

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT on

Machine Learning (23CS6PCMAL)

Submitted by

PRIYANSHU KUMAR (1BM22CS210)

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 **Priyanshu Kumar (1BM22CS210)**, who is a bonafide student of **B.M.S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Lab report has been approved as it satisfies the academic requirements in respect of an Machine Learning (23CS6PCMAL) work prescribed for the said degree.

Lab Faculty Incharge Sneha S Bagalkot 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
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	5
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	10
4	17-3-2025	Build Logistic Regression Model for a given dataset	15
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample	20
6	7-4-2025	Build KNN Classification model for a given dataset	25
7	21-4-2025	Build Support vector machine model for a given dataset	30
8	5-5-2025	Implement Random forest ensemble method on a given dataset	35
9	5-5-2025	Implement Boosting ensemble method on a given dataset	40
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file	45
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method	60

Github Link:

<https://github.com/pkcs210/6thSem-ML-Lab/tree/main>

▼ LAB 1: Using Pandas

Code

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

# Assuming 'data' is a DataFrame with stock data
reliance_data = data["RELIANCE.NS"]

print("Summary Statistics for Reliance Industries:")
print(reliance_data.describe())

# Calculate daily returns
reliance_data["Daily Return"] = reliance_data["Close"].pct_change()

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

# Plot Closing Price
plt.subplot(2, 1, 1)
reliance_data["Close"].plot(title="Reliance Industries Closing Price")

# Plot Daily Return
plt.subplot(2, 1, 2)
reliance_data["Daily Return"].plot(
    title="Reliance Industries Daily Return", color="orange"
)

plt.tight_layout()
plt.show()

import pandas as pd
import yfinance as yf
```

```

import matplotlib.pyplot as plt

tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
data = yf.download(
    tickers,
    start="2021-01-01",
    end="2021-12-30",
    group_by="ticker"
)

for ticker in tickers:
    try:
        bank_data = data[ticker]
        bank_data["Daily Return"] = bank_data["Close"].pct_change()

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

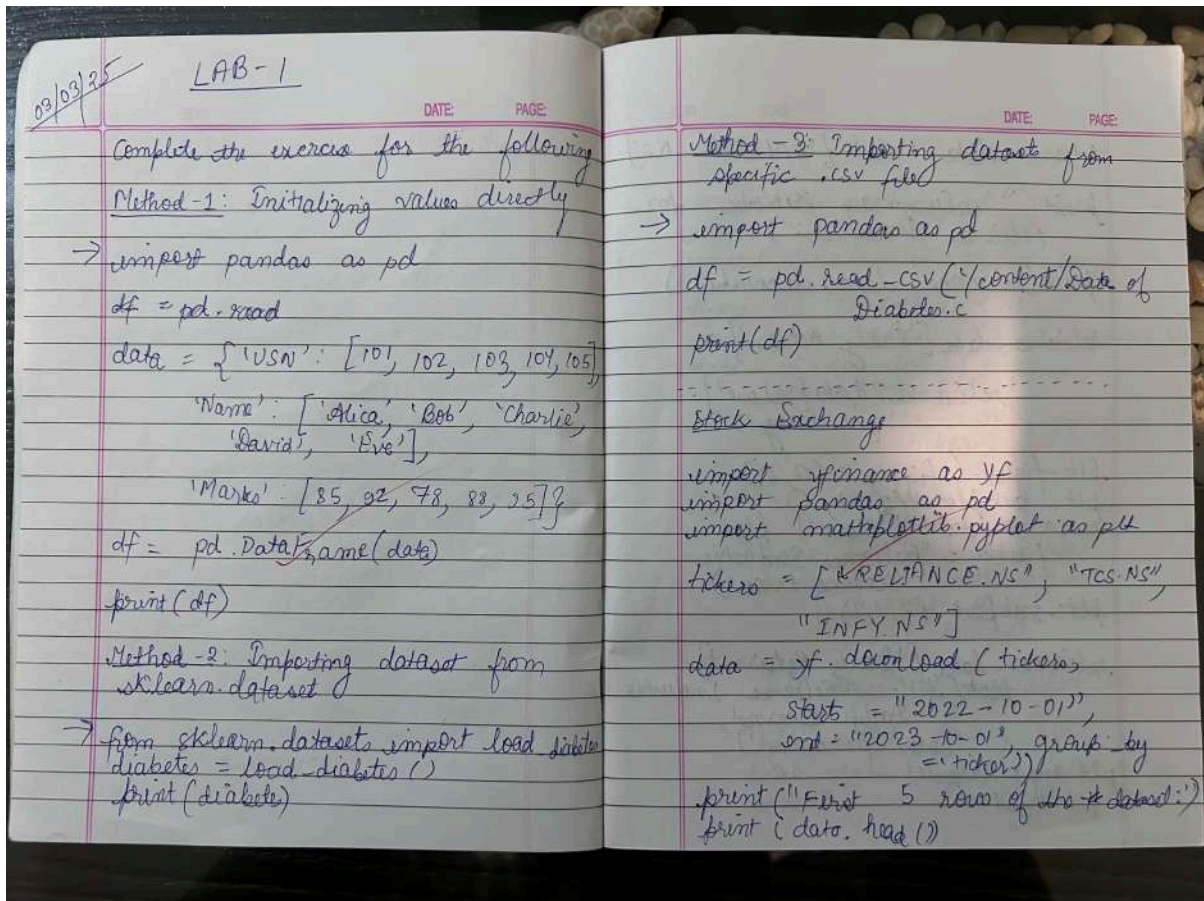
        # Plot Closing Price
        plt.subplot(2, 1, 1)
        bank_data["Close"].plot(title=f"{ticker} - Closing Price")

        # Plot Daily Return
        plt.subplot(2, 1, 2)
        bank_data["Daily Return"].plot(
            title=f"{ticker} - Daily Return", color="orange"
        )

