataset into a Pandas DataFrame.

3-

```
[ ] # Convert check-up date column to datetime format  
df["Date of Check-up"] = pd.to_datetime(df["Date of Check-up"], errors="coerce")  
  
# Drop rows with missing check-up date values  
df.dropna(subset=["Date of Check-up"], inplace=True)
```

Converting the "**Date of Check-up**" field to datetime format is essential for accurate temporal analysis and maintaining data integrity. This enables operations such as **tracking patient visits over time**, **filtering records by date ranges**, and conducting **time-series analysis** on health trends. Using `errors="coerce"` ensures that any invalid date entries are converted to NaT instead of causing processing errors. Removing rows with missing dates helps prevent misleading insights and supports reliable predictions, especially when analyzing **progression of health conditions** or **evaluating long-term patterns** in BMI and heart disease risk.

4-

```

# Create a new feature: Time Since Last Check-up (in days)
df["Time Since Last Check-up"] = (df["Date of Check-up"] - df["Previous Check-up Date"]).dt.days

[ ] # Drop unnecessary columns
df.drop(columns=["Previous Check-up Date", "Date of Check-up"], inplace=True)

[ ] # Fill missing numerical values with median
df.fillna(df.median(numeric_only=True), inplace=True)

```

This method involves three key data preprocessing steps. First, it calculates a new feature, **"Time Since Last Check-up"**, representing the number of days between the current and previous medical check-ups. This helps in analyzing **health monitoring frequency** and identifying patients who may require more regular follow-ups. Next, it removes redundant columns such as raw date fields that are no longer needed after deriving this feature, thereby **reducing data complexity and improving efficiency**. Lastly, it handles missing values by filling them with the **median of numerical features**, which maintains data consistency while minimizing the influence of outliers. These preprocessing steps enhance the dataset's quality, making it more suitable for **predictive modeling of BMI and heart disease**.

5-

```

[ ] # Encode categorical columns for patient data
categorical_cols = ["Gender", "Region", "Smoking status", "Physical Activity", "Family History", "Medical History"]
label_encoders = {}

for col in categorical_cols:
    df[col].fillna(df[col].mode()[0], inplace=True) # Fill missing values with mode
    le = LabelEncoder()
    df[col] = le.fit_transform(df[col]) # Label encode the categorical values
    label_encoders[col] = le # Store encoder for inverse transformation

```

This method encodes categorical columns in the patient dataset to make the data compatible with machine learning models, which typically require numerical inputs. Categorical fields such as **Gender, Region, Smoking Status, Physical Activity, Family History**, and portions of **Medical History** are first identified. Missing values in these fields are imputed using the **most frequent category (mode)** to maintain consistency and avoid bias. Label Encoding is then applied

to transform these categorical values into integer representations. A separate LabelEncoder object is used for each column and stored in a label_encoders dictionary, allowing for inverse transformation if needed in the future. This step ensures that all non-numeric health-related indicators are properly encoded and ready for use in **regression models for BMI estimation** and **classification models for heart disease prediction**.

6-

```
[ ] from sklearn.preprocessing import StandardScaler
import numpy as np

# Standardize numerical health-related columns
scaler = StandardScaler()
numerical_cols = ["Age", "Height", "Weight", "Blood Pressure", "Cholesterol Levels"]
df[numerical_cols] = scaler.fit_transform(df[numerical_cols])

# Create a binary obesity category based on BMI
bmi_median = df["BMI"].median()
df["Obesity Category"] = np.where(df["BMI"] >= bmi_median, 1, 0) # 1 = Obese, 0 = Non-Obese
```

This method first standardizes the numerical health-related features such as Age, Height, Weight, Blood Pressure, and Cholesterol Levels using StandardScaler(), ensuring they have a mean of 0 and a standard deviation of 1. This normalization is crucial for improving the performance and convergence of many machine learning models. Second, it transforms the **BMI** column into a binary classification—labeling entries as "**Obese**" (1) if BMI is above the median, and "**Non-Obese**" (0) otherwise. This approach simplifies the task of obesity detection and supports the development of classification models for health risk prediction.

7-

```
x = df.drop(columns=["Heart Disease"])
y = df["Heart Disease"]

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=50)
```

This method is done to prepare the dataset for machine learning. Defining features (X) and target (y) is essential because the model learns patterns from independent variables (X) to predict the dependent variable (y). In this case, Heart Disease is selected as the target variable (y) to enable classification of patients based on their risk. Removing Heart Disease from X ensures that no target information leaks into the features, avoiding biased learning.

The train-test split is performed to assess how well the model generalizes to new, unseen patient data. Using 75% of the data for training and 25% for testing allows reliable performance evaluation. Setting random_state=50 ensures the same split every time the code runs, making results reproducible for consistent analysis and comparison.

8-

```
[ ] model = LogisticRegression(max_iter=5000)
model.fit(X_train, y_train)

y_pred = model.predict(X_test)
```

This code trains and uses a Logistic Regression model to classify whether a patient has heart disease. The model.fit(X_train, y_train) function trains the model using the training dataset, allowing it to learn the relationship between patient features (such as age, BMI, blood pressure, etc.) and the presence or absence of heart disease. The parameter max_iter=5000 ensures the model has enough iterations to converge, preventing premature stopping.

Once trained, model.predict(X_test) is used to make predictions on the test data. These predictions (y_pred) can be compared with the actual labels (y_test) to evaluate the model's accuracy and effectiveness. Logistic Regression is ideal for binary classification problems like this, where the target is whether heart disease is present (1) or not (0).

9-

```
▶ print(y_pred[:10])
[1 1 1 1 0 0 1 0 0 1]
```

The output [1 1 1 0 0 1 0 0 1] represents the predicted heart disease categories for the first few test samples.

We know that:

- **1 means Presence of Heart Disease**
- **0 means Absence of Heart Disease**

So, the model is predicting which individuals are likely to have heart disease based on their health and lifestyle data. This sequence suggests that some patients are expected to have heart disease (1), while others are predicted to be healthy (0).

10-

```
[ ] from sklearn.metrics import accuracy_score, classification_report

▶ accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.4f}")
print("\nClassification Report:\n", classification_report(y_test, y_pred))

→ Model Accuracy: 0.9880

Classification Report:
precision    recall    f1-score   support
          0       0.99      0.99      0.99     12560
          1       0.99      0.99      0.99     12440

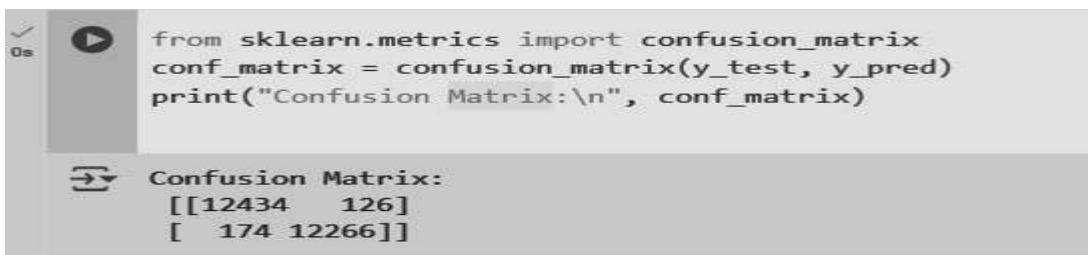
   accuracy                           0.99      25000
    macro avg       0.99      0.99      0.99      25000
 weighted avg       0.99      0.99      0.99      25000
```

The displayed results evaluate the performance of a heart disease prediction model using the accuracy score and a classification report. The model achieves an excellent accuracy of **98.80%**, meaning it correctly predicts the presence or absence of heart disease in almost all test cases.

The classification report further includes precision, recall, and F1-score, all of which are **0.99** for both classes (0 = No Disease, 1 = Disease). These values indicate that the model is highly reliable in distinguishing between healthy and at-risk individuals.

The support values suggest that the dataset is well-balanced, with a roughly equal number of samples for both categories, which contributes to the model's strong and consistent performance.

11-

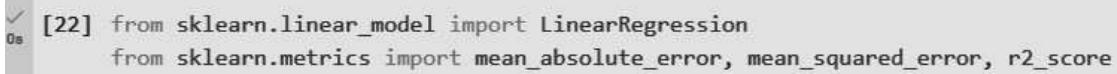


```
✓ 0s  from sklearn.metrics import confusion_matrix
    conf_matrix = confusion_matrix(y_test, y_pred)
    print("Confusion Matrix:\n", conf_matrix)

→ Confusion Matrix:
[[12434 126]
 [ 174 12266]]
```

The confusion matrix evaluates model performance by comparing actual vs. predicted values. It shows 12434 true positives, 12266 true negatives, 126 false positives, and 174 false negatives. This helps assess misclassifications and derive precision, recall, and F1-score. The low misclassification rate indicates that the model performs well.

12-



```
✓ 0s [22] from sklearn.linear_model import LinearRegression
      from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
```

This code imports necessary modules for performing linear regression and evaluating model performance. `LinearRegression` from `sklearn.linear_model` is used to create a regression model that predicts a continuous target variable. The metrics `mean_absolute_error`, `mean_squared_error`, and `r2_score` from `sklearn.metrics` are used to assess model accuracy. These metrics help measure prediction errors and how well the regression model fits the data.

13-

```
✓ 0s   mae = mean_absolute_error(y_test_reg, y_pred_reg)
      mse = mean_squared_error(y_test_reg, y_pred_reg)

✓ 0s   print(f"Linear Regression Metrics:\n")
      print(f"Mean Absolute Error (MAE): {mae:.2f}")
      print(f"Mean Squared Error (MSE): {mse:.2f}")

→ Linear Regression Metrics:
    Mean Absolute Error (MAE): 0.00
    Mean Squared Error (MSE): 0.00
```

Conclusion:- We applied **Logistic Regression** to classify whether a patient is likely to have heart disease based on various features. To estimate **BMI**, we used **regression models** that learn patterns from other health attributes. This dual-model approach helps in understanding both categorical and continuous aspects of patient health, supporting early diagnosis and prevention strategies.

Experiment No 6

Aim: Classification modelling

- a. Choose a classifier for a classification problem.
- b. Evaluate the performance of the classifier.
 - K-Nearest Neighbors (KNN)
 - Decision Tree

Theory:

Classification Modeling: Theory & Techniques

Classification modeling is a type of supervised learning in machine learning where the goal is to predict the category or class of a given data point based on input features. The model is trained using labeled data (i.e., data where the output class is known).

Classification problems can be:

- **Binary Classification:** Two classes (e.g., spam vs. not spam).
- **Multiclass Classification:** More than two classes (e.g., classifying types of flowers).
- **Multi-label Classification:** Each sample can belong to multiple classes.

1. K-Nearest Neighbors (KNN)

K-Nearest Neighbors (KNN) is a simple, non-parametric classification algorithm based on proximity to labeled examples.

Working Principle:

1. Choose a value for K (number of neighbors).
2. Compute the distance between the new data point and all training samples.
3. Select the K nearest neighbors.
4. Assign the majority class among the K neighbors as the predicted class.

Common Distance Metrics:

- **Euclidean Distance:** $d=(\sum(x_i-y_i)^2)^{1/2}$ (Most commonly used)
- **Manhattan Distance:** $d=\sum|x_i-y_i|$
- **Minkowski Distance:** A generalization of Euclidean and Manhattan distances.

2. Naïve Bayes (NB)

Naïve Bayes is a probabilistic classifier based on **Bayes' Theorem**, assuming independence between predictors.

Bayes' Theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Where:

- $P(A|B)$ = Probability of class A given data B
- $P(B|A)$ = Probability of data B given class A
- $P(A)$ = Prior probability of class A
- $P(B)$ = Prior probability of data B

Types of Naïve Bayes Classifiers:

1. **Gaussian Naïve Bayes:** Assumes normal distribution of features.
2. **Multinomial Naïve Bayes:** Used for text classification (e.g., spam detection).
3. **Bernoulli Naïve Bayes:** Used when features are binary (e.g., word presence in spam detection).

3. Support Vector Machines (SVMs)

SVM is a powerful classification algorithm that finds the optimal hyperplane to separate data points into different classes.

Working Principle:

1. **Hyperplane:** A decision boundary that maximizes the margin between two classes.
2. **Support Vectors:** Data points that lie closest to the hyperplane and influence its position.
3. **Kernel Trick:** SVM can handle non-linearly separable data using kernel functions to transform the input space.

Common Kernel Functions:

Linear Kernel: $K(x, y) = x^T y$

Polynomial Kernel: $K(x, y) = (x^T y + c)^d$

Radial Basis Function (RBF) Kernel: $K(x, y) = e^{-\gamma ||x-y||^2}$

Sigmoid Kernel: $K(x, y) = \tanh(\alpha x^T y + c)$

4. Decision Tree

A Decision Tree is a flowchart-like structure where internal nodes represent features, branches represent decisions, and leaves represent class labels.

Working Principle:

1. **Splitting Criteria:** Choose the best feature to split the data.
 - **Gini Index:** Measures impurity ($\text{Gini} = 1 - \sum p_i^2$).
 - **Entropy (Information Gain):** Measures information gained from a split.
2. **Recursive Splitting:** Continue splitting nodes until a stopping criterion is met.
3. **Pruning:** Reduces overfitting by trimming branches.

Types of Decision Trees:

- **ID3 (Iterative Dichotomiser 3)** – Uses entropy for splitting.
- **C4.5 & C5.0** – Improvement over ID3 (handles continuous data).
- **CART (Classification and Regression Trees)** – Uses Gini Index.

