

**VISVESVARAYA TECHNOLOGICAL
UNIVERSITY**

“JnanaSangama”, Belgaum -590014, Karnataka.



**LAB REPORT
on**

Machine Learning (23CS6PCMAL)

Submitted by

Sanjeet Prajwal Pandit (1BM22CS241)

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 **Sanjeet Prajwal Pandit (1BM22CS241)**, who is 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: Ms. Saritha 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	1-4
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	5-9
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	10-16
4	17-3-2025	Build Logistic Regression Model for a given dataset	17-23
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample	24-27
6	7-4-2025	Build KNN Classification model for a given dataset	28-31
7	21-4-2025	Build Support vector machine model for a given dataset	32-34
8	5-5-2025	Implement Random forest ensemble method on a given dataset	35-37
9	5-5-2025	Implement Boosting ensemble method on a given dataset	38-42
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file	43-46
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method	47-50

Github Link:

<https://github.com/Sanjeet-108/6thSem-ML-Lab>

Program 1

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

Screenshot:

```
'lab 0,1:  
# Stock Market Data for Analysis  
# Import required libraries  
import yfinance as yf  
import pandas as pd  
import matplotlib.pyplot as plt  
  
tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]  
  
data = yf.download(tickers, start="2024-01-01", end="2024-12-31")  
group_by='ticker'  
  
print(data.head())  
  
HDFC_data = data['HDFCBANK.NS']  
print(HDFC_data.describe())  
  
HDFC_data['Daily Return'] = HDFC_data['close'].pct_change()  
  
plt.figure(figsize=(12,6))  
plt.subplot(2,1,1)  
HDFC_data['close'].plot(title="HDFC-Closing Price")  
plt.subplot(2,1,2)  
HDFC_data['Daily Return'].plot(title="HDFC-Daily Returns", color='orange')  
plt.tight_layout()  
plt.show()  
HDFC_data.to_csv('HDFC-stock-data.csv')
```

```
IICCI-data = data['IICCI BANK.NS']
```

```
IICCI-data['Daily Return'] = IICCI-data['Close'].pct_change
```

```
plt.figure(figsize=(12,6))
```

```
IICCI-data['Close'].plot(title="IICCI - Closing Price")
```

```
IICCI-data['Daily Return'].plot(title="IICCI - Daily Returns", color='blue')
```

```
plt.show()
```

```
KOTAK-data = data['KOTAKBANK.NS']
```

```
KOTAK-data['Daily Return'] = KOTAK-data['Close'].pct_change
```

```
plt.figure(figsize=(12,6))
```

```
KOTAK-data['Close'].plot(title="KOTAK - Closing Price")
```

```
KOTAK-data['Daily Return'].plot(title="KOTAK - Daily Returns", color='blue')
```

```
plt.show()
```

Importing and Exporting Data

Method 1: Initializing values directly into DataFrame

import pandas as pd

data = {

'USN': ['IBN22CS001', 'IBN22CS002', 'IBN22CS003', 'IBN22CS004', 'IBN22CS005'],

'Names': ['A', 'B', 'C', 'D', 'E'],

'Marks': [96, 97, 98, 98, 99]

}

df = pd.DataFrame(data)

print(df.head())

DATE: _____
PAGE: _____

Method 2: Importing datasets from sklearn.datasets

```
import pandas as pd
from sklearn.datasets import load_diabetes
diabetes = pd.DataFrame(diabetes.data, columns=diabetes.feature_names)
df['target'] = diabetes.target
print(df.head())
```

Method 3: Importing datasets from specific .csv file

```
import pandas as pd
file_path = 'sampledata.csv'
df = pd.read_csv(file_path)
print(df.head())
```

Method 4: Downloading datasets from existing dataset repositories like Kaggle, UCI, Mendeley, KEEL etc

1) Download diabetes datasets from mendeley

```
df = pd.read_csv('diabetes.csv')
print(df.head())
```

Code:

```
import pandas as pd
df=pd.DataFrame({
    "USN": [21, 22, 23, 24, 25],
    "Name": ["Sam", "Ram", "Dam", "rahul", "rakshith"],
    "Marks": [45, 78, 24, 90, 60]
})
Df
```

```
from sklearn.datasets import load_diabetes
dia=load_diabetes()
df=pd.DataFrame(dia.data,columns=dia.feature_names)
df
```

```
filepath="/content/sample_sales_data.csv"
df=pd.read_csv(filepath)
df.head()
```