        plt.tight_layout()
        plt.show()
    except KeyError:
        print(f"Data not found for {ticker}")

```

Record Book



DATE PAGE
reliance_data = data['RELIANCE.NS']

print("\nSummary statistics for
Reliance Industries")

print(reliance_data.describe())

reliance_data['Daily Return'] =
reliance_data['close'].
pct_change()

plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['close'].plot
(title = "Reliance Industries
- Closing Prices")

plt.subplot(2, 1, 2)

reliance_data['Daily Return'].
plot(title = "Reliance Industries
- Daily Returns",
color = 'orange')

plt.tight_layout()
plt.show()

DATE PAGE
TODO:

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

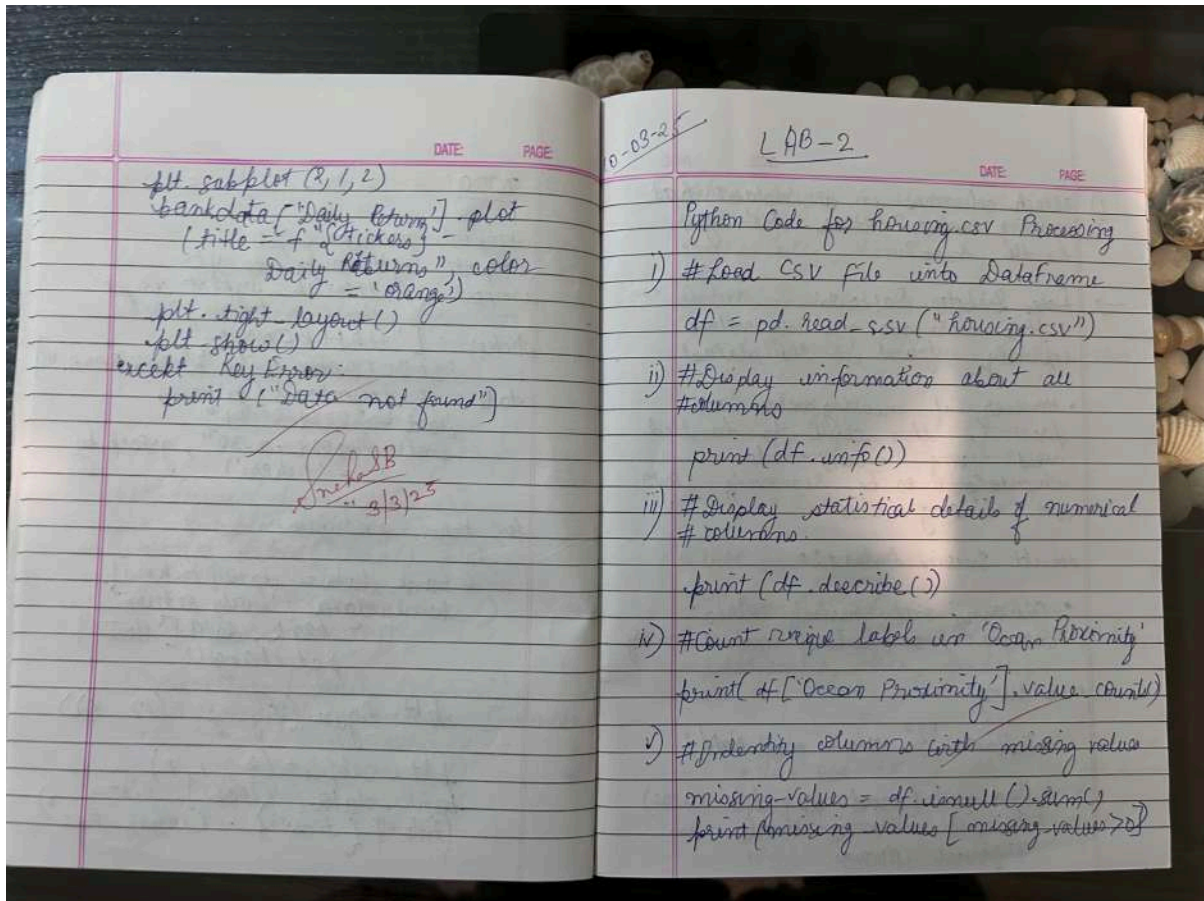
tickers = ["HDFEBANK.NS",
"ICICIBANK.NS", "KOTAKBANK.NS"]
data = yf.download(tickers,
start = "2024-01-01",
end = "2024-12-30", group_by
= 'tickers')

for ticker in tickers:

bank_data = data[ticker]
bank_data['Daily Return']
= bank_data['close'].
pct_change()

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

plt.subplot(2, 1, 1)
bank_data['close'].plot
(title = f"{ticker} - Closing Price")



▼ LAB 2: Data Preprocessing

Code

```
import pandas as pd
```

```
# i) Read CSV File into DataFrame
```

```
df = pd.read_csv("housing.csv")
```

```
# ii) Display information about all columns
```

```
print(df.info())
```

```
# iii) Display statistical details of numerical columns
```

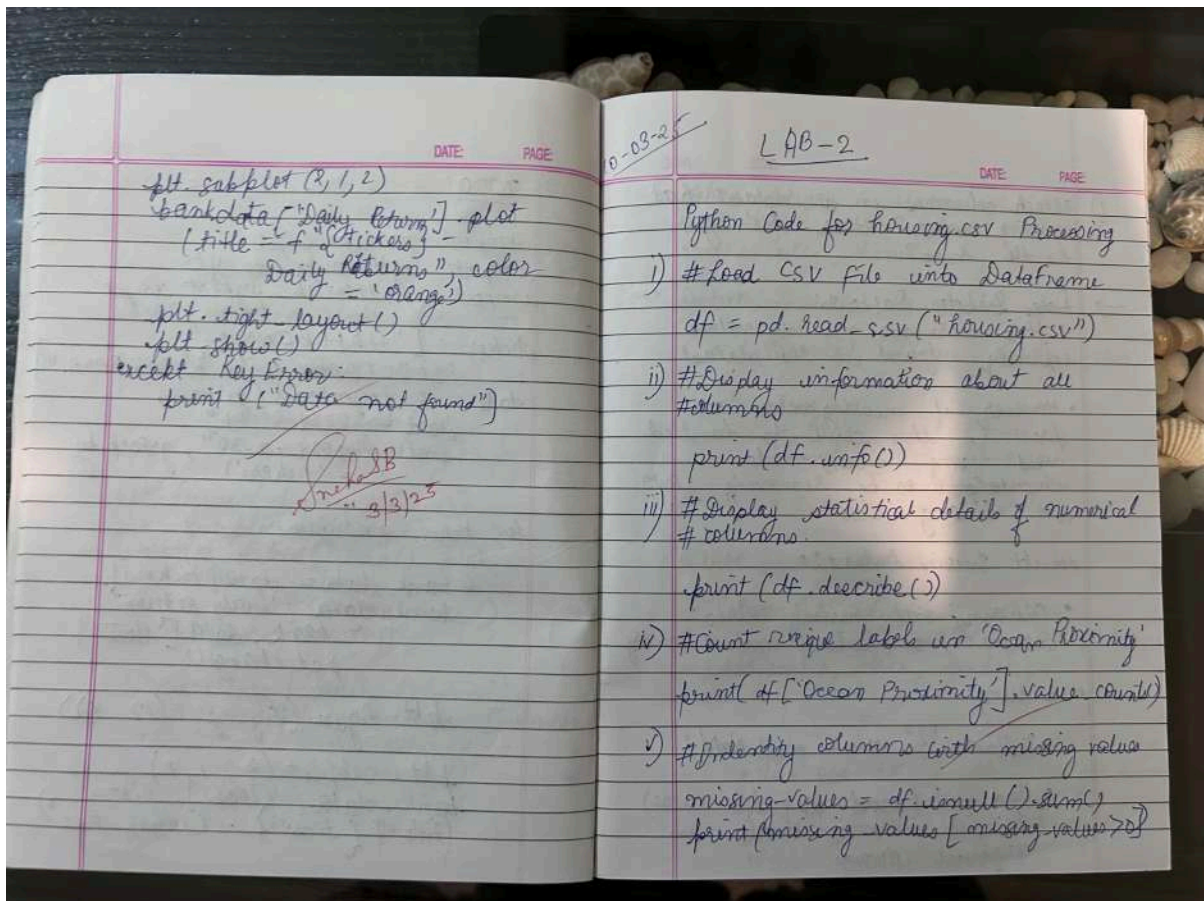
```
print(df.describe())
```

