

**VISVESVARAYA TECHNOLOGICAL  
UNIVERSITY**

“JnanaSangama”, Belgaum -590014, Karnataka.



**LAB REPORT**

**on**

**Machine Learning (23CS6PCMAL)**

*Submitted by*

**Setu Mishra (1BM22CS250)**

*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 **Setu Mishra (1BM22CS250)**, 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 a Machine Learning (23CS6PCMAL) work prescribed for the said degree.

<b>Lab Faculty Incharge</b>  Name: <b>Shravya AR</b>  Assistant Professor Department of CSE, BMSCE	<b>Dr. Kavitha Sooda</b> Professor & HOD Department of CSE, BMSCE
---	---

## Index

<b>Sl. No.</b>	<b>Date</b>	<b>Experiment Title</b>	<b>Page No.</b>
1	5-3-2025	Write a python program to import and export data using Pandas library functions	4
2	5-3-2025	Demonstrate various data pre-processing techniques for a given dataset	11
3	19-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	15
4	2-4-2025	Build Logistic Regression Model for a given dataset	21
5	12-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample	32
6	2-4-2025	Build KNN Classification model for a given dataset	39
7	9-4-2025	Build Support vector machine model for a given dataset	46
8	7-5-2025	Implement Random Forest ensemble method on a given dataset	48
9	7-5-2025	Implement Boosting ensemble method on a given dataset	52
10	7-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file	56
11	7-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method	60

