

VISVESVARAYA TECHNOLOGICAL UNIVERSITY
“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT
on
Machine Learning (23CS6PCMAL)

Submitted by

SUJNYAN ULHAS KINI (1BM22CS340)

in partial fulfillment for the award of the degree of

BACHELOR OF ENGINEERING
in
COMPUTER SCIENCE AND ENGINEERING



B.M.S. COLLEGE OF ENGINEERING
(Autonomous Institution under VTU)
BENGALURU-560019
Sep-2024 to Jan-2025

**B.M.S. College of Engineering,
Bull Temple Road, Bangalore 560019**
(Affiliated To Visvesvaraya Technological University, Belgaum)
Department of Computer Science and Engineering



CERTIFICATE

This is to certify that the Lab work entitled “Machine Learning (23CS6PCMAL)” carried out by **SUJNYAN KINI(1BM22CS340)**, who is a bonafide student of **B.M.S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Lab report has been approved as it satisfies the academic requirements in respect of an Machine Learning (23CS6PCMAL) work prescribed for the said degree.

Lab Faculty Incharge Name: Sarita A N Assistant Professor Department of CSE, BMSCE	Dr. Kavitha Sooda Professor & HOD Department of CSE, BMSCE
--	--

Index

Sl. No.	Date	Experiment Title	Page No.
1	21-2-2025	Write a python program to import and export data using Pandas library functions	5-11
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	12-19
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	20-26
4	17-3-2025	Build Logistic Regression Model for a given dataset	27-30
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	31-39
6	7-4-2025	Build KNN Classification model for a given dataset.	40-46
7	21-4-2025	Build Support vector machine model for a given dataset	47-52
8	5-5-2025	Implement Random forest ensemble method on a given dataset.	53-56
9	5-5-2025	Implement Boosting ensemble method on a given dataset.	57-60
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	61-64
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	65-78

GitHub ID:
<https://github.com/sujkini/ML-Lab>

INDEX

Name : Dijyan Kini Class :
Section : D Roll No. : 340 Subject : ML Lab - Observation

Program 1

Write a python program to import and export data using Pandas library functions.

Screenshot:

The screenshot shows a handwritten note on a lined notebook page. At the top left, there is a section labeled "To Do:" followed by several imports:

```
import yfinance as yf
import pandas as pd
import matplotlib.pyplot as plt
```

Below the imports, there is a list of tickers:

```
tickers = ["HDFC BANK.NS", "ICICIBANK.NS",  
          "KOTAKBANK.NS"]
```

Then, the code uses `yf.download` to download data for these tickers from January 1, 2024, to December 31, 2024, grouped by ticker:

```
data = yf.download(tickers, start="2024-01-01",  
                  end="2024-12-31",  
                  group_by='ticker')
```

Next, the code extracts the closing price data for HDFC Bank:

```
reliance_data = data['HDFC BANK']
```

It then calculates the daily return for HDFC Bank:

```
reliance_data['Daily Returns'] = reliance_data['Close'].pct_change()
```

The code then creates a plot for HDFC Bank's closing prices:

```
plt.figure(figsize=(12, 16))  
plt.subplot(2, 1, 1)  
reliance_data['Close'].plot(title="HDFC Industries -  
                           Closing Prices")
```

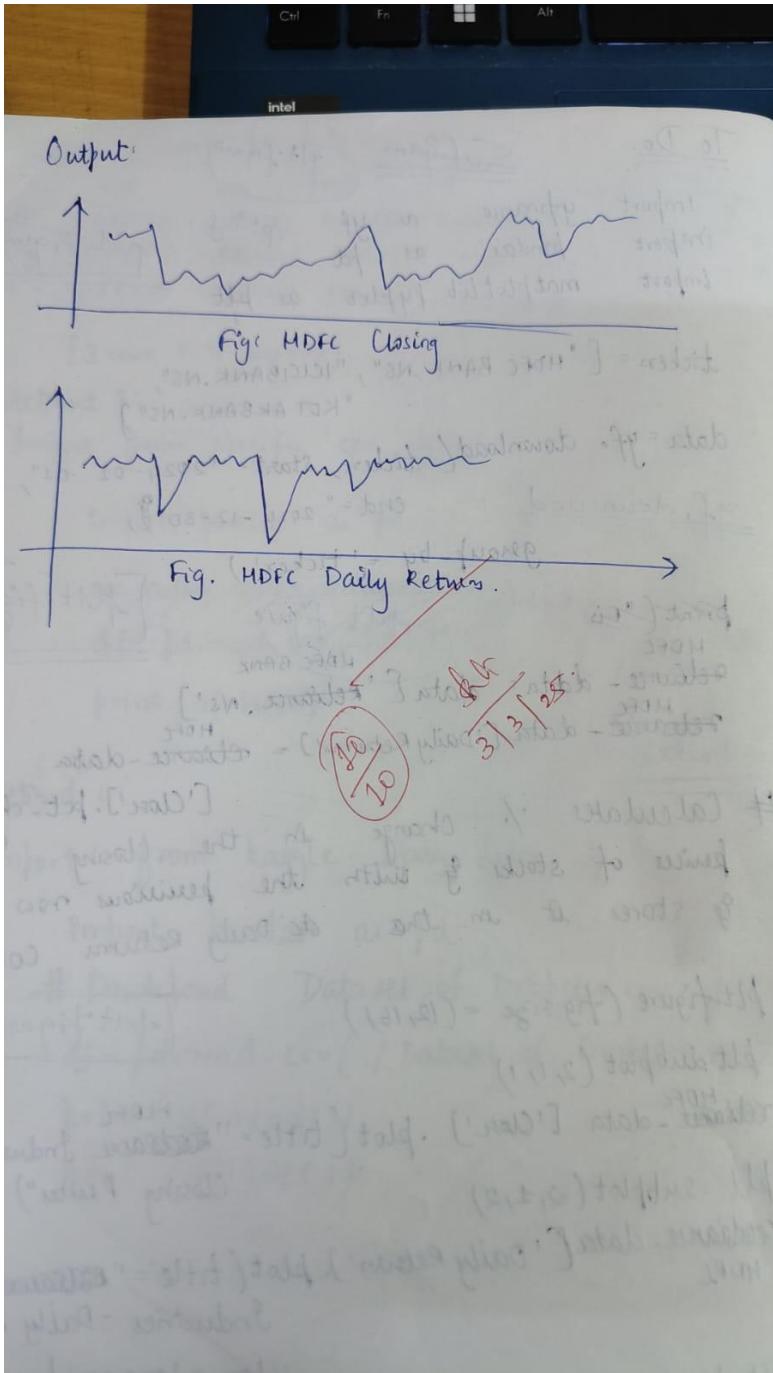
It also plots the daily returns for HDFC Bank:

```
plt.subplot(2, 1, 2)  
reliance_data['Daily Returns'].plot(title="HDFC  
                                      Industries - Daily Returns",  
                                      color='orange')
```

Finally, it saves the data to a CSV file:

```
reliance_data.to_csv('reliance-stock-data.csv')
```

Annotations in the margin include circled sections for `plt.figure`, `reliance_data`, and `reliance_data['Close'].pct_change()`.