```
# iv) Count unique labels in 'Ocean Proximity'
```

```
print(df["Ocean Proximity"].value_counts())
```

```
# v) Find columns with missing values
missing_values = df.isnull().sum()
print("Missing values:\n", missing_values[missing_values > 0])
```

Record Book



DATE PAGE
1) Which columns in the dataset had missing values? How did you handle them?

→ Diabetes Dataset:

- There were no missing values explicitly found in the dataset.

- However, if missing values were present, they could be handled using mean, mode imputation or by removing rows with missing values.

Adult Income Dataset:

- Columns with missing values: workclass, occupation, native-country

- Handling:

- Replace "?" values with NaN.

- Used mode (most frequent value) to fill missing entries in categorical columns.

DATE PAGE
2) Which categorical columns did you identify in the dataset? How did you encode them?

→ Diabetes Dataset:

- Categorical columns: Gender, CLASS

- Encoding:

- Gender: M → 1, F → 0

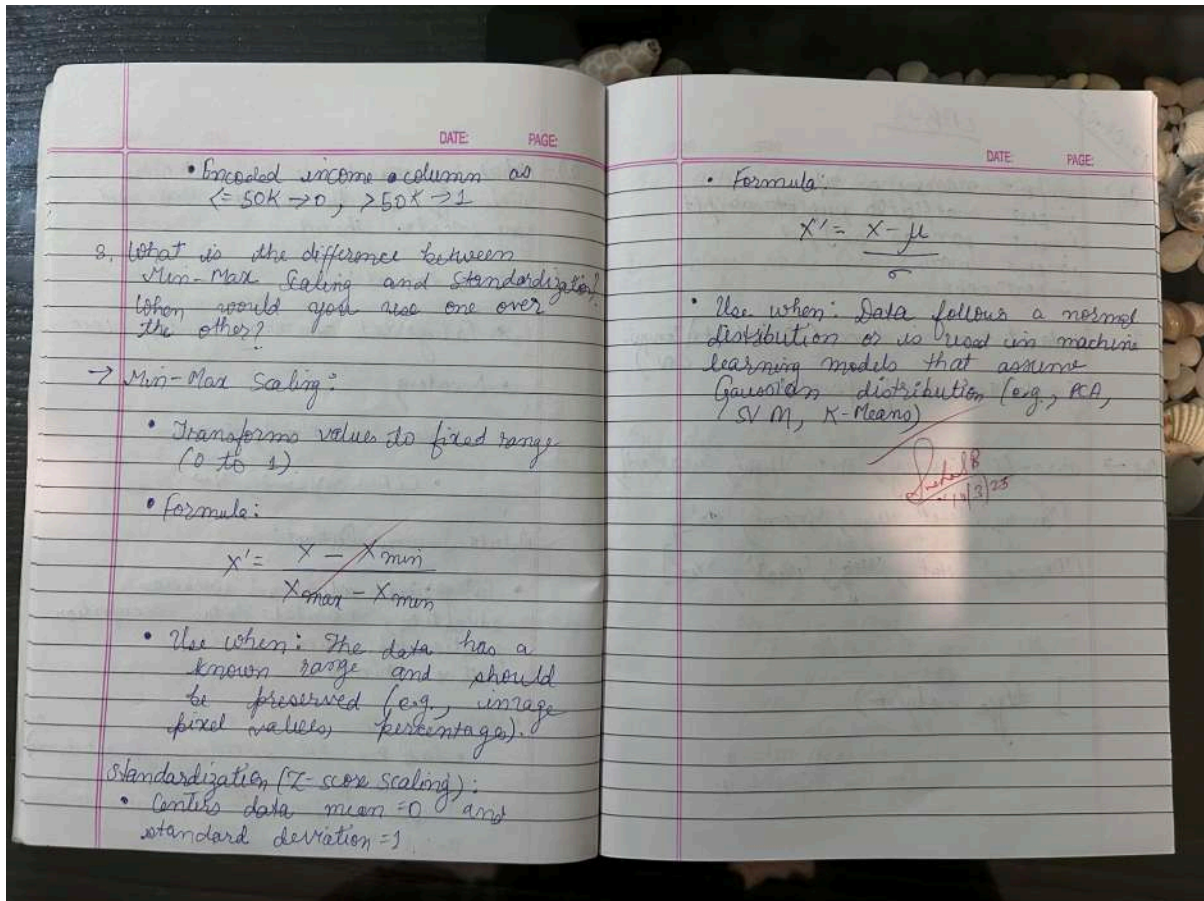
- CLASS: Y → 1, N → 0

Adult Income Dataset:

- Categorical columns: workclass, education, marital-status, occupation, relationship, race, gender, native-country

- Encoding:

- Used one-hot encoding (pd.get_dummies) for categorical columns, dropping the first category to avoid multicollinearity.



▼ LAB 3: ID3

Code

```
import math

def find_entropy(x, y):
    if x != 0 and y != 0:
        entropy = -1 * (x * math.log2(x) + y * math.log2(y))
        return entropy
    if x == 1:
        return 1
    if y == 1:
        return 0
    return 0

def find_max_gain(data, rows, columns):
```

```

max_gain = 0
retidx = -1
entropy_ans = find_entropy(data, rows)
if entropy_ans == 0:
    return max_gain, retidx, entropy_ans

for j in columns:
    mydict = {}
    ddd = 1
    for i in rows:
        key = data[i][j]
        if key not in mydict:
            mydict[key] = 1
        else:
            mydict[key] += 1

    gain = entropy_ans
    for key in mydict:
        yes = 0
        no = 0
        for k in rows:
            if data[k][j] == key:
                if data[k][-1] == 1: # Assuming last column is the label
                    yes += 1
                else:
                    no += 1
        x = yes / (yes + no)
        y = no / (yes + no)
        if x != 0 and y != 0:
            gain += (mydict[key] * (x * math.log2(x) + y * math.log2(y))) / 14 # '

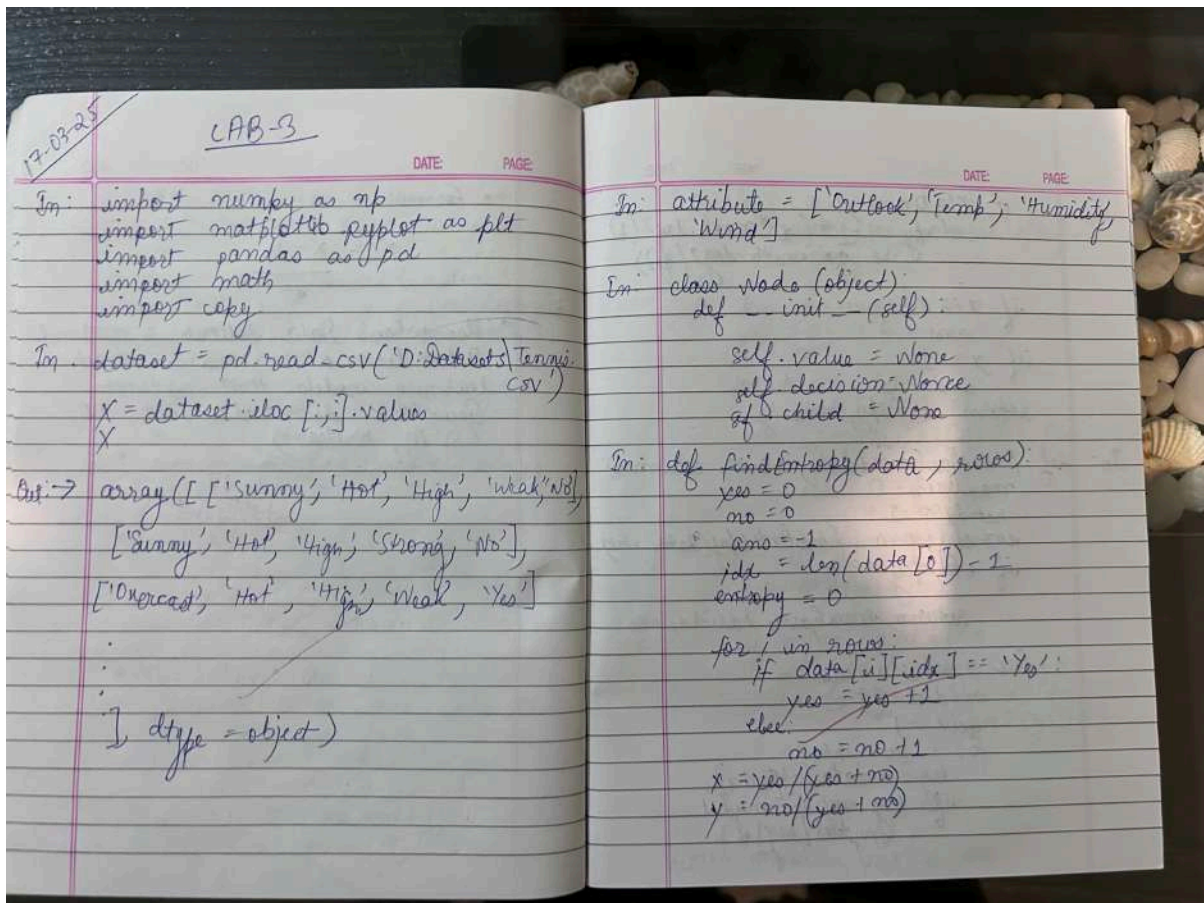
    if gain > max_gain:
        max_gain = gain
        retidx = j