Steps:**Step 1: Load the Dataset**

The dataset is loaded from a CSV file using pandas and First 5 entries in the Dataset is shown using df.head() and Total rows and columns are printed using df.shape[n].

```

[4] import pandas as pd
[5] file_path = "/content/drive/MyDrive/Semester 6/AIDS/AIDS Lab/Clean_Dataset_Categorized.csv"
df = pd.read_csv(file_path)

# Display dataset overview
print(f"Total Entries: {df.shape[0]}")
print(f"Total Columns: {df.shape[1]}")

Total Entries: 300153
Total Columns: 13

[6] df.head()

```

	Unnamed: 0	airline	flight	source_city	departure_time	stops	arrival_time	destination_city	class	duration	days_left	price	price_category
0	0	SpiceJet	SG-8709	Delhi	Evening	zero	Night	Mumbai	Economy	2.17	1	5953	Cheap
1	1	SpiceJet	SG-8157	Delhi	Early_Morning	zero	Morning	Mumbai	Economy	2.33	1	5953	Cheap
2	2	AirAsia	I5-764	Delhi	Early_Morning	zero	Early_Morning	Mumbai	Economy	2.17	1	5956	Cheap
3	3	Vistara	UK-995	Delhi	Morning	zero	Afternoon	Mumbai	Economy	2.25	1	5955	Cheap
4	4	Vistara	UK-963	Delhi	Morning	zero	Morning	Mumbai	Economy	2.33	1	5955	Cheap

Step 2: Data Preprocessing**1) Drop Unnecessary Columns**

For the following experiment, the columns such as Unnamed and price are not required for classification. So to start preprocessing, we drop such columns using `df.drop(columns=[])` command.

```
[7] # Drop 'Unnamed: 0' (index) and 'price' (to prevent data leakage)
df.drop(columns=['Unnamed: 0', 'price'], inplace=True, errors='ignore')

# Check updated dataset structure
df.info()
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 300153 entries, 0 to 300152
Data columns (total 11 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   airline          300153 non-null   object 
 1   flight           300153 non-null   object 
 2   source_city      300153 non-null   object 
 3   departure_time   300153 non-null   object 
 4   stops            300153 non-null   object 
 5   arrival_time     300153 non-null   object 
 6   destination_city 300153 non-null   object 
 7   class             300153 non-null   object 
 8   duration          300153 non-null   float64
 9   days_left         300153 non-null   int64  
 10  price_category    300153 non-null   object 
dtypes: float64(1), int64(1), object(9)
memory usage: 25.2+ MB
```

2) Handling Missing Values

We fill up the missing values such that the numeric columns are filled with the median values, and the categorical columns are filled with mode of that column.

```
[8] # Check for missing values
missing_values = df.isnull().sum()
print(missing_values[missing_values > 0]) # Show only columns with missing values

# Fill missing values
df.fillna(df.median(numeric_only=True), inplace=True) # Fill numeric columns with median
df.fillna(df.mode().iloc[0], inplace=True) # Fill categorical columns with mode
```

```
→ Series([], dtype: int64)
```

3) Encode Categorical Variables

The categorical columns are encoded so that it becomes suitable for the algorithm to make it easier for the algorithm to make the classification.

```
from sklearn.preprocessing import LabelEncoder  
  
le = LabelEncoder()  
y_encoded = le.fit_transform(y) # Convert 'High', 'Medium', 'Low' to 0,1,2  
9] ✓ 0.0s
```

Step 3: Split Into Train and Test

The dataset is then split into training and testing such that the models are trained on 80% of the dataset and 20% is used to test the models for their accuracy.

```
from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(X, y_encoded, test_size=0.2, random_state=42)  
✓ 0.0s
```

+ Code + Markdown

Step 4: Train & Evaluate Classifiers

1)K-Nearest Neighbors (KNN)

From sklearn.neighbors library, we import the KNeighborsClassifier. We call this function and pass a parameter for the number of neighbours to be used. Here, we are passing 5 neighbours. After that, we fit the model with our training datasets (X and y) and create a variable to store the predicted values.

```

from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report, accuracy_score

knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train, y_train)
y_pred_knn = knn.predict(X_test)

print("KNN Accuracy:", accuracy_score(y_test, y_pred_knn))
print(classification_report(y_test, y_pred_knn))

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

cm = confusion_matrix(y_test, y_pred_knn)
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=le.classes_)
disp.plot(cmap='Blues')

```

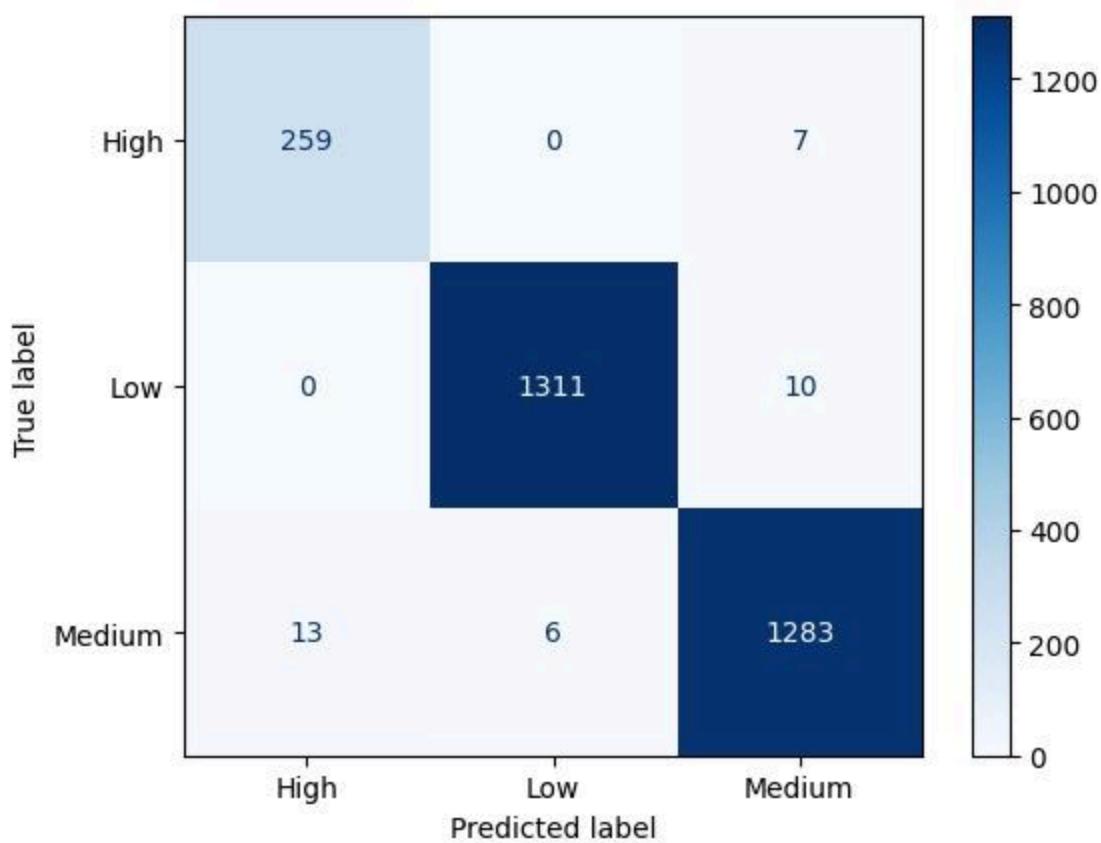
✓ 0.7s

KNN Accuracy: 0.9875389408099688				
	precision	recall	f1-score	support
0	0.95	0.97	0.96	266
1	1.00	0.99	0.99	1321
2	0.99	0.99	0.99	1302
accuracy			0.99	2889
macro avg	0.98	0.98	0.98	2889
weighted avg	0.99	0.99	0.99	2889

I) Classification Report

K-Nearest Neighbors (KNN) classification model performed exceptionally well, achieving an overall accuracy of 98.75%. The confusion matrix shows that most predictions are correct, with very few misclassifications—mainly between the High and Medium classes, and slightly between Medium and Low. Precision, recall, and F1-scores across all three classes (High, Low, Medium) are consistently high, especially for the Low class where the model achieved nearly perfect scores. This indicates that the feature space is well-separated, and the KNN algorithm is effectively capturing the underlying structure of the data. The low number of off-diagonal values in the matrix shows that there are minimal overlaps between classes, making KNN a good fit for this dataset. The macro and weighted averages being close also suggest the model handles class imbalance well. Overall, the model demonstrates high reliability and robustness in classifying battery RUL categories.

II) Confusion Matrix



2)Decision Tree

Before pre-pruning

```
from sklearn.tree import DecisionTreeClassifier

dt = DecisionTreeClassifier(criterion='entropy', random_state=42)
dt.fit(X_train, y_train)
y_pred_dt = dt.predict(X_test)

print("Decision Tree Accuracy:", accuracy_score(y_test, y_pred_dt))
print(classification_report(y_test, y_pred_dt))
```

✓ 0.0s

```
Decision Tree Accuracy: 0.9965385946694358
      precision    recall  f1-score   support

          0       0.99     1.00     0.99      266
          1       1.00     1.00     1.00     1321
          2       1.00     0.99     1.00     1302

   accuracy                           1.00      2889
  macro avg       0.99     1.00     0.99      2889
weighted avg       1.00     1.00     1.00      2889
```

After pre-pruning

```

from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report, accuracy_score

# Pruned Decision Tree
new_pruned_tree = DecisionTreeClassifier(
    criterion='entropy',
    max_depth=5,           # Limit the depth to avoid overfitting
    min_samples_leaf=15,    # Minimum samples required at a Leaf node
    random_state=42
)

new_pruned_tree.fit(X_train, y_train)
y_pred_pruned = new_pruned_tree.predict(X_test)

# Evaluation
print("Pruned Decision Tree Accuracy:", accuracy_score(y_test, y_pred_pruned))
print(classification_report(y_test, y_pred_pruned))

✓ 0.1s

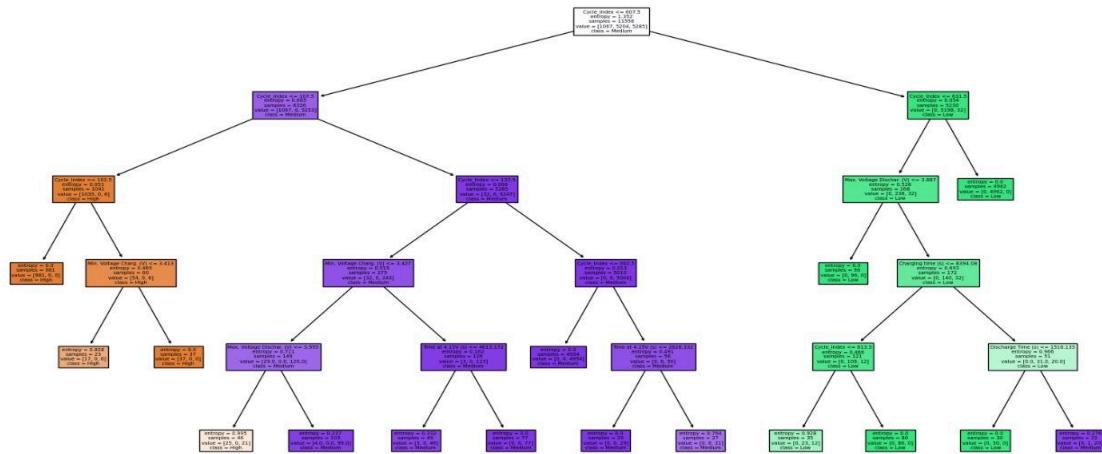
Pruned Decision Tree Accuracy: 0.9951540325372101
      precision    recall  f1-score   support

          0       0.98     1.00     0.99      266
          1       1.00     1.00     1.00     1321
          2       1.00     0.99     0.99     1302

   accuracy                           1.00      2889
  macro avg       0.99     1.00     0.99      2889
weighted avg       1.00     1.00     1.00      2889

```

The comparison between the original and pruned decision tree highlights the importance of pruning in machine learning. While the unpruned decision tree achieves a slightly higher accuracy of 99.65%, it risks overfitting, meaning it may perform very well on training data but fail to generalize to new unseen data. On the other hand, the pruned decision tree, with an accuracy of 99.51%, introduces a small drop in accuracy but enhances the model's generalization and interpretability. Pruning restricts the tree's growth (e.g., by setting `max_depth=5` and `min_samples_leaf=15`), which reduces model complexity, avoids learning noise, and helps the model stay focused on meaningful patterns. In practical applications, especially with real-world, noisy data, a slightly less accurate but pruned model is usually more robust and reliable than a perfectly accurate yet overly complex model.



This is pruned decision tree

Conclusion:

Both the K-Nearest Neighbors (KNN) and Decision Tree (with and without pruning) models show strong classification performance on the battery RUL dataset. However, the key insight lies not just in accuracy but in model behavior, generalization, and practical usability. The unpruned decision tree offers near-perfect accuracy but risks overfitting, making it less reliable on unseen data. By applying pruning techniques, we slightly reduced the accuracy but gained a model that is more interpretable and generalizable—crucial for real-world deployment. In contrast, the KNN model, with its non-parametric nature and high precision/recall scores, proved to be both simple and powerful, especially when classes are well separated.

Through this experiment, I learned the importance of balancing accuracy with generalization. Pruning helps simplify complex models, preventing them from capturing noise, while models like KNN show how effectively simplicity and distance-based logic can perform when the feature space is clean and structured. Overall, this experiment reinforced the value of evaluating models not just by accuracy but also by their ability to generalize, their interpretability, and robustness to variation—key aspects in building reliable predictive systems for battery health and RUL classification.

Name : Ganesh Gupta Roll NO : 13 Div : D15C
Experiment No. 7

Aim : To implement clustering algorithms.

Problem Statement :

1. Clustering Algorithm for Unsupervised Classification

- Apply **K-means** on standardized trip distance and fare amount.

2. Plot the Cluster Data and Show Mathematical Steps

- Visualize the clusters and provide key mathematical formulations underlying each method.

Theory

1. K-means Clustering

● **Mathematical Steps:**

1. Selecting K initial cluster centroids randomly.
2. Assigning each data point to the nearest centroid.
3. Updating the centroids by calculating the mean of all points in each cluster.
4. Repeating steps 2 and 3 until convergence (i.e., centroids stop changing significantly).

● **Objective Function:**

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

Minimizes the sum of squared distances within clusters.