Code:

```
import pandas as pd
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David'],
    'Age': [25, 30, 35, 40],
    'City': ['New York', 'Los Angeles', 'Chicago', 'Houston']
}
df = pd.DataFrame(data)
print("Sample data:")
print(df.head())

from sklearn.datasets import load_iris
iris = load_iris()
df = pd.DataFrame(iris.data, columns=iris.feature_names)
df['target'] = iris.target
print("Sample data:")
print(df.head())

from sklearn.datasets import load_iris
iris = load_iris()
df = pd.DataFrame(iris.data, columns=iris.feature_names)
df['target'] = iris.target
print("Sample data:")
print(df.head())

file_path = 'mobiles-dataset-2025.csv'
df = pd.read_csv(file_path, encoding='latin-1') # or 'cp1252' or other suitable encoding
print("Sample data:")
print(df.head())

import pandas as pd
data = {
```

```

'USN': ['IS001','IS002','IS003','IS004','IS005'],
'Name': ['Alice', 'Bob', 'Charlie', 'David','Eve'],
'Marks': [25, 30, 35, 40,45]
}

df=pd.DataFrame(data)
print("Sample data:")
print(df.head())

file_path = 'sample_sales_data.csv'
df = pd.read_csv(file_path)
print("Sample data:")
print(df.head())
print("\n")

df = pd.read_csv("/content/dataset-of-diabetes .csv",encoding='latin-1')
print("Sample data:")
print(df.head())
print("\n")

df =pd.read_csv('sample_sales_data.csv')
print("Sample data:")
print(df.head())
df.to_csv('output.csv',index=False)
print("Data saved to output.csv")

sales_df =pd.read_csv('sample_sales_data.csv')
print("Sample data:")
print(sales_df.head())
sales_by_region =sales_df.groupby('Region')['Sales'].sum()
print("\nTotal sales by region:")
print(sales_by_region)

best_selling_products =sales_df.groupby('Product')['Quantity'].sum().sort_values(ascending=False)

```

```
print("\nBest-selling products by quantity:")
print(best_selling_products)
sales_by_region.to_csv('sales_by_region.csv')
best_selling_products.to_csv('best_selling_products.csv')
print("Data saved to sales_by_region.csv and best_selling_products.csv")
```

```
import yfinance as yf
import matplotlib.pyplot as plt
tickers = ["RELIANCE.NS", "TCS.NS", "INFY.NS"]
data = yf.download(tickers, start="2022-10-01", end="2023-10-01",
                   group_by='ticker')
print("First 5 rows of the dataset:")
print(data.head())
print("\nShape of the dataset:")
print(data.shape)
print("\nColumn names:")
print(data.columns)
print("\n")
reliance_data = data['RELIANCE.NS']
print("\nSummary statistics for Reliance Industries:")
print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['Close'].pct_change()
print("\n")
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="Reliance Industries - Closing Price")
plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="Reliance Industries - Daily Returns", color='orange')
plt.tight_layout()
plt.show()
reliance_data.to_csv('reliance_stock_data.csv')
```

```

tickers = ["HDFCBANK.NS", "ICICI.NS", "KOTAKBANK.NS"]
data = yf.download(tickers, start="2024-01-01", end="2024-12-30",
                   group_by='ticker')
print("First 5 rows of the dataset:")
print(data.head())
print("\nShape of the dataset:")
print(data.shape)
print("\nColumn names:")
print(data.columns)
print("\n")
reliance_data = data['HDFCBANK.NS']
print("\nSummary statistics for Reliance Industries:")
print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['Close'].pct_change()
print("\n")
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="HDFC Industries - Closing Price")
plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="HDFC Industries - Daily Returns", color='red')
plt.tight_layout()
plt.show()
reliance_data.to_csv('hdfc_stock_data.csv')
print("\nhdfc stock data saved to 'hdfc_stock_data.csv'.")

```

Program 2

Demonstrate various data pre-processing techniques for a given dataset.

Screenshot:

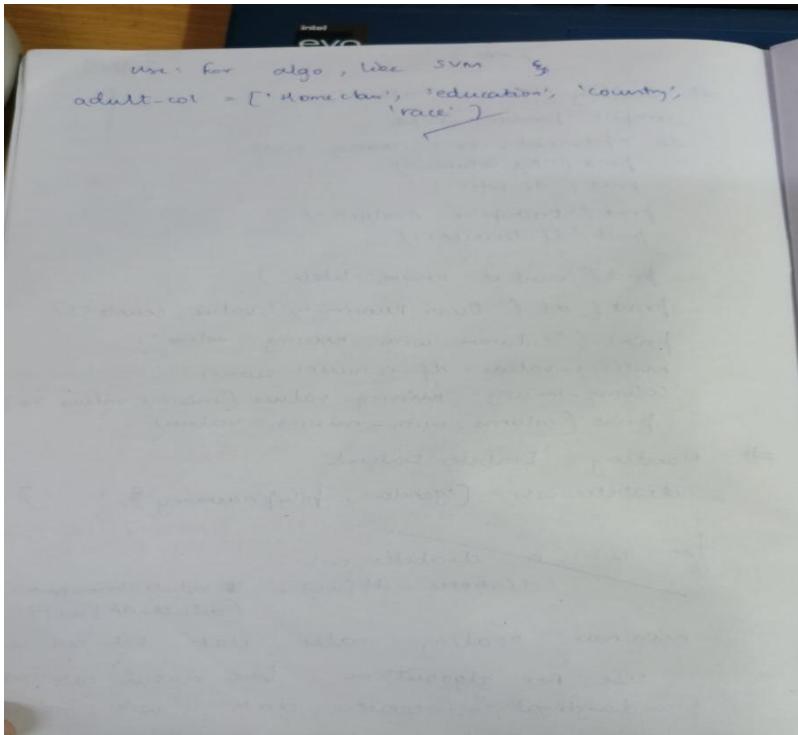
The image shows a laptop screen with handwritten notes in blue ink. At the top left, it says "LAB-2.". The code is organized into two main sections: "# Handling 'Housing Dataset'" and "# Handling 'Diabetes Dataset'".

```
# Handling 'Housing Dataset'
import pandas as pd
df = pd.read_csv('housing.csv')
print("All Columns")
print(df.info())
print("Descriptive Analysis")
print(df.describe())
print("Count of unique labels")
print(df['Ocean Proximity'].value_counts())
print("Columns with missing values")
missing_values = df.isnull().sum()
column_missing = missing_values[missing_values > 0]
print(column_missing)

# Handling "Diabetes Dataset"
diabetes_cols = ['gender', 'polypharmacy', ...]

for col in diabetes_cols:
    diabetes_df[col] = adult_transform(adult_df[col])

min-max scaling: scales data between 0 and 1.
use: for algorithms like neural network
standardization: center data with mean
```



Code:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.impute import SimpleImputer

try:
    diabetes_df = pd.read_csv('diabetes.csv')
    adult_df = pd.read_csv('adult.csv')
except FileNotFoundError:
    print("Error: Please upload 'diabetes.csv' and 'adult.csv' to your Google Colab environment.")
    exit()

diabetes_df.head(10)
adult_df.head(10)

diabetes_df.shape
adult_df.shape

#Handling Missing Values
diabetes_numeric_cols = diabetes_df.select_dtypes(include=[np.number]).columns
diabetes_categorical_cols = diabetes_df.select_dtypes(exclude=[np.number]).columns

adult_numeric_cols = adult_df.select_dtypes(include=[np.number]).columns
adult_categorical_cols = adult_df.select_dtypes(exclude=[np.number]).columns

diabetes_numeric_imputer = SimpleImputer(strategy='mean')
adult_numeric_imputer = SimpleImputer(strategy='mean')
diabetes_df[diabetes_numeric_cols] =
diabetes_numeric_imputer.fit_transform(diabetes_df[diabetes_numeric_cols])
```

```
adult_df[adult_numeric_cols] = adult_numeric_imputer.fit_transform(adult_df[adult_numeric_cols])
```