return max_gain, retidx, entropy_ans

```

```
def build_tree(X, rows, columns):
    max_gain, best_col, entropy = find_max_gain(X, rows, columns)
    if best_col == -1 or entropy == 0:
        return
```

Record Book



DATE: PAGE:

```

if x!=0 and y!=0:
    entropy = -1*(x*math.log2(x)
               + y*math.log2(y))

if x==1:
    ans = 1
if y==1:
    ans = 0
return entropy, ans

In: def find_max_gain(data, rows, columns)
    max_gain = 0
    retidx = -1
    entropy, ans = find_entropy(data, rows)
    if entropy == 0:
        return max_gain, retidx, ans

    for j in columns:
        mydict = {}
        total = 0
        for i in rows:
            key = data[i][j]
            if key not in mydict:
                mydict[key] = 1

```

DATE: PAGE:

```

else:
    mydict[key] = mydict[key] + 1
    gain = entropy

for key in mydict:
    yes = 0
    no = 0
    for k in rows:
        if data[k][j] == key:
            if data[k][1] == 'yes':
                yes = yes + 1
            else:
                no = no + 1
    x = yes/(yes + no)
    y = no/(yes + no)

    if x!=0 and y!=0:
        gain = (mydict[key]*
                (x*math.log2(x) +
                 y*math.log2(y)))/4

    if gain > max_gain:
        max_gain = gain
        retidx = j

return max_gain, retidx, ans

```


DATE: PAGE:
In: def buildTree(data, rows, columns):

maxGain, idx, ans = findMaxGain
(X, rows, columns)
root = Node()
root.child = []

if maxGain == 0:
if ans == 1:
root.value == 'Yes'
else:
root.value == 'No'

return root
root.value = attribute[idx]
mydict = {}

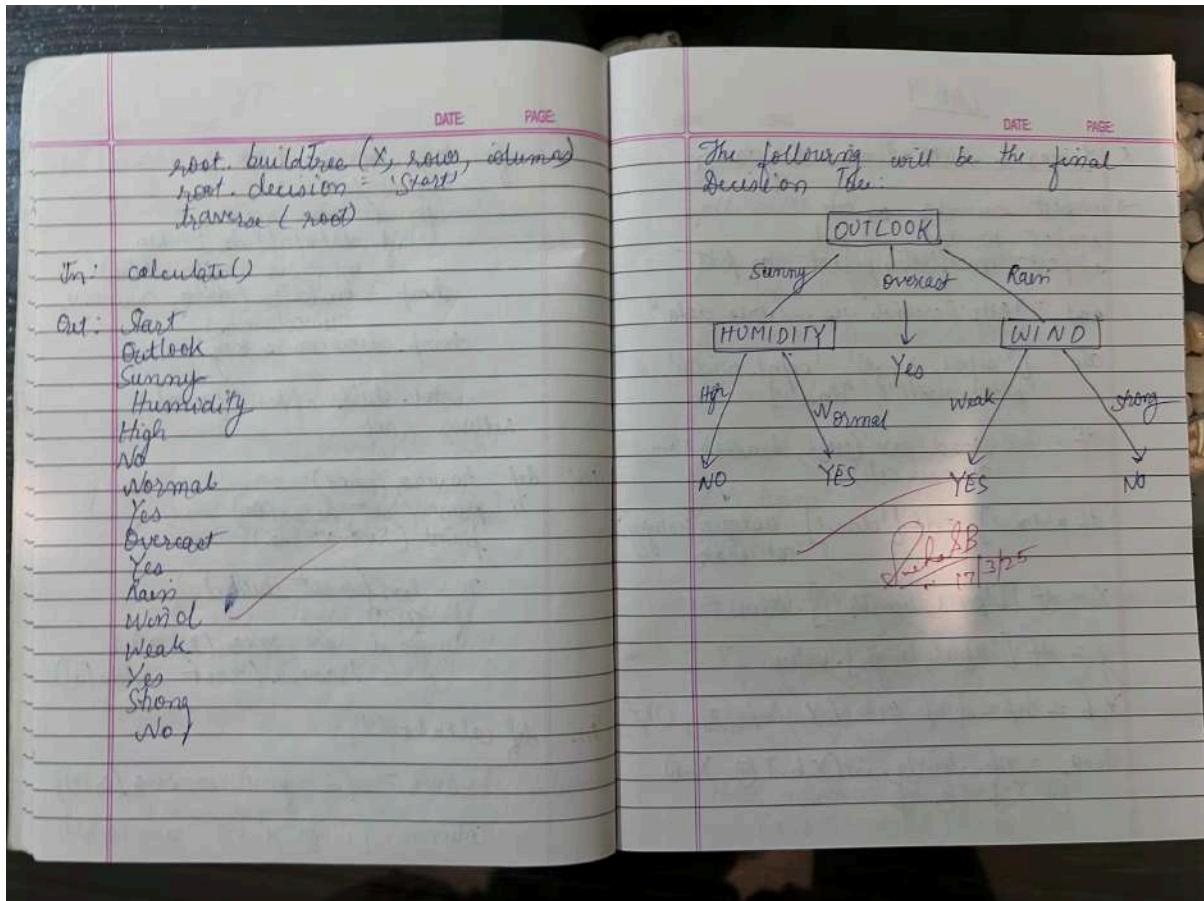
for j in rows:
key = data[j][idx]
if key not in mydict:
mydict[key] = 1
else:
mydict[key] += 1

newcolumns = copy.deepcopy(columns)
newcolumns.remove(idx)

DATE: PAGE:
for key in mydict:
newrows = []
for i in rows:
if data[i][idx] == key:
newrows.append(i)
temp = buildTree(data, newrows,
newcolumns)
temp.decision = key
root.childs.append(temp)
return root

In: def traverse(root):
print(root.decision)
print(root.value)
n = len(root.childs)
if n > 0:
for i in range(0, n):
traverse(root.childs[i])

In: def calculate():
rows = [i for i in range(0, 14)]
columns = [i for i in range(0, 4)]



▼ LAB 4: Linear Regression, Multiple Regression, Logical Regression

Code

```

import numpy as np
import pandas as pd

# Sample data
data = {
    "Temp": [30, 28, 25, 27, 32],
    "Humidity": [70, 65, 80, 85, 75],
    "Windspeed": [5, 10, 12, 8, 7],
    "Rain": [1, 0, 1, 0, 1]
}
  
```

```

df = pd.DataFrame(data)