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

```
# Define the tickers and date range
tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
```

```

start_date = "2024-01-01"
end_date = "2024-12-30"

# Download the data
data = yf.download(tickers, start=start_date, end=end_date)

# Check the structure of the downloaded data
print(data.head()) # Print the first few rows to see the column names

# Extract closing prices (if 'Adj Close' doesn't exist, fall back to 'Close')
closing_prices = data["Close"]

# Handle missing values by forward filling
closing_prices = closing_prices.fillna()

# Calculate daily returns
daily_returns = closing_prices.pct_change()

# Calculate 50-day moving average
moving_average = closing_prices.rolling(window=50).mean()

# Plotting
plt.figure(figsize=(12, 6))

# Closing Prices and 50-day Moving Average
plt.subplot(2, 1, 1)
for ticker in tickers:
    plt.plot(closing_prices.index, closing_prices[ticker], label=f"{ticker} Closing Price")
    plt.plot(moving_average.index, moving_average[ticker], label=f"{ticker} 50-Day MA", linestyle="--")
plt.title("Closing Prices of HDFC Bank, ICICI Bank, and Kotak Mahindra Bank with 50-Day Moving Average")
plt.xlabel("Date")
plt.ylabel("Closing Price")
plt.legend()

# Daily Returns
plt.subplot(2, 1, 2)
for ticker in tickers:
    plt.plot(daily_returns.index, daily_returns[ticker], label=ticker)
plt.title("Daily Returns of HDFC Bank, ICICI Bank, and Kotak Mahindra Bank")
plt.xlabel("Date")
plt.ylabel("Daily Return")
plt.legend()

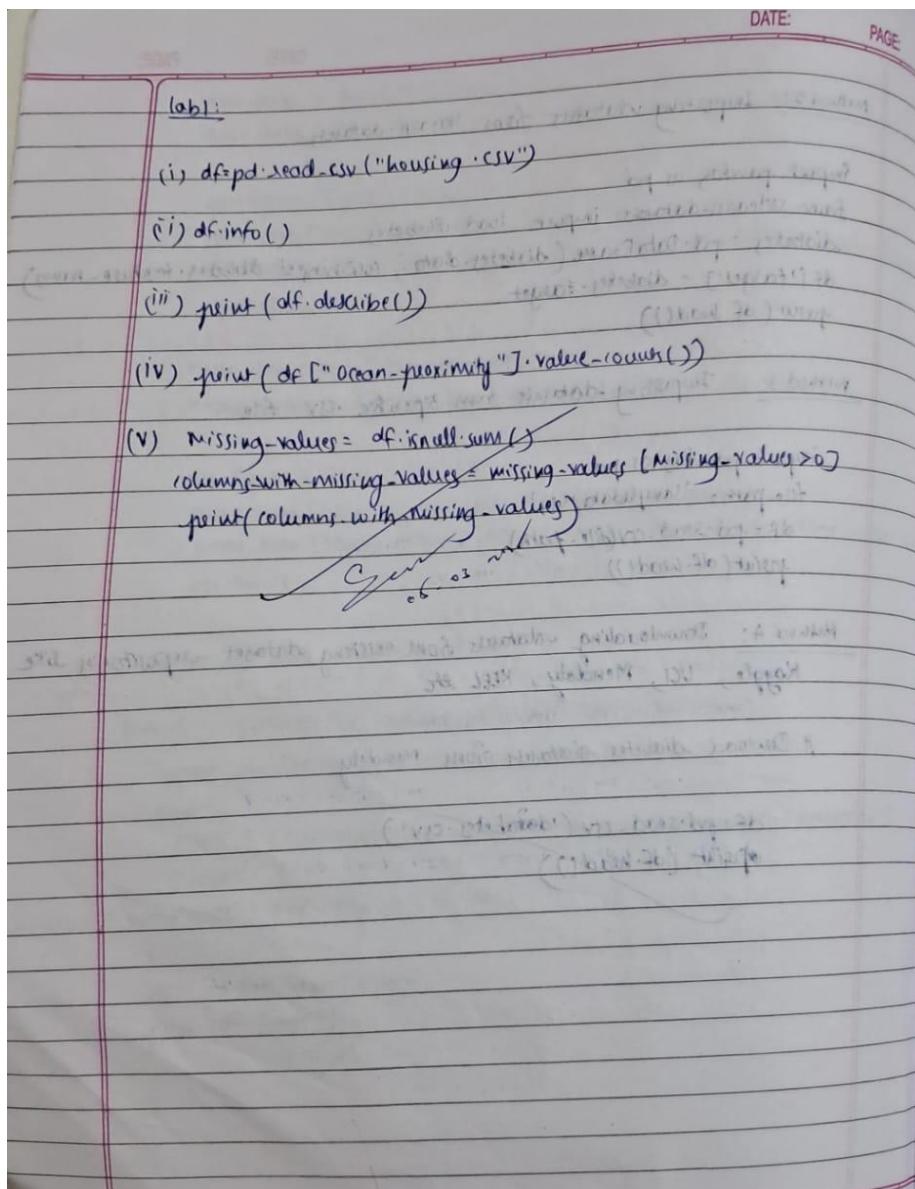
plt.tight_layout()
plt.show()

```

Program 2

Demonstrate various data pre-processing techniques for a given dataset

Screenshot:



```

label:
import pandas as pd
import numpy as np
//data = { "Day": [---], --- }
df = pd.read_csv("weather.csv") // Day, Outlook, Temp, Humidity, Wind, Decision
data = pd.read_csv("weatherdata.csv"), df = pd.DataFrame(data)
//function to calculate entropy
def entropy(target):
    class_counts = target.value_counts()
    probabilities = class_counts / len(target)
    return -np.sum(probabilities * np.log2(probabilities))

//function to calculate info gain
def information_gain(data, feature, target):
    entropy_before = entropy(target)
    feature_values = data[feature].unique()
    weighted_entropy = 0
    for value in feature_values:
        subset = target[data[feature] == value]
        weighted_entropy += (len(subset) / len(target)) * entropy(subset)
    return entropy_before - weighted_entropy

def print_entropy_and_gain(data, features, target):
    print("\nEntropy and Information gain for each feature:")
    for feature in features:
        gain = information_gain(data, feature, target)
        ent = entropy(target)
        print(f"Feature: {feature} Entropy: {ent:.4f} | Information Gain: {gain:.4f}")

```

Code:

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.preprocessing import MinMaxScaler, StandardScaler, OneHotEncoder
```

```
from sklearn.impute import SimpleImputer
```

```
from sklearn.compose import ColumnTransformer
```

```
from sklearn.pipeline import Pipeline
```

```
# Load Diabetes dataset (update with your file path in Google Colab)
```

```
diabetes_df = pd.read_csv("/content/Dataset of Diabetes .csv")
```

```
# Load Adult Income dataset (update with your file path in Google Colab)
```

```
adult_income_df = pd.read_csv("/content/adult.csv")
```

```
# Function for data cleaning and transformations
```

```
def preprocess_data(df):
```

```
# 1. Data Cleaning: Handle Missing Values
```

```
# Identify columns with missing values
```

```
missing_values = df.isnull().sum()
```

```
print("Columns with missing values:\n", missing_values[missing_values > 0])
```

```
# Handle missing values: Use median for numerical columns and most frequent for categorical columns
```

```
numerical_cols = df.select_dtypes(include=["float64", "int64"]).columns
```

```
categorical_cols = df.select_dtypes(include=["object"]).columns
```

```
# Impute numerical columns with median, and categorical columns with the most frequent value
```

```
imputer = SimpleImputer(strategy='most_frequent')
```

```
df[categorical_cols] = imputer.fit_transform(df[categorical_cols])
```

```
imputer = SimpleImputer(strategy='median')
```

```
df[numerical_cols] = imputer.fit_transform(df[numerical_cols])
```

```
# 2. Handling Categorical Data: One-Hot Encoding
```

```
print("\nCategorical Columns:\n", categorical_cols)
```

```
# Apply OneHotEncoder to categorical columns
```

```

encoder = OneHotEncoder(drop='first', sparse_output=False)

encoded_categorical_data = encoder.fit_transform(df[categorical_cols])

# Creating a DataFrame for encoded categorical data

encoded_categorical_df = pd.DataFrame(encoded_categorical_data,
columns=encoder.get_feature_names_out(categorical_cols))