```

diabetes_categorical_imputer = SimpleImputer(strategy='most_frequent')
adult_categorical_imputer = SimpleImputer(strategy='most_frequent')
diabetes_df[diabetes_categorical_cols] =
diabetes_categorical_imputer.fit_transform(diabetes_df[diabetes_categorical_cols])
adult_df[adult_categorical_cols] =
adult_categorical_imputer.fit_transform(adult_df[adult_categorical_cols])

print("Missing values in Diabetes dataset after imputation:")
print(diabetes_df.isnull().sum())
print("Missing values in Adult Income dataset after imputation:")
print(adult_df.isnull().sum())

adult_df.replace("?", np.nan, inplace=True)
print("Missing values in Adult Income dataset after replacing '?'")
print(adult_df.isnull().sum())

from sklearn.impute import SimpleImputer

# Identify numeric and categorical columns
adult_numeric_cols = adult_df.select_dtypes(include=[np.number]).columns
adult_categorical_cols = adult_df.select_dtypes(exclude=[np.number]).columns

# Handle missing values in numeric columns using mean imputation
adult_numeric_imputer = SimpleImputer(strategy='mean')
adult_df[adult_numeric_cols] = adult_numeric_imputer.fit_transform(adult_df[adult_numeric_cols])

# Handle missing values in categorical columns using most frequent imputation
adult_categorical_imputer = SimpleImputer(strategy='most_frequent')
adult_df[adult_categorical_cols] =
adult_categorical_imputer.fit_transform(adult_df[adult_categorical_cols])
print("Missing values in Adult Income dataset after imputation:")

```

```

print(adult_df.isnull().sum())

from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()

# Encode categorical columns in Diabetes dataset
for col in diabetes_categorical_cols:
    diabetes_df[col] = label_encoder.fit_transform(diabetes_df[col])

# Encode categorical columns in Adult Income dataset
for col in adult_categorical_cols:
    adult_df[col] = label_encoder.fit_transform(adult_df[col])

print("Encoded columns in Diabetes dataset:")
print(diabetes_df.head())
print("Encoded columns in Adult Income dataset:")
print(adult_df.head())

#Handling outliers
def remove_outliers(df):
    Q1 = df.quantile(0.25)
    Q3 = df.quantile(0.75)
    IQR = Q3 - Q1
    df_no_outliers = df[~((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))).any(axis=1)]
    return df_no_outliers

diabetes_df_no_outliers = remove_outliers(diabetes_df)
adult_df_no_outliers = remove_outliers(adult_df)

print("Diabetes dataset shape after removing outliers:", diabetes_df_no_outliers.shape)
print("Adult Income dataset shape after removing outliers:", adult_df_no_outliers.shape)

#Min-max scaling
from sklearn.preprocessing import MinMaxScaler

```

```

min_max_scaler = MinMaxScaler()
diabetes_scaled_minmax = pd.DataFrame(min_max_scaler.fit_transform(diabetes_df_no_outliers),
columns=diabetes_df_no_outliers.columns)
adult_scaled_minmax = pd.DataFrame(min_max_scaler.fit_transform(adult_df_no_outliers),
columns=adult_df_no_outliers.columns)
print("Diabetes dataset after Min-Max scaling:")
print(diabetes_scaled_minmax.head())
print("Adult Income dataset after Min-Max scaling:")
print(adult_scaled_minmax.head())

# Initialize Standard Scaler
from sklearn.preprocessing import StandardScaler
standard_scaler = StandardScaler()
diabetes_scaled_standard = pd.DataFrame(standard_scaler.fit_transform(diabetes_df_no_outliers),
columns=diabetes_df_no_outliers.columns)
adult_scaled_standard = pd.DataFrame(standard_scaler.fit_transform(adult_df_no_outliers),
columns=adult_df_no_outliers.columns)
print("Diabetes dataset after Standard scaling:")
print(diabetes_scaled_standard.head())
print("Adult Income dataset after Standard scaling:")
print(adult_scaled_standard.head())

```

Program 3

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset.

Screenshot:

LAB-3

Simple Linear Regression

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

def estimate_coeff(x, y):
    n = np.size(x)
    m_x = np.mean(x)
    m_y = np.mean(y)
    ss_xy = np.sum((x - m_x) * (y - m_y))
    ss_xx = np.sum((x - m_x) ** 2)
    b_1 = ss_xy / ss_xx
    b_0 = m_y - b_1 * m_x
    return (b_0, b_1)

def plot(x, y, b):
    plt.scatter(x, y, color='m')
    y_pred = b[0] + b[1] * x
    plt.plot(x, y_pred, color='g')
    plt.xlabel('x')
    plt.ylabel('y')
    plt.show()

file_path = input("Enter path")
df = pd.read_csv(file_path)
x = df.iloc[:, 0].values
y = df.iloc[:, 1].values
b = estimate_coeff(x, y)
plot(x, y, b)

Output:
Enter path: /content/
estimated coefficients:
b_0 = 7.032
b_1 = 0.047
```

Multiple Linear Regression
 $\text{data} = \{$
 "Feature 1": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
 "Feature 2": [2, 3, 5, 7, 13, 17, 19, 23, 29],
 "Feature 3": [3, 6, 9, 12, 15, 18, 21, 24, 27, 30],
 "Feature 4": [5, 9, 15, 22, 31, 41, 53, 66, 80, 96]
 $\}$
 y
 $df = pd.DataFrame(\text{data})$
 $x = df.drop(\text{columns} = ["Target"], axis=1)$
 $y = df["Target"]$
 $x = df.drop("Target", axis=1)$
 $X = np.vstack((np.ones((X.shape[0], 1)), X))$
 $\beta = np.linalg.solve(X.T @ X + 0.01 * np.identity(X.shape[1]), X.T @ y)$
 $y_pred = X @ \beta$
 $mse = np.mean((y - y_pred) ** 2)$
 $tot_var = np.sum((y - np.mean(y)) ** 2)$
 $exp_var = np.sum((y - y_pred - np.mean(y)) ** 2)$
 $r^2 = exp_var / tot_var$
Output:
 Model Coefficients: [0.040, 3.313, 0.120]
 Intercept: -18.1599
 Mean Squared Error: 5.362
 R-squared Score: 0.993

Code:

```
import pandas as pd  
import numpy as np  
from sklearn import linear_model  
import matplotlib.pyplot as plt
```

```

df=pd.read_csv('housing_area_price.csv')
plt.xlabel('area')
plt.ylabel('price')
plt.scatter(df.area,df.price,color='red',marker='+')
new_df = df.drop('price',axis='columns')
new_df
price = df.price
reg = linear_model.LinearRegression()
reg.fit(new_df,price)

```

#(1) Predict price of a home with area = 3300 sqrft

```

reg.predict([[3300]])
reg.coef_
reg.intercept_
3300*135.78767123 + 180616.43835616432

```

#(2) Predict price of a home with area = 5000 sqrft

```
reg.predict([[5000]])
```

```

df=pd.read_csv('homeprices_Multiple_LR.csv')
df.bedrooms.median()
df.bedrooms = df.bedrooms.fillna(df.bedrooms.median())
reg = linear_model.LinearRegression()
reg.fit(df.drop('price',axis='columns'),df.price)
reg.coef_
reg.intercept_

```

#Find price of home with 3000 sqr ft area, 3 bedrooms, 40 year old

```

reg.predict([[3000, 3, 40]])
112.06244194*3000 + 23388.88007794*3 + -3231.71790863*40 + 221323.00186540384

