

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



**LAB REPORT  
on**

**Machine Learning (23CS6PCMAL)**

*Submitted by*

**Manav Kumar (1BM22CS348)**

*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**

**Feb-2025 to June-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 **Manav Kumar (1BM22CS348)**, 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.

Rajeshwari Madli Assistant Professor Department of CSE, BMSCE	Dr. Kavitha Sooda Professor & HOD Department of CSE, BMSCE
---	--

## Index

<b>Sl. No.</b>	<b>Date</b>	<b>Experiment Title</b>	<b>Page No.</b>
1	21-2-2025	Write a python program to import and export data using Pandas library functions	1-2
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	3-5
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	6-17
4	17-3-2025	Build Logistic Regression Model for a given dataset	18-21
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	22-24
6	7-4-2025	Build KNN Classification model for a given dataset.	25-27
7	21-4-2025	Build Support vector machine model for a given dataset	28-30
8	5-5-2025	Implement Random forest ensemble method on a given dataset.	31-33
9	5-5-2025	Implement Boosting ensemble method on a given dataset.	34-37
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	38-41
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	42-46

Github Link:

<https://github.com/Manav-Kumar123/6thSem-ML-L>

## Program 1

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

### **Code:**

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OrdinalEncoder, OneHotEncoder
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from scipy import stats
df=pd.read_csv('/content/Dataset of Diabetes .csv')
df.head()
df.shape
print(df.info())
# Summary statistics
print(df.describe())
missing_values=df.isnull().sum()
categorical_cols = df.select_dtypes(include=['object']).columns
print("Categorical columns identified:", categorical_cols)
if len(categorical_cols) > 0:
    df = pd.get_dummies(df, columns=categorical_cols, drop_first=True)
    print("\nDataFrame after one-hot encoding:")
    print(df.head())
else:
    print("\nNo categorical columns found in the dataset.")
from sklearn.preprocessing import MinMaxScaler, StandardScaler
import pandas as pd

numerical_cols = df.select_dtypes(include=['number']).columns

scaler = MinMaxScaler()
df_minmax = df.copy() # Create a copy to avoid modifying the original
df_minmax[numerical_cols] = scaler.fit_transform(df[numerical_cols])

scaler = StandardScaler()
df_standard = df.copy()
df_standard[numerical_cols] = scaler.fit_transform(df[numerical_cols])
print("\nDataFrame after Min-Max Scaling:")
print(df_minmax.head())
print("\nDataFrame after Standardization:")
print(df_standard.head())
df1=pd.read_csv('/content/adult.csv')
df1.head()
```

```
from sklearn.preprocessing import MinMaxScaler, StandardScaler
import pandas as pd

numerical_cols = df1.select_dtypes(include=['number']).columns

scaler = MinMaxScaler()
df_minmax = df1.copy() # Create a copy to avoid modifying the original
df_minmax[numerical_cols] = scaler.fit_transform(df1[numerical_cols])

scaler = StandardScaler()
df_standard = df1.copy()
df_standard[numerical_cols] = scaler.fit_transform(df1[numerical_cols])
print("\nDataFrame after Min-Max Scaling:")
print(df_minmax.head())
print("\nDataFrame after Standardization:")
print(df_standard.head())
```

## Program 2

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

### Observation:

LAB-2 DATA	
	import pandas as pd (i) df = pd.read_csv("housing.csv") (ii) print(df.info()) (iii) print(df.describe()) (iv) print(df["Ocean Proximity"].value_counts()) (v) missing = df.isnull().sum() print(missing[missing > 0])
	Questions: 1. As per analysis no columns in the dataset had the missing values. The missing values were handled by filling any numerical columns with median values and categorical columns with mode.
	2. The identified categorical columns in Diabetes dataset are Gender, class. Categorical columns in Income dataset are 'workclass', education, marital-status, occupation, relationship, race, gender, native-country, income. Label encoding was used to encode them.
	3. Min-max scaling scales values between 0 and 1. Standardization (Z-score transformation) transforms data to have mean = 0 and standard deviation = 1. Min-max is used when preserving the original distribution is important. Standardization is used when data has different scales or follows a normal distribution.

Figure 2.1

**Code:**

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder
from google.colab import files
```

```
# Upload Files Manually in Google Colab
uploaded = files.upload()
```

```
# Load the datasets (replace filenames accordingly after uploading)
diabetes_df = pd.read_csv("diabetes.csv")
adult_df = pd.read_csv("adult.csv")
```

```
# --- Data Cleaning ---
```

```
# Handling Missing Values: Fill numerical columns with median, categorical with mode
for df in [diabetes_df, adult_df]:
    for col in df.columns:
        if df[col].isnull().sum() > 0:
            if df[col].dtype == "object":
                df[col].fillna(df[col].mode()[0], inplace=True)
            else:
                df[col].fillna(df[col].median(), inplace=True)
```

```
# Handling Outliers: Capping values beyond 1.5*IQR
```

```
for df in [diabetes_df, adult_df]:
    for col in df.select_dtypes(include=np.number).columns:
        Q1, Q3 = df[col].quantile(0.25), df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower, upper = Q1 - 1.5 * IQR, Q3 + 1.5 * IQR
        df[col] = np.clip(df[col], lower, upper)
```

```
# --- Handling Categorical Data ---
```

```
for df in [diabetes_df, adult_df]:
    categorical_cols = df.select_dtypes(include="object").columns
    for col in categorical_cols:
        df[col] = LabelEncoder().fit_transform(df[col])
```

```
# --- Data Transformations ---
```

```
scaler_minmax = MinMaxScaler()
scaler_standard = StandardScaler()

for df in [diabetes_df, adult_df]:
    numerical_cols = df.select_dtypes(include=np.number).columns
    df[numerical_cols] = scaler_minmax.fit_transform(df[numerical_cols])
    df[numerical_cols] = scaler_standard.fit_transform(df[numerical_cols])
```

```

# Save processed datasets
diabetes_df.to_csv("processed_diabetes.csv", index=False)
adult_df.to_csv("processed_adult.csv", index=False)

# Download processed files
files.download("processed_diabetes.csv")
files.download("processed_adult.csv")
import pandas as pd
from google.colab import files

# Upload Files Manually in Google Colab
uploaded = files.upload()

# Load the datasets
diabetes_df = pd.read_csv("diabetes.csv")
adult_df = pd.read_csv("adult.csv")

# Check for missing values
missing_diabetes = diabetes_df.isnull().sum()
missing_adult = adult_df.isnull().sum()

# Display columns with missing values
print("Missing values in Diabetes Dataset:")
print(missing_diabetes[missing_diabetes > 0])

print("\nMissing values in Adult Income Dataset:")
print(missing_adult[missing_adult > 0])

print("Missing Values Count in Diabetes Dataset:")
print(missing_diabetes)

print("\nMissing Values Count in Adult Income Dataset:")
print(missing_adult)

categorical_diabetes = diabetes_df.select_dtypes(include="object").columns.tolist()
categorical_adult = adult_df.select_dtypes(include="object").columns.tolist()

# Display categorical columns
print("Categorical Columns in Diabetes Dataset:", categorical_diabetes)
print("\nCategorical Columns in Adult Income Dataset:", categorical_adult)

```

### Program 3

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

#### Observation:

		<u>LAB-3</u>	06.		
17/03/2025:-					
Q. $y_i$ (Week)		$y_i$ (sales in thousands)			
1		2			
2		4			
3		5			
4		9			
$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} 1 & n_1 \\ 1 & n_2 \\ 1 & n_3 \\ 1 & n_4 \end{bmatrix} \begin{bmatrix} B_0 \\ B_1 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}$					
$\beta = ((x^T x)^{-1} x^T) Y$ $x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}$ $x^T x = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$					
$(x^T x)^{-1} x^T = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$ $\approx \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$ $\approx \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$ $\therefore \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$					
$Y = 2.2n - 0.5$					