# Drop original categorical columns and concatenate the encoded columns

df = df.drop(categorical_cols, axis=1)

df = pd.concat([df, encoded_categorical_df], axis=1)

# 3. Handling Outliers: Using Z-score (Optional - can be added if needed)

from scipy import stats

z_scores = np.abs(stats.zscore(df[numerical_cols]))

df = df[(z_scores < 3).all(axis=1)] # Remove rows with outliers

# 4. Data Transformations: Min-Max Scaling

min_max_scaler = MinMaxScaler()

df[numerical_cols] = min_max_scaler.fit_transform(df[numerical_cols])

# 5. Data Transformations: Standardization

standard_scaler = StandardScaler()

df[numerical_cols] = standard_scaler.fit_transform(df[numerical_cols])

return df

```

```
# Apply preprocessing for Diabetes dataset

diabetes_df = preprocess_data(diabetes_df)

print("\nPreprocessed Diabetes Dataset:")

print(diabetes_df.head())

# Apply preprocessing for Adult Income dataset

adult_income_df = preprocess_data(adult_income_df)

print("\nPreprocessed Adult Income Dataset:")

print(adult_income_df.head())
```

Program 3

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

Screenshot:

The image shows handwritten Python code on lined paper. The code is a recursive function to build a decision tree. It starts with importing pandas and numpy, then defining a function `build-tree` that takes `data`, `target`, and `features` as arguments. The function checks if the target has unique values or if there are no features left. If either condition is met, it returns the target's mode. Otherwise, it calculates information gain for each feature, finds the best feature, and splits the data and target based on that feature's value. It then recursively builds subtrees for each subset of data and stores them in a dictionary under the best feature's key.

```
import pandas as pd
import numpy as np

//function to build decision tree recursively
def build-tree (data, target, features):
    if len(target.unique()) == 1:
        return target.mode()[0]
    if len(features) == 0:
        return target.mode()[0]
    gains = {feature: information-gain(data, feature, target) for feature in features}
    best-feature = max(gains, key=gains.get)

    tree = {best-feature: {}}
    feature_value = data[best-feature].unique()

    for value in feature_value:
        subset_data = data[data[best-feature] == value]
        subset_target = target[data[best-feature] == value]

        remaining_features = [f for f in features if f != best-feature]
        subtree = build-tree(subset_data, subset_target, remaining_features)
        tree[best-feature][value] = subtree

    return tree
```

```

def print-tree(tree, indent=""):
    if not isinstance(tree, dict):
        for feature, branches in tree.items():
            print(f"{' ' * indent}{feature}:")
            for value, subtree in branches.items():
                print(f"{' ' * (indent + 2)}{value} →")
                print-tree(subtree, indent + 2)
    else:
        print(f"{' ' * indent}{tree}")

```

target = df['Decision']

features = ['outlook', 'Temperature', 'Humidity', 'Wind']

print_entropy_and_gain(df, features, target)

tree = build-tree(df, target, features)

print("\nDecision Tree: ")

print-tree(tree, indent=" ")

Output:

Entropy and Information gain for each feature:

Feature: outlook | Entropy: 0.9403 | Information gain: 0.2967

Feature: Temperature | Entropy: 0.9403 | Information gain: 0.0292

Feature: Humidity | Entropy: 0.9403 | Information gain: 0.1578

Feature: Wind | Entropy: 0.9403 | Information gain: 0.0481

Decision Tree :

~~outlook:~~

Sunny →	Humidity:
High →	No
Normal →	Yes
Overcast →	Yes
Rainy →	Wind:
Wet →	Yes

Q1) Consider binary classification problem where we want to predict whether a student will pass or fail based on their Study hours. Logistic regression model has been trained, the learned parameters are $a_0 = -5$ (intercept) & $a_1 = 0.8$ (coeff for study hrs)

- Write logistic regression equation for this problem
- Calculate probability that student who studies for 7 hours will pass.
- Determine predicted class (pass or fail) for this student based on threshold of 0.5.

(2) Consider $\mathbf{z} = [2, 1, 0]$ for three classes. Apply softmax function to find probability values of three classes.

④ Linear Regression

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

```
dataset = pd.read_csv('contew/Sales.csv')
```

```
print(dataset.head())
```

Sales.csv	
Week	Sales(ink)
1	1.2
2	1.8
3	2.6
4	3.2
5	3.8

```
weeks = dataset['x1(week)'].values
```

```
sales = dataset['y1(sale in k)'].values
```

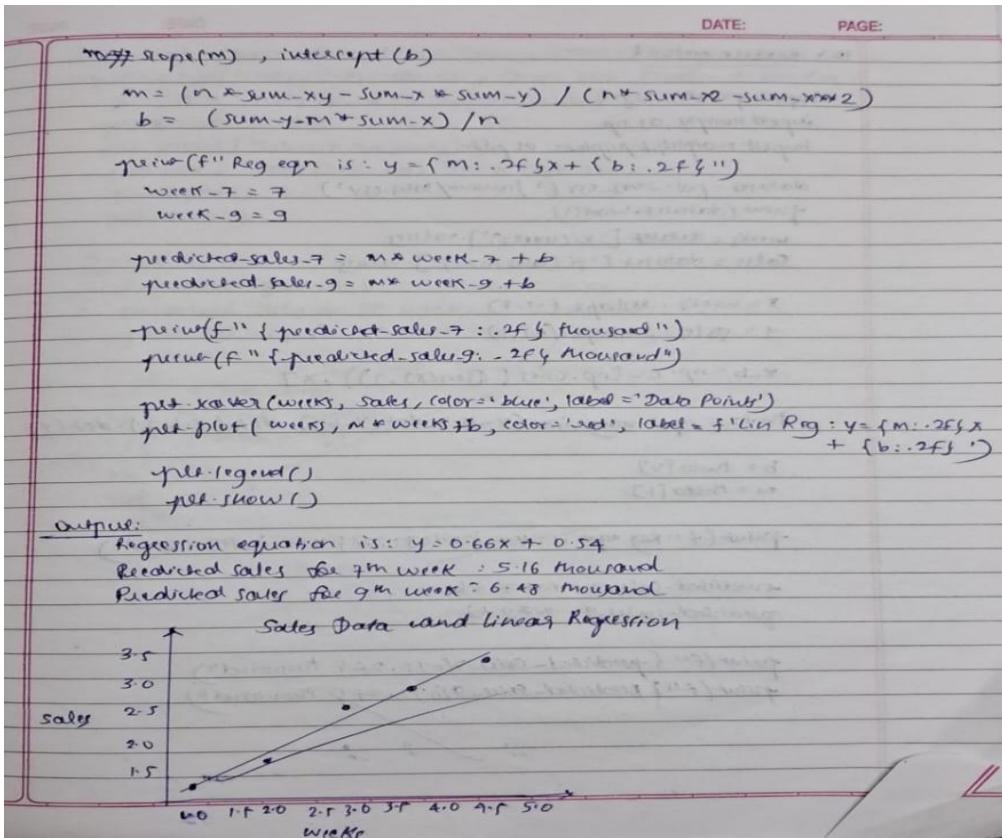
```
n = len(weeks)
```

$$\text{sum-}x = \text{np.sum}(weeks)$$

$$\text{sum-}y = \text{np.sum}(sales)$$

$$\text{sum-}x^2 = \text{np.sum}(weeks**2)$$

$$\text{sum-}xy = \text{np.sum}(weeks * sales)$$



Code:

```
import numpy as np

import matplotlib.pyplot as plt

# Given data

x = np.array([1, 2, 3, 4]) # weeks

y = np.array([1, 3, 4, 8]) # corresponding y values

X = np.vstack([np.ones(len(x)), x]).T
```

```

theta = np.linalg.inv(X.T @ X) @ X.T @ y

print(f"Intercept (theta_0): {theta[0]}")
print(f"Slope (theta_1): {theta[1]}")

week_to_predict = 5

y_pred = theta[0] + theta[1] * week_to_predict

print(f"Predicted value for week {week_to_predict}: {y_pred}")

plt.scatter(x, y, color='blue', label='Data points') # scatter plot of the data
plt.plot(x, theta[0] + theta[1] * x, color='red', label='Regression line') # plot the regression line

plt.scatter(week_to_predict, y_pred, color='green', label=f'Prediction for week {week_to_predict}')

plt.xlabel('Weeks')
plt.ylabel('Y values')
plt.title('Linear Regression: Week vs Y values')

# Show legend
plt.legend()

```

```
# Show the plot
```

```
plt.show()
```

2. Approach:

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

Given data

```
x = np.array([1, 2, 3, 4]) # weeks
```

```
y = np.array([1, 3, 4, 8]) # corresponding y values
```

Calculate means

```
x_mean = np.mean(x)
```

```
y_mean = np.mean(y)
```

Calculate slope (a_1) using the given formula

```
a1 = (np.mean(x * y) - x_mean * y_mean) / (np.mean(x**2) - x_mean**2)
```

Calculate intercept (a_0) using the given formula

```
a0 = y_mean - a1 * x_mean
```

Print the result

```
print(f"Intercept (a0): {a0}")
```

```

print(f"Slope (a1): {a1}")

# Predict for week 5

week_to_predict = 5

y_pred = a0 + a1 * week_to_predict

print(f"Predicted value for week {week_to_predict}: {y_pred}")


# Plotting the data points and the regression line

plt.scatter(x, y, color='blue', label='Data points') # scatter plot of the data

plt.plot(x, a0 + a1 * x, color='red', label='Regression line') # plot the regression line


# Mark the prediction point for week 5

plt.scatter(week_to_predict, y_pred, color='green', label=f'Prediction for week {week_to_predict}')


# Adding labels and title

plt.xlabel('Weeks')

plt.ylabel('Y values')

plt.title('Linear Regression: Week vs Y values')


# Show legend

plt.legend()


# Show the plot

plt.show()

```

Program 4:

Build Logistic Regression Model for a given dataset

Screenshot:

DATE: _____ PAGE: _____

2/4/25

Q Consider binary classification problem where we want to predict whether a student will pass or fail based on their study hours. Logistic regression model has been trained and learned parameters are $\alpha_0 = -5$ (intercept) & $\alpha_1 = 0.8$ (coefficient of study hours).

(a) Write logistic regression eqn for this problem

$$p(y=1|x) = \frac{1}{1 + e^{-(\alpha_0 + \alpha_1 x)}}$$

(b) calculate probability that student who studies for 7 hrs will pass

$$p(\text{pass}) = \frac{1}{1 + e^{-(-5 + 0.8 \cdot 7)}} = 0.6479 \quad (\because x = 7)$$

$$\therefore \alpha_0 = -5 + 0.8(7) = 0.6$$

(c) Determine predicted class for this student based on threshold of 0.5. If $p(\text{pass}) > 0.5$, Student will pass else he will fail.

(d) Consider $x = (2, 1, 0)$ for their classes, apply softmax fun to find prob of values three classes

$$\text{softmax}(\alpha_0) = \frac{e^{\alpha_0}}{\sum_{i=1}^k e^{\alpha_i}}$$

$$p(1) = \frac{e^2}{e^2 + e^1 + e^0} = 0.665$$

$$p(2) = \frac{e^1}{e^2 + e^1 + e^0} = 0.295$$

$$p(3) = \frac{e^0}{e^2 + e^1 + e^0} = 0.09$$

Q After building logistic regression models, answer following questions:

- (i) For dataset file "HR_comma_sep.csv"
 - (i) Which variables did you identify as having an direct and clear impact on employee retention? Why?
 - (ii) What was accuracy of your logistic regression model? Do you think this is good accuracy? Why or why not?

(2) For zoo dataset

- (i) Did you perform any data preprocessing steps? If yes, what were they, and why those necessary?
- (ii) Were there any missing or inconsistent values in dataset? How did you handle them?
- (iii) What does confusion matrix tell you about performance of your model
- (iv) Which class types were most frequently misclassified? Why do you think this happened?

Answer:

(1) (i) Satisfaction level

→ Time spent in company

→ Number of projects

→ Salary

These variables were chosen based on trends in data visualization

(ii) The accuracy of logistic regression was 78%. This accuracy is fairly good, it suggests that model captures most of properties affecting employee retention.

(a)	
(i)	Yes → Dropped 'animal-name' column → checked for missing values → converted categorical variable if needed
(ii)	No missing values found in dataset. In case of inconsistency, we could have used mean/mode imputation or remove them.
(iii)	It shows how well the model predicted different class types. A high no. success prediction, along diagonal of matrix indicates good performance.
(iv)	→ Amphibians, Reptiles Reason → Since they have similar features Logistic regression assumes linear decision boundaries which may not work well for complex class separation.

Code:

```
#LogisticRegression_Multiclass.ipynb
```

```
# Import necessary libraries

import pandas as pd

from sklearn.datasets import load_iris

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import accuracy_score

from sklearn import metrics

import matplotlib.pyplot as plt
```

```

# Load the Iris dataset

iris = pd.read_csv("/content/iris (2).csv")

iris.head()

X=iris.drop('species',axis='columns')# Features (sepal length, sepal width, petal length, petal width)

y = iris.species # Target labels (0: Setosa, 1: Versicolor, 2: Virginica)

# Split the dataset 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)

# Initialize the Multinomial Logistic Regression model

# Use 'multinomial' for multi-class classification and 'lbfgs' solver

model = LogisticRegression(multi_class='multinomial')

# Train the model on the training data

model.fit(X_train, y_train)

# Make predictions on the test data

y_pred = model.predict(X_test)

# Calculate the accuracy of the model on the test data

accuracy = accuracy_score(y_test, y_pred)

# Display the accuracy

```

```
print(f"Accuracy of the Multinomial Logistic Regression model on the test set: {accuracy:.2f}")

confusion_matrix = metrics.confusion_matrix(y_test, y_pred)

cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix = confusion_matrix, display_labels = ["Setosa", "Versicolor", "Virginica"])

cm_display.plot()

plt.show()
```

Binary Classification:

#*LogisticRegression_Binary.ipynb*

```
# Commented out IPython magic to ensure Python compatibility.

import pandas as pd

from matplotlib import pyplot as plt

# %matplotlib inline

# "%matplotlib inline" will make your plot outputs appear and be stored within the notebook.
```

```
df = pd.read_csv("/content/insurance_data (1).csv")

df.head()
```

```
plt.scatter(df.age, df.bought_insurance, marker='+', color='red')
```

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test =
```

```
train_test_split(df[['age']],df.bought_insurance,train_size=0.9,random_state=10)
```

```
X_train.shape
```

```
X_test
```

```
from sklearn.linear_model import LogisticRegression
```

```
model = LogisticRegression()
```

```
model.fit(X_train, y_train)
```

```
X_test
```

```
y_test
```

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

```
y_predicted
```

```
model.score(X_test,y_test)
```

```
model.predict_proba(X_test)
```

```
y_predicted = model.predict([[60]])
```

```
y_predicted
```

*#model.coef_ indicates value of m in $y=m*x + b$ equation*

```
model.coef_
```

*#model.intercept_ indicates value of b in $y=m*x + b$ equation*

```
model.intercept_
```

#Lets defined sigmoid function now and do the math with hand

```
import math
```

```
def sigmoid(x):
```

```
    return 1 / (1 + math.exp(-x))
```

```
def prediction_function(age):
```

```
    z = 0.127 * age - 4.973 # 0.12740563 ~ 0.0127 and -4.97335111 ~ -4.97
```

```
    y = sigmoid(z)
```

```
    return y
```

```
age = 35
```

```
prediction_function(age)
```

"""0.37 is less than 0.5 which means person with 35 will not buy the insurance"""

Program 5

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

Screenshot:

```
task:  
import pandas as pd  
import numpy as np  
// data = {  
    'Day': [...] }  
  
// df = pd.read_csv("weather.csv") // Day, Outlook, Temp, Humidity, Wind, Decision  
data = pd.read_csv("weatherdata.csv"), df = pd.DataFrame(data)  
// Function to calculate entropy  
def entropy(target):  
    class_counts = target.value_counts()  
    probabilities = class_counts / len(target)  
    return -np.sum(probabilities * np.log2(probabilities))  
  
// Function to calculate info gain  
def information_gain(data, feature, target):  
    entropy_before = entropy(target)  
    feature_values = data[feature].unique()  
  
    weighted_entropy = 0  
    for value in feature_values:  
        subset = target[data[feature] == value]  
        weighted_entropy += (len(subset) / len(target)) * entropy(subset)  
  
    return entropy_before - weighted_entropy  
  
def print_entropy_and_gain(data, features, target):  
    print("\nEntropy and Information Gain for each feature: ")  
    for feature in features:  
        gain = information_gain(data, feature, target)  
        ent = entropy(target)  
        print(f"Feature: {feature} Entropy: {ent:.4f}, Information Gain: {gain:.4f}")
```

DATE: _____
PAGE: _____

```

//function to build decision tree recursively
def build_tree(data, target, features):
    if len(target.unique()) == 1:
        return target.mode()[0]
    if len(features) == 0:
        return target.mode()[0]

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

    tree = {best_feature: {}}
    feature_value = data[best_feature].unique()

    for value in feature_value:
        subset_data = data[data[best_feature] == value]
        subset_target = target[data[best_feature] == value]

        remaining_features = [f for f in features if f != best_feature]
        subtree = build_tree(subset_data, subset_target, remaining_features)
        tree[best_feature][value] = subtree

    return tree

```

Code:

```

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import LabelEncoder

from sklearn.tree import DecisionTreeClassifier

```

```
from sklearn import tree
import matplotlib.pyplot as plt

# Load the newly uploaded dataset
data = pd.read_csv('/content/seattle-weather.csv')

# Data Preparation
data_cleaned = data.drop('date', axis=1)
label_encoder = LabelEncoder()
data_cleaned['weather'] = label_encoder.fit_transform(data_cleaned['weather'])

# Split data into features and target
X = data_cleaned.drop('weather', axis=1)
y = data_cleaned['weather']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the Decision Tree Classifier with improved clarity
id3_classifier = DecisionTreeClassifier(criterion='entropy', max_depth=4, random_state=42)
id3_classifier.fit(X_train, y_train)

# Visualize the Decision Tree with larger font size for clarity
plt.figure(figsize=(20, 12))
```

```
tree.plot_tree(id3_classifier, feature_names=X.columns, class_names=list(label_encoder.classes_),  
filled=True, fontsize=10)  
  
plt.title("ID3 Decision Tree Classifier (Enhanced Clarity)")  
  
plt.show()
```

Program 6

Build KNN Classification model for a given dataset

Screenshot:

KNN (K-Nearest-Neighbours)					DATE:	PAGE:
Person	Age	Salary	Dist	Rank	target	
A	18	50	52.8			
B	23	55	46.6			
C	24	70	31.9	2	N	
D	41	60	40.4	3	Y	
E	43	70	31.1	1	Y	
F	38	40	60.1			