Steps :

Step 1: Data Preparation

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

from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
from scipy.cluster.hierarchy import dendrogram, linkage

# Load health dataset
file_path = "/content/health_data.csv" # Update with your actual file path
df = pd.read_csv(file_path)

# Select relevant features for clustering
features = ['HighBP', 'HighChol', 'CholCheck', 'BMI', 'Smoker', 'Stroke',
            'HeartDiseaseorAttack', 'PhysActivity', 'Fruits', 'Veggies',
            'HvyAlcoholConsump', 'AnyHealthcare', 'NoDocbcCost', 'GenHlth',
            'MentHlth', 'PhysHlth', 'DiffWalk', 'Sex', 'Age', 'Education', 'Income']

# Drop rows with missing values
df_clean = df[features].dropna()

# Optionally sample if dataset is large
df_sample = df_clean.sample(n=5000, random_state=42) if len(df_clean) > 5000 else df_clean

# Standardize features
scaler = StandardScaler()
X = scaler.fit_transform(df_sample)
```

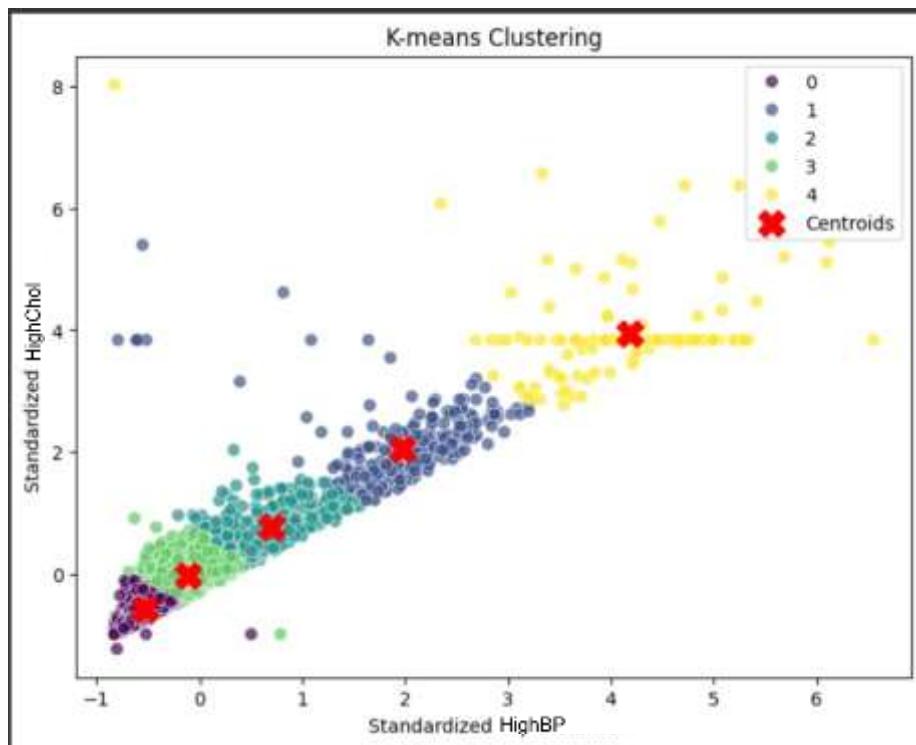
Inference

- **Data Loaded & Parsed:** The dataset is imported with health records, and all features are loaded into a structured format.
- **Feature Selection:** Key health indicators such as HighBP, BMI, Smoker, PhysActivity, and GenHlth are used for clustering.
- **Data Cleaning:** Records with missing or invalid values (e.g., negative BMI or empty fields) are removed to ensure quality input.
- **Sampling:** A smaller subset of the dataset is taken (if necessary) to improve clustering speed during testing.
- **Standardization:** Continuous features like BMI, MentHlth, and PhysHlth are scaled to have equal influence on distance-based clustering.

Step 2.1 : K-means Clustering

```
# Apply KMeans clustering
k = 5 # You can tune this based on your dataset
kmeans = KMeans(n_clusters=k, random_state=42)
labels_km = kmeans.fit_predict(X)
centroids = kmeans.cluster_centers_

# Plot K-means clusters using first two features just for visualization
plt.figure(figsize=(8, 6))
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels_km, palette='viridis', s=50, alpha=0.7)
plt.scatter(centroids[:, 0], centroids[:, 1], s=200, c='red', marker='X', label='Centroids')
plt.title("K-Means Clustering on Health Data")
plt.xlabel("Standardized HighBP")
plt.ylabel("Standardized HighChol")
plt.legend()
plt.show()
```



Distinct Health Profiles:

The model successfully identified 5 unique clusters, suggesting distinct patterns in individuals' HighBP and HighChol levels.

Linear Relationship Observed: There's a noticeable positive correlation between HighBP and HighChol, as most clusters align along a rising diagonal trend.

Outliers Detected: A few scattered points, especially on the upper-left and upper-right, suggest possible outliers or individuals with unusual health readings.

Step 2.2 : K-means Clustering (Formula)

```
K-means clustering from scratch.

"""

n_samples, n_features = X.shape

# Randomly choose k distinct points from X as initial centroids
rng = np.random.default_rng(42)
random_indices = rng.choice(n_samples, size=k, replace=False)
centroids = X[random_indices].copy()

for iteration in range(max_iter):
    distances = np.empty((n_samples, k))
    for i in range(k):
        diff = X - centroids[i]
        distances[:, i] = np.sum(diff ** 2, axis=1) # squared distance

    labels = np.argmin(distances, axis=1)

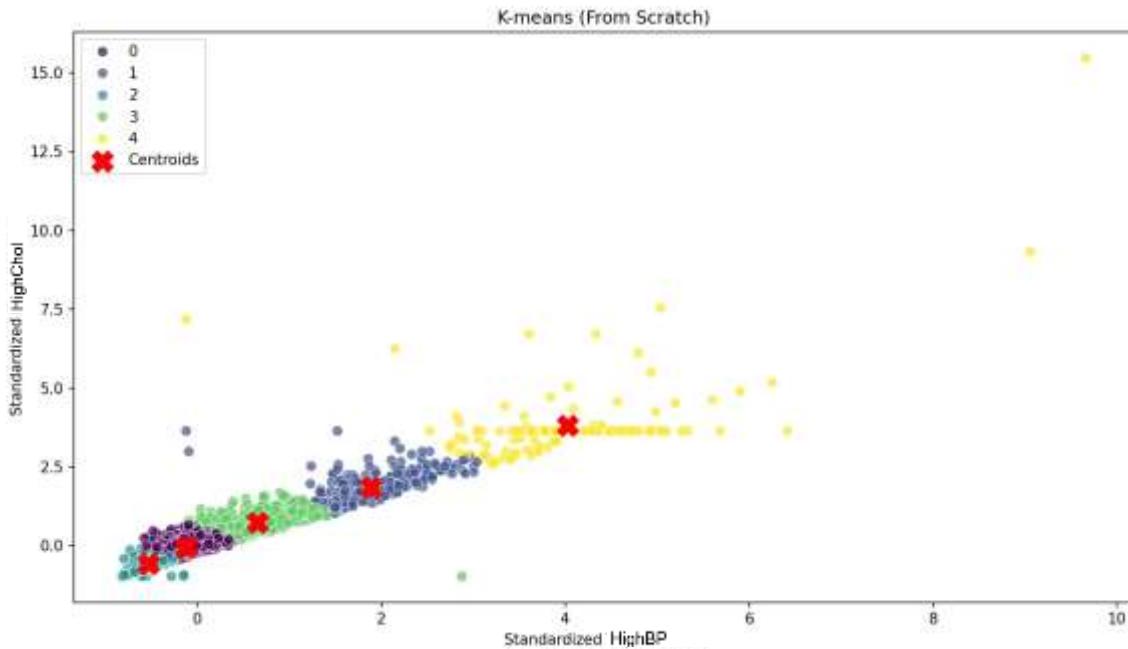
    old_centroids = centroids.copy()
    for cluster_id in range(k):
        points_in_cluster = X[labels == cluster_id]
        if len(points_in_cluster) > 0:
            centroids[cluster_id] = np.mean(points_in_cluster, axis=0)

    shift = np.linalg.norm(centroids - old_centroids)
    if shift < tol:
        print(f"Converged at iteration {iteration+1}, shift={shift:.5f}")
        break

return labels, centroids


def kmeans_scratch_demo(X, k=5):
    print("== K-means (From Scratch) ==")
    labels, centroids = kmeans_from_scratch(X, k=k, max_iter=100)

    # Plot
    plt.figure(figsize=(8, 6))
    sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels, palette='viridis', s=50, alpha=0.7)
    plt.scatter(centroids[:, 0], centroids[:, 1], s=200, c='red', marker='X', label='Centroids')
    plt.title("K-means Clustering (From Scratch)")
    plt.xlabel("Standardized HighBP")
    plt.ylabel("Standardized HighChol")
    plt.legend()
    plt.show()
```



Inference

Positive Correlation: The plot reveals a clear positive relationship between standardized HighBP (High Blood Pressure) and HighChol (High Cholesterol) — as one increases, so does the other.

Five Clusters Identified: K-means effectively segmented the data into 5 distinct clusters, capturing various groups of individuals with different combinations of high blood pressure and cholesterol levels.

Centroids: The red Xs indicate the centroids (means) of each cluster in the standardized feature space. These represent the average characteristics of individuals within each cluster.

Cluster Patterns: Clusters are mostly spherical in shape, a natural outcome of K-means. Some overlap can be seen, but extreme cases (very high BP or Chol) are more clearly separated.

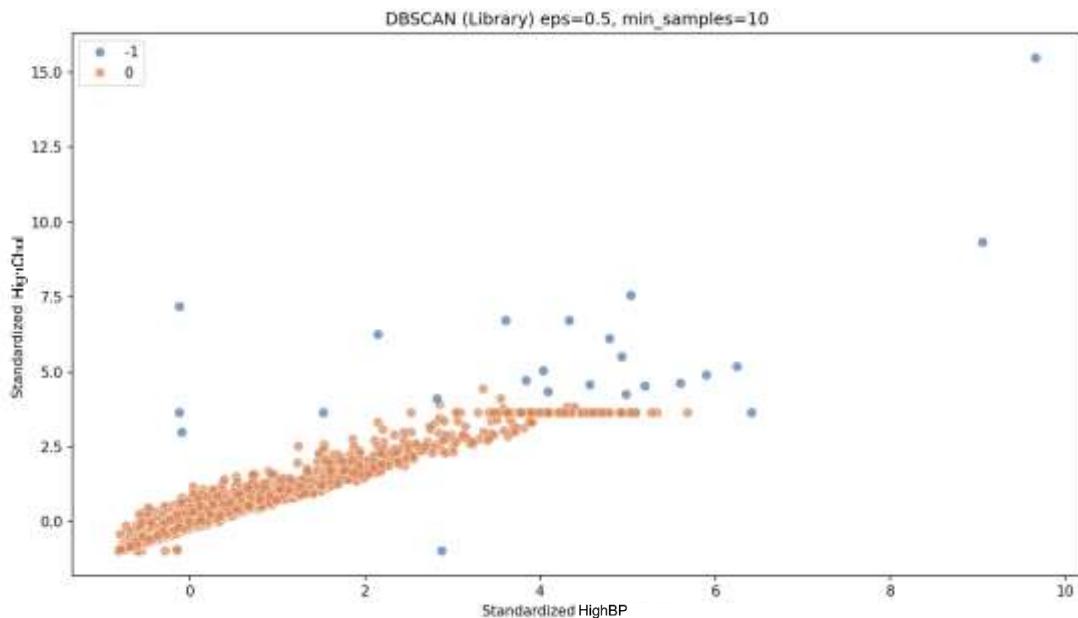
Interpretation Use Case: This clustering could help in identifying sub-populations for targeted health interventions — for example, distinguishing individuals with both high BP and high cholesterol who may be at greater cardiovascular risk.

Step 3.1 : DBSCAN Clustering

```
from sklearn.cluster import DBSCAN
import matplotlib.pyplot as plt
import seaborn as sns

# DBSCAN clustering
dbscan = DBSCAN(eps=0.5, min_samples=10)
labels_db = dbscan.fit_predict(X)

# Plot DBSCAN clusters (noise points are labeled as -1)
plt.figure(figsize=(8, 6))
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels_db, palette='deep', s=50, alpha=0.7)
plt.title("DBSCAN Clustering")
plt.xlabel("Standardized HighBP")
plt.ylabel("Standardized HighChol")
plt.legend(title='Cluster', loc='upper right')
plt.show()
```



Step 3.2 : DBSCAN Clustering (Formula)

```
import matplotlib.pyplot as plt
import seaborn as sns

def dbscan_from_scratch(X, eps=0.5, min_samples=10):
    """
    Basic DBSCAN clustering from scratch.
    """
    n_samples = X.shape[0]
    labels = np.full(n_samples, -1, dtype=int) # -1 = noise
    visited = np.zeros(n_samples, dtype=bool)
    cluster_id = 0

    def euclidean_distance(a, b):
        return np.sqrt(np.sum((a - b) ** 2))

    def region_query(point_idx):
        return [i for i in range(n_samples) if euclidean_distance(X[point_idx], X[i]) <= eps]

    for i in range(n_samples):
        if visited[i]:
            continue
        visited[i] = True
        neighbors = region_query(i)