# Features and target
X = df[["Temp", "Humidity", "Windspeed"]].values
y = df["Rain"].values

# Feature scaling (Standardization)
X = (X - np.mean(X, axis=0)) / np.std(X, axis=0)

# Add bias (intercept) term
X = np.c_[np.ones(X.shape[0]), X]

# Sigmoid function
def sigmoid(z):
    return 1 / (1 + np.exp(-z))

# Cost function
def compute_cost(X, y, weights):
    m = len(y)
    h = sigmoid(np.dot(X, weights))
    cost = -(1/m) * np.sum(y * np.log(h) + (1 - y) * np.log(1 - h))
    return cost

# Gradient Descent
def gradient_descent(X, y, weights, alpha, iterations):
    m = len(y)
    cost_history = []
    for i in range(iterations):
        h = sigmoid(np.dot(X, weights))
        gradient = np.dot(X.T, (h - y)) / m
        weights -= alpha * gradient
        cost_history.append(compute_cost(X, y, weights))
    return weights, cost_history

# Initialize weights, learning rate, and iterations

```

```

weights = np.zeros(X.shape[1])
alpha = 0.1
iterations = 1000

# Train model
weights, cost_history = gradient_descent(X, y, weights, alpha, iterations)
print("Final parameters:", weights)

# Prediction function
def predict(X, weights):
    return sigmoid(np.dot(X, weights)) >= 0.5

# Make predictions
predictions = predict(X, weights).astype(int)
print("Predicted labels:", predictions)

# Calculate accuracy
accuracy = np.mean(predictions == y)
print("Accuracy:", accuracy * 100, "%")

import numpy as np
import pandas as pd

# Example data
data = {
    "Area": [2104, 1600, 2400, 1416, 3000],
    "Bedrooms": [3, 3, 3, 2, 4],
    "Price": [399900, 329900, 369000, 232000, 539900]
}
df = pd.DataFrame(data)

X = df[["Area", "Bedrooms"]].values
y = df["Price"].values

# Feature scaling

```

```

X = (X - np.mean(X, axis=0)) / np.std(X, axis=0)
X = np.c_[np.ones(X.shape[0]), X] # Add bias

theta = np.zeros(X.shape[1])
alpha = 0.01
iterations = 1000

def compute_cost(X, y, theta):
    m = len(y)
    predictions = X.dot(theta)
    cost = (1/(2*m)) * np.sum((predictions - y) ** 2)
    return cost

def gradient_descent(X, y, theta, alpha, iterations):
    m = len(y)
    cost_history = []
    for i in range(iterations):
        predictions = X.dot(theta)
        gradient = (1/m) * X.T.dot(predictions - y)
        theta -= alpha * gradient
        cost_history.append(compute_cost(X, y, theta))
    return theta, cost_history

theta, cost_history = gradient_descent(X, y, theta, alpha, iterations)
print("Final parameters:", theta)

import numpy as np
import pandas as pd

# Sample data
data = {
    "Temp": [30, 28, 25, 27, 32],
    "Humidity": [70, 65, 80, 85, 75],
    "Windspeed": [5, 10, 12, 8, 7],

```

```

    "Rain": [1, 0, 1, 0, 1]
}

df = pd.DataFrame(data)

# Features and target
X = df[["Temp", "Humidity", "Windspeed"]].values
y = df["Rain"].values

# Feature scaling (Standardization)
X = (X - np.mean(X, axis=0)) / np.std(X, axis=0)

# Add bias (intercept) term
X = np.c_[np.ones(X.shape[0]), X]

# Sigmoid function
def sigmoid(z):
    return 1 / (1 + np.exp(-z))

# Cost function
def compute_cost(X, y, weights):
    m = len(y)
    h = sigmoid(np.dot(X, weights))
    cost = -(1/m) * np.sum(y * np.log(h) + (1 - y) * np.log(1 - h))
    return cost

# Gradient Descent
def gradient_descent(X, y, weights, alpha, iterations):
    m = len(y)
    cost_history = []
    for i in range(iterations):
        h = sigmoid(np.dot(X, weights))
        gradient = np.dot(X.T, (h - y)) / m
        weights -= alpha * gradient
        cost_history.append(compute_cost(X, y, weights))
    return weights, cost_history

```

```
# Initialize weights, learning rate, and iterations
weights = np.zeros(X.shape[1])
alpha = 0.1
iterations = 1000

# Train model
weights, cost_history = gradient_descent(X, y, weights, alpha, iterations)
print("Final parameters:", weights)

# Prediction function
def predict(X, weights):
    return sigmoid(np.dot(X, weights)) >= 0.5

# Make predictions
predictions = predict(X, weights).astype(int)
print("Predicted labels:", predictions)

# Calculate accuracy
accuracy = np.mean(predictions == y)
print("Accuracy:", accuracy * 100, "%")
```

Record Book

LAB-4

DATE

PAGE

Linear & Multi-linear Regression

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

url = "http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
cols = ['sepal-length', 'sepal-width', 'petal-width', 'class']

df = pd.read_csv(url, header=None, names=cols)

df['class'] = df['class'].astype('category').cat.codes

X = df[['petal-length']].values
y = df['sepal-length'].values

X_b = np.c_[np.ones((X.shape[0], 1))], X

theta = np.linalg.inv(X_b.T @ X_b) @ X_b.T @ y
```

DATE

PAGE

$$y_{pred} = X_b @ \theta$$

```
plt.scatter(X, y, label='Actual')

plt.plot(X, y_pred, color='red', label='Predicted')

plt.title("Linear Regression")
plt.xlabel("Petal Length")
plt.ylabel("Sepal Length")
plt.legend()
plt.show()
```

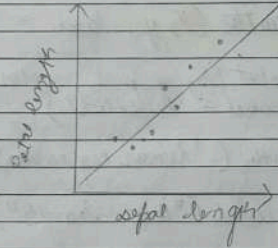
Multi-linear: Use petal length + petal width

```
X_multi = df[['petal-length', 'petal-width']].values
X_b = np.c_[np.ones((X_multi.shape[0], 1))], X_multi

theta_multi = np.linalg.inv(X_b.T @ X_b) @ X_b.T @ y
```


$y_{pred_multi} = X \cdot b @ theta_multi$

print



Multiple linear regression

import numpy as np
import pandas as pd

df = pd.DataFrame(data)

x = df[['Temperature', 'Humidity',
 'Windspeed']].values

y = df[['Rainfall']].values.reshape
 (-1, 1)

def normalize_features(x):
 mean = np.mean(x, axis=0)

def normalize

std = np.std(x, axis=0)

return (x - mean) / std, mean, std

x, mean_x, std_x = normalize_features(x)

X = np.hstack((np.ones((X.shape[0],
 1))), X))

theta = np.zeros(X.shape[1])

alpha = 0.01

iterations = 1000

def hypothesis(x, theta):

 return np.dot(X, theta)

def compute_cost(X, y, theta):

 m = len(y)

 predictions = hypothesis(X, theta)

 cost = (1/2 * m) * np.sum
 (predictions - y)**2

 return cost

DATE PAGE:

```

def gradient_descent(X, y, theta,
                    alpha, iterations):
    m = len(y)
    cost_history = []
    for i in range(iterations):
        gradient = (1/m) * np.dot(X.T,
        (hypothesis(X, theta) - y))
        theta = theta - alpha * gradient
        cost_history.append(compute_cost(X, y, theta))
    if i % 100 == 0:
        print(f"Iteration {i}:")
        cost_history.append(i)
    return theta, cost_history

theta, cost_history = gradient_descent(X, y,
theta, alpha, iterations)

print("Final parameters:", print(theta))