Step 1: $DIST(d) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$(x_2, y_2) = (35, 100)$

$d_1 = \sqrt{(35-18)^2 + (100-50)^2} = 52.8$

$d_2 = \sqrt{(35-23)^2 + (100-70)^2} = 46.6$

KNN (K-Nearest-Neighbors)

Consider following dataset, for $K=3$ and test data $(x, 35, 100)$ at $(\text{person}, \text{Age}, \text{Salary})$ and predict the target

Person	Age	Salary	Dist	Rank	target
A	18	50	52.8		
B	23	55	46.6		
C	24	70	31.9	2	N
D	41	60	40.4	3	y
E	43	70	31.1	1	y
F	38	40	60.1		

$$\text{Step 1: } D(x_2, x_1) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$(x_2, x_1) = (35, 100)$$

$$d_1 = \sqrt{(35-18)^2 + (100-50)^2} = 52.8$$

$$d_2 = \sqrt{(35-23)^2 + (100-55)^2} = 46.6$$

Code:

```
import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

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

# Load the Iris dataset
iris = pd.read_csv("/content/iris (2).csv")

# Prepare the data
X = iris.drop('species', axis='columns')
y = iris.species

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize the KNN classifier with k=3
knn = KNeighborsClassifier(n_neighbors=3) # You can experiment with different k values

# Train the classifier
knn.fit(X_train, y_train)

# Make predictions
y_pred = knn.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
```

```
print(f"Accuracy: {accuracy}")

print("\nClassification Report:\n", classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)
print("\nConfusion Matrix:\n", cm)

# Visualize the confusion matrix (optional)

plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=["Setosa", "Versicolor", "Virginica"],
            yticklabels=["Setosa", "Versicolor", "Virginica"])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()
```

Program 7

Demonstrate various data pre-processing techniques for a given dataset

Screenshot:

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

dataset = pd.read_csv('content/sales.csv')
print(dataset.head())
weeks = dataset['x1(weeks)'].values
Sales = dataset['y1(sale in m)'].values

X = weeks.reshape(-1,1)
y = Sales.reshape(-1,1)

X_b = np.c_[np.ones((len(X), 1)), X]

theta = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)

b = theta[0]
m = theta[1]

print(f"Reg equation: y = {m}x + {b}")

predicted_sales7 = m * 7 + b
predicted_sales9 = m * 9 + b

print(f"Predicted sales at 7 weeks: {predicted_sales7} thousand")
print(f"Predicted sales at 9 weeks: {predicted_sales9} thousand")
```

Code:

```
import pandas as pd

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, classification_report, confusion_matrix

# Load the dataset from CSV

df = pd.read_csv("/content/iris (1).csv")

# Features and target

X = df[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']]

y = df['species']

# Encode species labels (setosa → 0, versicolor → 1, virginica → 2)

label_encoder = LabelEncoder()

y_encoded = label_encoder.fit_transform(y)

# Split into train and test sets

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

# Feature scaling

scaler = StandardScaler()

X_train_scaled = scaler.fit_transform(X_train)
```

```
X_test_scaled = scaler.transform(X_test)

# Train the SVM model

svm_model = SVC(kernel='rbf', C=1.0, gamma='scale')
svm_model.fit(X_train_scaled, y_train)

# Predictions

y_pred = svm_model.predict(X_test_scaled)