Figure 3.1

$\Rightarrow$ 

Normal form:-

$$\beta_1 = (\bar{u}\bar{y}) - (\bar{u})(\bar{y})$$

$$(\bar{u}^2) - (\bar{u})^2$$

$$\beta_0 = \bar{y} - \beta_1 \times \bar{u}$$

$$\begin{array}{|c|c|c|c|} \hline u & y & u^2 & uy \\ \hline \end{array}$$

$$1 & 2 & 1 & 2 \\ \hline$$

$$2 & 4 & 4 & 8 \\ \hline$$

$$3 & 5 & 9 & 15 \\ \hline$$

$$4 & 9 & 16 & 36 \\ \hline \end{array}$$

$$\bar{u} = 10/4 = 2.5 \quad \bar{u}^2 = 7.5$$

$$\bar{y} = 20/4 = 5 \quad \bar{u}y = 15.25$$

$$(\bar{u})^2 = 6.25$$

$$\beta_1 = (\bar{u}\bar{y}) - (\bar{u})(\bar{y}) = 15.25 - (2.5)(5)$$

$$(\bar{u}^2)^2 - (\bar{u})^2 = 42.25 - 6.25$$

$$= 2.75 \Rightarrow \boxed{\beta_1 = 2.2}$$

$$\beta_0 = (\bar{y}) - \beta_1 \times \bar{u}$$

$$= 5 - 2.2(2.5)$$

$$\boxed{\beta_0 = -0.5}$$

$$\text{Now } y = \beta_0 + \beta_1 u$$

$$\boxed{y = -0.5 + 2.2u}$$

$$\begin{array}{|c|c|c|} \hline \bar{u} & \bar{y} & \beta_1 \\ \hline 1 & 2 & 2.2 \\ \hline 2 & 4 & 2.2 \\ \hline 3 & 5 & 2.2 \\ \hline 4 & 9 & 2.2 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline \bar{u} & \bar{y} & \beta_1 \\ \hline 1 & 2 & 2.2 \\ \hline 2 & 4 & 2.2 \\ \hline 3 & 5 & 2.2 \\ \hline 4 & 9 & 2.2 \\ \hline \end{array}$$

Figure 3.2

Q.	Diameter in inch (x)	Price in \$ (y)
	8	$(P) = 10$
	10	$(P) = 13$
	12	$(P) = 16$

Price of 20 inch pizza = ?

$\Rightarrow$  Matrix Method:

$$X = \begin{bmatrix} 1 & 8 \\ 1 & 10 \\ 1 & 12 \end{bmatrix}; Y = \begin{bmatrix} 10 \\ 13 \\ 16 \end{bmatrix}; B = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$B = ((X^T X)^{-1} X^T) Y = \begin{bmatrix} -2 \\ 1.5 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$\beta_0 = -2$$

$$\beta_1 = 1.5$$

$$Y = \beta_0 + \beta_1 x$$

$$Y = -2 + 1.5 x$$

Now,  $y = ?$  when  $x = 20$

$$y = -2 + 1.5(20)$$

$$y = 28 \$$$

$\Rightarrow$  Normal Method:

x	y	xy	$x^2$
8	10	80	64
10	13	130	100
12	16	192	144

$$\bar{x} = 10$$

$$\bar{y} = 13$$

$$\bar{xy} = 308/3$$

$$\bar{x^2} = 10^2 = 100$$

$$\bar{y^2} = 402/3 = 134$$

$$\beta_1 = \frac{\bar{xy} - (\bar{x})(\bar{y})}{(\bar{x}^2) - (\bar{x})^2} = \frac{134 - 130}{100 - 100} = \frac{308/3 - 300/3}{8} = \frac{3(4)}{8} = 1.5$$

Figure 3.3

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\beta_1 = 1.5$$

$$\begin{aligned}\beta_0 &= (\bar{y}) - \beta_1(\bar{n}) \\ &= 13 - 1.5(10)\end{aligned}$$

$$\boxed{\beta_0 = -2}$$

$$\therefore y = \beta_0 + \beta_1 n$$

$$y = -2 + 1.5n$$

at  $n = 20$ ;

$$y = -2 + 30$$

$$\boxed{y = 28\$}$$

Q. Did you perform any data preprocessing steps?

→ Yes, I handled missing values by fitting them with the column mean. I also applied label encoding to categorical columns & scaled numerical features for `1000-companies.csv` to normalize the data.

Q. Did you visualize the regression line for `canada-per-capita-income.csv`?

→ Yes, the regression line was plotted. The plot shows a strong linear relationship between year & per capita income, meaning that as the year increases, per capita income also increases.

Q.

→ The predicted salary is printed in the script & depends on the trained model's coefficients.

Figure 3.4

- Q. Did you encode categorical variables for 1000-companies.csv?

→ Yes, the state column was encoded using LabelEncoder().

- Q. Did you scale the features? *as - it has*

→ Yes, because R&D spend, Administration & marketing spend have different units, feature scaling was applied to compare model performance.

Figure 3.5

### Code:

```
import numpy as np

# Given data
# x: Week numbers
# y: Sales in thousands
x = np.array([1, 2, 3, 4])
y = np.array([2, 4, 5, 9])

# Construct the design matrix X by adding a column of ones (for the intercept)
X = np.column_stack((np.ones(x.shape[0]), x))

# Compute the coefficients using the formula: beta = (X^T X)^(-1) X^T y
XtX = X.T.dot(X)      # Compute X^T X
XtX_inv = np.linalg.inv(XtX) # Invert X^T X
XtY = X.T.dot(y)       # Compute X^T y

beta = XtX_inv.dot(XtY) # Compute beta

# Display the computed coefficients
print("Computed coefficients (beta):", beta)

import matplotlib.pyplot as plt

# ... (previous code)

# Generate points for the regression line
x_line = np.linspace(x.min(), x.max(), 100) # Create 100 points for a smooth line
y_line = beta[0] + beta[1] * x_line      # Calculate y-values for the line

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

# Plot the regression line
plt.plot(x_line, y_line, label='Linear Regression', color='red')

# Customize the plot
plt.xlabel('Week Number (x)')
plt.ylabel('Sales (thousands) (y)')
plt.title('Linear Regression Plot')
plt.legend() # Show the legend
plt.grid(True) # Show the grid

# Display the plot
plt.show()
```

```

import numpy as np

# Given data
x = np.array([8, 10, 12])
y = np.array([10, 13, 16])

# Construct the design matrix X (adding a column of ones for the intercept)
X = np.column_stack((np.ones(x.shape[0]), x))

# Compute beta using the normal equation: beta = (X^T X)^(-1) X^T y
XtX = X.T.dot(X)
XtX_inv = np.linalg.inv(XtX)
XtY = X.T.dot(y)
beta = XtX_inv.dot(XtY)

# Extract coefficients
beta0, beta1 = beta
print("Intercept (beta0):", beta0)
print("Slope (beta1):", beta1)

# Predict the price for a 20-inch pizza
x_new = 20
y_pred = beta0 + beta1 * x_new
print("Predicted price for a 20-inch pizza: $", y_pred)

import pandas as pd
from sklearn.linear_model import LinearRegression
# Load the data
income_data = pd.read_csv("canada_per_capita_income.csv")
# Assumed data columns: 'Year' and 'PerCapitaIncome'
print("Canada Income Data Head:")
print(income_data.head())
# Prepare feature and target
X_income = income_data[["year"]] # Predictor variable: Year
y_income = income_data["per capita income (US$)"] # Target variable: Per capita income

# Build and train the linear regression model
model_income = LinearRegression()
model_income.fit(X_income, y_income)
# Predict per capita income for the year 2020
predicted_income = model_income.predict([[2020]])
print("\nPredicted per capita income for Canada in 2020:", predicted_income[0])

import matplotlib.pyplot as plt

# ... (previous code)

# Predict per capita income for the year 2020