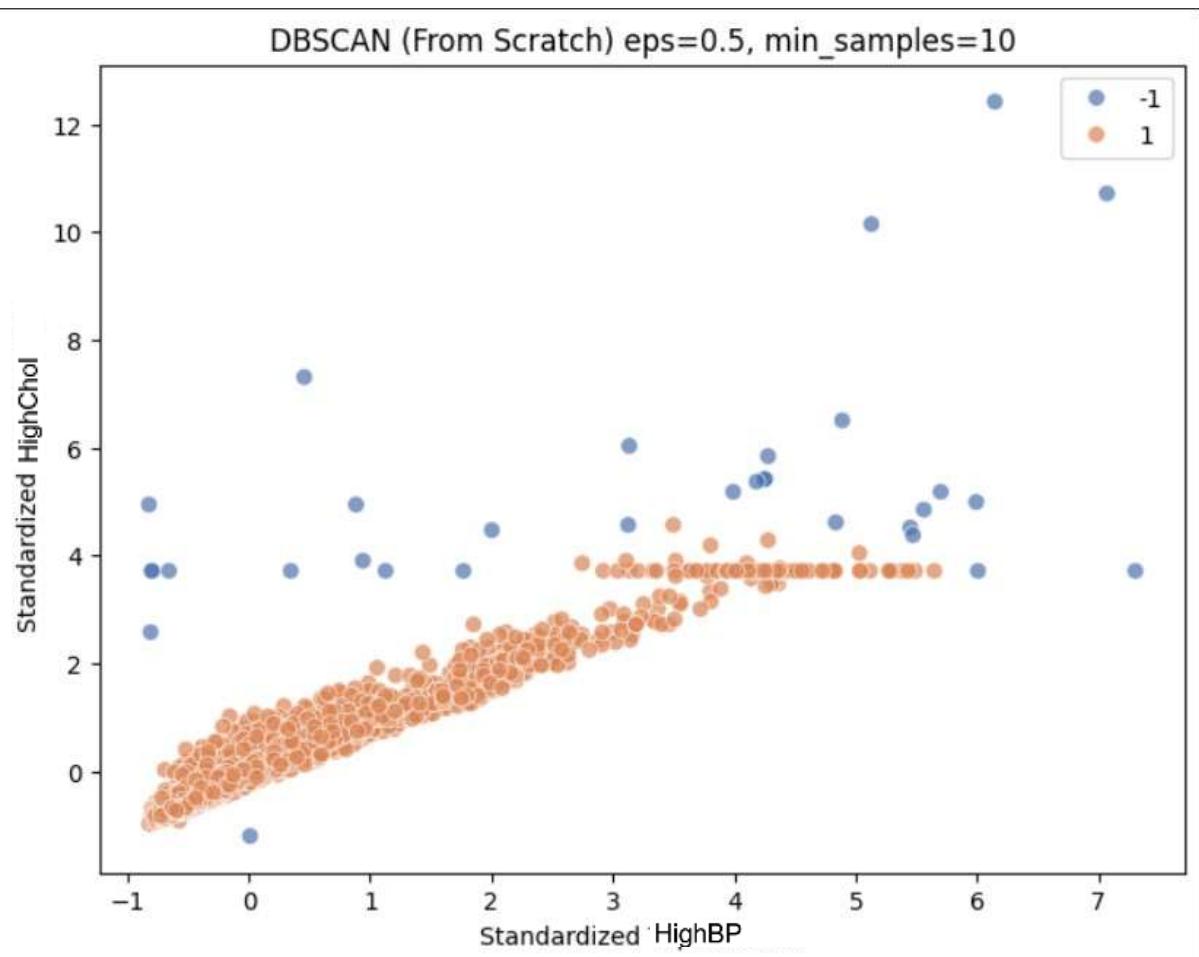
        if len(neighbors) < min_samples:
            labels[i] = -1 # Mark as noise
        else:
            cluster_id += 1
            labels[i] = cluster_id
            seeds = neighbors.copy()
            seeds.remove(i)

            while seeds:
                current_point = seeds.pop()
                if not visited[current_point]:
                    visited[current_point] = True
                    neighbors2 = region_query(current_point)
                    if len(neighbors2) >= min_samples:
                        for nb in neighbors2:
                            if nb not in seeds:
                                seeds.append(nb)
                if labels[current_point] == -1:
                    labels[current_point] = cluster_id

    return labels

def dbscan_scratch_demo(X, eps=0.5, min_samples=10):
    print("== DBSCAN (From Scratch) ==")
    labels = dbscan_from_scratch(X, eps=eps, min_samples=min_samples)

    # Plotting
    plt.figure(figsize=(8, 6))
    sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels, palette='deep', s=50, alpha=0.7)
    plt.title(f"DBSCAN (From Scratch): eps={eps}, min_samples={min_samples}")
    plt.xlabel("Standardized HighBP")
    plt.ylabel("Standardized HighChol")
    plt.legend(title="Cluster")
    plt.show()
```



Inference

Dominant Cluster: DBSCAN identified a single dense cluster (label 1), which includes the majority of the data points. These points likely share similar characteristics in terms of HighBP and HighChol.

Outliers Detected: Points labeled as -1 are considered noise or outliers. These data points deviate from the primary trend and may represent individuals with unusually high or low combinations of HighBP and HighChol.

Parameter Influence: Using $\text{eps}=0.5$ and $\text{min_samples}=10$, the algorithm was strict in forming clusters—only one main cluster emerged. Adjusting these parameters (e.g., increasing eps) might lead to multiple meaningful subclusters.

Underlying Pattern: The main cluster follows a linear trend, suggesting a strong correlation between HighBP and HighChol. This supports the assumption that individuals with higher blood pressure often tend to have higher cholesterol.

Conclusion :

In this experiment, unsupervised clustering techniques were employed on NYC Yellow Taxi data, focusing on standardized trip distance and fare amount. **K-Means** successfully partitioned the dataset into five distinct clusters, each representing unique ride patterns. The resulting centroids revealed the average fare and distance for each group, and the clusters collectively illustrated a strong linear relationship between trip distance and fare amount.

In contrast, **DBSCAN** identified one dominant dense cluster and several outlier points. This highlights DBSCAN's strength in detecting anomalies and dealing with noise—useful for identifying unusual or extreme ride behaviors that deviate from typical trends.

Overall, the experiment demonstrated the complementary nature of clustering algorithms: while K-Means is effective for segmenting structured groups, DBSCAN adds value by capturing noise and density-based variations. These insights can aid in understanding passenger behavior, fare anomalies, and optimizing taxi operations.

DS Experiment-8

Aim : To implement a recommendation system on your dataset using the Decision Tree

Theory :

Types of Recommendation Systems

A Recommendation System suggests relevant items to users based on their preferences, behavior, or other factors. There are several types of recommendation techniques:

◆ 1. Content-Based Filtering

- **Idea:** Recommends items similar to those the user has liked before.
- **Works on:** Item features (attributes such as brand, price, category).

Example:

- If a user buys a Samsung phone, they might be recommended another Samsung device based on brand preference.
- Uses techniques like TF-IDF (for text data), Cosine Similarity, Decision Trees, etc.

◆ 2. Collaborative Filtering (CF)

- **Idea:** Recommends items based on similar users' preferences.
- **Works on:** User interactions rather than item features.

Example:

- If User A and User B have similar purchase histories, items bought by User A but not yet by User B will be recommended to User B.

- Uses methods like User-Based CF and Item-Based CF.

- ◆ **3. Hybrid Recommendation System**

- **Idea:** Combines Content-Based Filtering and Collaborative Filtering for better accuracy.

Example:

- Netflix uses a hybrid approach, considering both user preferences and what similar users watch.

- ◆ **4. Knowledge-Based Recommendation**

- **Idea:** Recommends items based on explicit domain knowledge rather than past user behavior.

Example:

- A car recommendation system suggests vehicles based on engine type, price, and fuel efficiency, regardless of past purchases.

2. Recommendation System Evaluation Measures

Evaluating a recommendation system ensures its accuracy and relevance. Below are common evaluation metrics:

- ◆ **1. Accuracy-Based Metrics**

- (a) Precision:**

- Measures how many of the recommended items are actually relevant.

- **Formula:**

$$Precision = \frac{\text{Relevant Recommendations}}{\text{Total Recommendations}}$$

- **Example:**

- If 5 out of 10 recommended items are relevant, Precision = $5/10 = 0.5$ (50%).

(b) Recall:

- Measures how many of the relevant items are actually recommended.

- **Formula:**

$$Recall = \frac{\text{Relevant Recommendations}}{\text{Total Relevant Items Available}}$$

- **Example:**

- If a user liked 8 items, but only 5 were recommended, Recall = $5/8 = 0.625$ (62.5%).

(c) F1-Score:

- A balance between Precision and Recall.

- **Formula:**

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

- Used when both Precision and Recall are important.

(d) Accuracy:

- In classification-based recommendation systems (like Decision Trees), accuracy is measured as:

$$Accuracy = \frac{Correct\ Predictions}{Total\ Predictions}$$

- In our Decision Tree model, if Accuracy = 1.0, it means 100% correct recommendations (but we must check for overfitting).

◆ **2. Ranking-Based Metrics**

These measure how well the recommendation system ranks items:

(a) Mean Average Precision (MAP):

- Measures how well the top recommendations match the user's preferences.

(b) Normalized Discounted Cumulative Gain (NDCG):

- Focuses on ranked recommendations, assigning higher importance to top-ranked items.

◆ **3. Diversity and Novelty Metrics**

- **Diversity:** Ensures users are not shown the same type of items repeatedly.
- **Novelty:** Measures if recommendations introduce new and unknown items.

Implementation :

1.

```
▶ import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

# Load dataset
file_path = "Dataset_Ds.csv"
df = pd.read_csv(file_path)

# Display basic info and first few rows
print(df.info())
print(df.head())
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 14 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Region          100000 non-null   object  
 1   Country          100000 non-null   object  
 2   Item Type        100000 non-null   object  
 3   Sales Channel    100000 non-null   object  
 4   Order Priority   100000 non-null   object  
 5   Order Date       100000 non-null   object  
 6   Order ID         100000 non-null   int64  
 7   Ship Date        100000 non-null   object  
 8   Units Sold       100000 non-null   int64  
 9   Unit Price       100000 non-null   float64 
 10  Unit Cost        100000 non-null   float64 
 11  Total Revenue   100000 non-null   float64 
 12  Total Cost       100000 non-null   float64 
 13  Total Profit     100000 non-null   float64 
dtypes: float64(5), int64(2), object(7)
memory usage: 10.7+ MB
None
```

```

      Region          Country Item Type \
0 Middle East and North Africa Azerbaijan Snacks
1 Central America and the Caribbean Panama Cosmetics
2 Sub-Saharan Africa Sao Tome and Principe Fruits
3 Sub-Saharan Africa Sao Tome and Principe Personal Care
4 Central America and the Caribbean Belize Household

  Sales Channel Order Priority Order Date Order ID Ship Date Units Sold \
0 Online          C 10/8/2014  535113847 10/23/2014        934
1 Offline         L 2/22/2015  874708545  2/27/2015       4551
2 Offline         M 12/9/2015  854349935  1/18/2016       9986
3 Online          M 9/17/2014  892836844 10/12/2014       9118
4 Offline         H 2/4/2010   129280602  3/5/2010       5858

  Unit Price Unit Cost Total Revenue Total Cost Total Profit
0 152.58    97.44   142509.72  91008.96  51500.76
1 437.20    263.33  1989697.20 1198414.83 791282.37
2 9.33      6.92    93169.38   69103.12  24066.26
3 81.73     56.67   745214.14  516717.06 228497.08
4 668.27    502.54  3914725.66 2943879.32 970846.34

```

The code loads the dataset Dataset_Ds.csv using Pandas and displays basic information about it. The df.info() function provides an overview of the dataset, including the number of entries, column names, data types, and memory usage. The df.head() function prints the first few rows to understand the dataset structure, which contains categorical (Region, Country, Item Type, etc.) and numerical (Unit Price, Total Revenue, Total Profit, etc.) features.

2.

```

[2] # Convert Categorical Data to Numeric using Label Encoding

# Selecting categorical columns
categorical_cols = ['Region', 'Country', 'Item Type', 'Sales Channel', 'Order Priority']

# Apply Label Encoding
label_encoders = {}
for col in categorical_cols:
    le = LabelEncoder()
    df[col] = le.fit_transform(df[col])
    label_encoders[col] = le

```

This code converts categorical data into numerical values using Label Encoding. First, it selects categorical columns: 'Region', 'Country', 'Item Type', 'Sales Channel', and 'Order Priority'. Then, it applies LabelEncoder() to each column, transforming categorical values into unique numerical labels. The encoded values allow machine learning models to process the data efficiently.

3.

```
✓ [3] # Convert Dates to Numerical Format
Ds
# Convert date columns to datetime format
df['Order Date'] = pd.to_datetime(df['Order Date'], format='%m/%d/%Y')
df['Ship Date'] = pd.to_datetime(df['Ship Date'], format='%m/%d/%Y')

# Feature: Days taken to ship the order
df['Shipping Days'] = (df['Ship Date'] - df['Order Date']).dt.days

# Drop unnecessary columns
df.drop(['Order Date', 'Ship Date', 'Order ID'], axis=1, inplace=True)

# Display processed data
print(df.head())
```

	Region	Country	Item Type	Sales Channel	Order Priority	Units Sold	\
0	4	9	10	1	0	934	
1	2	124	4	0	2	4551	
2	6	139	5	0	3	9986	
3	6	139	9	1	3	9118	
4	2	15	6	0	1	5858	

	Unit Price	Unit Cost	Total Revenue	Total Cost	Total Profit	\
0	152.58	97.44	142509.72	91088.96	51500.76	
1	437.20	263.33	1989697.20	1198414.83	791282.37	
2	9.33	6.92	93169.38	69103.12	24066.26	
3	81.73	56.67	745214.14	516717.06	228497.08	
4	668.27	502.54	3914725.66	2943879.32	970846.34	

	Shipping Days
0	15
1	5
2	40
3	25
4	29

This code converts date columns into a numerical format for better processing. The Order Date and Ship Date columns are first converted into datetime format. Then, a new feature Shipping Days is created by calculating the difference between Ship Date and Order Date. Unnecessary columns (Order Date, Ship Date, Order ID) are dropped to reduce redundancy. Finally, the processed dataset is displayed, showing encoded categorical values and numerical data ready for machine learning.

4.

```
[4] # Split Data into Training & Testing Sets
#
# Define features (X) and target variable (y)
X = df.drop(columns=['Item Type']) # All columns except 'Item Type'
y = df['Item Type'] # Target variable

# Split data into 80% training and 20% testing
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

This code splits a dataset into training and testing sets. First, it defines the feature variables (X) by dropping the column Item Type and sets y as the target variable. Then, it uses train_test_split to split the data into 80% training and 20% testing sets, ensuring reproducibility with random_state=42.

5.

```
✓ [5] # Train the Decision Tree Model  
0s  
# Initialize Decision Tree model  
dt_model = DecisionTreeClassifier(max_depth=5, random_state=42)  
  
# Train the model  
dt_model.fit(X_train, y_train)  
  
print("Model training complete!")
```

→ Model training complete!

This code trains a Decision Tree model. It initializes a DecisionTreeClassifier with a maximum depth of 5 and random_state=42 for reproducibility. The model is then trained using the fit method on the training data (X_train, y_train). A message confirms successful training.

6.

```
✓ [6] # Make Predictions & Evaluate Accuracy  
0s  
# Make predictions  
y_pred = dt_model.predict(X_test)  
  
# Calculate accuracy  
accuracy = accuracy_score(y_test, y_pred)  
print(f"Model Accuracy: {accuracy:.2f}")
```

→ Model Accuracy: 0.83

This code makes predictions and evaluates model accuracy. It uses the trained DecisionTreeClassifier to predict labels for X_test. The accuracy is then calculated

using accuracy_score(y_test, y_pred). The result obtained is 0.83 which means the model has an accuracy of 83%.

7.

```
✓ [7] # Extracting a real row from dataset (row=3) and checking it to predict item type
0s
new_order = df.iloc[2, :-1].values.reshape(1, -1)
recommended_item = dt_model.predict(new_order)
print(f'Recommended Item: {label_encoders['Item Type'].inverse_transform(recommended_item)[0]}')

➡ Recommended Item: Fruits
/usr/local/lib/python3.11/dist-packages/sklearn/utils/validation.py:2739: UserWarning: X does not have valid feature names, but DecisionTreeClassifier expects them
warnings.warn('X does not have valid feature names, but DecisionTreeClassifier expects them')
```

A row from the dataset (excluding the "Item Type") is fed into the trained Decision Tree model. The model predicts the "Item Type", which is then decoded from its numerical representation. The output confirms the model accurately predicts "Fruits" for the given row.