```

DATE PAGE:

```

def:
Final parameters
[ 0.60584029]
[-0.67652179]
[ 0.68426521]
[-0.42303921]

Logistic Regression

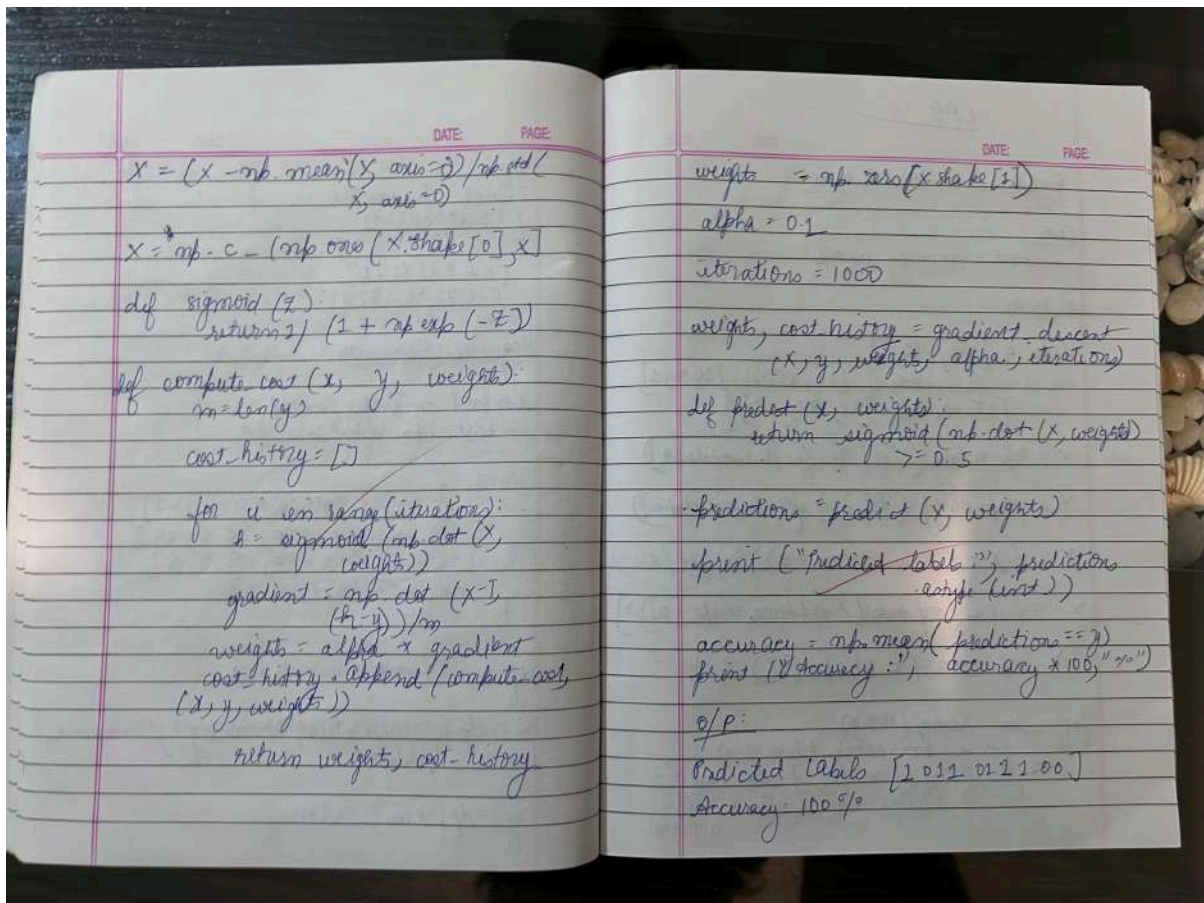
import numpy as np
import pandas as pd

data = {
    "Temp": [30, 20.5, 27, 32],
    "Humidity": [70, 65, 80, 85],
    "Windspeed": [5, 10, 12, 8],
    "Rain": [1, 0, 1, 1]
}

df = pd.DataFrame(data)

X = df[["Temp", "Humidity", "Windspeed"]].values
y = df[["Rain"]].values