```

```

df = pd.read_csv('canada_per_capita_income.csv')
print(df.head())
X = df[['year']]
y = df['per capita income (US$)']
reg = LinearRegression()
reg.fit(X, y)
predicted_income_2020 = reg.predict([[2020]])
print(f"Predicted per capita income for Canada in 2020: {predicted_income_2020[0]:.2f}")

```

```

plt.scatter(X, y, color='blue')
plt.plot(X, reg.predict(X), color='red')
plt.xlabel('Year')
plt.ylabel('Per Capita Income')
plt.title('Per Capita Income in Canada Over the Years')
plt.show()

```

```

df = pd.read_csv('salary.csv')
print(df.head())
print("Missing values in the dataset:")
print(df.isnull().sum())

```

```

df['YearsExperience'] = df['YearsExperience'].fillna(df['YearsExperience'].median())
print("\nMissing values after filling:")
print(df.isnull().sum())
X = df[['YearsExperience']]
y = df['Salary']
reg = LinearRegression()
reg.fit(X, y)
predicted_salary_12_years = reg.predict([[12]])
print(f"\nPredicted salary for an employee with 12 years of experience:
${predicted_salary_12_years[0]:,.2f}")

```

```

plt.scatter(X, y, color='blue')
plt.plot(X, reg.predict(X), color='red')
plt.xlabel('Years of Experience')
plt.ylabel('Salary')
plt.title('Salary vs. Years of Experience')
plt.show()

def convert_to_numeric(value):
    word_to_num = {
        'zero': 0, 'one': 1, 'two': 2, 'three': 3, 'four': 4, 'five': 5,
        'six': 6, 'seven': 7, 'eight': 8, 'nine': 9, 'ten': 10,
        'eleven': 11, 'twelve': 12, 'thirteen': 13, 'fourteen': 14,
        'fifteen': 15
    }
    return word_to_num.get(value.lower(), value) if isinstance(value, str) else value

df_hiring = pd.read_csv('hirings.csv')
print(df.head())
df_hiring['experience'] = df_hiring['experience'].apply(convert_to_numeric)

df_hiring['experience'].fillna(0, inplace=True)
df_hiring['test_score(out of 10)'].fillna(df_hiring['test_score(out of 10)'].median(), inplace=True)
df_hiring['interview_score(out of 10)'].fillna(df_hiring['interview_score(out of 10)'].median(),
inplace=True)

X_hiring = df_hiring[['experience', 'test_score(out of 10)', 'interview_score(out of 10)']]
y_hiring = df_hiring['salary($)']
reg_hiring = LinearRegression()
reg_hiring.fit(X_hiring, y_hiring)
candidates = np.array([[2, 9, 6], [12, 10, 10]])
predicted_salaries = reg_hiring.predict(candidates)

for i, candidate in enumerate(candidates):

```

```

print(f"\nPredicted salary for candidate with {candidate[0]} yrs experience, {candidate[1]} test score,
{candidate[2]} interview score: {predicted_salaries[i]:.2f} USD")

plt.scatter(y_hiring, reg_hiring.predict(X_hiring), color='blue', label='Predicted vs Actual')
plt.xlabel("Actual Salary")
plt.ylabel("Predicted Salary")
plt.title("Actual vs Predicted Salary")
plt.legend()
plt.show()

df_companies = pd.read_csv('1000_Companies.csv')
print(df.head())
label_encoder = LabelEncoder()
df_companies['State'] = label_encoder.fit_transform(df_companies['State'])
X_companies = df_companies[['R&D Spend', 'Administration', 'Marketing Spend', 'State']]
y_companies = df_companies['Profit']
df_companies.fillna(df_companies.median(), inplace=True)
reg_companies = LinearRegression()
reg_companies.fit(X_companies, y_companies)

input_data = np.array([[91694.48, 515841.3, 11931.24, label_encoder.transform(['Florida'])[0]]])
predicted_profit = reg_companies.predict(input_data)
print(f"Predicted profit: {predicted_profit[0]:.2f} USD")

plt.scatter(y_companies, reg_companies.predict(X_companies), color='blue', label='Predicted vs
Actual')
plt.xlabel("Actual Profit")
plt.ylabel("Predicted Profit")
plt.title("Actual vs Predicted Profit")
plt.legend()
plt.show()

```

Program 4

Build Logistic Regression Model for a given dataset.

Screenshot:

```
Logistic Regression
def sigmoid(z):
    return 1 / (1 + np.exp(-z))

def Compute(m, y, theta):
    m = len(y)
    h = sigmoid(m @ theta)
    cost = (-1/m) * np.sum(y * np.log(h) + (1 - y) * np.log(1 - h))
    return cost

def grad_descent(X, y, theta, alpha, iterations):
    m = len(y)
    cost_history = []
    for i in range(iterations):
        gradient = (-1/m) * X.T @ (sigmoid(X @ theta) - y)
        theta -= alpha * gradient
    return theta, cost_history

def predict(X, theta):
    return (sigmoid(X @ theta) >= 0.5).astype(int)

accuracy = np.mean(y == predict(X, theta))
print(f'Accuracy: {accuracy * 100} %')

Output:
Accuracy: 1
```

Code:

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

df=pd.read_csv("HR_comma_sep.csv")
print(df.info())
numericCols = df.select_dtypes(include=['float64', 'int64']).columns
plt.figure(figsize=(10, 8))
sns.heatmap(df[numericCols].corr(), annot=True, cmap='coolwarm', fmt='.2f')
plt.title("Correlation Matrix (Numeric Features)")
plt.show()

plt.figure(figsize=(8, 6))
sns.countplot(x='salary', hue='left', data=df)
plt.title("Impact of Salary on Employee Retention")
plt.xlabel("Salary Level")
plt.ylabel("Employee Count")
plt.show()

import pandas as pd
df=pd.read_csv("zoo-data.csv")
print(df.info())
print(df.head())
print(df.isnull().sum())
df.drop(columns=['animal_name'], inplace=True)
X = df.drop(columns=['class_type'])
y = df['class_type']

from sklearn.model_selection import train_test_split
```

```

from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
logreg = LogisticRegression(max_iter=200, multi_class='multinomial', solver='lbfgs')
logreg.fit(X_train, y_train)

from sklearn.metrics import accuracy_score
y_pred = logreg.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.2f}")

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
import matplotlib.pyplot as plt
cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=logreg.classes_)
disp.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix for Zoo Animal Classification")
plt.show()

y_pred = logreg.predict(X_test)
pred_classes = [class_mapping[pred] for pred in y_pred]
print("Predicted Classes:", pred_classes)

import seaborn as sns
import matplotlib.pyplot as plt
sns.countplot(x='class_type', data=df)
plt.title("Class Distribution of Animals in Zoo Dataset")
plt.xlabel("Class Type")

```

```
plt.ylabel("Count")
plt.show()

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
cm = confusion_matrix(y_test, y_pred)
class_labels = [class_mapping[num] for num in logreg.classes_]
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=class_labels)
disp.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix with Class Names")
plt.show()
```

Program 5

Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.

Screenshot:

```
LAB-3: (ID3)
import pandas as pd
import numpy as np
from collections import Counter
def Entropy(y):
    counter = Counter(y)
    total = len(y)
    return -sum((c/total) * np.log2(x/total) for
                c in counter.values())
def Information_gain(data, feature, target):
    total_entropy = Entropy(data[target])
    values = data[feature].unique()
    weighted_entropy = 0
    for val in values:
        subset = data[data[feature] == val]
        weight = len(subset) / len(data)
        weighted_entropy += weight * Entropy(subset)
    return total_entropy - weighted_entropy
def id3(data, features, target):
    labels = data[target]
    if len(labels.unique()) == 1:
        return labels.iloc[0]
    if not features:
        return labels.mode()[0]
    gains = {f: information_gain(data, f, target) for f in
            features}
    best_feature = max(gains, key=gains.get)
    tree = {best_feature: {}}
    gains.pop(best_feature)
    for value in data[best_feature].unique():
        subtree = id3(data[data[best_feature] == value], [f for f in
                                                          features if f != best_feature],
                      target)
        tree[best_feature][value] = subtree
    return tree
Output:
{'Outlook': {'Overcast': 'Yes',
              'Rain': {'Wind': {'Weak': 'Yes',
                               'Strong': 'No'}}}}
```

Code:

```
import pandas as pd
import numpy as np
from collections import Counter
from pprint import pprint

# -----
# Step 1: Entropy calculation
# -----
def entropy(y):
    counts = Counter(y)
    total = len(y)
    return -sum((c/total) * np.log2(c/total) for c in counts.values())

# -----
# Step 2: Information Gain
# -----
def information_gain(data, feature, target):
    total_entropy = entropy(data[target])
    values = data[feature].unique()
    weighted_entropy = 0

    for val in values:
        subset = data[data[feature] == val]
        weight = len(subset) / len(data)
        weighted_entropy += weight * entropy(subset[target])

    return total_entropy - weighted_entropy

# -----
# Step 3: ID3 Algorithm
# -----
def id3(data, features, target):
    labels = data[target]

    # Base cases
    if len(labels.unique()) == 1:
        return labels.iloc[0]
    if not features:
        return labels.mode()[0]

    # Choose best feature
    gains = {f: information_gain(data, f, target) for f in features}
    best_feature = max(gains, key=gains.get)
```

```

tree = {best_feature: {}}

for value in data[best_feature].unique():
    subset = data[data[best_feature] == value].drop(columns=[best_feature])
    subtree = id3(subset, [f for f in features if f != best_feature], target)
    tree[best_feature][value] = subtree

return tree

# -----
# Step 4: Predict function
# -----
def predict(tree, sample):
    if not isinstance(tree, dict):
        return tree
    root = next(iter(tree))
    value = sample.get(root)
    subtree = tree[root].get(value)
    if subtree is None:
        return None # unknown value
    return predict(subtree, sample)

# -----
# Step 5: Example dataset
# -----
data = pd.DataFrame({
    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain', 'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast', 'Rain'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'High'],
    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Strong'],
    'Play': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
})

features = [col for col in data.columns if col != 'Play']
target = 'Play'

# -----
# Step 6: Build and print tree
# -----
tree = id3(data, features, target)
pprint(tree)

```

Program 6

Build KNN Classification model for a given dataset.

Screenshot:

Lab 4

① KNN Algo -

```
import math
from collections import Counter
import matplotlib.pyplot as plt

# data = [
#     ([65, 100], 'Female'),
#     ([70, 150], 'Male'),
#     ([60, 120], 'Female'),
#     ...
# ]
#
# label_colors = {'Male': 'blue', 'Female': 'red'}
```

def euclidean_distance(point1, point2):
 return math.sqrt(sum((a - b)**2 for a, b in zip(point1, point2)))

def knn(data, query, k):
 distances = []
 for features, label in data:
 dist = euclidean_distance(features, query)
 distances.append((dist, label, features))
 distances.sort(key=lambda x: x[0])
 k_nearest = distances[:k]
 k_nearest_labels = [label for _, label, _ in k_nearest]
 return Counter(k_nearest_labels).most_common(1)[0][0]

plt.Scatter(features[0], features[1], color=label, s=100, edgecolor='black')

```

plt.scatter(query[0], query[1], color='green', label='query',
           s=150, marker='x')

for i, j, neighbor_label in K_nearest:
    plt.plot([query[0], neighbor_label[0]], [query[1], neighbor_label[1]],
              neighbor_label[0], color='gray', linestyle='--')

handles, labels = plt.gca().get_legend_handles_labels()
by_label = dict(zip(handles, labels))
plt.legend(by_label, by_label.keys())

plt.xlabel("Height (in')")
plt.ylabel('Weight')
plt.title(f'kNN Classification (k={k})')
plt.grid(True)
plt.show()

prediction = Counter(K_nearest_labels).most_common(1)[0][0]
return prediction

```

query-point = [69, 130]
~~k=3~~
result = knn(data, query-point, k)
print(f"Predicted Gender: {result}")

Code:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import seaborn as sns
import matplotlib.pyplot as plt

iris_df = pd.read_csv('iris.csv')
le = LabelEncoder()
iris_df['species'] = le.fit_transform(iris_df['species'])
X = iris_df.drop('species', axis=1)
y = iris_df['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

error_rates = []
accuracies = []
k_values = range(1, 10)
for k in k_values:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred_k = knn.predict(X_test)
    error = 1 - accuracy_score(y_test, y_pred_k)
    error_rates.append(error)
    accuracies.append(accuracy_score(y_test, y_pred_k))

plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
plt.plot(k_values, accuracies, marker='o', color='blue')
plt.title("Accuracy vs K")
```

```

plt.xlabel("K Value")
plt.ylabel("Accuracy")

plt.subplot(1, 2, 2)
plt.plot(k_values, error_rates, marker='o', color='red')
plt.title("Error Rate vs K")
plt.xlabel("K Value")
plt.ylabel("Error Rate")
plt.tight_layout()
plt.show()

best_k = k_values[accuracies.index(max(accuracies))]
print(f"Best K: {best_k} with Accuracy: {max(accuracies):.2f}")

```

```

knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

```

Evaluation

```

print("\n==== Final Evaluation on IRIS Dataset ====")
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nClassification Report:")
print(classification_report(y_test, y_pred, labels=[0, 1, 2], target_names=le.classes_))

```

Confusion Matrix

```

cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
            xticklabels=le.classes_, yticklabels=le.classes_)
plt.title("Confusion Matrix - IRIS")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

```

df = pd.read_csv('diabetes.csv')
X = df.drop('Outcome', axis=1)
y = df['Outcome']
scaler      =      StandardScaler()
X_scaled = scaler.fit_transform(X)
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.2, random_state=42, stratify=y
)
accuracy_scores = []
k_range = range(1, 21)

for k in k_range:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred_k = knn.predict(X_test)
    acc = accuracy_score(y_test, y_pred_k)
    accuracy_scores.append(acc)

plt.figure(figsize=(8, 5))
plt.plot(k_range, accuracy_scores, marker='o', color='purple')
plt.title("Accuracy vs K (Diabetes Dataset)")
plt.xlabel("K Value")
plt.ylabel("Accuracy")
plt.xticks(k_range)
plt.grid()
plt.show()

best_k = k_range[accuracy_scores.index(max(accuracy_scores))]
print(f"Best K: {best_k} with Accuracy: {max(accuracy_scores):.2f}")

knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

```

```

print("==== Final Evaluation (Diabetes Dataset) ====")
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nClassification Report:")
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Purples', xticklabels=['No Diabetes', 'Diabetes'],
            yticklabels=['No Diabetes', 'Diabetes'])
plt.title("Confusion Matrix - Diabetes")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

```

heart_df = pd.read_csv('heart.csv')
X = heart_df.drop('target', axis=1)
y = heart_df['target']
scaler      =      StandardScaler()
X_scaled = scaler.fit_transform(X)
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.2, random_state=42, stratify=y
)
accuracy_scores = []
k_range = range(1, 21)