Conclusion

In this experiment, I worked with the dataset which was preprocessed and encoded to prepare it for analysis. After splitting the data into features (X) and labels (Y), I applied a Decision Tree model to predict product categories. Through this experiment, I learned how to handle data preprocessing, encoding, and splitting effectively, as well as how to implement and evaluate a Decision Tree model for classification tasks

Experiment 9

Aim: To perform exploratory data analysis using Apache Spark and Pandas.

Theory:

1. What is Apache Spark and how does it work?

Ans) Apache Spark is an open-source data processing engine designed for speed and scalability. It supports batch processing, real-time streaming, machine learning, and graph analytics in one unified platform.

Unlike Hadoop MapReduce, Spark processes data in memory, making it much faster for iterative and complex tasks. It supports languages like Python, Scala, Java, and SQL, and can run on various cluster managers like YARN, Kubernetes, or standalone. Spark is widely used in big data environments for its performance, ease of use, and flexibility.

Core Components of Apache Spark are:

- **Spark Core:** Core engine (scheduling, memory, etc.)
- **Spark SQL:** For SQL queries and DataFrames
- **Structured Streaming:** Real-time processing
- **Mlib:** Machine learning
- **GraphX:** Graph processing

It works in the following ways:

1. **Driver Program:** Starts the app, defines the workflow (DAG), and manages coordination.
2. **Cluster Manager:** Allocates resources (Standalone, YARN, Mesos, Kubernetes).
3. **Executors:** Run tasks on worker nodes and store data.
4. **In-Memory Computing:** Keeps data in memory for fast processing.
5. **APIs:** Uses RDDs for low-level tasks, DataFrames/Datasets for optimized, SQL-like operations.

2. How is data exploration done in Apache Spark? Explain with steps.

Ans) Data exploration in Apache Spark is done as follows:

1. **Initialize Spark Session:** Start a Spark session to enable interaction with Spark's features and APIs.
2. **Load Data:** Import your dataset into Spark from sources like CSV, JSON, databases, or Parquet files.
3. **View Data Sample:** Look at the first few rows to get a quick sense of the data content and format.
4. **Check Schema:** Examine the structure of the dataset, including column names and data types.
5. **Summary Statistics:** Generate basic statistics like mean, count, min, and max for numerical columns to understand distributions.
6. **Check for Missing Values:** Identify columns that have null or missing values to assess data quality.
7. **Value Counts / Grouping:** Group data by specific columns to analyze category frequencies or distributions.
8. **Filter and Query Data:** Apply filters or queries to focus on specific subsets of the data for deeper analysis.

Conclusion: In this experiment, we learned how to perform exploratory data analysis using Apache Spark and Pandas. By examining the data's structure, statistics, missing values, and groupings, we gain key insights to guide cleaning, feature selection, and modeling. Spark's speed and scalability make it ideal for exploring large datasets efficiently.

Experiment 10

Aim: To perform Batch and Streamed Data Analysis using Apache Spark.

Theory:

1. What is streaming? Explain batch and stream data.

Ans) Streaming is the process of transmitting and processing data continuously and in real-time. Instead of waiting for all the data to be collected (like in batch processing), streaming handles data as soon as it's generated. This allows for immediate analysis and actions — useful in scenarios like fraud detection, live video, and real-time analytics.

A. Batch Data Analysis:

- Processed in groups after collecting over time.
- Not real-time, higher latency.
- Examples: Monthly billing, daily sales reports.
- Tools: Hadoop, Hive.

B. Streamed Data Analysis:

- Processed in real-time as it's generated.
- Low latency, continuous flow.
- Examples: Live chat, stock market feeds.
- Tools: Kafka, Spark Streaming.

The main difference is that in batch processing, data is processed later in chunks, while in stream processing, data is processed continuously in real-time as it arrives.

2. How does data streaming take place in Apache Spark?

Ans) In Apache Spark, data streaming is handled by Spark Streaming, which allows for the processing of real-time data in a distributed and fault-tolerant manner. Here's how the process works:

1. **Data Ingestion:** Spark Streaming can ingest real-time data from various sources like Kafka, Flume, Sockets, or file systems like HDFS or Amazon S3. The data is continuously received in small chunks.
2. **Micro-batches:** Spark Streaming divides the incoming data stream into micro-batches, which are small, manageable batches that are processed in intervals (typically in milliseconds or seconds). This approach allows Spark to handle streaming data while leveraging its batch processing engine.
3. **Processing:** Each micro-batch is processed using Spark's standard transformations (like map, filter, etc.) and actions (like reduce, count, etc.). These operations are applied on the data in the batch, similar to how traditional Spark processes batch data.
4. **Output:** After processing, the results of the micro-batches are typically written to external storage systems (such as HDFS, databases, or dashboards) or used to trigger further actions.
5. **Fault Tolerance:** Spark Streaming ensures fault tolerance by replicating the data across multiple nodes and periodically checkpointing the data, so that the system can recover from failures without losing progress.

Conclusion: In this experiment, we learned how Apache Spark Streaming processes real-time data using micro-batches. It enables efficient, scalable, and fault-tolerant real-time analytics, making it ideal for applications like monitoring and fraud detection. Spark Streaming provides a powerful solution for large-scale data processing.

Chapter 1: Introduction

1.1. Introduction

Lending for real estate, consumer, mortgage, and business loans is central to most banks' business models, with loan profits forming a major part of their assets. However, lending to unfit borrowers remains a key credit risk. Banks also face challenges in identifying intentional defaulters and biased internal practices influenced by defaulting companies. While modern verification processes help, selecting truly creditworthy applicants remains uncertain. Housing finance companies can use machine learning to identify ideal customers and automate loan approval by predicting applicant reliability.

1.2. Objectives

- Lending is key to bank's profits, but credit risk from unfit borrowers and internal bias remains a challenge. Despite better verification, identifying reliable applicants is uncertain.
- Machine learning helps housing finance firms target and approve trustworthy borrowers.

1.3. Motivation

Rising loan volumes require faster, more accurate evaluation. Manual reviews are slow and prone to bias, prompting banks to use machine learning for automated, consistent loan eligibility predictions and quicker customer feedback.

1.4. Scope of the Work

- Uses supervised learning for binary classification of loan approvals.
- Applies ML models to a public dataset with applicant and loan details.
- Tests algorithms like Logistic Regression, Random Forest, KNN, Naive Bayes, SVM, and Gradient Boosting.
- Visualizes insights with graphs, confusion matrices, and performance metrics.

1.5. Feasibility Study

Technical Feasibility: Utilizes Python along with libraries like Scikit-learn, Pandas, and Matplotlib, which are widely supported and suitable for implementing machine learning models.

Operational Feasibility: The prediction system can be integrated into bank software or customer-facing portals to provide real-time loan eligibility checks.

Economic Feasibility: Cost-effective due to the use of open-source tools and publicly available datasets, requiring no additional financial investment

Chapter 2: Literature Survey

2.1. Introduction

Loan approval systems have gained significant importance in the banking and financial services industry due to increasing demand for fast, reliable, and risk-aware credit disbursement. Traditional loan processing methods are often manual, time-consuming, and prone to bias or error. With the evolution of machine learning technologies, predictive models are now being used to automate and enhance the decision-making process. This literature survey presents a comparative study of two significant research papers in the domain of loan eligibility prediction using machine learning.

2.2. Problem Definition

The goal of this literature review is to explore how machine learning techniques can be applied to predict loan approval status based on applicant information. It also aims to examine the challenges associated with data quality, model selection, and evaluation criteria when applying ML to real-world financial scenarios.

2.3. Review of Literature Survey

In the research paper **“Loan Approval Prediction using Machine Learning Algorithms” by Aditi Sharma et al.**, the authors outline a structured pipeline for predicting whether a loan should be approved or not. The study is divided into four phases: data collection, model comparison, training, and testing. Various machine learning models including Logistic Regression, Decision Tree, and Support Vector Machine were evaluated based on their accuracy and precision. Among these, Logistic Regression showed promising results due to its simplicity and effectiveness in binary classification tasks. The study emphasized the importance of data preprocessing steps such as handling missing values, encoding categorical data, and normalizing numeric features to improve model performance.

Another notable contribution is the paper **“Risk Reduction in Loan Lending using Machine Learning” by Priya R. et al.**, which focuses on reducing non-performing assets (NPAs) in banks and NBFCs by improving the loan applicant screening process. The paper highlights how feeding historical loan data into trained machine learning models can aid in identifying trustworthy borrowers. The authors experimented with models such as Random Forest and Gradient Boosting, which delivered better accuracy and robustness compared to simpler models. The study also sheds light on real-world challenges such as imbalanced datasets and the need for interpretability in financial decision-making. It concludes that incorporating machine learning into the lending process can significantly improve risk management and operational efficiency

Chapter 3: Design and Implementation

3.1. Introduction

This chapter outlines the design and implementation process for the Loan Eligibility Prediction System. The project adopts a structured machine learning pipeline to ensure accurate prediction of whether a customer is eligible for a loan based on various personal and financial attributes. The workflow consists of data preprocessing, feature selection, model training, validation, and performance evaluation. Several supervised machine learning algorithms are employed, including Logistic Regression, Random Forest, Support Vector Machines (SVM), and Gradient Boosting, with the goal of identifying the most effective model for real-time deployment.

3.2. Requirement Gathering Hardware

Requirements:

- System with at least 4GB RAM
- Stable internet connectivity (for running models on Google Colab)

Software Requirements:

- Google Colab or Jupyter Notebook (for implementation and testing)
- Python 3.x (programming environment)

3.3. Proposed Design

The system design follows the standard CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology, ensuring a systematic approach to problem-solving and model development:

1. Data Collection: The dataset includes applicant details like gender, income, credit history, and property area. The target variable is `Loan_Status`, showing loan approval.
2. Data preprocessing included imputing missing values, encoding categorical features, scaling with `StandardScaler`, and handling outliers through visual and statistical methods..
3. New features like `Total_Income` and `Loan_Amount_Term_Years` were added, while redundant or correlated features were removed to boost accuracy.
4. Model Building:
 - Logistic regression
 - Random forest classifier
 - K-Nearest Neighbors (KNN)
 - Gaussian Naive Bayes
 - Support Vector Machine (SVM)
 - Gradient Boosting Classifier
5. Model Evaluation: Models were assessed using accuracy, precision, recall, F1-score, confusion matrix, and cross-validation, with the best-balanced model chosen for deployment.
6. Visualization: Data and model insights were visualized using count plots, bar charts, heatmaps, confusion matrices, and ROC curves.

3.4. Proposed Algorithm

Logistic Regression

Logistic Regression is a statistical model used for binary classification. It estimates the probability that a given input point belongs to a particular category. Due to its simplicity and interpretability, it was used as a baseline model. It performs well when the features have a linear relationship with the target variable.

Random Forest Classifier

Random Forest is an ensemble learning method that builds multiple decision trees and merges their predictions to improve accuracy and reduce overfitting. Each tree is trained on a random subset of the data and features. This model handles both categorical and numerical data efficiently and provides feature importance metrics for analysis.

K-Nearest Neighbors (KNN)

KNN is a non-parametric algorithm that classifies new data points based on the majority label of their 'k' nearest neighbors in the feature space. It is simple to implement but sensitive to the scale of the data and the choice of 'k'. Therefore, preprocessing steps like normalization were applied before training.

Gaussian Naive Bayes

Naive Bayes is a probabilistic classifier based on Bayes' theorem with the assumption of independence between predictors. The Gaussian variant assumes that the continuous features follow a normal distribution. It is efficient and performs well with high-dimensional data, especially when the assumption of independence roughly holds.

Support Vector Machine (SVM)

SVM constructs a hyperplane in a high-dimensional space to separate classes with maximum margin. It is particularly useful for datasets with clear class boundaries. The RBF (Radial Basis Function) kernel was used to capture non-linear patterns in the loan data. Hyperparameter tuning was performed on C and gamma to balance bias-variance trade-off.

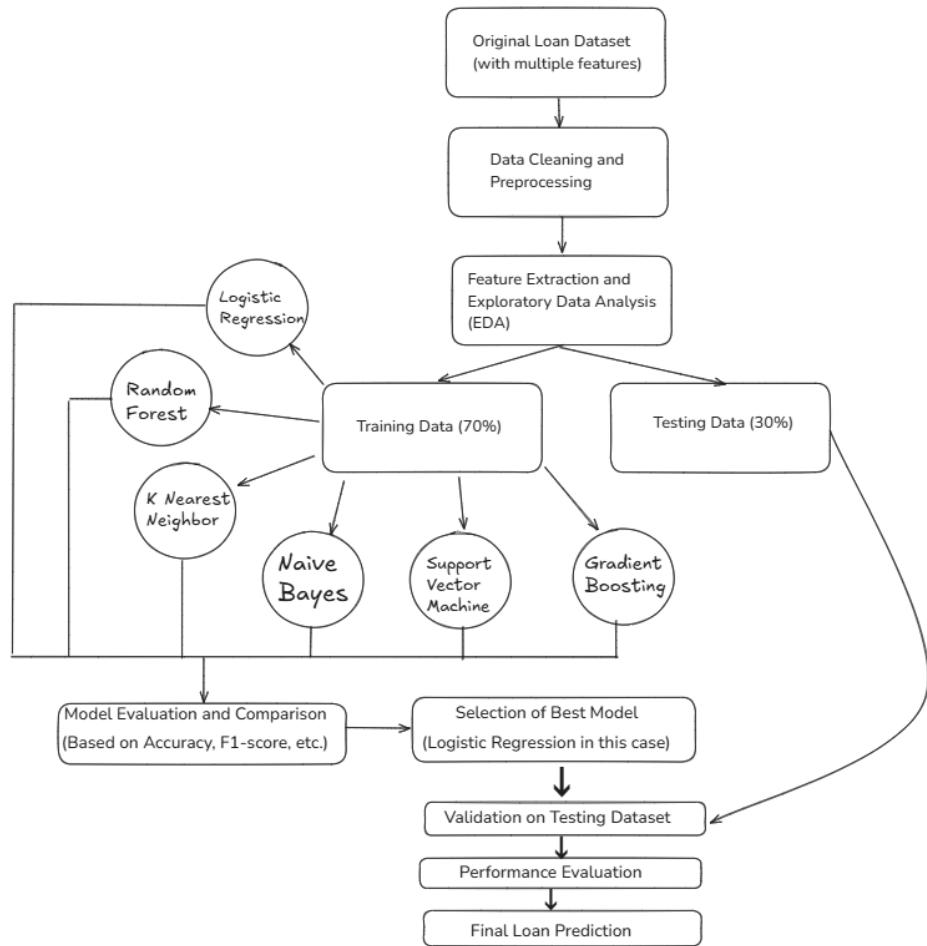
Gradient Boosting Classifier

Gradient Boosting is an advanced ensemble technique that builds models sequentially, each one trying to correct the errors made by the previous. It is highly effective for structured data and often provides better performance than other models, although it may require more training time. Learning rate, number of estimators, and tree depth were tuned to optimize its performance.