```

```

predicted_income = model_income.predict([[2020]])
print("\nPredicted per capita income for Canada in 2020:", predicted_income[0])

# Plot the data points and the regression line
plt.scatter(X_income, y_income, color='blue', label='Actual Data')
plt.plot(X_income, model_income.predict(X_income), color='red', label='Regression Line')

# Plot the prediction for 2020
plt.scatter(2020, predicted_income[0], color='green', label='Prediction for 2020')

# Customize the plot
plt.xlabel('Year')
plt.ylabel('Per Capita Income (US$)')
plt.title('Canada Per Capita Income Prediction')
plt.legend()
plt.grid(True)

# Display the plot
plt.show()

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

# Load the salary data
salary_data = pd.read_csv("salary.csv")
print(salary_data.head())
# Check for null values and handle them (e.g., imputation or removal)
if salary_data.isnull().values.any():
    print("Null values found in the salary dataset. Handling null values...")
    # Example: Fill null values with the mean of the 'YearsExperience' column
    salary_data['YearsExperience'].fillna(salary_data['YearsExperience'].mean(), inplace=True)
    # Other options: Remove rows with nulls or use more sophisticated imputation methods

# Prepare feature and target
X_salary = salary_data[['YearsExperience']] # Predictor variable: Years of Experience
y_salary = salary_data['Salary'] # Target variable: Salary
# Build and train the linear regression model
model_salary = LinearRegression()
model_salary.fit(X_salary, y_salary)
# Predict salary for an employee with 12 years of experience
predicted_salary = model_salary.predict([[12]])
print("\nPredicted salary for an employee with 12 years of experience:", predicted_salary[0])

import matplotlib.pyplot as plt
# Plot the data points and the regression line
plt.scatter(X_salary, y_salary, color='blue', label='Actual Data')

```

```

plt.plot(X_salary, model_salary.predict(X_salary), color='red', label='Regression Line')

# Plot the prediction for 12 years of experience
plt.scatter(12, predicted_salary[0], color='green', label='Prediction for 12 years')

# Customize the plot
plt.xlabel('Years of Experience')
plt.ylabel('Salary')
plt.title('Salary Prediction based on Experience')
plt.legend()
plt.grid(True)

# Display the plot
plt.show()

import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression

# Read the CSV file (ensure the file is uploaded in your Colab environment)
df = pd.read_csv("hiring.csv")
# Rename columns for convenience
df.columns = ['experience', 'test_score', 'interview_score', 'salary']

print("Original Data:")
print(df)
# Define a mapping for text to numeric conversion for the 'experience' column
num_map = {
    "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
}

# Function to convert experience values to numeric
def convert_experience(x):
    try:
        return float(x)
    except:

```

```

x_lower = str(x).strip().lower()
return num_map.get(x_lower, np.nan)

# Convert the 'experience' column using the mapping
df['experience'] = df['experience'].apply(convert_experience)

# Convert 'test_score', 'interview_score', and 'salary' to numeric (coerce errors to NaN)
df['test_score'] = pd.to_numeric(df['test_score'], errors='coerce')
df['interview_score'] = pd.to_numeric(df['interview_score'], errors='coerce')
df['salary'] = pd.to_numeric(df['salary'], errors='coerce')

print("\nData After Conversion:")
print(df)
# Fill missing values in numeric columns using the column mean
df['experience'].fillna(df['experience'].mean(), inplace=True)
df['test_score'].fillna(df['test_score'].mean(), inplace=True)
df['interview_score'].fillna(df['interview_score'].mean(), inplace=True)

print("\nData After Filling Missing Values:")
print(df)
# Prepare the feature matrix X and target vector y
X = df[['experience', 'test_score', 'interview_score']]
y = df['salary']

# Build and train the Multiple Linear Regression model
model = LinearRegression()
model.fit(X, y)
# Predict salaries for the given candidate profiles
# Candidate 1: 2 years of experience, 9 test score, 6 interview score
candidate1 = np.array([[2, 9, 6]])
predicted_salary1 = model.predict(candidate1)

# Candidate 2: 12 years of experience, 10 test score, 10 interview score
candidate2 = np.array([[12, 10, 10]])
predicted_salary2 = model.predict(candidate2)
import matplotlib.pyplot as plt

# Create the plot
plt.figure(figsize=(10, 6)) # Adjust figure size for better visualization
plt.scatter(df['experience'], y, color='blue', label='Actual Salary') #Plot actual salary against years of experience

# Plot the regression line (this is an approximation since it's a multi-variable regression)
# You can visualize a single feature against the predicted salary
plt.plot(df['experience'], model.predict(X), color='red', label='Regression Line')

# Highlight predictions
plt.scatter(candidate1[0, 0], predicted_salary1, color='green', label='Candidate 1 Prediction')

```

```

plt.scatter(candidate2[0, 0], predicted_salary2, color='purple', label='Candidate 2 Prediction')

# Add labels and title
plt.xlabel("Years of Experience")
plt.ylabel("Salary")
plt.title("Salary Prediction based on Experience, Test Score, Interview Score")

# Add a legend
plt.legend()
plt.grid(True)
plt.show()

import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression

# Read the CSV file (ensure the file is uploaded in your Colab environment)
df = pd.read_csv("1000_Companies.csv")
# Display the first few rows
print("Original Data:")
print(df.head())
# --- Data Preprocessing ---

# For numeric columns, fill missing values with the column mean
numeric_cols = ["R&D Spend", "Administration", "Marketing Spend", "Profit"]
for col in numeric_cols:
    df[col].fillna(df[col].mean(), inplace=True)

# For the categorical column 'State', fill missing values with a placeholder
df["State"].fillna("Unknown", inplace=True)

# Confirm that missing values are handled
print("\nMissing Values After Processing:")
print(df.isnull().sum())
# Separate the features and target variable
features = ["R&D Spend", "Administration", "Marketing Spend"] + \
    [col for col in df_encoded.columns if col.startswith("State_")]
X = df_encoded[features]
y = df_encoded["Profit"]
# --- Prediction for a New Company ---

# Given sample data:
# R&D Spend = 91694.48, Administration = 515841.3, Marketing Spend = 11931.24, State = 'Florida'
new_company = pd.DataFrame({
    "R&D Spend": [91694.48],
    "Administration": [515841.3],
    "Marketing Spend": [11931.24],
    "State": ["Florida"]}

```

```

})
# One-hot encode the 'State' column using the same strategy as training data
new_company_encoded = pd.get_dummies(new_company, columns=["State"], drop_first=True)

# Align the new data's columns with the training features (fill missing columns with 0)
new_company_encoded = new_company_encoded.reindex(columns=X.columns, fill_value=0)

# Predict the profit using the trained model
predicted_profit = model.predict(new_company_encoded)
print("\nPredicted Profit for the New Company: ${}", round(predicted_profit[0], 2))

import matplotlib.pyplot as plt

# Assuming 'df_encoded', 'features', 'X', 'y', 'model', 'new_company_encoded', and 'predicted_profit' are
# defined from the previous code

# Create the plot
plt.figure(figsize=(10, 6))

# Scatter plot of actual profits vs. R&D Spend
plt.scatter(df_encoded["R&D Spend"], y, color='blue', label='Actual Profit')

# Plot the regression line (approximation for visualization)
plt.plot(df_encoded["R&D Spend"], model.predict(X), color='red', label='Regression Line')

# Highlight the new company's prediction
plt.scatter(new_company_encoded["R&D Spend"], predicted_profit, color='green', label='New
Company Prediction')

# Add labels and title
plt.xlabel("R&D Spend")
plt.ylabel("Profit")
plt.title("Profit Prediction based on R&D Spend")

# Add a legend
plt.legend()
plt.grid(True)
plt.show()

```

## Program 4

Build Logistic Regression Model for a given dataset.