```

for k in k_range:

```

knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X_train, y_train)
y_pred_k = knn.predict(X_test)
acc = accuracy_score(y_test, y_pred_k)
accuracy_scores.append(acc)

```

```

plt.figure(figsize=(8, 5))
plt.plot(k_range, accuracy_scores, marker='o', color='red')
plt.title("Accuracy vs K (Heart Dataset)")
plt.xlabel("K Value")
plt.ylabel("Accuracy")
plt.xticks(k_range)
plt.grid()
plt.show()

best_k = k_range[accuracy_scores.index(max(accuracy_scores))]
print(f"Best K: {best_k} with Accuracy: {max(accuracy_scores):.2f}")

knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

print("==== Final Evaluation (Heart Dataset) ====")
print("\nAccuracy Score:", accuracy_score(y_test, y_pred))
print("\nClassification Report:")
print(classification_report(y_test, y_pred, target_names=['No Disease', 'Disease']))
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Reds', xticklabels=['No Disease', 'Disease'],
            yticklabels=['No Disease', 'Disease'])
plt.title("Confusion Matrix - Heart Disease")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

Program 7

Build Support vector machine model for a given dataset.

Screenshot

```
def SVM():
    import numpy as np
    import matplotlib.lib.pyplot as plt
    X = np.array([[1, 2, 3], [3, 5, 7], [2, 2], [6, 6], [7, 8]])
    y = np.array([1, 1, -1, -1, -1])
    learning_rate = 0.001
    lambda_param = 0.01
    epochs = 1000
    w = np.zeros(X.shape[1])
    b = 0
    for epoch in range(epochs):
        for i in enumerate(X):
            condition = y[i] * (np.dot(X[i], w) + b) >= 1
            if not condition:
                w -= learning_rate * (2 * lambda_param * w + y[i] * X[i])
                b -= learning_rate * y[i]
    def predict(x):
        return np.sign(np.dot(x, w) + b)
    def plot_svm(X, y, w, b):
        plt.figure(figsize=(8, 6))
```

```

for i in range (len(x)):
    if y[i] == 1:
        plt.scatter (x[i][0], x[i][1], c="red", marker='o')
        label = 'class +1' if i == 0 else ""
    else:
        plt.scatter (x[i][0], x[i][1], c="blue", marker='s')
        label = 'class -1', if i == 3 else ""

plot_SUM(x,y,w,b)
test = np.array ([[3,3],[7,7]])
print ("Prediction:", predict(test))

```

Code:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, confusion_matrix, roc_auc_score, roc_curve
from sklearn.preprocessing import label_binarize
import matplotlib.pyplot as plt
import seaborn as sns

iris = pd.read_csv("iris.csv")
label_encoder = LabelEncoder()
iris['species'] = label_encoder.fit_transform(iris['species'])
class_names_iris = label_encoder.classes_
X_iris = iris.drop('species', axis=1)
y_iris = iris['species']
X_train_iris, X_test_iris, y_train_iris, y_test_iris = train_test_split(X_iris, y_iris, test_size=0.2,
random_state=42)

scaler = StandardScaler()
X_train_iris = scaler.fit_transform(X_train_iris)
X_test_iris = scaler.transform(X_test_iris)
svm_linear = SVC(kernel='linear')
svm_linear.fit(X_train_iris, y_train_iris)
y_pred_linear = svm_linear.predict(X_test_iris)
acc_linear = accuracy_score(y_test_iris, y_pred_linear)
cm_linear = confusion_matrix(y_test_iris, y_pred_linear)

plt.figure(figsize=(6,4))
```

```

sns.heatmap(cm_linear, annot=True, fmt='d', cmap='Blues', xticklabels=class_names_iris,
yticklabels=class_names_iris)

plt.title(f'IRIS SVM Linear Kernel\nAccuracy: {acc_linear:.2f}')

plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.tight_layout()
plt.show()

```

```

svm_rbf = SVC(kernel='rbf')
svm_rbf.fit(X_train_iris, y_train_iris)
y_pred_rbf = svm_rbf.predict(X_test_iris)
acc_rbf = accuracy_score(y_test_iris, y_pred_rbf)
cm_rbf = confusion_matrix(y_test_iris, y_pred_rbf)

```

```

plt.figure(figsize=(6,4))

sns.heatmap(cm_rbf, annot=True, fmt='d', cmap='Greens', xticklabels=class_names_iris,
yticklabels=class_names_iris)

plt.title(f'IRIS SVM RBF Kernel\nAccuracy: {acc_rbf:.2f}')

plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.tight_layout()
plt.show()

```

```

letters = pd.read_csv("letter-recognition.csv")
X_letters = letters.drop('letter', axis=1)
y_letters = letters['letter']
label_encoder_letters = LabelEncoder()
y_letters_encoded = label_encoder_letters.fit_transform(y_letters)
class_names_letters = label_encoder_letters.classes_

X_train_letters, X_test_letters, y_train_letters, y_test_letters = train_test_split(
    X_letters, y_letters_encoded, test_size=0.2, random_state=42)

```

```

scaler_letters = StandardScaler()
X_train_letters = scaler_letters.fit_transform(X_train_letters)
X_test_letters = scaler_letters.transform(X_test_letters)
svm_letters = SVC(kernel='rbf', probability=True)
svm_letters.fit(X_train_letters, y_train_letters)

y_pred_letters = svm_letters.predict(X_test_letters)
acc_letters = accuracy_score(y_test_letters, y_pred_letters)
cm_letters = confusion_matrix(y_test_letters, y_pred_letters)

plt.figure(figsize=(14, 12))
sns.heatmap(cm_letters, annot=True, fmt='d', cmap='Purples',
            xticklabels=class_names_letters,
            yticklabels=class_names_letters,
            annot_kws={"size": 8},
            cbar=True)
plt.title(f'Letter Recognition - SVM RBF Kernel\nAccuracy: {acc_letters*100:.2f}%', fontsize=16)
plt.xlabel("Predicted Label", fontsize=14)
plt.ylabel("True Label", fontsize=14)
plt.xticks(rotation=45)
plt.yticks(rotation=0)
plt.tight_layout()
plt.show()

y_test_binarized = label_binarize(y_test_letters, classes=np.arange(len(class_names_letters)))
y_score = svm_letters.predict_proba(X_test_letters)
auc_score = roc_auc_score(y_test_binarized, y_score, average='macro')

fpr = dict()
tpr = dict()
for i in range(len(class_names_letters)):

```

```
fpr[i], tpr[i], _ = roc_curve(y_test_binarized[:, i], y_score[:, i])

plt.figure(figsize=(8, 6))
for i in range(0, len(class_names_letters), 4): # Plot every 4th class
    plt.plot(fpr[i], tpr[i], lw=1.5, label=f'Class {class_names_letters[i]}')

plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title(f"Multi-Class ROC Curve (Macro AUC = {auc_score:.6f})")
plt.legend(loc="lower right", fontsize='small')
plt.grid()
plt.tight_layout()
plt.show()

print(f"Exact AUC Score = {auc_score}")
```

Program 8

Implement Random forest ensemble method on a given dataset.