Model Selection and Comparison

Each of the above algorithms was trained and evaluated on the same preprocessed dataset. Accuracy, precision, recall, F1-score, and ROC-AUC were used as evaluation metrics. The model that achieved the best balance across all metrics was selected for deployment. In this project, the Gradient Boosting Classifier and Random Forest showed superior performance, making them strong candidates for the final system.

3.5. Architectural Diagrams



4: Results and Discussion

4.1. Introduction

To assess the performance of the developed loan eligibility prediction system, multiple machine learning models were trained and evaluated using a common dataset. Each model was tested using standard performance metrics such as **Accuracy**, which reflects the percentage of correctly predicted outcomes. These metrics provided insight into the strengths and limitations of each model in predicting whether a loan should be granted or not.

4.2 Cost Estimation

The system was developed using open-source Python libraries such as Scikit-learn, Pandas, NumPy, and Matplotlib, and executed in a Google Colab environment. This eliminated the need for expensive local infrastructure, making the solution highly cost-effective and suitable for academic research or small-scale enterprise deployment.

4.3 Feasibility Study

Given its low hardware footprint and ability to run effectively in cloud environments, the system demonstrates strong feasibility for integration into online loan application portals. The real-time prediction capability and scalability of the chosen models make this system well-suited for financial institutions looking to automate and streamline their loan approval workflow.

4.4 Results of Implementation

Model	Accuracy
Logistic Regression	97.00
Naive Bayes	95.92
Gradient Boosting Classifier	93.84
Random Forest	89.76
Support Vector Machine (SVM)	71.89
K Nearest Neighbor	61.08

5: Conclusion

5.1. Conclusion

In this project titled "*Loan Eligibility Prediction using Machine Learning*", we built a robust and automated system to predict whether a customer is eligible for a home loan based on various demographic and financial factors. We explored multiple machine learning algorithms, including Logistic Regression, Random Forest, K-Nearest Neighbors, Naive Bayes, SVM, and Gradient Boosting, and compared their performance to identify the most accurate model.

We began by collecting and preprocessing the data, which involved handling missing values, encoding categorical features, and normalizing inputs where necessary. We then split the data into training and testing sets to evaluate model performance fairly. We trained each classification model and evaluated them using accuracy as the primary metric.

We found that Logistic Regression achieved the highest accuracy (~97%), making it the most effective model for this classification problem. Other models such as Naive Bayes and Gradient Boosting also performed reasonably well, while SVM and KNN showed comparatively lower accuracy.

We concluded that logistic regression not only offered good predictive power but also maintained interpretability, making it a suitable choice for real-time loan eligibility checks. We implemented the solution using Python and open-source libraries on Google Colab, which kept the project cost-effective and accessible.

Overall, we demonstrated that machine learning can effectively support financial institutions in automating the loan eligibility process and targeting the right customer segments, thereby improving decision-making and operational efficiency.

5.2. Future Scope

- While the current system successfully predicts loan eligibility using multiple machine learning models, there are several directions in which this work can be extended to enhance its applicability and performance:
- **Real-time Deployment:** Integrating the trained models into production using RESTful APIs can enable real-time eligibility checks during the loan application process, improving user experience and operational efficiency.
- **Feature Expansion:** Including more detailed financial and behavioral attributes such as previous loan history, spending habits, credit card usage, and repayment patterns could provide deeper insights and improve prediction accuracy.
- **Ensemble and Stacking Models:** Future versions could explore stacked ensembles combining the strengths of top-performing models to boost overall accuracy and generalization.
- **Continuous Model Training:** Incorporating mechanisms for continuous learning and retraining with new customer data would help the system remain current with changing customer behaviors and financial trends.
- **Mobile and Web Integration:** Embedding the predictive system into mobile or web loan application portals can make the process seamless for users, helping financial institutions provide instant decisions and reduce dropout rates.

5.3. Societal Impact

The development and deployment of loan eligibility prediction systems have significant implications for society, particularly in the context of financial inclusion and responsible lending practices.

- **Financial Inclusion:** Automated loan eligibility prediction helps provide equal access to credit by reducing bias in decision-making. By using data-driven models, individuals from diverse socioeconomic backgrounds can be evaluated fairly, improving access to financial services for underserved communities.
- **Responsible Lending:** Predictive models assist financial institutions in identifying customers who are more likely to repay their loans, reducing the risks of over-lending and defaults. This promotes responsible lending practices and contributes to the stability of financial markets.

- **Empowering Individuals:** By automating loan eligibility assessments, individuals can receive faster decisions, reducing the waiting time and uncertainty associated with loan applications. This improves customer experience and empowers individuals to make informed financial decisions.
- **Data Privacy and Protection:** Predicting loan eligibility involves handling sensitive personal and financial data. Ensuring robust data privacy measures and compliance with regulations like GDPR is critical to safeguarding user information and building trust in automated financial systems.
- **Educational and Research Value:** The techniques used in this project provide valuable learning opportunities for students and professionals in the fields of machine learning and data science. It promotes a deeper understanding of financial analytics and predictive modeling, contributing to the advancement of knowledge in financial technologies.

AIDS Assignment

Ques

- Q) What is AI? Considering the COVID-19 pandemic situation, how AI helped to survive & renovated our way of life with different applications?

→

AI is a branch of computer science that focuses on creating intelligent systems capable of performing tasks that typically requires human intelligence. AI systems uses techniques such as machine learning, rule, etc.

Q)

AI can be categorised into:-

1. Narrow AI (weak AI):-
2. General AI
3. Super AI.

- Role in surviving & renovating life during covid..
- i) Healthcare & medical diagnosis.
 - CT Scan analysis & X-ray diagnosis helped in rapid detection of covid, cases.
 - ii) Virus spread control.
 - Social distance monitoring tools helped to monitor & enforce laws in public places
 - iii) remote work & education.
 - block from home optimization & AI in online education help people connect with their work.
2. What are AI agents terminology & explain with example
- - Environment:-
 - Everything which surrounds the agent & influences its action. It can be complete or partially observable.

Applying the laws of motion to
current

- This is how does magnetism work
- Electromagnet

Another source of
magnetism

- Electromagnet

- Take out the permanent magnet and
keep it one side. It has two poles
both a north pole and south pole alternating
- South

- The final stage is to connect the bulb
to your DC source and turn on the current

- Now if we rotate the wire used to connect the bulb
continuously

- The 3rd stage situation is a state where each
rotation of the wire for every rotation of the
magnet will generate some current
- This is how we can get a small
generator and this is how it works

- This is how it is connected to the bulb -

- This is how the circuit will look like

- This - Induction of current in the coil
- This - Induced current in the coil

1) Induced current

- G.P. - Generator Principle

- Now when there is no current in
the circuit then there is no current

- If current in the circuit then there is current

but it is not

Initial State :- 1 2 3
5 6 0
4 7 8

Goal State :- 1 2 3
4 5 6
7 8 0

- compute heuristic of each possible move.
- expand the state with the lowest f(n) & repeat

4. What is PEAS description? give PEAS description for following?

→ Performance Measure :- How success of agent is evaluated

Environment :- Surrounding in which agent operates.

Actuators :- Component that allows agent to take actions

Sensors :- Component that allows agent to perceive the environment

- i) Taxi driver agent:-

Performance measure	Environment	Actuators	Sensors
-safe driving	-Traffic Signal	-Steering wheel	-camera,
-travel time.	-roads.		-GPS
-traffic rules	-Weather	-accelerator	-fuel gauge,
		-brakes	

- ii) Medical Diagnostic Agents:-

perf. measure	Environment	Actuator	Sensors
-health of patient	-patient data	display screen	heart rate
-accuracy of diagnosis	-symptoms	alarm system	monitor
-recommended treatment	-test reports	robotic arms	lab results

- iii) Music component Agent:-

perf. measure	Environment	Actuator	Sensors
-originality	-music db	-Speaker	-microphones,
-listener engagement	-user preferences	-digital music	-user feedback
-quality		FOR EDUCATIONAL USE Intention	-recognition level

iv) Aircraft auto lander:

pref. measure	Environment	Actuator	Sensors
- smooth landing	- runway	- landing gear	- altimeter
- accuracy in touchdown	- wind conditions	- flap	- GPS
	- air traffic	- air brakes	- camera

v) Essay evaluator:

pref. measure	Environment	Actuator	Sensors
- grading	- plagiarism database	- display	- optical
- grammar	- rubric criteria	- screen	- character
- paradigm check.		- text to recognition	(OCR), Speech system.

vi) Robotic sentry gun:

pref. measure	Environment	Actuators	Sensors
- neutralize threats	- lab area	- gun mech	- camera
- target tracking	- potential	- anism.	- thermal
- false alarms	- intruders	- alarm siren	Sensors
		- tracking system	

5. Categorize a shopping bot for a shopping bot for an offline bookstore according to the following system.

- - observability: Partially observable - Relies on limited sensor input.
- Deterministic or stochastic: Stochastic: customer pref. is unpredictable
- Episodic v/s sequential: Sequential decisions affect future actions.
- static v/s dynamic: Dynamic, customer behaviour is always evolving
- Discrete v/s continuous: Finite no. of choices. Such as books

6. Differentiate between model based & utility based agent.

→ Model Based Agent Utility Based Agent
 - Agent that maintains the internal model of the env. based on the utility function to understand its current aiming to maximize long term state to predict future states satisfaction or benefit

- Model updates its information measures how desirable about the environment different states are.

- less complex • more complex.

- Doesn't concern long term re- • focuses on long term words reward.

- e.g:- self driving cars. e.g:- shopping recommendation system

7.

Explain the architecture of knowledge based agent & learning agent.

→ Knowledge based agent:- stores knowledge & reason based on enginee inference
 - knowledge base:- stores facts, rules & heuristic function about the environment

- inference engine:- uses logical reasoning techniques like forward & backward chaining.
 - perception:- gathers data from the environment
 - actions:- executes actions based on inferred knowledge.

- knowledge update mechanism:- updates itself as new facts are learned.

- Learning Based agents:- agent that improves performance over time by learning from experience, also provides feedback.

- Learning element:- responsible for improving agents perf by analyzing past experiences using ml techniques.
- critic :- provides feedback on agents actions by evaluating success or failure
- problem generation:- supports new experiences for learning & explanation.

→ v. convert the following to predicates.

Anita travels by car if available otherwise travels by bus.
travels(x, y) → person x travels by y .
Available(x) → x vehicle is available
(availability, x) → vehicle goes via x .
purch(x, y) → x purchases y .

available(x) → Travels(x , bus).

b. Bus goes via andheri & goregaon.
(Goregaon, bus, andheri)
Goes via C bus, goregaon.

c. car has a puncture, so it's not available.
punctured(x)
~ Available(x).

1Q What do you mean by depth limit search?

Explain iterative deepening search with example.



→ Depth Limit Search(DLS):-

DLS is a DESS variant with a fixed depth L preventing infinite loop & saving memory.

Advantages:-

- Avoids infinite recursion.
- Memory efficient.

Disadvantages:-

- May miss less deeper solutions
- Need for good L choice (limit).

Iterative Deepening Search (IDS):-
IDS runs DLS repeatedly, increasing L until the goal is found.

Advantages:-

- Completes the whole search.
- Gives most optimized answer.
- Memory efficient

Disadvantages:-

- Redundant computation.
- Higher time cost.

example:- searches level by level until the goal appears

(b)

Explain hill climbing & its drawback in detail with example. Also state limitations of steep hill climbing.

It is an optimization problem that moves toward higher values (between solutions), until a peak (local optimum) is reached.

Algorithm:-

1. Start with an initial state.
2. Move to the best neighbouring state.
3. Repeat until no strictly better neighbour exists.

Example:- 8-Queens problem:-

- Adjust queen's position to minimize conflicts.
- Stop when no improvements are possible

Drawbacks:-

- Local maxima:- stuck at suboptimal peaks.
- Plateau:- no directions for improvement.
- Ridges:- needs special move to progress.
- Vanishing solution
- Steepest-ascent:- evaluates all neighbors but still get stuck.

(ii)

Explain simulated annealing & write its example. It improves hill climbing by allowing occasionally bad moves to escape local maxima, inspired by metal annealing.

Algorithm:-

1. Start with an initial solution & temperature T.

2. pick a random neighbour & compute ΔE .

3. Accept S is better; otherwise with probability

$$P = e^{-C - \Delta E / T}$$

4. Reduce T until stopping condition

Advantages:-

- 1) Escapes local minima.
- 2) Handles large problems
- 3) Near-optimal solutions

Disadvantages:-

- 1) Tricky cooling schedule.
- 2) No guarantee of best solution.

Q2) Explain A* algorithm with an example.

~~Algorithm~~

A* is a best-first search algorithm for pathfinding containing:

- 1) Uniform cost search (cheapest path).
- 2) Greedy best-first search (heuristic-based speed)

key-formula:-

$$f(c_n) = g(c_n) + h(c_n)$$

• $f(c_n)$: cost from start to n .

• $g(c_n)$: cost from n to goal. (estimated).

Steps:-

1. Start with the initial node, compute $f(n)$.
2. expand the node with lowest $f(n)$.
3. if goal is reached, return the path; else, update & continue.

Advantages:-

- 1) Optimal paths.
- 2) Efficient in AI applications.

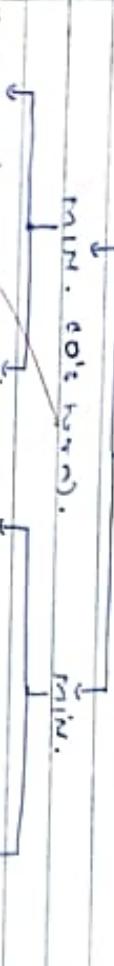
Disadvantages:-

- 1) High memory usage.
- 2) Depends on Heuristic.