```



▼ LAB 5: SVM & KNN

Code

```
import numpy as np

# Data: [Temp, Humidity, Label]
data = [
    [30, 80, 1], [25, 70, 1], [27, 65, 1], [20, 90, 1],
    [25, 40, 0], [35, 30, 0]
]

def distance(a, b):
    return sum((x - y) ** 2 for x, y in zip(a, b)) ** 0.5

def KNN_predict(test_point, k=3):
    # Compute distances
```

```

    dists = sorted(data, key=lambda row: distance(row[:2], test_point))
    # Get labels of k nearest neighbors
    labels = [row[2] for row in dists[:k]]
    # Majority vote
    return 1 if labels.count(1) > labels.count(0) else 0

test_point = [26, 60]
prediction = KNN_predict(test_point)
weather = "Rain" if prediction == 1 else "No Rain"
print(f"Test input: Temperature = {test_point[0]}°C, Humidity = {test_point[1]}")
print(f"Predicted weather: {weather}")

import numpy as np

# Data: [Temp, Humidity, Label]
data = [
    [30, 80, 1], [25, 70, 1], [27, 65, 1], [20, 90, 1],
    [25, 40, 0], [35, 30, 0]
]

def distance(a, b):
    return sum((x - y) ** 2 for x, y in zip(a, b)) ** 0.5

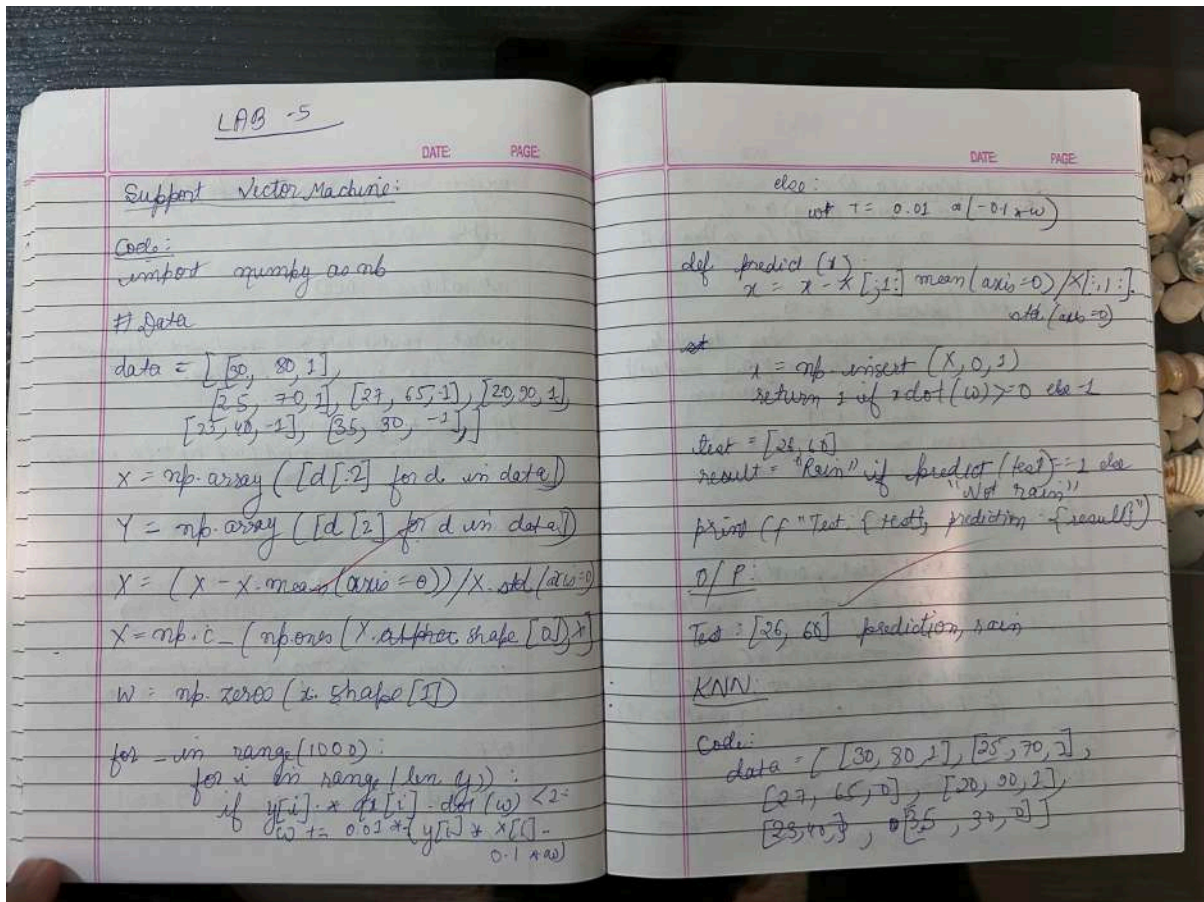
def KNN_predict(test_point, k=3):
    # Compute distances
    dists = sorted(data, key=lambda row: distance(row[:2], test_point))
    # Get labels of k nearest neighbors
    labels = [row[2] for row in dists[:k]]
    # Majority vote
    return 1 if labels.count(1) > labels.count(0) else 0

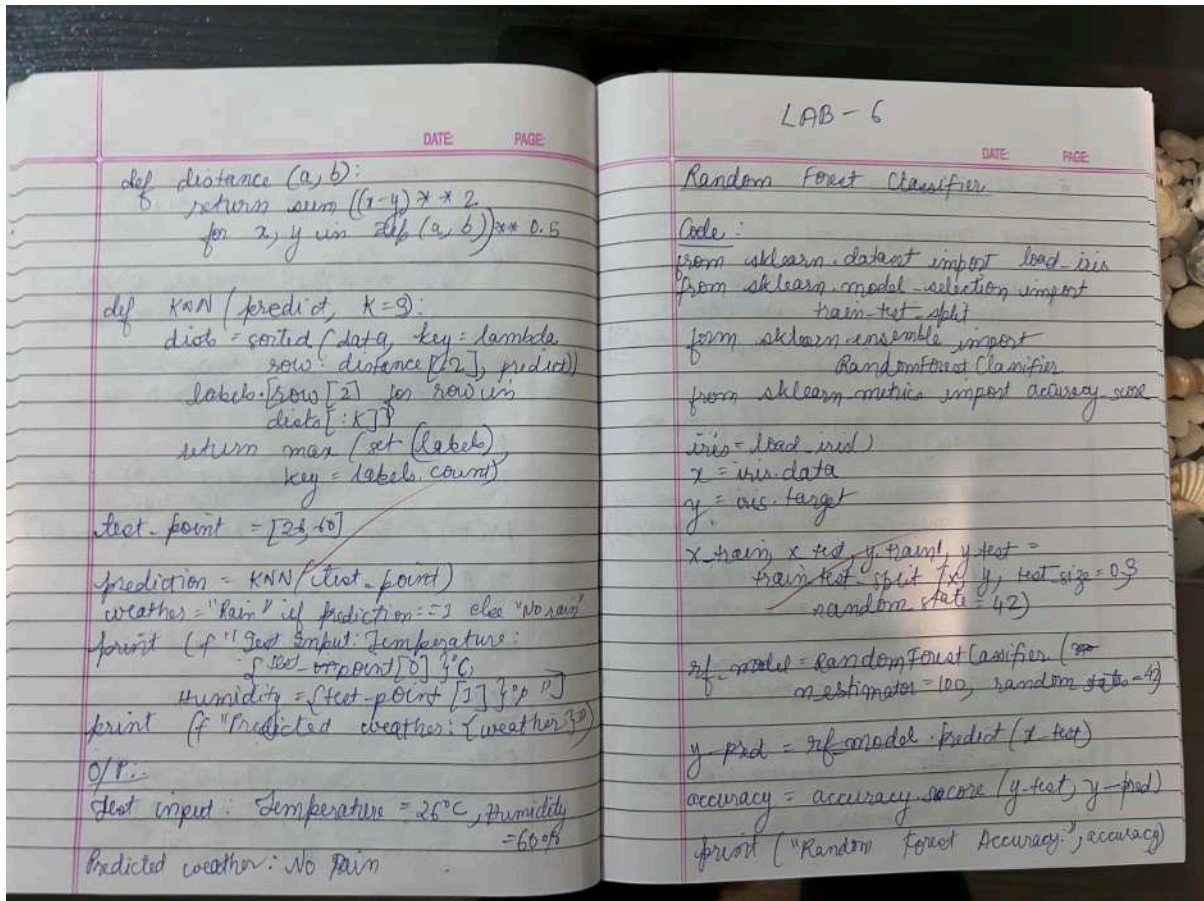
test_point = [26, 60]
prediction = KNN_predict(test_point)
weather = "Rain" if prediction == 1 else "No Rain"
print(f"Test input: Temperature = {test_point[0]}°C, Humidity = {test_point[1]}")

```

```
print(f"Predicted weather: {weather}")
```

Record Book





▼ LAB 6: Random Forest, AdaBoost, K-Means Clustering

Code

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
```

```
# Load data
iris = load_iris()
X = iris.data
y = iris.target
```

```
# Split data
```

```

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42
)

# Train Random Forest
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

# Predict
y_pred = rf_model.predict(X_test)

# Accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Random Forest Accuracy:", accuracy)

import numpy as np

class DecisionStump:
    def __init__(self):
        self.feature_index = None
        self.threshold = None
        self.polarity = 1

    def predict(self, X):
        n_samples = X.shape[0]
        predictions = np.ones(n_samples)
        if self.polarity == 1:
            predictions[X[:, self.feature_index] < self.threshold] = -1
        else:
            predictions[X[:, self.feature_index] > self.threshold] = -1
        return predictions

def adaboost(X, y, n_clf=5):
    n_samples, n_features = X.shape
    w = np.full(n_samples, (1 / n_samples))
    models = []

```

```

alphas = []

for _ in range(n_clf):
    clf = DecisionStump()
    min_error = float('inf')

    # Find best stump
    for feature_i in range(n_features):
        feature_values = np.unique(X[:, feature_i])
        for threshold in feature_values:
            for polarity in [1, -1]:
                predictions = np.ones(n_samples)
                if polarity == 1:
                    predictions[X[:, feature_i] < threshold] = -1
                else:
                    predictions[X[:, feature_i] > threshold] = -1
                error = np.sum(w[y != predictions])
                if error < min_error:
                    min_error = error
                    clf.polarity = polarity
                    clf.threshold = threshold
                    clf.feature_index = feature_i