Screenshot:

```
Random forest
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from google.colab import files
uploaded = files.upload()

for filename in uploaded.keys():
    df = pd.read_csv(filename)
    print(f"Data loaded from {filename}.")
    display(df.head(1))

X = df.iloc[:, :-1]
y = df.iloc[:, -1]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

rf_model = RandomForestClassifier(n_estimators=100,
                                  random_state=42)
rf_model.fit(X_train, y_train)

y_pred = rf_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy of Random Forest Model : {accuracy * 100 :.2f}%")
```


Code:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import confusion_matrix

iris_df = pd.read_csv("iris.csv")
X = iris_df.drop('species', axis=1)
y = iris_df['species']
le = LabelEncoder()
y_encoded = le.fit_transform(y)

X_train, X_test, y_train, y_test = train_test_split(X, y_encoded, test_size=0.3, random_state=42)
rf_model = RandomForestClassifier(n_estimators=10, random_state=42)
rf_model.fit(X_train, y_train)
y_pred = rf_model.predict(X_test)
print("Random Forest Accuracy with 10 trees:", accuracy_score(y_test, y_pred))

scores = []
n_range = range(1, 101)
best_model = None
best_preds = None

for n in n_range:
    model = RandomForestClassifier(n_estimators=n, random_state=42)
    model.fit(X_train, y_train)
```

```
preds = model.predict(X_test)
```

```

acc = accuracy_score(y_test, preds)
scores.append(acc)
if acc == max(scores):
    best_model = model
    best_preds = preds

best_score = max(scores)
best_n = n_range[scores.index(best_score)]
print(f"Best Random Forest Accuracy: {best_score:.4f} with {best_n} trees")

plt.figure(figsize=(10, 5))
plt.plot(n_range, scores, marker='o', linestyle='-', color='blue')
plt.title('Random Forest Accuracy vs Number of Trees (Iris Dataset)')
plt.xlabel('Number of Trees')
plt.ylabel('Accuracy')
plt.grid(True)
plt.show()

cm = confusion_matrix(y_test, best_preds)
plt.figure(figsize=(6, 5))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=le.classes_, yticklabels=le.classes_)
plt.title(f"Confusion Matrix for Best Random Forest Model ({best_n} Trees)")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

Program 9

Implement Boosting ensemble method on a given dataset.

Screenshot:

#B-4.
Random Forest Boosting

```
import numpy as np
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import make_classification
from sklearn.metrics import accuracy_score

class AdaBoost:
    def __init__(self, n_estimators = 50):
        self.n_estimators = n_estimators
        self.alphas = []
        self.models = []

    def fit(self, X, y):
        n_samples, n_features = X.shape
        w = np.ones(n_samples)

        for estimator in range(self.n_estimators):
            #Train weak classifier (decision stump)
            model = DecisionTreeClassifier(max_depth=1)
            model.fit(X, y, sample_weight=w)
            y_pred = model.predict(X)

            #Calculate error rate
            err = np.sum(w * (y_pred != y)) / np.sum(w)
            alpha = 0.5 + np.log((1 - err) / err) if err < 1 else 0
            self.alphas.append(alpha)
            self.models.append(model)

            w = w * np.exp(-alpha * y * y_pred)
            w = w / np.sum(w)
```

```

def __init__(self):
    def predict(self, x):
        final_pred += alpha * model.predict(x)
        for model, alpha in zip(self.models,
                               self.alphas):
            final_pred += alpha * model.predict(x)
        return np.sign(final_pred)

    def score(self, X, y):
        return accuracy_score(y, self.predict(X))

```

$X, y = \text{make_classification}(\text{n_samples}=500, \text{n_features}:$
 $\text{n_redundant}=0, \text{n_classes}=2, \text{random_state}=4)$

Plot error over iterations:

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

```
plt.plot(range(1, adaboost.n_estimators+1),
```

adaboost, error = 0, k = 10

`plt.xlabel('Number of Estimators')`

~~pt. yodel (number of errors)~~

~~It. title (error vs Number of estimators)~~

`plt.grid (True)`

plt.show()

Code:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import AdaBoostClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import confusion_matrix

income_df = pd.read_csv("income.csv")
X_income = income_df.drop('income_level', axis=1)
y_income = income_df['income_level']
X_train_i, X_test_i, y_train_i, y_test_i = train_test_split(X_income, y_income, test_size=0.3,
random_state=42)

ada_model = AdaBoostClassifier(n_estimators=10, random_state=42)
ada_model.fit(X_train_i, y_train_i)
y_pred_i = ada_model.predict(X_test_i)
print("AdaBoost Accuracy with 10 estimators:", accuracy_score(y_test_i, y_pred_i))

scores_ada = []
n_range_ada = range(1, 51)
best_model_ada = None
best_preds_ada = None

for n in n_range_ada:
    model = AdaBoostClassifier(n_estimators=n, random_state=42)
    model.fit(X_train_i, y_train_i)
    preds = model.predict(X_test_i)
```

```

scores_ada.append(acc)

if acc == max(scores_ada):
    best_model_ada = model
    best_preds_ada = preds

best_score_ada = max(scores_ada)
best_n_ada = n_range_ada[scores_ada.index(best_score_ada)]
print(f"Best AdaBoost Accuracy: {best_score_ada:.4f} with {best_n_ada} estimators")

```

```

plt.figure(figsize=(10, 5))
plt.plot(n_range_ada, scores_ada, marker='o', linestyle='-', color='orange')
plt.title('AdaBoost Accuracy vs Number of Estimators (Income Dataset)')
plt.xlabel('Number of Estimators')
plt.ylabel('Accuracy')
plt.grid(True)
plt.show()

```

```

cm_ada = confusion_matrix(y_test_i, best_preds_ada)
plt.figure(figsize=(6, 5))
sns.heatmap(cm_ada, annot=True, fmt='d', cmap='Oranges', xticklabels=[0, 1], yticklabels=[0, 1])
plt.title(f"Confusion Matrix for Best AdaBoost Model ({best_n_ada} Estimators)")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

Program 10

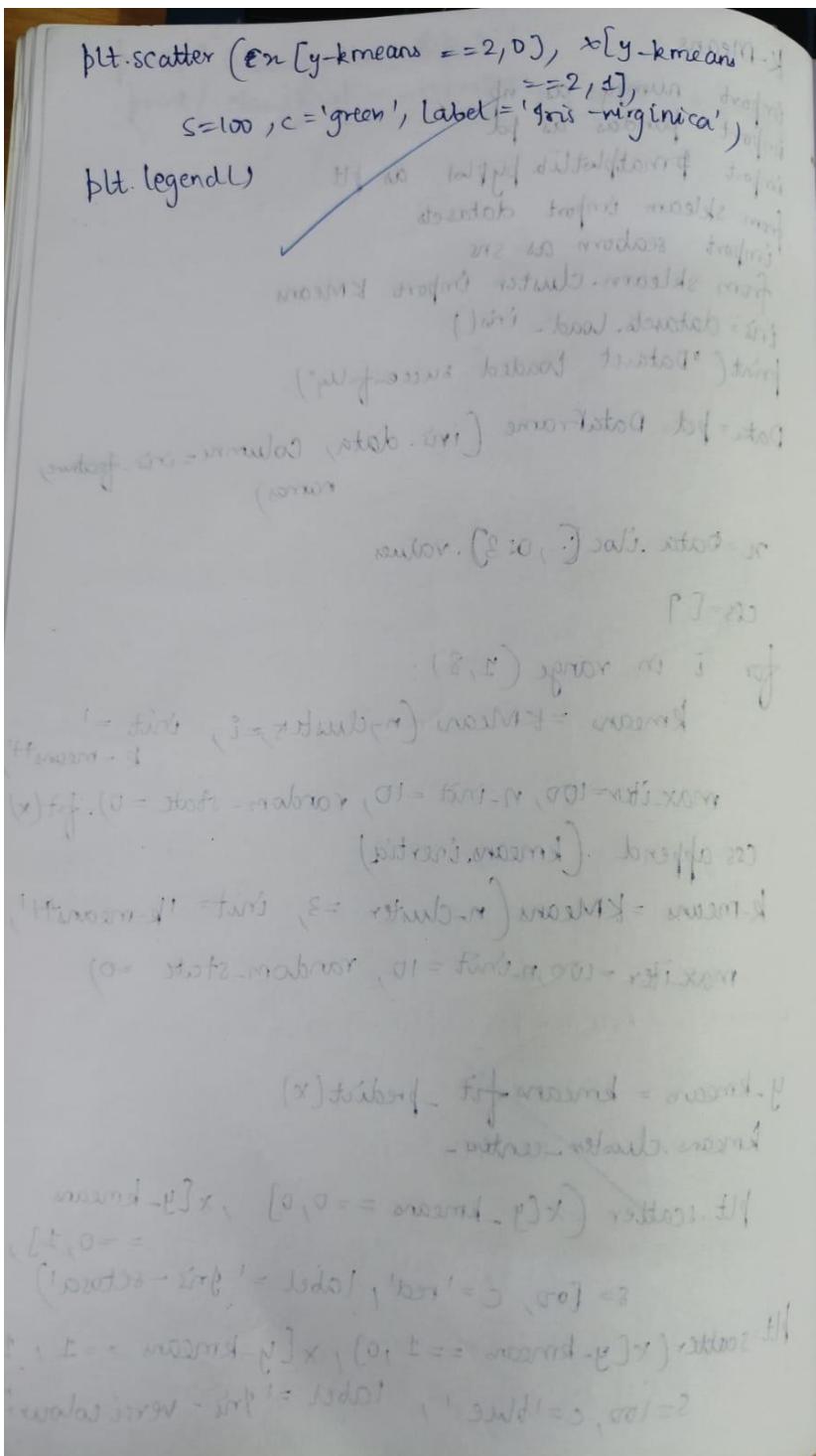
Build k-Means algorithm to cluster a set of data stored in a .CSV file.