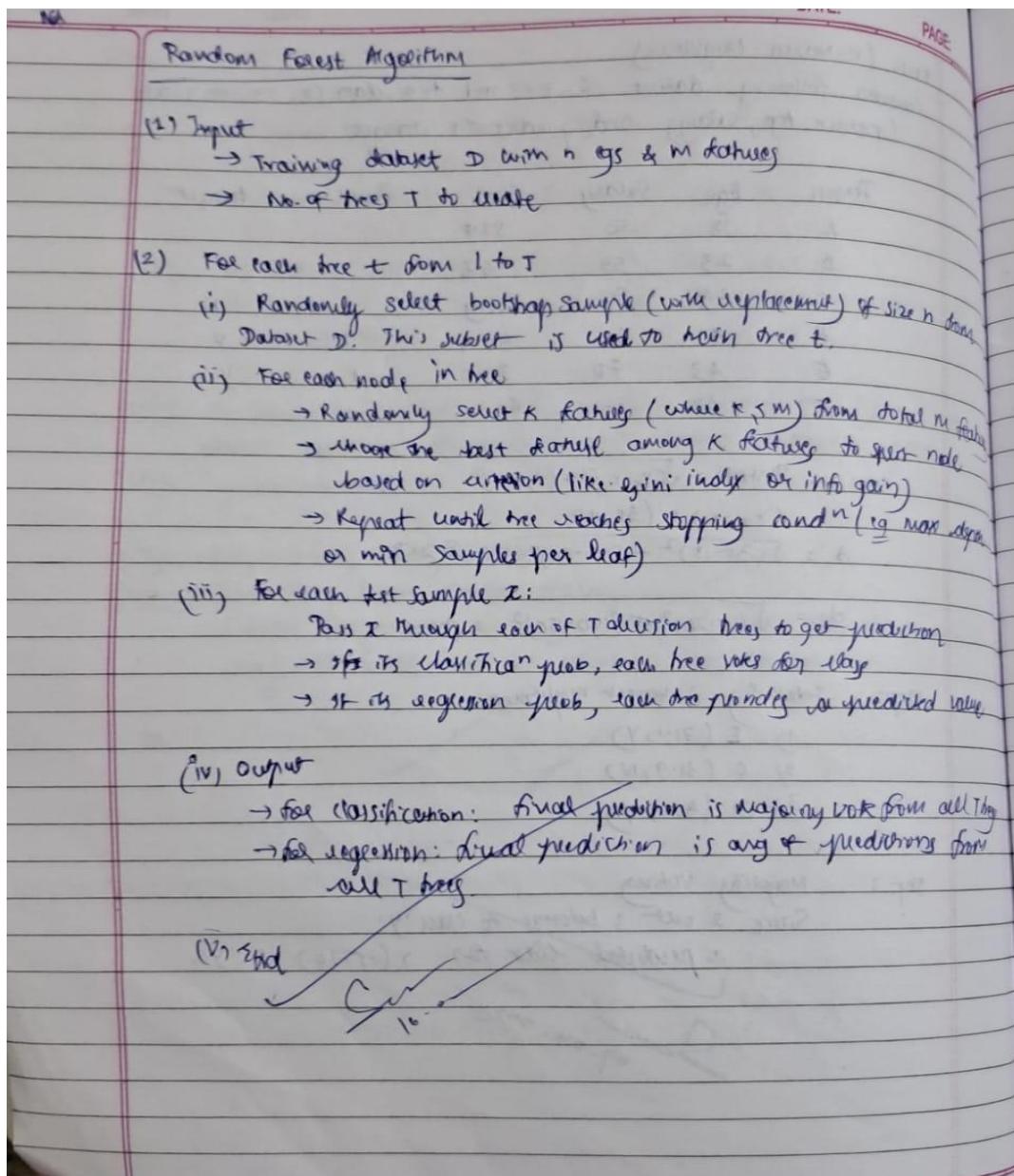
# Evaluation

print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred,
target_names=label_encoder.classes_))
```

Program 8

Implement Random forest ensemble method on a given dataset

Screenshot:



SNO.	CGPA	Interactivity	Communication	Practical Knowledge	Job
1	≥ 9	Yes	Good	Good	Yes
2	< 9	No	Moderate	Very Good	Yes
3	≥ 9	No	Moderate	Very	No
4	≥ 9	No	Moderate	Very	No
5	≥ 9	Yes	Moderate	Good	Yes

(1)

```

graph TD
    CGPA[CGPA] --<9-->|Yes| Interactivity1{Interactivity}
    CGPA -->|≥9| Interactivity1
    Interactivity1 --Yes-->|Yes| Leaf1[Yes]
    Interactivity1 --No-->|No| Leaf2[No]
  
```

(2)

```

graph TD
    Interactivity2{Interactivity} --Yes-->|Yes| Leaf3[Yes]
    Interactivity2 --No-->|No| PracticalKnowledge1{Practical Knowledge}
    PracticalKnowledge1 --Good-->|Yes| Leaf4[Yes]
    PracticalKnowledge1 --Very-->|No| Leaf5[No]
  
```

Code:

```

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score
  
```

```
df = pd.read_csv('/content/iris (3).csv')
```

```
# Prepare the data
```

```
X = df.drop('species', axis=1)
```

```
y = df['species']
```

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

```
# Build and evaluate the Random Forest classifier with default n_estimators
```

```
rf_classifier_default = RandomForestClassifier(random_state=42)
```

```
rf_classifier_default.fit(X_train, y_train)
```

```
y_pred_default = rf_classifier_default.predict(X_test)
```

```
default_accuracy = accuracy_score(y_test, y_pred_default)
```

```
print(f"Accuracy with default n_estimators (10): {default_accuracy}")
```

```
# Fine-tune the model by varying the number of trees
```

```
best_accuracy = 0
```

```
best_n_estimators = 0
```

```
for n_estimators in range(1, 101): # Try different number of trees
```

```
    rf_classifier = RandomForestClassifier(n_estimators=n_estimators, random_state=42)
```

```
    rf_classifier.fit(X_train, y_train)
```

```
    y_pred = rf_classifier.predict(X_test)
```

```
    accuracy = accuracy_score(y_test, y_pred)
```

```
    if accuracy > best_accuracy:
```

```
        best_accuracy = accuracy
```

```
        best_n_estimators = n_estimators
```

```
print(f"\nBest Accuracy: {best_accuracy} achieved with n_estimators = {best_n_estimators}").
```

Program 9

Implement Boosting ensemble method on a given dataset

Screenshot:

Adaboost						
CGPA	Interactions	P. knowledge	Comm	Job /in chile	weight	
≥9	Yes	good	Good	Yes	1/6	
<9	No	bad	Moderate	Yes	1/6	
≥9	No	bad	Moderate	No	1/6	
<9	No	bad	Good	No	1/6	
≥9	yes	bad	Moderate	Yes	1/6	
≥9	yes	bad	Moderate	Yes	1/6	

Initial wt = $1/6$

If CGPA $\geq g \rightarrow$ Job profile = Yes

$$\varepsilon = \sum w_i = \frac{2}{6} = 0.333$$

CGPA $< g \rightarrow$ Job profile = No

CGPA Actual Product

$$\alpha = \frac{1}{2} \ln\left(\frac{1-\varepsilon}{\varepsilon}\right) = \frac{1}{2} 0.347$$

$\geq g$ Yes Yes

$< g$ Yes No

$$2cgpa = \frac{1}{6} \times 4 \times e^{-0.347}$$

$\geq g$ No Yes

$$+ \frac{1}{6} \times 2 \times e^{0.347}$$

$< g$ No No

$$= 0.9438$$

$\geq g$ Yes Yes

$\geq g$ Yes Yes

$$w_{obj,14} = \frac{1}{6} \times e^{-0.347} = 0.125$$

$$0.9438$$

$$wt(d_i)_{14} = 0.2508$$

But accuracy score $\times 83.3\%$

Code:

```
import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import LabelEncoder

from sklearn.ensemble import AdaBoostClassifier

from sklearn.metrics import accuracy_score

import matplotlib.pyplot as plt

# Load the dataset

df = pd.read_csv("/content/income.csv")

# Encode the target column (income_level)

le = LabelEncoder()

df['income_level'] = le.fit_transform(df['income_level']) # e.g., <=50K → 0, >50K → 1

# Split features and target

X = df.drop("income_level", axis=1)

y = df["income_level"]

# Train-test split

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

# 1 Train AdaBoostClassifier with n_estimators = 10

model = AdaBoostClassifier(n_estimators=10, random_state=42)
```

```

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

print(f"Accuracy with 10 estimators: {accuracy:.4f}")

# 2 Fine-tune the model by changing n_estimators

scores = []

n_values = range(10, 101, 10)

for n in n_values:

    model = AdaBoostClassifier(n_estimators=n, random_state=42)

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)

    acc = accuracy_score(y_test, y_pred)

    scores.append(acc)

    print(f"n_estimators={n} => Accuracy: {acc:.4f}")

# 3 Find the best result

best_n = n_values[scores.index(max(scores))]

best_score = max(scores)

print(f"\n Best accuracy: {best_score:.4f} with {best_n} estimators")