- Q) Explain minimax algorithm & draw game tree for tic-tac-toe game.
 → minimax is a game strategy for 2 player games like tic-tac-toe.

How - it works?

- maximize(x) aims to increase the score (to win).
- minimizer(o) aims for lowest score (to lose).
- explore all possible moves, assign score & picking the best one.
- Game tree:- max. (x 's turn), min. (o 's turn).



Advantage:-

- i) Always finds the best moves
- ii) Disadvantages:-

- i) Slow for deep tree (Alpha-beta pruning helps).

Q) Explain alpha beta pruning for adversarial search with example.

→ Alpha Beta pruning optimizes the minimax algorithm by skipping unnecessary branches, making it faster without affecting the result.

- Explanation
- AlphaBeta: Best max value found so far.

2) BestFirst:- Best min value found so far

3) If a move worse than α or β is found then the further exploration is stopped. (pruned).

Example (simplified game tree):-

max.

\nearrow
min min

3. 5. 2. Pruned.

Here, if min finds move worse than 5, it stops exploring that branch.

Advantages:-

1) Speeds up minimax by ignoring bad choices.

2) Same result as minimax but faster

Q) Explain Wumpus world environment giving its PEAS description. Explain how percept sequence is generated?

→ Wumpus world is a grid based game environment where an agent navigates a cave to find gold while avoiding pits & the wumpus controller.

PEAS description:-

P (performance measure):- +1000(Gold), -1000(Wumpus), -100(Pit), -1(Move).

E(Environment):- Grid world with Wumpus, Gold, Pits & Agents.

A(Actuators):- Move, Grab(Gold), Shoot(Arrow), Climb.
S(Sensors):- a) Breeze(near pit) b) Stench(near wumpus) c) Gittern(near Gold).

Except sequence generation:-

The agent receives sensory input at each step based on its current location.

2. Example:- If agent move next to pit, it perceive Breeze.

3. Using percept history it finds safe path & avoid danger.

Q) Solve the following crypto arithmetic problem

1. ~~SEND + MORE = MONEY~~

\rightarrow each letter represents a unique digit (0-9).

$$\text{Step 1:- } \text{Evaluation setup} \\ (1000S + 100O + 10N + O) + (1000M + 100O + 10R + E) \\ = (10000M + 10000 + 1000N + 100E + Y).$$

Step 2:- Constraints.

- $M = 1$. C since MONEY is a 5 digit number.
- $S \neq 0$ (it's the first digit in SEND).

~~All letters have unique values.~~

Step 3:- Assigning Digits.

Letters	Digits
S	9
E	5
N	6
D	7
M	1
O	0
R	8
Y	2

19)

Consider the following axioms.
 All people who are graduating are happy
 all happy people are smiling.
 someone is graduating.

Explain the following:-

- 1) Represent these axioms in first order predicate
- 2) Convert each formula to clause form
- 3) Prove that "is someone smiling?" with resolution technique. Draw the resolution tree

\rightarrow 1) FOL:

$L(n) \rightarrow m : n \text{ is graduating}, H(m) \rightarrow m : n \text{ is happy}, S(m) \rightarrow m : n \text{ is smiling}.$

Axioms :- 1. $\forall n (L(n) \rightarrow H(n))$

2. $\forall n (H(n) \rightarrow S(n))$.

3. convert to clause form

1. $\neg L(n) \vee H(n)$.

2. $\neg H(n) \vee S(n)$.

3. prove "is someone smiling?"

a. $L(n)$ given.

b. $\neg L(n) \vee H(n) \rightarrow$ axiom 1

c. $\neg H(n) \vee S(n) \rightarrow$ axiom 2

Since we derived $S(n)$, the proof confirms that someone is smiling.

Resolution tree

$L(n) \quad \neg L(n) \vee H(n)$

$H(n)$

$\neg S(n)$.

20) Explain modus ponens with suitable example.

→ Modus ponens is a fundamental rule of inference in logic. It states:

If $P \rightarrow Q$ (if P, then Q) is true.
 $\& P$ is true, then Q must be true.

Symbolically:-

$$P \rightarrow Q, P \vdash Q$$

example:-

- 1) If it rains, ground will be wet.
- 2) It is raining. (P)
- 3) Therefore, the ground is wet. (Q).

This rule is widely used in mathematical proofs & AI reasoning systems.

21.0) Explain Forward & Backward chaining with an example.

→ There are inference techniques used in rule-based systems & AI - reasoning.

1. Forward chaining (Data-Driven):-

Starts with known facts & applies rules to infer new facts until the goal is reached.
 Works from cause to effect (bottom-up approach).

Ex:-

1. If it is raining then the ground is wet ($R \rightarrow Q$)
2. If ground is wet then traffic is slow ($W \rightarrow T$)

1. Starts with the goal & goes backwards to find

supporting facts.

• works from effect to cause. (top-down approach).

Ex:- $w \rightarrow r : \text{Is the ground wet}$

$r \rightarrow w : \text{Is it raining}$

AIDS-I Assignment No: 2

Q.1: Use the following data set for question 1

82, 66, 70, 59, 90, 78, 76, 95, 99, 84, 88, 76, 82, 81, 91, 64, 79, 76, 85, 90

1. Find the Mean (10 pts)

To determine the mean (average), sum all the data points and divide by the number of values present in the dataset.

- Total sum of all values:
 $82 + 66 + 70 + 59 + 90 + 78 + 76 + 95 + 99 + 84 + 88 + 76 + 82 + 81 + 91 + 64 + 79 + 76 + 85 + 90 = \mathbf{1691}$
- Number of values = **20**
- Mean = Total sum \div Number of values = $1691 \div 20 = \mathbf{84.55}$

The mean of the dataset is 84.55.

2. Find the Median (10 pts)

The median represents the middle value of a sorted dataset. Since the dataset has an even number of values, the median is the average of the two central values.

- Ordered dataset:
59, 64, 66, 70, 76, 76, 76, 78, 79, 81, 82, 82, 84, 85, 88, 90, 90, 91, 95, 99
- The 10th and 11th values are 81 and 82.
- Median = $(81 + 82) \div 2 = \mathbf{81.5}$

The median of the dataset is 81.5.

3. Find the Mode (10 pts)

The mode is the value(s) that appear most frequently in a dataset. We count the frequency of each number to find the mode.

- From the sorted list, we observe:
 - 76 appears 3 times.
 - All other numbers occur less frequently.

The mode of the dataset is 76.

4. Find the Interquartile Range (IQR) (20 pts)

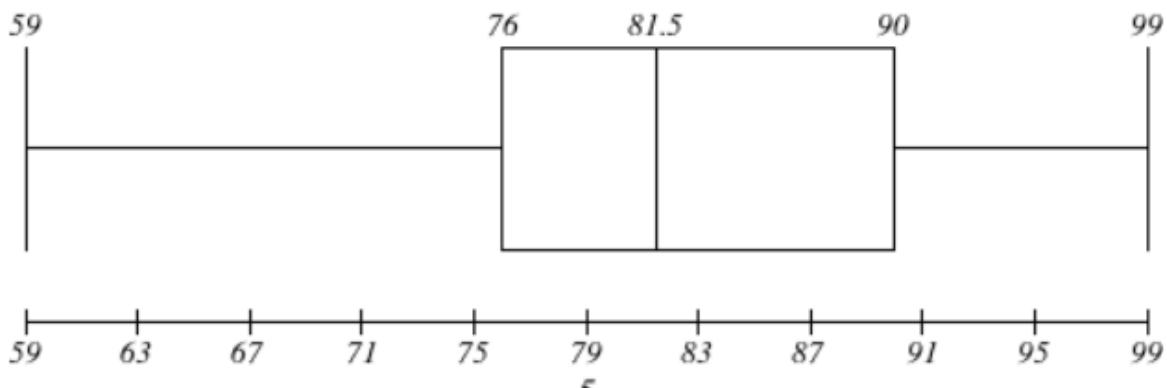
The interquartile range measures the spread of the middle 50% of the data and is calculated as:

$$\text{IQR} = \text{Q3} - \text{Q1}$$

- First, divide the dataset into two halves:
 - Lower half (first 10 values): 59, 64, 66, 70, 76, 76, 76, 78, 79, 81
 - Upper half (last 10 values): 82, 82, 84, 85, 88, 90, 90, 91, 95, 99
- Q1 is the median of the lower half $= (5\text{th} + 6\text{th}) \div 2 = (76 + 76) \div 2 = \mathbf{76}$
- Q3 is the median of the upper half $= (15\text{th} + 16\text{th}) \div 2 = (90 + 90) \div 2 = \mathbf{90}$
- $\text{IQR} = \text{Q3} - \text{Q1} = \mathbf{90} - \mathbf{76} = \mathbf{14}$

The interquartile range is 14.

Box Plot



Q.2 1) Machine Learning for Kids 2) Teachable Machine

1. For each tool listed above:

- Identify the target audience
- Discuss the use of this tool by the target audience
- Identify the tool's benefits and drawbacks

2. 1) Machine Learning for Kids

- **Target Audience:** The main audience for Machine Learning for Kids includes young learners and educators aiming to teach the basics of machine learning in a simple and approachable manner. It is tailored for beginners with minimal or no experience in coding or machine learning.
- **Use of the Tool:** This platform enables kids to build machine learning models through easy-to-use, kid-friendly interfaces. They can develop models to identify text, images, audio, or numerical data. Once trained, these models can be incorporated into Scratch or Python projects, allowing children to design interactive games or applications that react to the inputs they've taught the model. This practical, hands-on method helps young learners grasp the fundamentals of machine learning through active participation.

Benefits:

- **Ease of Use:** The tool breaks down complex machine learning ideas into simpler concepts that are easy for children to grasp.

- **Engaging Interface:** Integrating with Scratch and Python adds an element of fun and interactivity to the learning experience.
- **Educational Value:** It provides a practical introduction to AI and machine learning, fostering computational thinking skills.
- **Drawbacks:**
 - **Limited Complexity:** Due to its simplicity, it may not be suitable for advanced machine learning projects.
 - **Dependency on External Platforms:** Relies on Scratch or Python for project integration, which might require some additional setup.

3. 2) Teachable Machine

- **Target Audience:** Teachable Machine is geared towards a broader audience, including students, artists, educators, and makers. It's designed for anyone who wants to quickly and easily create machine learning models without writing code.
- **Use of the Tool:** Teachable Machine simplifies the process of training machine learning models for image, audio, and pose recognition. Users can train models directly in their browser by providing examples through a webcam, microphone, or file uploads. The trained models can then be exported and used in various applications, websites, or projects. This tool is excellent for rapid prototyping and integrating machine learning into creative projects.
- **Benefits:**
 - **No Coding Required:** It's very accessible, as it doesn't require any coding knowledge.
 - **Fast Prototyping:** Allows for quick creation and testing of machine learning models.
 - **Versatile Export Options:** Models can be exported in various formats, making them usable in different environments.
- **Drawbacks:**
 - **Limited Control:** Offers less fine-grained control over the model training process compared to more advanced tools.
 - **Browser-Based Limitations:** Performance and capabilities can be limited by the browser's capabilities.

4. From the two choices listed below, how would you describe each tool listed above? Why did you choose the answer?

- Predictive analytic
- Descriptive analytic

5. **1) Machine Learning for Kids:** Predictive analytic

- **Why?** Machine Learning for Kids enables users to train models that can make predictions or classifications based on the input data. For instance, a child can train a model to predict whether an image is a cat or a dog. This is a clear example of predictive analytics, as the model is used to predict a future outcome or classify an input.

6. **2) Teachable Machine:** Predictive analytic

- **Why?** Teachable Machine is also a predictive analytic tool. It allows users to train models that predict outputs (like image, audio, or pose classifications) from new input data. The core function is to predict or classify, which aligns with the definition of predictive analytics.

7. From the three choices listed below, how would you describe each tool listed above? Why did you choose the answer?

- Supervised learning
- Unsupervised learning
- Reinforcement learning
-

8. **1) Machine Learning for Kids:** Supervised learning

- **Why?** In Machine Learning for Kids, children provide labeled examples to train the models. For example, they show the model pictures labeled "cat" or "dog." The model learns from these labeled examples to make predictions on new, unseen images. This process of learning from labeled data is the fundamental characteristic of supervised learning.

9. **2) Teachable Machine:** Supervised learning

- **Why?** Teachable Machine also operates on the principle of supervised learning. Users provide labeled data (e.g., images labeled with categories) to train the model. The model then learns to map the input data to the provided labels, enabling it to classify new inputs. The need for labeled data to train the model classifies it as supervised learning.

Q.3 Data Visualization: Read the following two short articles:

Read the article Kakande, Arthur. February 12. "What's in a chart? A Step-by-Step guide to Identifying Misinformation in Data Visualization." Medium

Read the short web page Foley, Katherine Ellen. June 25, 2020. "How bad Covid-19 data visualizations mislead the public." Quartz

Research a current event which highlights the results of misinformation based on data visualization. Explain how the data visualization method failed in presenting accurate information. Use newspaper articles, magazines, online news websites or any other legitimate and valid source to cite this example. Cite the news source that you found.

Current Event Highlighting Misinformation Through Data Visualization: Misrepresented NOAA Temperature Graph

Event Overview

A graph sourced from the National Oceanic and Atmospheric Administration (NOAA) was selectively cropped and presented on social media to falsely claim that the Earth has been experiencing a cooling trend, contradicting the established scientific consensus on human-caused global warming.¹ The graph focused solely on temperature data from 2015 to 2022, omitting long-term historical data that clearly shows a warming trend. This misrepresentation was fact-checked by the Associated Press (AP).

How the Data Visualization Method Failed

- **Cherry-Picking Data:** The visualization deliberately showcased only a limited timeframe (2015-2022), which was chosen to exploit natural climate variability (El Niño and La Niña events) and create a false impression of cooling. This selective presentation ignored the broader 140-year temperature record, a clear example of cherry-picking data to support a predetermined narrative.
- **Lack of Context:** The graph failed to provide essential context by omitting the long-term temperature data from NOAA. This omission prevented viewers from understanding the true extent of global warming and the significance of the short-term fluctuations shown in the graph. Legitimate data visualizations should always include sufficient context to ensure accurate interpretation.
- **Misleading Trendline:** A black line was added to the graph to emphasize a minor downward trend within the selected timeframe. This visual element drew attention to a statistically insignificant fluctuation, further reinforcing the false narrative of global cooling and obscuring the overall warming trend.
- **Exploitation of Authority:** The NOAA logo was prominently displayed on the graph, lending a false sense of authority and scientific validity to the misleading data. This tactic exploited the credibility of a reputable scientific institution to promote misinformation.