Screenshot:

The image shows handwritten Python code for a K-Means clustering program. The code imports numpy, pandas, and matplotlib.pyplot, and loads the Iris dataset from sklearn. It then defines a function to print the loaded data and creates a DataFrame. The code initializes a list 'css' with 7 elements and iterates through values 1 to 8. For each value, it creates a KMeans object with n_clusters=i, init='k-means++', max_iter=100, n_init=10, random_state=0, and fits it to the data 'x'. The inertia value is appended to 'css'. A second KMeans object is created with n_clusters=3, init='k-means++', max_iter=100, n_init=10, and random_state=0. The predicted cluster labels are stored in 'y_kmeans' and the cluster centers are printed. Finally, two scatter plots are shown: one for clusters 0 and 1, and another for clusters 1 and 2. The x-axis for both plots is 'y_kmeans == 0, 1' and the y-axis is 'y_kmeans == 0, 1'. The first plot uses red points for cluster 0 and blue points for cluster 1. The second plot uses blue points for cluster 1 and green points for cluster 2.

```
K-Means
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn import datasets
import seaborn as sns
from sklearn.cluster import KMeans
iris = datasets.load_iris()
print("Dataset loaded successfully")
Data = pd.DataFrame(iris.data, columns=iris.feature_names)
n = Data.iloc[:, 0:2].values
css = [7]
for i in range(1, 8):
    kmeans = KMeans(n_clusters=i, init='k-means++',
                     max_iter=100, n_init=10, random_state=0).fit(x)
    css.append(kmeans.inertia_)
kmeans = KMeans(n_clusters=3, init='k-means++',
                max_iter=100, n_init=10, random_state=0)

y_kmeans = kmeans.fit_predict(x)
kmeans.cluster_centers_
plt.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1],
            s=100, c='red', label='Iris-setosa')
plt.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1],
            s=100, c='blue', label='Iris-versicolour')
plt.show()
```



Code:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from scipy import stats
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

df1=pd.read_csv("iris.csv")
df1.head()
df=df1.drop(['sepal_length','sepal_width','species'],axis=1)
scaler = StandardScaler()
scaled_df=scaler.fit_transform(df)
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    kmeans.fit(scaled_df)
    wcss.append(kmeans.inertia_)
plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()

kmeans = KMeans(n_clusters=3, init='k-means++', max_iter=300, n_init=10, random_state=0)
pred_y = kmeans.fit_predict(scaled_df)
```

```
plt.scatter(df['petal_length'], df['petal_width'], c=df['cluster'])  
plt.title('Clusters of Iris Flowers')  
plt.xlabel('Petal Length')  
plt.ylabel('Petal Width')  
plt.show()
```

Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method.

Screenshot:

The screenshot shows a Jupyter Notebook cell with handwritten Python code. The code is for performing Principal Component Analysis (PCA) on a dataset. It includes importing pandas, numpy, and various sklearn modules, reading files from Google Colab, standardizing data, applying PCA, visualizing the first two components, and calculating explained variance ratios. A large portion of the code is crossed out with a red marker.

```
Principal Component Analysis:  
import pandas as pd  
import numpy as np  
from sklearn.decomposition import PCA  
from sklearn.preprocessing import StandardScaler  
import matplotlib.pyplot as plt  
from google.colab import files  
uploaded = files.upload()  
  
for filename in uploaded.keys():  
    df = pd.read_csv(filename)  
    print(f"Uploaded {filename}")  
    display(df.head())  
  
numeric_df = df.select_dtypes(include=[np.number])  
  
## Standardize Data  
  
## Apply PCA  
pca = PCA(n_components=2)  
principal_components = pca.fit_transform(x_scaled)  
  
## Visualize first two Principle Components  
## Explained variance ration  
  
print(f"Model accuracy : {accuracy}%")  
  
S  
S  
S  
S  
S  
at 0.4%
```

Code:

```
#PCA
```

```
# STEP 1: Import packages
```

```
import pandas as pd  
import numpy as np  
from sklearn.decomposition import PCA  
from sklearn.preprocessing import StandardScaler  
import matplotlib.pyplot as plt  
from google.colab import files
```

```
# STEP 2: Upload the CSV file
```

```
uploaded = files.upload()
```

```
# STEP 3: Load the dataset
```

```
for filename in uploaded.keys():  
    df = pd.read_csv(filename)  
    print(f" ✅ Uploaded: {filename}")  
    display(df.head())
```

```
# STEP 4: Select numerical columns
```

```
numeric_df = df.select_dtypes(include=[np.number])  
  
print(" 📈 Numerical features found:",  
list(numeric_df.columns))
```

```
# OPTIONAL: Manually select columns if needed
```

```
# selected_features = ['feature1', 'feature2', ...]
```

```
selected_features = numeric_df.columns # use all numeric  
features for now
```

STEP 5: Standardize data

```
X = numeric_df[selected_features].dropna()  
X_scaled = StandardScaler().fit_transform(X)
```

STEP 6: Apply PCA

```
pca = PCA(n_components=2)  
principal_components = pca.fit_transform(X_scaled)
```

STEP 7: Create DataFrame for components

```
pca_df = pd.DataFrame(data=principal_components,  
columns=['PC1', 'PC2'])
```

STEP 8: Visualize the first two principal components

```
plt.figure(figsize=(8,6))  
plt.scatter(pca_df['PC1'], pca_df['PC2'], alpha=0.7)  
plt.xlabel('Principal Component 1')  
plt.ylabel('Principal Component 2')  
plt.title('2D PCA Visualization')  
plt.grid(True)  
plt.show()
```

STEP 9: Explained variance ratio

```
print("📈 Explained Variance Ratio:",  
pca.explained_variance_ratio_)
```