# Optional: Plot the result

plt.plot(n_values, scores, marker='o')

```

```
plt.title("Accuracy vs Number of Estimators")  
plt.xlabel("n_estimators (number of trees)")  
plt.ylabel("Accuracy")  
plt.grid(True)  
plt.show()
```

Program 10

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

Screenshot:

DATE: PAGE:

1ab9

K-means Algorithm
for given data, compute two clusters using K-means alg for clustering
where initial cluster centers $(1,1)$ & $(5,7)$ execute two iterations

Record NO.	A	B	$c_1 = (1,1) \Rightarrow c_2 = (5,7)$
R1	1.0	1.0	$d = \sqrt{(x_2 - y_2)^2 + (y_2 - y_1)^2}$
R2	1.5	2.0	<u>Record (A,B)</u> $d(1,1)$ $d(1,2)$ C1
R3	3.0	4.0	<u>R1</u> (1,2) 0.0 7.2 C1
R4	5.0	7.0	<u>R2</u> (1,2) 1.1 6.1 C1
R5	3.5	5.0	<u>R3</u> (3,4) 3.6 4.2 C1
R6	4.5	5.0	<u>R4</u> (5,7) 2.2 0 C2
R7	3.5	4.5	<u>R5</u> (3,5) 1.0 2.5 C2
			<u>R6</u> (1,5) 3.32 2.2 C2
			<u>R7</u> (3,5) 4.3 3.2 C2
$c_1 = \left(\frac{1+1.5+3}{3}, \frac{1+2+4}{3} \right) = (1.83, 2.33)$			
$c_2 = \left(\frac{5+3.5+4.5+3.5}{4}, \frac{7+5+5+4.5}{4} \right) = (4.125, 5.325)$			

$d(1.83, 2.33)$	$d(4.12, 5.77)$	Clusters
1.57	9.62	C1
0.47	4.12	C1
2.12	1.63	C2
5.71	1.85	C2
3.53	0.72	C2
3.92	0.53	C2
3.07	1.01	C2

Code:

```

import pandas as pd

import matplotlib.pyplot as plt

from sklearn.cluster import KMeans

from sklearn.preprocessing import StandardScaler

# 1. Load dataset

df = pd.read_csv("/content/iris.csv")

# 2. Drop Sepal features, keep only Petal length and width

X = df[['petal_length', 'petal_width']]

# 3. Scaling (K-Means is distance-based, scaling helps!)

```

```

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 4. Elbow Method to find optimal k

inertia = []
k_range = range(1, 11)

for k in k_range:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertia.append(kmeans.inertia_)

# Plot elbow curve

plt.figure(figsize=(8,5))
plt.plot(k_range, inertia, 'bo-')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia (Within-cluster sum of squares)')
plt.title('Elbow Method for Optimal k')
plt.grid(True)
plt.show()

# 5. Choose optimal k (usually at 'elbow' point, say k=3)

optimal_k = 3
kmeans = KMeans(n_clusters=optimal_k, random_state=42)

```

```
clusters = kmeans.fit_predict(X_scaled)

# 6. Add clusters to original dataframe for visualization

df['Cluster'] = clusters

# Plot clusters

plt.figure(figsize=(8,5))

for i in range(optimal_k):

    cluster_data = df[df['Cluster'] == i]

    plt.scatter(cluster_data['petal_length'], cluster_data['petal_width'], label=f'Cluster {i}')

    plt.xlabel('Petal Length')
    plt.ylabel('Petal Width')
    plt.title('K-Means Clustering of Iris (Petal Features)')
    plt.legend()
    plt.grid(True)
    plt.show()
```

Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

Screenshot:

→ Dimensionality Reduction using principle component analysis
Reduce dim from 2 to 1 using PCA

Feature	Ex-1	Ex-2	Ex-3	Ex-4
x_1	9	8	13	7
x_2	11	4	5	14

eigen values $\lambda_1 = 30.3849$ $\lambda_2 = 6.6151$
eigen vectors $e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$ $e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$

Mean of $x_1 = 8$
Mean of $x_2 = 8.5$

$X_{\text{centered}} = \begin{bmatrix} 4-8 & 8-8 & 13-8 & 7-8 \\ 11-8.5 & 4-8.5 & 5-8.5 & 14-8.5 \end{bmatrix} = \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & 3.5 & 11 \end{bmatrix}$

about eigen value = λ_1
(corresponding vector = $e_1 = \begin{bmatrix} -0.5574 \\ -0.8303 \end{bmatrix}$)

DATE _____ PAGE _____

$$Z = C^T \cdot X_{\text{reduced}}$$

$$Z = [0.5574 \ -0.8303] \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & 3.5 & 5.5 \end{bmatrix}$$

$$Z_1 = (0.5574)(-4) + (-0.8303)(2.5)$$

$$Z_1 = 1.5385$$

Sum of squares

Code:

```

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import LabelEncoder, StandardScaler

from sklearn.svm import SVC

from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score

from sklearn.decomposition import PCA

```

Load the dataset

```
df = pd.read_csv('/content/heart (2).csv')
```

Detect and encode categorical columns

```
text_cols = df.select_dtypes(include=['object']).columns.tolist()
```

```
for col in text_cols:
```

```
    le = LabelEncoder()
```

```
    df[col] = le.fit_transform(df[col].astype(str))
```

```
# Use correct target column
```

```
target_col = 'HeartDisease'
```

```
X = df.drop(columns=[target_col])
```

```
y = df[target_col]
```

```
# Scale the features
```

```
scaler = StandardScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

```
# Train-test split
```

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

```
# Models
```

```
models = {
```

```
'SVM': SVC(),
```

```
'Logistic Regression': LogisticRegression(),
```

```
'Random Forest': RandomForestClassifier()
```

```
}
```

```
# Evaluate without PCA
```

```
results = {}

for name, model in models.items():

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)

    results[name] = accuracy_score(y_test, y_pred)

    print(f'{name} Accuracy: {results[name]:.4f}'")
```

Apply PCA

```
pca = PCA(n_components=0.95)

X_train_pca = pca.fit_transform(X_train)

X_test_pca = pca.transform(X_test)
```

Evaluate after PCA

```
pca_results = {}

for name, model in models.items():

    model.fit(X_train_pca, y_train)

    y_pred = model.predict(X_test_pca)

    pca_results[name] = accuracy_score(y_test, y_pred)

    print(f'{name} Accuracy after PCA: {pca_results[name]:.4f}'")
```

Final comparison

```
print("\nComparison (Original vs PCA):")

for name in models:

    print(f'{name}: Original = {results[name]:.4f}, PCA = {pca_results[name]:.4f}'")
```