### Observation:

Q.	LAB-4.	11.
	Date: 24/03/2025:-	Date _____ Page _____
	Binary Regression.	
Given:	$a_0 = -5$ (Intercept) $a_1 = 0.8$ (co-efficient for study hours)	
(a)	Logistic regression equation: $P(y=1) = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$ $P(y=1) = \frac{1}{1 + e^{-( -5 + 0.8x)}} ; a_0 = -5$ $a_1 = 0.8$	
(b)	Probability calculation: $x = 7$ (study hours) $P(y=1) = \frac{1}{1 + e^{-( -5 + 0.8(7))}}$ $= \frac{1}{1 + e^{-( -5 + 5.6)}}$ $= \frac{1}{1 + e^{-0.6}}$ Now, $e^{-0.6} \approx 0.54$ $\text{so, } P(y=1) = \frac{1}{1 + 0.5488}$ $= \frac{1}{1.5488} \approx 0.6457$ The probability that student will pass is 64.57%.	
(c)	Predicted class (threshold = 0.5) If; $P(y=1) \geq 0.5$ , $\rightarrow$ Pass $P(y=1) < 0.5$ , $\rightarrow$ fail. Since, $0.6457 > 0.5$ $\Rightarrow$ The student will pass.	

Figure 4.1

**Code:**

**Hr.csv**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

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

# Basic Info
print("Dataset Info:")
print(df.info())
print("\nFirst few rows:")
print(df.head())
plt.figure(figsize=(8, 6))
# sns.countplot(x='salary', hue='left', data=df)
sns.barplot(x='Department', y='satisfaction_level', data=df)

# plt.title('Salary vs Employee Retention')
plt.xlabel('Departments')
plt.ylabel('Satisfaction level')

plt.show()

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

# Encode categorical variables (drop_first avoids dummy variable trap)
df_encoded = pd.get_dummies(df, columns=['salary', 'Department'], drop_first=True)

plt.figure(figsize=(15, 8))
sns.heatmap(df_encoded.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()

plt.figure(figsize=(8, 5))
sns.countplot(x='salary', hue='left', data=df, order=['low', 'medium', 'high'])

plt.title('Impact of Salary on Employee Retention')
```

```

plt.xlabel('Salary Level')
plt.ylabel('Number of Employees')
plt.legend(title='Left', labels=['Stayed', 'Left'])
plt.show()
df_encoded = pd.get_dummies(df, columns=['Department', 'salary'], drop_first=True)

# Calculate the correlation matrix
correlation_matrix = df_encoded.corr()

# Extract the correlation with 'left' (employee retention)
correlation_with_left = correlation_matrix['left'].sort_values(ascending=False)

# Display the correlation
print(correlation_with_left)
plt.figure(figsize=(12, 6))
sns.countplot(x='Department', hue='left', data=df)

# Title and labels
plt.title('Impact of Department on Employee Retention')
plt.xlabel('Department')
plt.ylabel('Number of Employees')
plt.legend(title='Left', labels=['Stayed', 'Left'])
plt.xticks(rotation=45) # Rotate department names for readability
plt.show()

# Step 1: Preprocess the data
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report

# Load the dataset
df = pd.read_csv('HR_comma_sep.csv')

# Select important features and encode categorical variable
df_encoded = pd.get_dummies(df, columns=['salary'], drop_first=True) # This encodes salary (low -> low salary column)

# Step 2: Define features (X) and target (y)
X = df_encoded[['satisfaction_level', 'time_spend_company', 'salary_low']] # Using low salary as a feature
y = df_encoded['left'] # Target variable (whether the employee left or stayed)

# Step 3: Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Step 4: Build and train the logistic regression model
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

```

```
# Step 5: Make predictions
y_pred = model.predict(X_test)

# Step 6: Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy * 100:.2f}%')
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
```

### Program 5

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

#### Observation:

Decision Trees:			
Instance	$a_2$	$a_3$	class
1	Hot	High	No
2	Hot	High	No
6	Cool	High	No
7	Hot	High	No
8	Hot	Normal	Yes

Identify whether the splitting node should be  $a_2$  or  $a_3$  attribute.

$$\text{entropy}(S) = - \sum_{i=1}^C p_i \log_2 p_i$$

$$p(\text{No}) = \frac{4}{5} \quad p(\text{Yes}) = \frac{1}{5}$$

$$\text{Entropy}(S) = -\left(\frac{4}{5} \log_2 \left(\frac{4}{5}\right) + \frac{1}{5} \log_2 \left(\frac{1}{5}\right)\right)$$

$$= 0.7219$$

Attribute:

$a_2 \Rightarrow$  values : Hot, Cool

$$\text{Entropy}(\text{Hot}) = -\left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4}\right)$$

$$= 0.8112$$

$$\text{Entropy}(\text{Cool}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) = 0$$

$$I_G(a_2) = \text{Entropy}(S) - \left(\frac{4}{5} (0.8112) + \frac{1}{5} (0)\right)$$

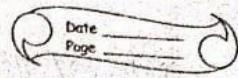
$$= 0.7219 - 0.64896$$

$$= 0.07294$$

$a_3 \Rightarrow$  values  $\rightarrow$  (High, Normal)

$$\text{Entropy}(\text{High}) = -\left(\frac{4}{7} \log_2 \left(\frac{4}{7}\right) + \frac{3}{7} \log_2 \left(\frac{3}{7}\right)\right) = 0$$

Figure 5.1



$$\text{Entropy (Normal)} = -( \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} ) + 0 = 0.$$

$$IG(a_3) = 0.7219 - 0 = 0.7219$$

Maximum  $IG \Rightarrow a_3 = 0.7219$

∴  $a_3$  attribute should be the splitting Node.

✓

24/3

Figure 5.2

### Code:

```
import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report
from sklearn.tree import plot_tree
import matplotlib.pyplot as plt
data = {
    'a1': [True, True, False, False, True, True, True, False, False],
    'a2': ['Hot', 'Hot', 'Hot', 'Cool', 'Cool', 'Hot', 'Hot', 'Cool', 'Cool'],
    'a3': ['High', 'High', 'High', 'Normal', 'Normal', 'High', 'High', 'Normal', 'Normal'],
    'Classification': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'No', 'Yes', 'Yes']}
}
df = pd.DataFrame(data)
df.head()
label_encoders = {}
for column in df.columns:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
    label_encoders[column] = le
df.head()
X = df.drop('Classification', axis=1)
y = df['Classification']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
clf = DecisionTreeClassifier(criterion='entropy')
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
accuracy
```

## Program 6

Build KNN Classification model for a given dataset.

### Observation:

Date: 07/04/2025:- LAB-05							
KNN							
(#)	Given, $K=3$ test data ( $x, 35, 100$ )	Person	Age	Salary/N	Target	Distance	Rank
A	18	50	N	52.81	2		
B	23	55	N	46.57	3		
C	24	70	N	31.95			
D	41	60	Y	40.45			
E	43	70	Y	31.05			
F	38	40	X	60.07	1		

∴ Possibilities = N, N, Y  
So, test data ( $x, 35, 100$ ) belongs to N class.

(#) \* For the iris dataset, first the dataset is split into training and testing dataset.  
\* A range of  $K$ -values are selected. ex. (1, 3)  
\* For each  $K$ , value, a model is trained.  
\* For each value, accuracy and error rate is calculated.

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

Error rate =  $1 - \text{Accuracy}$

\*  $K$ -value for highest accuracy and/or lowest error rate is chosen.  
\* For iris dataset, best value found  
 $K=3$ .

Figure 6.1

Date \_\_\_\_\_  
Page \_\_\_\_\_

(1) \* Feature scaling is the process of normalizing the scale of values of independent features.

\* Need:

- (1) Improved performance & efficiency as scales are same or comparable.
- (2) Prevent dominance of large scale of feature.

For diabetes dataset, Blood pressure feature ranges b/w 50-150, while diabetes pedigree function is in the range 0-1.

\* Different methods:-

(1) Min-Max (Normalization)

for each value,

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

(2) Standardization (Z-score)

$$\bar{x} = \frac{x - \mu}{\sigma}$$

∴ for diabetes dataset, for BP column.

Ex.

$$x' = \frac{72 - 0}{122 - 0} = 0.59.$$

By  
27/11/25

Figure 6.2

### Code:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import OrdinalEncoder, StandardScaler
data = pd.read_csv("diabetes.csv")
data.head()
X = data.iloc[:, :-1]
y = data.iloc[:, -1]
ss = StandardScaler()
X[["Pregnancies"]] = ss.fit_transform(X[["Pregnancies"]])
X[["Glucose"]] = ss.fit_transform(X[["Glucose"]])
X[["BloodPressure"]] = ss.fit_transform(X[["BloodPressure"]])
X[["SkinThickness"]] = ss.fit_transform(X[["SkinThickness"]])
X[["Insulin"]] = ss.fit_transform(X[["Insulin"]])
X[["BMI"]] = ss.fit_transform(X[["BMI"]])
X[["DiabetesPedigreeFunction"]] = ss.fit_transform(X[["DiabetesPedigreeFunction"]])
X[["Age"]] = ss.fit_transform(X[["Age"]])
X.head()
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
knn = KNeighborsClassifier()
param_grid = {"n_neighbors": [1, 3, 5, 7, 9]}
grid = GridSearchCV(estimator = knn, param_grid = param_grid, cv = 5, scoring = "accuracy")
grid.fit(X_train, y_train)
grid.best_params_
best = grid.best_estimator_
best
y_pred = best.predict(X_test)
accuracy_score(y_test, y_pred)
```

## Program 7

Build Support vector machine model for a given dataset.

### Observation:

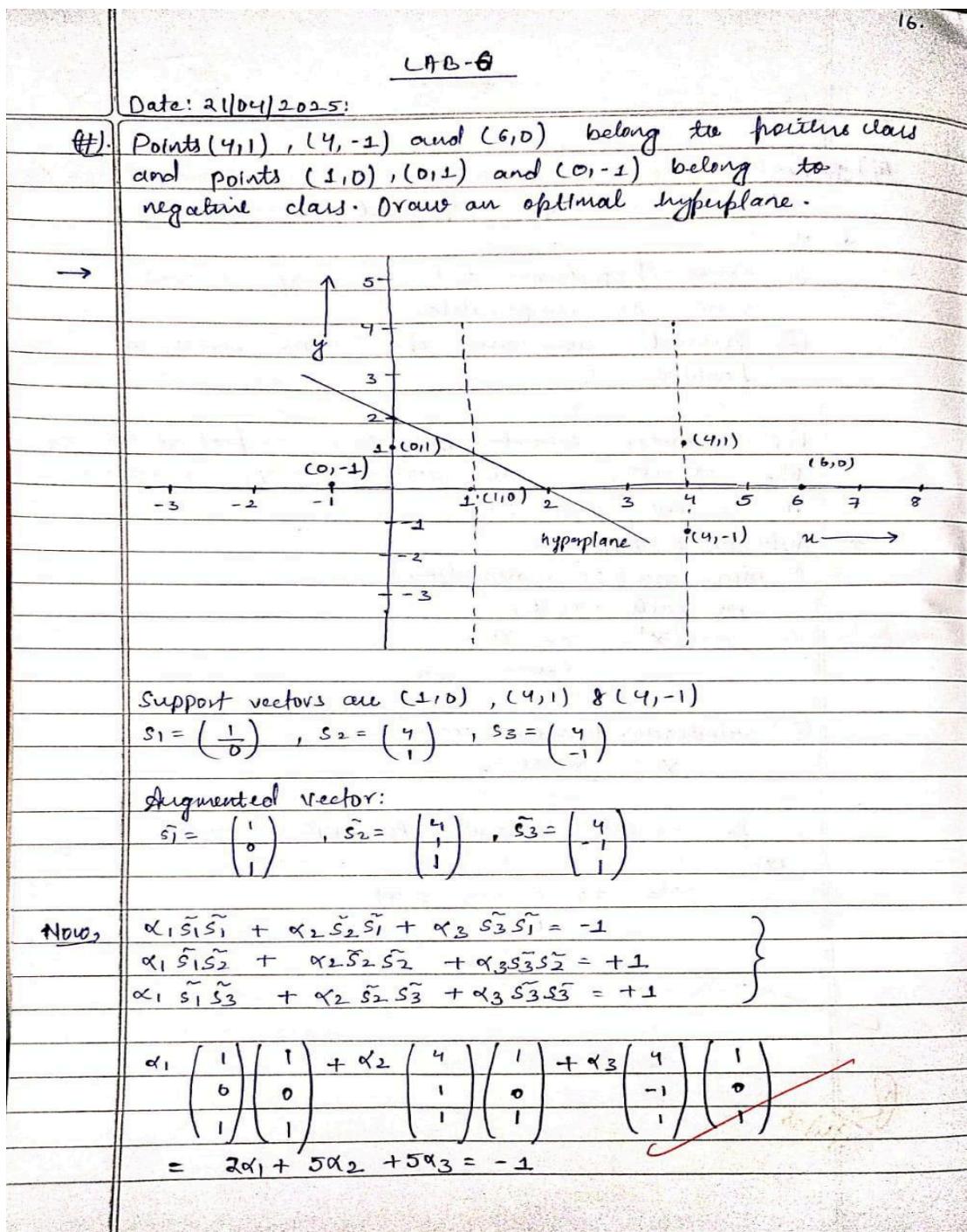


Figure 7.1

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix}$$

$$= 5\alpha_1 + 18\alpha_2 + 16\alpha_3 = +1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix}$$

$$= 5\alpha_1 + 16\alpha_2 + 18\alpha_3 = +1$$

Solving these simultaneous equations:

$$\alpha_1 = -3; \alpha_2 = +1; \alpha_3 = 0$$

The optimal hyperplane is given as:

$$\omega = \sum_{i=1}^3 \alpha_i x_i \tilde{s}_i$$

$$= -3 \times \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 1 \times \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} + 0 \times \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ -2 \end{pmatrix}$$

The hyperplane is  $(1, -1)$  with an offset -2.

This dataset:

Q1.)

Linear Kernel:

Accuracy: 1.0000

confusion Matrix:

$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

RBF Kernel:

Accuracy: 1.0000

confusion Matrix:

$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

Both RBF & Linear kernel performed equally.

Q2.)