Impact of the Misinformation

- The misleading visualization contributed to the spread of climate change denial, undermining public understanding of the severity and urgency of global warming.
- It fostered skepticism towards climate science and reputable scientific institutions like NOAA, potentially eroding public trust in evidence-based information.
- The misrepresentation could diminish public support for policies and actions aimed at mitigating climate change.

Lessons for Ethical Data Visualization

- **Provide Complete Context:** Always include sufficient background information and long-term data to ensure accurate interpretation.
- **Avoid Cherry-Picking:** Present a comprehensive view of the data, avoiding selective presentation that supports a specific narrative.²
- **Ensure Data Accuracy:** Verify the accuracy and reliability of the data sources and methodologies used.
- **Transparency:** Clearly label and explain any data manipulations or selections.

News Source

- Associated Press (AP) Fact Check: "Temperature graph misrepresented to deny climate change," authored by Sophia Tulp, published on January 19, 2023.

Q. 4 Train Classification Model and visualize the prediction performance of trained model

Required information

Data File: Classification data.csv

Class Label: Last Column

Use any Machine Learning model (SVM, Naïve Base Classifier)

Requirements to satisfy

Programming Language: Python

Class imbalance should be resolved

Data Pre-processing must be used

Hyper parameter tuning must be used

Train, Validation and Test Split should be 70/20/10

Train and Test split must be randomly done

Classification Accuracy should be maximized

Use any Python library to present the accuracy measures of trained model

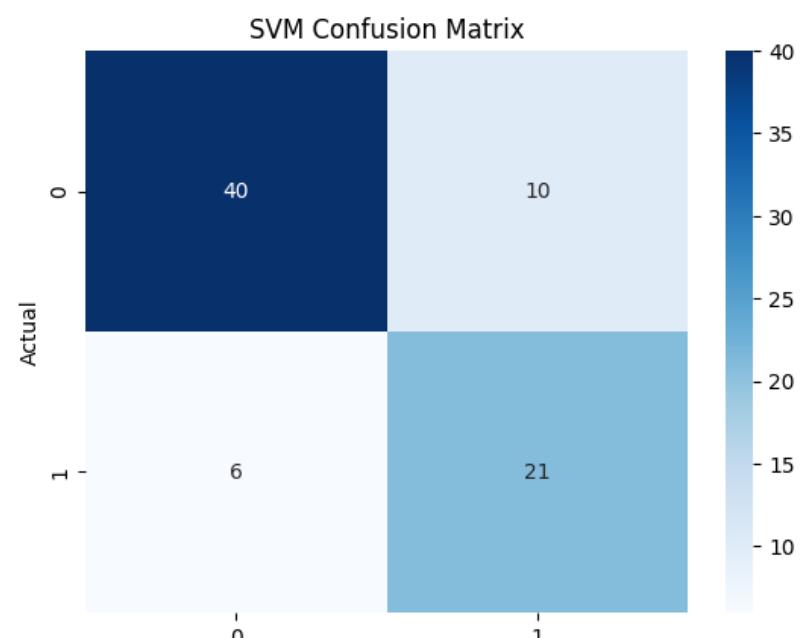
[Pima Indians Diabetes Database](#)

Explanation:

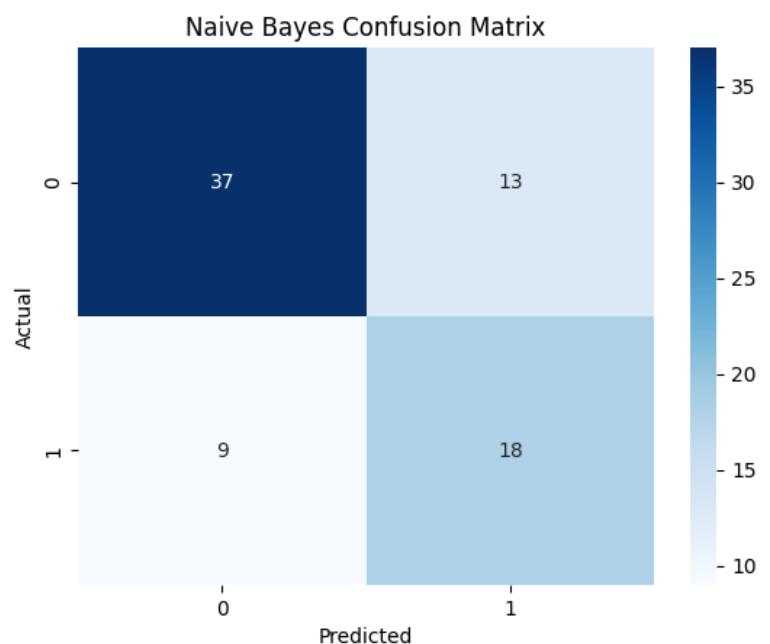
1. **Import Libraries:** Necessary libraries like numpy, pandas, sklearn, and matplotlib are imported.
2. **Load Data:** The "diabetes.csv" file is loaded using pandas.

3. **Data Separation:** The features (X) and the target variable (y) are separated. It's assumed the last column is the class label.
4. **Class Imbalance Resolution:** SMOTE (Synthetic Minority Over-sampling Technique) is used to handle class imbalance. This generates synthetic samples for the minority class.
5. **Data Pre-processing:** StandardScaler is applied to scale the features. This is crucial for models like SVM and Naive Bayes.
6. **Train-Validation-Test Split:** The data is split into 70% train, 20% validation, and 10% test sets, with stratification to maintain class ratios.
7. **Hyperparameter Tuning:** GridSearchCV is used with StratifiedKFold cross-validation to find the best hyperparameters for the Gaussian Naive Bayes model.
8. **Model Evaluation:** The model's performance is evaluated on the test set, and a classification report and confusion matrix are printed.
9. **Visualization:** The confusion matrix is visualized using Seaborn and Matplotlib.

--- SVM Validation Report ---				
	precision	recall	f1-score	support
0	0.81	0.73	0.77	100
1	0.58	0.69	0.63	54
accuracy			0.71	154
macro avg	0.69	0.71	0.70	154
weighted avg	0.73	0.71	0.72	154
--- SVM Test Report ---				
	precision	recall	f1-score	support
0	0.87	0.80	0.83	50
1	0.68	0.78	0.72	27
accuracy			0.79	77
macro avg	0.77	0.79	0.78	77
weighted avg	0.80	0.79	0.80	77



--- Naive Bayes Validation Report ---				
	precision	recall	f1-score	support
0	0.86	0.77	0.81	100
1	0.64	0.76	0.69	54
accuracy			0.77	154
macro avg	0.75	0.76	0.75	154
weighted avg	0.78	0.77	0.77	154
--- Naive Bayes Test Report ---				
	precision	recall	f1-score	support
0	0.80	0.74	0.77	50
1	0.58	0.67	0.62	27
accuracy			0.71	77
macro avg	0.69	0.70	0.70	77
weighted avg	0.73	0.71	0.72	77



Q.5 Train Regression Model and visualize the prediction performance of trained model

Data File: Regression data.csv

Independent Variable: 1st Column

Dependent variables: Column 2 to 5

Use any Regression model to predict the values of all Dependent variables using values of Ist column.

Requirements to satisfy:

Programming Language: Python

OOP approach must be followed

Hyper parameter tuning must be used

Train and Test Split should be 70/30

Train and Test split must be randomly done

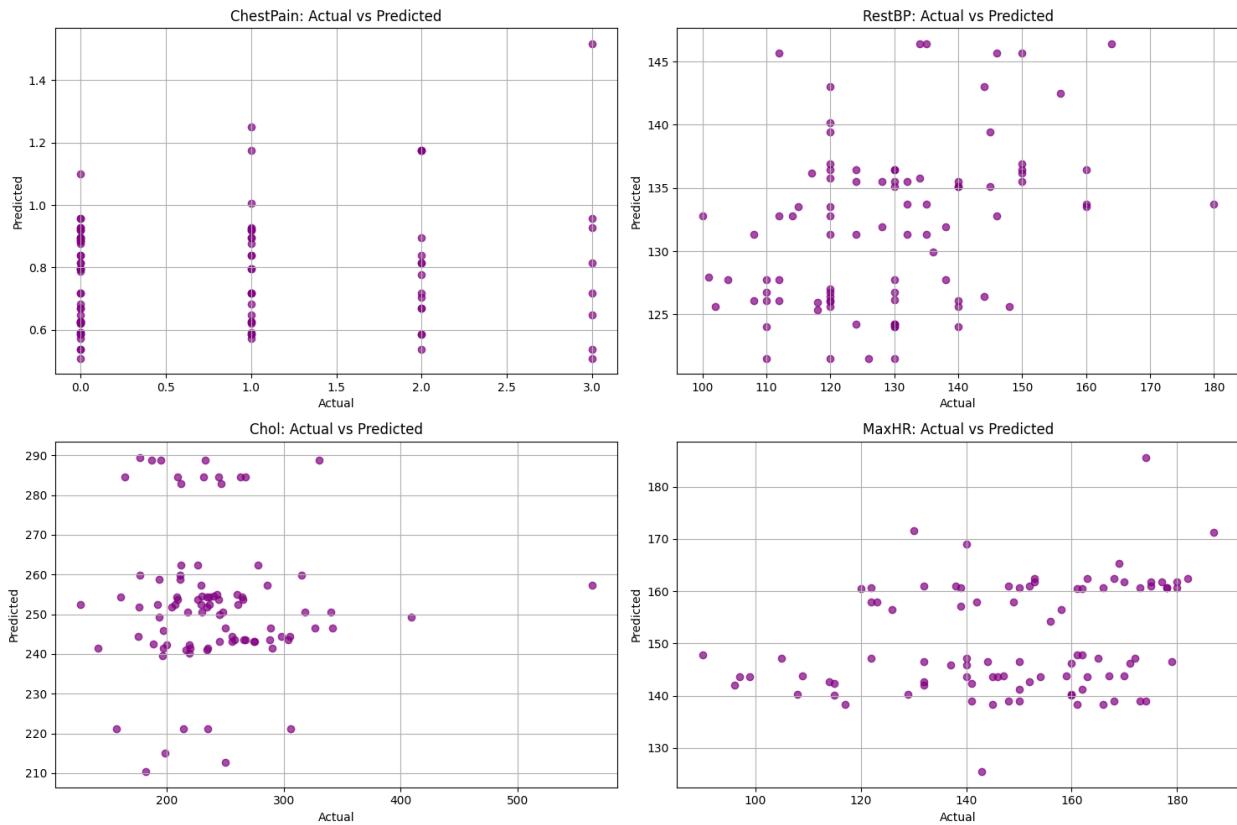
Adjusted R2 score should more than 0.99

Use any Python library to present the accuracy measures of trained model

```
Adjusted R2      : -0.0056
RestBP:
R2 Score       : 0.0992
Adjusted R2     : 0.0890

Chol:
R2 Score       : -0.1291
Adjusted R2     : -0.1418

MaxHR:
R2 Score       : 0.0389
Adjusted R2     : 0.0281
```



Q.6 What are the key features of the wine quality data set? Discuss the importance of each feature in predicting the quality of wine? How did you handle missing data in the wine quality data set during the feature engineering process? Discuss the advantages and disadvantages of different imputation techniques. (Refer dataset from Kaggle).

The wine quality data set from Kaggle primarily contains physicochemical measurements of wines and a quality score that reflects sensory evaluation. The key features (predictor variables) include:

- **Fixed Acidity:** Represents non-volatile acids that contribute to the wine's sour taste; it affects balance and structure.
- **Volatile Acidity:** Measures the presence of acetic acid; high levels can impart an unpleasant vinegar taste, negatively impacting quality.
- **Citric Acid:** Adds freshness and helps balance the wine's overall acidity; moderate amounts can enhance flavor.
- **Residual Sugar:** Indicates the unfermented sugar left in the wine; it plays a role in sweetness and overall flavor harmony, especially in white wines.

- **Chlorides:** Reflects the salt content; excessive levels can indicate poor sanitation or imbalance, thus reducing quality.
- **Free Sulfur Dioxide:** Acts as an antimicrobial and antioxidant; helps preserve wine flavor and stability, though excessive amounts may create off flavors.
- **Total Sulfur Dioxide:** Represents the overall amount used for preservation; important for shelf life but must be within safe sensory thresholds.
- **Density:** Closely related to the sugar content; contributes to the body and mouthfeel of the wine.
- **pH:** Provides an indication of wine acidity; optimal pH levels are necessary for microbial stability and overall flavor balance.
- **Sulphates:** Contribute to the wine's aroma and taste; they enhance complexity and act as an additional preservative measure.
- **Alcohol:** Affects body, viscosity, and flavor intensity; higher alcohol can balance high acidity in robust wines but may be overpowering if in excess.

Each of these features is vital because they collectively inform a model that predicts quality by capturing taste, balance, preservation, and overall sensory attributes. For example, while fixed and volatile acidity set the baseline for taste, residual sugar and citric acid can moderate harsh flavors; density and pH add to the body and stability, and alcohol contributes to both the structural and aromatic dimensions.

Missing Data Handling & Imputation Techniques:

In the wine quality data set, missing values may arise during data collection or preprocessing. Although many versions of this popular data set are complete, handling missing data is an essential step in feature engineering. The common imputation techniques include:

1. Mean/Median Imputation:

- *Advantages:* Simple to implement and preserves the dataset size.
- *Disadvantages:* Can distort distribution properties (especially with skewed data) and reduce variability.

2. K-Nearest Neighbors (KNN) Imputation:

- *Advantages:* Uses local similarities to estimate missing values, which often leads to more accurate imputations.

- *Disadvantages:* Computationally intensive and sensitive to the choice of distance metric and K value.

3. Regression Imputation:

- *Advantages:* Leverages relationships among variables to predict missing values, potentially capturing complex dependencies.
- *Disadvantages:* Can overfit and underestimate variability, as predicted values may cluster too tightly.

Each technique carries trade-offs; the best choice depends on the extent and pattern of missingness as well as the underlying data distribution. In practice, one might evaluate the “missing completely at random” (MCAR) assumption to decide if a simple mean/median imputation is sufficient, or if more sophisticated methods like KNN or regression imputation are required.

Overall, understanding the contribution of each feature helps build more accurate predictive models for wine quality, while careful handling of missing data ensures that the integrity of the model is maintained.