    # Compute alpha
    EPS = 1e-10
    alpha = 0.5 * np.log((1 - min_error + EPS) / (min_error + EPS))
    predictions = clf.predict(X)
    w *= np.exp(-alpha * y * predictions)
    w /= np.sum(w)

    models.append(clf)
    alphas.append(alpha)

return models, alphas

def predict(X, models, alphas):

```

```
clf_preds = [alpha * clf.predict(X) for clf, alpha in zip(models, alphas)]
y_pred = np.sign(np.sum(clf_preds, axis=0))
return y_pred
```

Example usage:

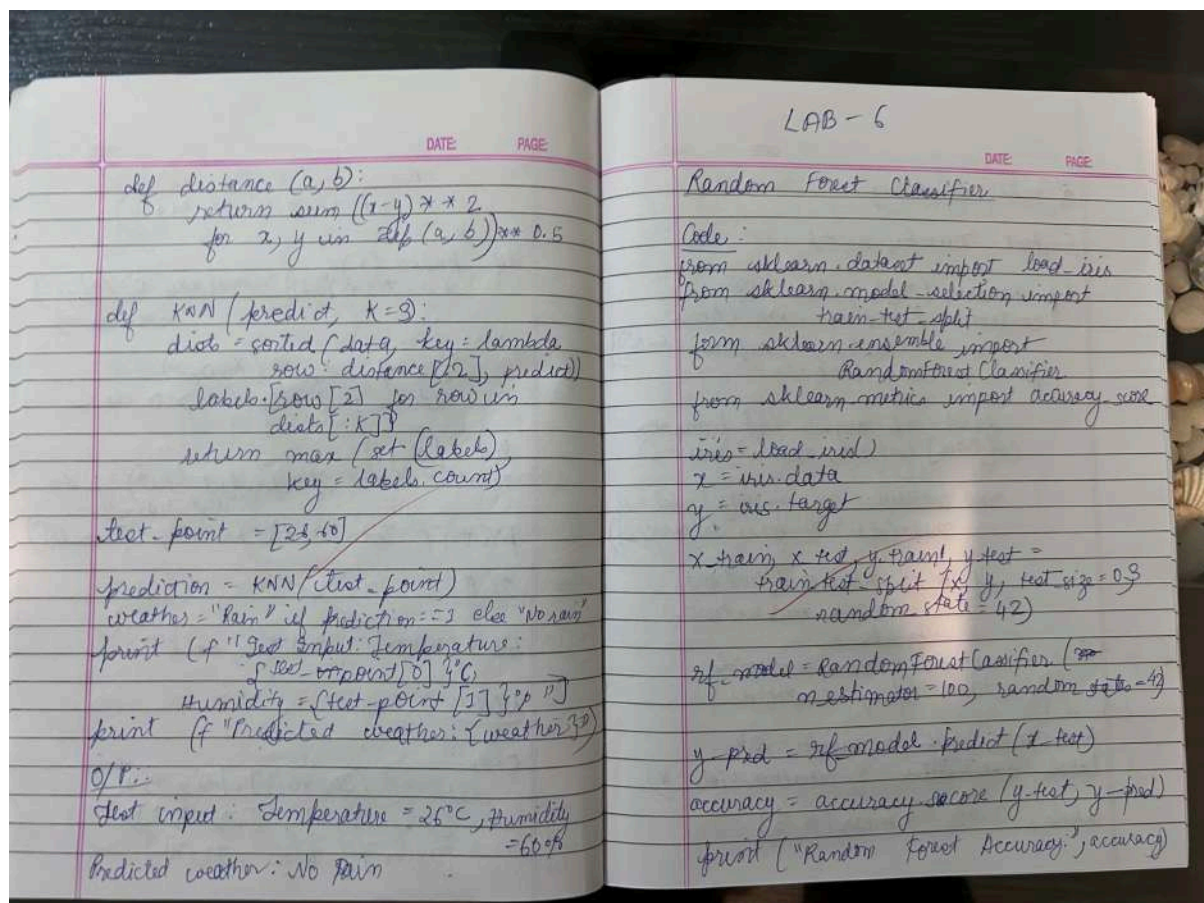
X = np.array([[...], ...]) # shape (n_samples, n_features)

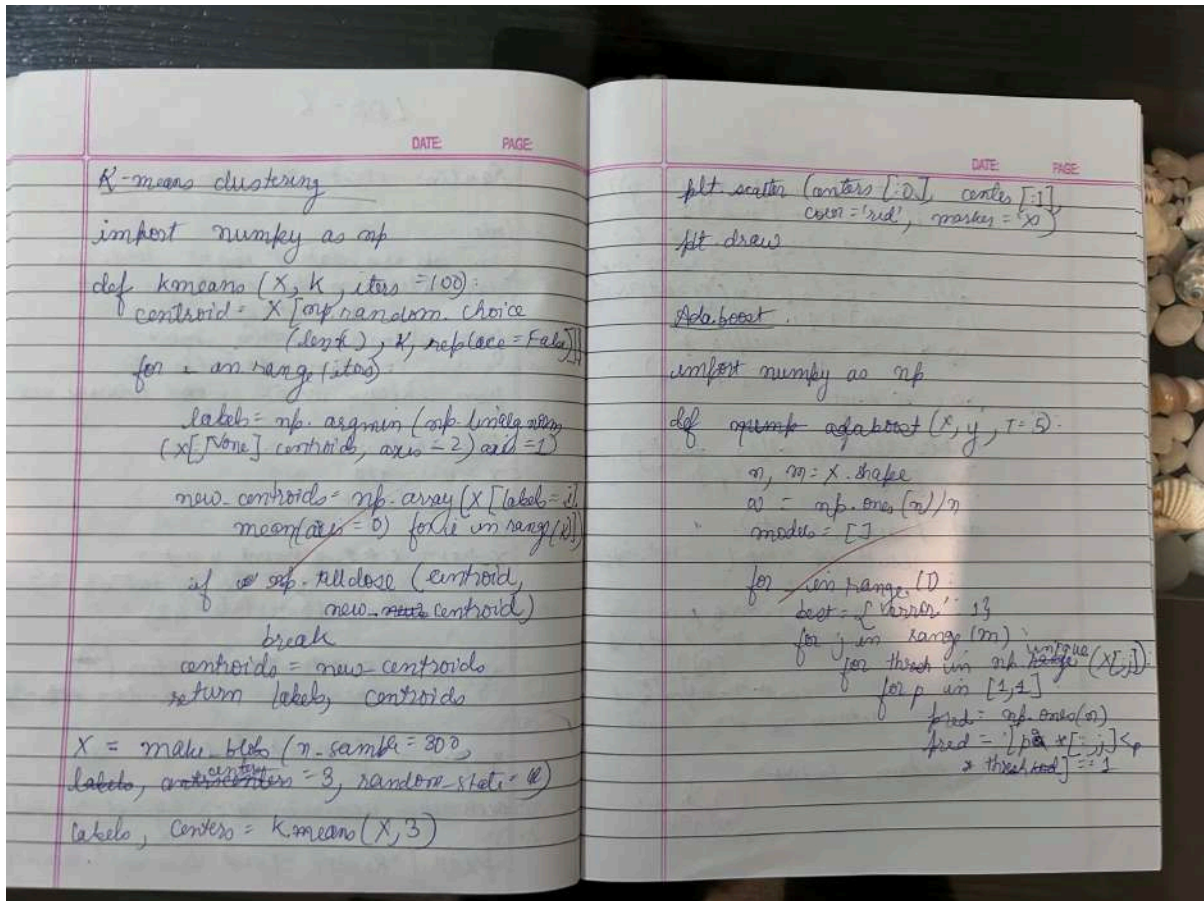
y = np.array([...]) # shape (n_samples,), labels must be -1 or 1

models, alphas = adaboost(X, y, n_clf=5)

y_pred = predict(X, models, alphas)

Record Book





DATE

PAGE

```

err = np.sum((y - pred) ** 2)
if err < best['error']:
    best.update({'error': err, 'theta': theta, 'p': p, 'pred': pred})
    alpha = 0.5 * np.log((1 - best['error']) / best['error'] + 1) * 10
    w = np.exp(-alpha * best['pred'])
    w /= w.sum()

```

```

model.append([alpha, best['p'], best['theta'], best['pred']])

```

```

def predict(x_test):
    result = np.zeros(x_test.shape)
    for alpha, j, theta, p in model:
        pred = np.exp(x_test * theta)
        result += alpha * pred
    return np.sign(result)

```

```

return predict

```

19/5/25

DATE

PAGE