Letley-recognition dataset:

Accuracy score: 93.05%

AUC score: 0.9985

- '4' confused with '0'.
- 'P' confused with 'F/F' '4'.

Figure 7.2

## Code:

### Iris.csv

```
import pandas as pd
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.preprocessing import OrdinalEncoder
data = pd.read_csv("iris (1).csv")
data.head()
oe = OrdinalEncoder()
data[["species"]] = oe.fit_transform(data[["species"]])
data.head()
y = data.iloc[:, -1]
X = data.iloc[:, :-1]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
rbf_model = SVC(kernel='rbf')
rbf_model.fit(X_train, y_train)
rbf_model.score(X_test,y_test)
y_pred = rbf_model.predict(X_test)
print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))
linear_model = SVC(kernel='linear')
linear_model.fit(X_train,y_train)
linear_model.score(X_test,y_test)
y_pred = rbf_model.predict(X_test)
print(confusion_matrix(y_test, y_pred))
```

### Digits.csv

```
import pandas as pd
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
digits = load_digits()
digits.target
dir(digits)
X_train, X_test, y_train, y_test = train_test_split(df.drop('target',axis='columns'), df.target,
test_size=0.3)
rbf_model = SVC(kernel='rbf')
rbf_model.fit(X_train, y_train)
linear_model = SVC(kernel='linear')
linear_model.fit(X_train,y_train)
```

## Program 8

Implement Random forest ensemble method on a given dataset.

### Observation:

		18.
	<u>LAB-07</u>	
	05/05/2025:	
(H).	Implement Random forest ensemble method on a given dataset.	
1.	$Y_{ini}(S) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 1 - 0.36 - 0.16$ $= 0.48$	
2.	Split with C4PA $Y_{ini} = \frac{1}{5}x_0 + \frac{4}{5}x_0 \cdot 0.375 = 0.300$ $\Delta Y_{ini} = 0.48 - 0.300 = 0.18$	
3.	Growing C4PA $\Rightarrow$ model based on interactions $\rightarrow Y_{ini} = \frac{2}{4}x_0 + \frac{2}{4}x_0 \cdot 0.5 = 0.25$ $\Delta Y_{ini} = 0.375 - 0.25 = 0.125$	
	$\rightarrow$ "based" on practical knowledge $Y_{ini} = \frac{3}{4}x_0 + \frac{1}{4}x_0 = 0$ $\Delta Y_{ini} = 0.375 - 0 = 0.375$ (best)	
	final cut $[C4PA = ?]$ <pre>           /   \          No    Yes         / \   / \        No  Yes No  Yes         [Practical knowledge]         yes   No       </pre> <p style="text-align: center;"><del>Algorithm with pure splits</del></p> <p style="text-align: center;"><del>Best split</del></p> <p style="text-align: center;"><del>Random forest</del></p>	

Figure 8.1

Date \_\_\_\_\_  
Page \_\_\_\_\_

1.  $Gini = 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 1 - 0.64 - 0.04 = 0.32$

2.  $Gini = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.444 = 0.2664$   
 $\Delta Gini = 0.32 - 0.2664 = 0.0536$

3. Growing the 'no' branch with other features.

CUPA:

Weighted Gini =  $\frac{1}{3} \times 0 + \frac{2}{3} \times 0.50 = 0.333$

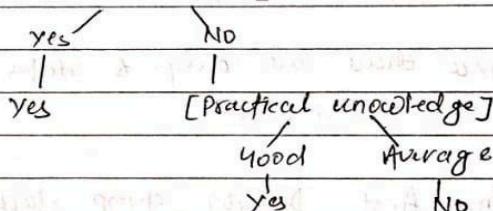
$\Delta Gini = 0.111$

→ Practical knowledge

Unweighted Gini: 0

$Gini = 0.444$  (best)

[Interactivity]



RF	n-estimators	mean-accuracy
1	10	0.9667
2	50	0.9667
3	100	0.9667
4	150	0.9667
5	200	0.9667

The best can be chosen with n-estimator

Best accuracy: 1.000

Best confusion Matrix:  $\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$

Figure 8.2

### Code:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
import matplotlib.pyplot as plt
from sklearn.preprocessing import OrdinalEncoder
data = pd.read_csv("iris (2).csv")
data.head()
oe = OrdinalEncoder()
data[["species"]] = oe.fit_transform(data[["species"]])
data.head()
y = data.iloc[:, -1]
X = data.iloc[:, :-1]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
rf = RandomForestClassifier(n_estimators=10, random_state=42)
rf.fit(X_train, y_train)
y_pred = rf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
accuracy
n_estimators_list = [10, 50, 100, 200, 500, 1000]
accuracies = []

for n in n_estimators_list:
    rf = RandomForestClassifier(n_estimators=n, random_state=42)
    rf.fit(X_train, y_train)
    y_pred = rf.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    accuracies.append(accuracy)
    print(f'Accuracy with n_estimators={n}: {accuracy:.4f}')
plt.plot(n_estimators_list, accuracies, marker='o')
plt.xlabel('Number of Trees (n_estimators)')
plt.ylabel('Accuracy')
plt.title('Random Forest Accuracy vs Number of Trees')
plt.show()
optimal_n_estimators = n_estimators_list[np.argmax(accuracies)]
print(f"Best accuracy is obtained with n_estimators={optimal_n_estimators}")
```

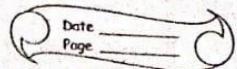
## Program 9

Implement Boosting ensemble method on a given dataset.

### Observation:

LAB-8					
05/05/2024 .					
Considering AdaBoost Algorithm. For the following sample data, show the decision stump calculation steps for the attribute CGPA.					
CGPA	Interactionness	Practical knowledge	Communication skill	Job profile	
$\geq 9$	Yes	Good	Good	Yes	
$< 9$	No	Good	Moderate	Yes	
$\geq 9$	No	Average	Moderate	No	
$< 9$	No	Average	Good	No	
$\geq 9$	Yes	Good	Moderate	Yes	
$\geq 9$	Yes	Good	Moderate	Yes	
$\Rightarrow$ Step 1: Initial weight assigned to each item = $1/6$ Step 2: Iterate for each weak classifier.					
<u>Decision stump for CGPA:</u> <ul style="list-style-type: none"> <li>• Train the Decision Stump <math>H_{CGPA}</math> with a random bootstrap sample from the training dataset <math>T</math>.</li> <li>• Since there are only 6 data instances, use the full training dataset.</li> <li>• The first Decision Stump classifies the instances based on the CGPA attribute.</li> <li>• If <math>CGPA \geq 9</math>, the data instance is predicted to have 'Job offer' as 'Yes' else 'No'.</li> </ul>					

Figure 9.1



2(b)

Compute the weighted error  $\epsilon_{C4PA}$  of H<sub>C4PA</sub> on current training dataset T:

$$\epsilon_t = \sum_{j=1}^N H(d_j) \text{wt}(d_j)$$

where,

$$H(d_j) = 0 \text{ for correct prediction}$$

$$H(d_j) = 1 \text{ for wrong prediction}$$

$$\epsilon_{C4PA} = 2 \times 1/6 = 0.333$$

2(c)

$$\alpha_{C4PA} = \frac{1}{2} \ln(1 - \epsilon_{C4PA})$$

$$= \frac{1}{2} \ln(1 - 0.333)$$

$$\alpha_{C4PA} = 0.347$$

2(d)

$$Z_{C4PA} = \frac{1}{6} \times 4 \times e^{-0.347} + \frac{1}{6} \times 2 \times e^{0.347}$$

$$= 0.9428$$

2(e)

$$\text{wt}(d_j)_{i+1} = \text{wt}(d_j) \text{ C4PA of correct instance} \times e^{-\alpha_{C4PA}}$$

$$= \frac{1}{6} \times e^{-0.347} = 0.1249$$

$$0.9428$$

$$\text{wt}(d_j)_{i+1} = \text{wt}(d_j) \text{ C4PA of incorrect instance} \times e^{\alpha_{C4PA}}$$

$$= \frac{1}{6} \times e^{0.347} = 0.1501$$

$$0.9428$$

Figure 9.2

CCPA	Predicted Job Offer.	Actual Job Offer.	Weight
$\geq 9$	Yes	Yes	0.1249
$< 9$	No	Yes	0.2501
$> 9$	Yes	No	0.2501
$< 9$	No	No	0.1249
$\geq 9$	Yes	Yes	0.1249
$\geq 9$	Yes	Yes	0.1249

Default Adaboost ( $n\text{-estimator} = 10$ )  $\Rightarrow$  Mean 5-fold CV accuracy  
0.8202

Tuning results:

n-estimators	mean-accuracy
10	0.820237
50	0.830023
100	0.832050
150	0.832296
200	0.832624

Best performance: 0.8326 accuracy using 200 Neets.

Run  
5/5/25

Figure 9.3

### Code:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import AdaBoostClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
import matplotlib.pyplot as plt
from sklearn.preprocessing import OrdinalEncoder
data = pd.read_csv("income.csv")
data.head()
y = data.iloc[:, -1]
X = data.iloc[:, :-1]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
rf = AdaBoostClassifier(n_estimators=1000, random_state=42)
rf.fit(X_train, y_train)
y_pred = rf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
accuracy
n_estimators_list = [10, 50, 100, 200, 500, 1000]
accuracies = []

for n in n_estimators_list:
    rf = AdaBoostClassifier(n_estimators=n, random_state=42)
    rf.fit(X_train, y_train)
    y_pred = rf.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    accuracies.append(accuracy)
    print(f"Accuracy with n_estimators={n}: {accuracy:.4f}")

plt.plot(n_estimators_list, accuracies, marker='o')
plt.xlabel('Number of Trees (n_estimators)')
plt.ylabel('Accuracy')
plt.title('Random Forest Accuracy vs Number of Trees')
plt.show()
optimal_n_estimators = n_estimators_list[np.argmax(accuracies)]
print(f"Best accuracy is obtained with n_estimators={optimal_n_estimators}")
```

## Program 10

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

### Observation:

		Lab - 9		Date: 12/05/2025	Date _____ Page _____		
$\rightarrow$ For the given data compute two clusters using k-Means algorithm for clustering where initial cluster centres are $(1.0, 1.0)$ and $(5.0, 7.0)$ . Execute for two iterations.					20.		
$\rightarrow$ No. of clusters ( $K=2$ ) Centroid for cluster 1 ( $C_1$ ) = $(1.0, 1.0)$ Centroid for cluster 2 ( $C_2$ ) = $(5.0, 7.0)$							
Record Number		A	B				
R <sub>1</sub>		1.0	1.0				
R <sub>2</sub>		1.5	2.0				
R <sub>3</sub>		3.0	4.0				
R <sub>4</sub>		5.0	7.0				
R <sub>5</sub>		3.5	5.0				
R <sub>6</sub>		4.5	5.0				
R <sub>7</sub>		3.5	4.5				
Iteration 1:							
Record Number		dist from C <sub>1</sub> (1.0, 1.0) (euclidean dist.)	dist from C <sub>2</sub> (5.0, 7.0)	cluster			
R <sub>1</sub>	(1.0, 1.0)	0.0	7.21	1			
R <sub>2</sub>	(1.5, 2.0)	1.12	6.12	1			
R <sub>3</sub>	(3.0, 4.0)	3.61	3.61	1			
R <sub>4</sub>	(5.0, 7.0)	7.21	0.0	2			
R <sub>5</sub>	(3.5, 5.0)	4.12	2.5	2			
R <sub>6</sub>	(4.5, 5.0)	5.31	2.06	2			
R <sub>7</sub>	(3.5, 4.5)	4.30	2.92	2			
Cluster 1: {R <sub>1</sub> , R <sub>2</sub> , R <sub>3</sub> } Cluster 2: {R <sub>4</sub> , R <sub>5</sub> , R <sub>6</sub> , R <sub>7</sub> }							

Figure 10.1

New centroids are

$$c_1 = (1.0 + 1.5 + 3.0)/3, (1.0 + 2.0 + 4.0)/3 \\ = 5.5/3, 7.0/3 \\ = (1.83, 2.33)$$

$$c_2 = (5.0 + 3.5 + 4.5 + 3.5)/4, (7 + 5 + 5 + 4.5)/4 \\ = 16.5/4, 21.5/4 \\ = (4.12, 5.37)$$

Iteration 2:

	Round Number	dist(c <sub>1</sub> , (1.83, 2.33))	dist(c <sub>2</sub> , (4.12, 5.37))	cluster
R <sub>1</sub>	(1.0, 1.0)	1.57	5.37	1
R <sub>2</sub>	(1.5, 2.0)	0.47	4.27	1
R <sub>3</sub>	(3.0, 4.0)	2.04	1.77	2
R <sub>4</sub>	(5.0, 7.0)	5.64	1.85	2
R <sub>5</sub>	(3.5, 5.0)	3.15	0.72	2
R <sub>6</sub>	(4.5, 5.0)	3.78	0.53	2
R <sub>7</sub>	(3.5, 4.5)	2.74	1.07	2

cluster 1: {R<sub>1</sub>, R<sub>2</sub>}

cluster 2: {R<sub>3</sub>, R<sub>4</sub>, R<sub>5</sub>, R<sub>6</sub>, R<sub>7</sub>}

New centroids are:

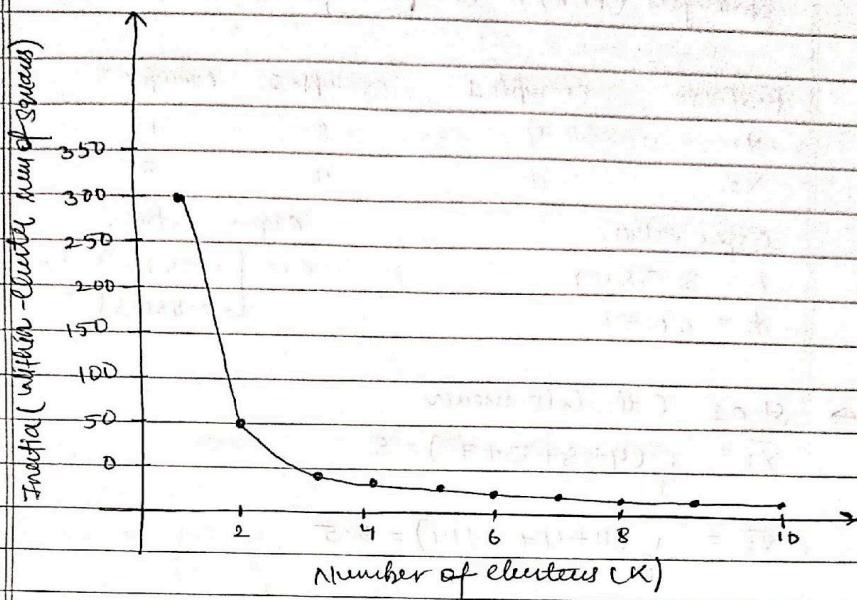
$$c_1 = (1.0 + 1.5)/2, (1.0 + 2.0)/2 \\ = 2.5/2, 3.0/2 \\ = 1.25, 1.5$$

$$c_2 = (3.0 + 5.0 + 3.5 + 4.5 + 3.5)/5, (7 + 5 + 5 + 4.5)/5 \\ = 19.5/5, 22.5/5 \\ = 3.9, 5.1$$

Figure 10.2

Date \_\_\_\_\_  
Page \_\_\_\_\_

Elbow plot:-



Optimal Number of Clusters (K): 3

## **==Code:**

```
# Import required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

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

# Use only petal_length and petal_width
X = df[["petal_length", "petal_width"]]

# Scale the features (helps with KMeans)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Elbow method to determine 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 the elbow curve
plt.figure(figsize=(8, 5))
plt.plot(k_range, inertia, marker='o')
plt.title("Elbow Method for Optimal k")
plt.xlabel("Number of clusters (k)")
plt.ylabel("Inertia (Within-Cluster Sum of Squares)")
plt.grid(True)
plt.show()

# Find optimal k using "elbow" (visually)
optimal_k = 3 # for IRIS, elbow is usually at 3
print(f"Optimal number of clusters (k): {optimal_k}")
```

## Program 11

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

### Observation:

Lab-10				
Date: 12/05/2025:-				
→	Given the data table, reduce the dimension from 2 to 1 using the principal component Analysis (PCA). Compute for first principal component.			
feature	example 1	example 2	example 3	example 4
$x_1$	4	8	13	7
$x_2$	11	4	5	14
Eigen values:		Eigen Vectors:		
$\lambda_1 = 30.3849$		$e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$	$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$	
$\lambda_2 = 6.6151$				
Step 1: Calculate mean				
$\bar{x}_1 = \frac{1}{4}(4+8+13+7) = 8$				
$\bar{x}_2 = \frac{1}{4}(11+4+5+14) = 8.5$				
Step 2: calculate covariance matrix.				
$S = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix}$				
$\text{cov}(x_1, x_1) = \frac{1}{N-1} \sum_{n=1}^N (x_{1n} - \bar{x}_1)(x_{1n} - \bar{x}_1)$				
$= \frac{1}{3} \left\{ (4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2 \right\}$				
$= 14$				
$\text{cov}(x_1, x_2) = \frac{1}{N-1} \sum_{n=1}^N (x_{1n} - \bar{x}_1)(x_{2n} - \bar{x}_2)$				
$= \frac{1}{3} \left\{ (4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) \right. \right. \\ \left. \left. (7-8)(14-8.5) \right\}$				
$= -11$				

Figure 11.1

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\text{Cov}(X_2, X_1) = \text{Cov}(X_1, X_2)$$

$$= -11$$

$$\begin{aligned}\text{Cov}(X_2, X_2) &= \frac{1}{N-1} \sum_{k=1}^N (X_{2k} - \bar{X}_2)(X_{2k} - \bar{X}_2) \\ &= \frac{1}{3} \left\{ (11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2 \right\} \\ &= 23\end{aligned}$$

$$S = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Step 3: Eigen values of the covariance matrix.

$$\begin{aligned}0 &= \det(S - \lambda I) \\ &= \begin{vmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{vmatrix} \\ &= (14-\lambda)(23-\lambda) - (-11) \times (-11) \\ &= \lambda^2 - 37\lambda + 201\end{aligned}$$

$$\begin{aligned}\therefore \lambda &= \frac{1}{2} (37 \pm \sqrt{565}) \\ &= 30.3849, 6.6151 \\ &= \lambda_1, \lambda_2 \text{ (say)}.\end{aligned}$$

Step 4: Computation of the eigen vectors.

Given:  $e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$

$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$

Figure 11.2

steps: computation of first Principal components

$$\mathbf{e}_1^T \begin{bmatrix} x_{1n} - \bar{x}_1 \\ x_{2n} - \bar{x}_2 \end{bmatrix}$$

$$= [0.5574 \quad -0.8303] \begin{bmatrix} x_{11} - \bar{x}_1 \\ x_{21} - \bar{x}_2 \end{bmatrix}$$

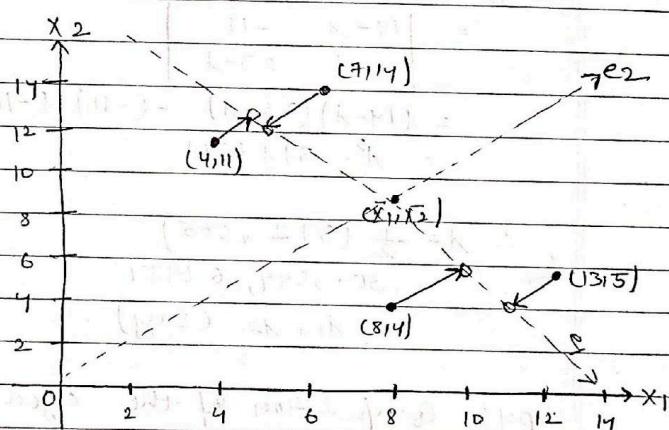
$$= 0.5574(x_{11} - \bar{x}_1) - 0.8303(x_{21} - \bar{x}_2)$$

$$= 0.5574(4 - 8) - 0.8303(11 - 8.5)$$

$$= -4.30535$$

Similarly computing for all other points:

feature	$e_{x1}$	$e_{x2}$	$e_{x3}$	$e_{x4}$
$x_1$	4	8	11	7
$x_2$	11	4	5	14
first principal component	-4.3052	3.7361	5.6928	-5.1238



Model Accuracy Comparison (Before vs. After PCA):

Model	Before PCA	After PCA
logistic Regression	0.8533	0.8533
SVM	0.8804	0.8804
Random forest	0.8750	0.8533

Figure 11.3

## Code:

```
# Importing necessary libraries
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split
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 dataset
df = pd.read_csv("heart.csv")

# Separate features and target
X = df.drop("HeartDisease", axis=1)
y = df["HeartDisease"]

# Identify categorical columns
cat_cols = X.select_dtypes(include=['object']).columns.tolist()

# Label Encode binary categorical columns
label_enc = LabelEncoder()
for col in cat_cols:
    if X[col].nunique() == 2:
        X[col] = label_enc.fit_transform(X[col])
        cat_cols.remove(col)

# One-hot encode remaining categorical columns
X = pd.get_dummies(X, columns=cat_cols)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, 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)

# Initialize models
models = {
    "Logistic Regression": LogisticRegression(max_iter=1000),
    "SVM": SVC(),
    "Random Forest": RandomForestClassifier()
}
```

```

# Store accuracy scores
accuracy_before_pca = {}
accuracy_after_pca = {}

# Training and evaluating models before PCA
for name, model in models.items():
    model.fit(X_train_scaled, y_train)
    y_pred = model.predict(X_test_scaled)
    acc = accuracy_score(y_test, y_pred)
    accuracy_before_pca[name] = acc

# Apply PCA
pca = PCA(n_components=0.95) # retain 95% variance
X_train_pca = pca.fit_transform(X_train_scaled)
X_test_pca = pca.transform(X_test_scaled)

# Training and evaluating models after PCA
for name, model in models.items():
    model.fit(X_train_pca, y_train)
    y_pred = model.predict(X_test_pca)
    acc = accuracy_score(y_test, y_pred)
    accuracy_after_pca[name] = acc

# Print accuracy comparison
print("Model Accuracy Comparison (Before vs After PCA):")
print(f"{'Model':<20} {'Before PCA':<15} {'After PCA':<15}")
for name in models.keys():
    print(f"{'name':<20} {accuracy_before_pca[name]:<15.4f} {accuracy_after_pca[name]:<15.4f}")

```