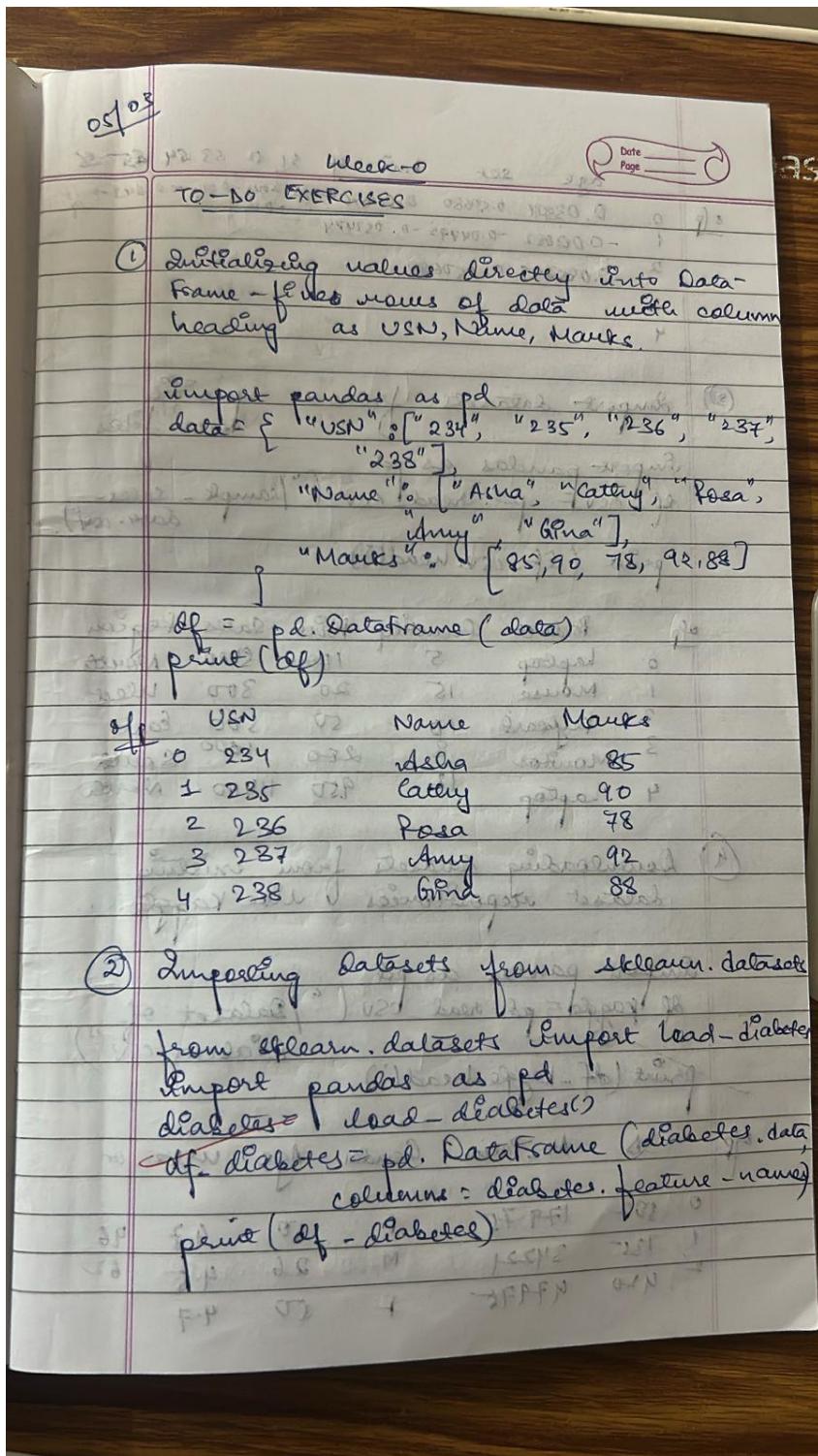
**Github Link:**

<https://github.com/setu-mishra/ML>

## Program 1

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

Screenshot



age sex bmi s1 s2 s3 s4 s5 s6

0 0.03807 0.50680 0.061696 -0.0248 -0.0161 -0.034 -0.043 -0.051 0.017  
1 -0.001882 -0.04472 -0.051474 -0.08 0.09

2 0.085299 0.050680 0.061696 -0.0248 -0.0161 -0.034 -0.043 -0.051 0.017

3 0.085299 0.050680 0.061696 -0.0248 -0.0161 -0.034 -0.043 -0.051 0.017

4 0.085299 0.050680 0.061696 -0.0248 -0.0161 -0.034 -0.043 -0.051 0.017

② Import datasets from a specific

Import pandas as pd

df = pd.read\_csv("sample-sales-data.csv")

print(df.head(1))

	Product	Quantity	Price	Sales	Region
0	Laptop	5	1000	5000	North
1	Mouse	15	20	300	West
2	Keyboard	10	50	500	East
3	Monitor	8	200	1600	South
4	Laptop	12	950	11400	North

④ Downloading datasets from existing  
dataset repositories like Kaggle...

Import pandas as pd

df\_kaggle = pd.read\_csv("Dataset of  
Diabetes.csv")

print(df\_kaggle.head(1))

⑤ No\_Patno Gender Age urea cr

0 502 179 F 50 4.7 46

1 735 34221 M 26 4.5 62

2 420 47975 F 50 4.7 46

## BMI CLASS

1 24.0 N  
 2 23.0 N  
 3 24.0 N  
 4 21.0 N

(no exercise required,  
 leg up exercise required,  
 leg up to shoulder height required)

"20. XMAS 121" "24. 2434" 27.354

["MATHS"]

"(maths) measured in = m/s  
 - 3m/s "10s = m/s  
 "and group "5s = m/s"  
 ("possible and if always 2 diff")  
 ((@) math. std 3) living

("except for speed n") living

(maths. std) living

("common words n") living

(common std) living

["20. XMAS 211"] std = standard

## TO-DO

### Stock Market Data Analysis

1. ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
  2. start date and end date
  3. Plot the closing price and daily returns for all three banks
- Import yfinance as yf  
Import pandas as pd  
Import matplotlib.pyplot as plt
- tickers = ["HDFC.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]

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

```
print("First 5 rows of the dataset")
```

```
print(data.head(5))
```

```
print("In shape of dataset")
```

```
print(data.shape)
```

```
print("In column names")
```

```
print(data.columns)
```

```
hdfc_data = data["HDFCBANK.NS"]
```

print ("In Summary statistics for  
HDFC Industries")

print (hdfc\_data.describe())

hdfc\_data['Daily Return'] = hdfc\_data

[['close'].pct\_change(1)]

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

plt.subplot(2, 1, 1)

hdfc\_data['Daily Return'] = plt

(+title="HDFC Industries Daily  
Returns", color='orange')

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

plt.show()

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

0.1014  
0.2831  
0.3222  
0.2021

Code:

```
import pandas as pd

data={
    'USN':['1BM22CS001','1BM22CS002','1BM22CS003','1BM22CS004','1BM22CS005'],
    'Name':['Ankita','Anita','Amit','Anish','Arun'],
    'Marks':[99,56,96,85,45]
}

df=pd.DataFrame(data)
print(df)

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

path=r"/content/sample_sales_data.csv"
df=pd.read_csv(path)
print(df)

path=r"/content/Dataset of Diabetes .csv"
df=pd.read_csv(path)
print(df.head())

import yfinance as yf
import matplotlib.pyplot as plt

tickers=['HDFCBANK.NS','ICICIBANK.NS','KOTAKBANK.NS']

data=yf.download(tickers,start="2024-01-01",end="2024-12-30",group_by=tickers)
print(data)

#HDFCBANK
HDFC=data['HDFCBANK.NS']
HDFC['Daily Return']=HDFC['Close'].pct_change()
print(HDFC)

plt.figure(figsize=(12,6))
plt.subplot(2,1,1)
HDFC['Close'].plot(title='HDFC BANK - Closing Price')
plt.subplot(2,1,2)
HDFC['Daily Return'].plot(title='HDFC BANK - Daily Return',color='orange')
plt.tight_layout()
```

```
plt.show()

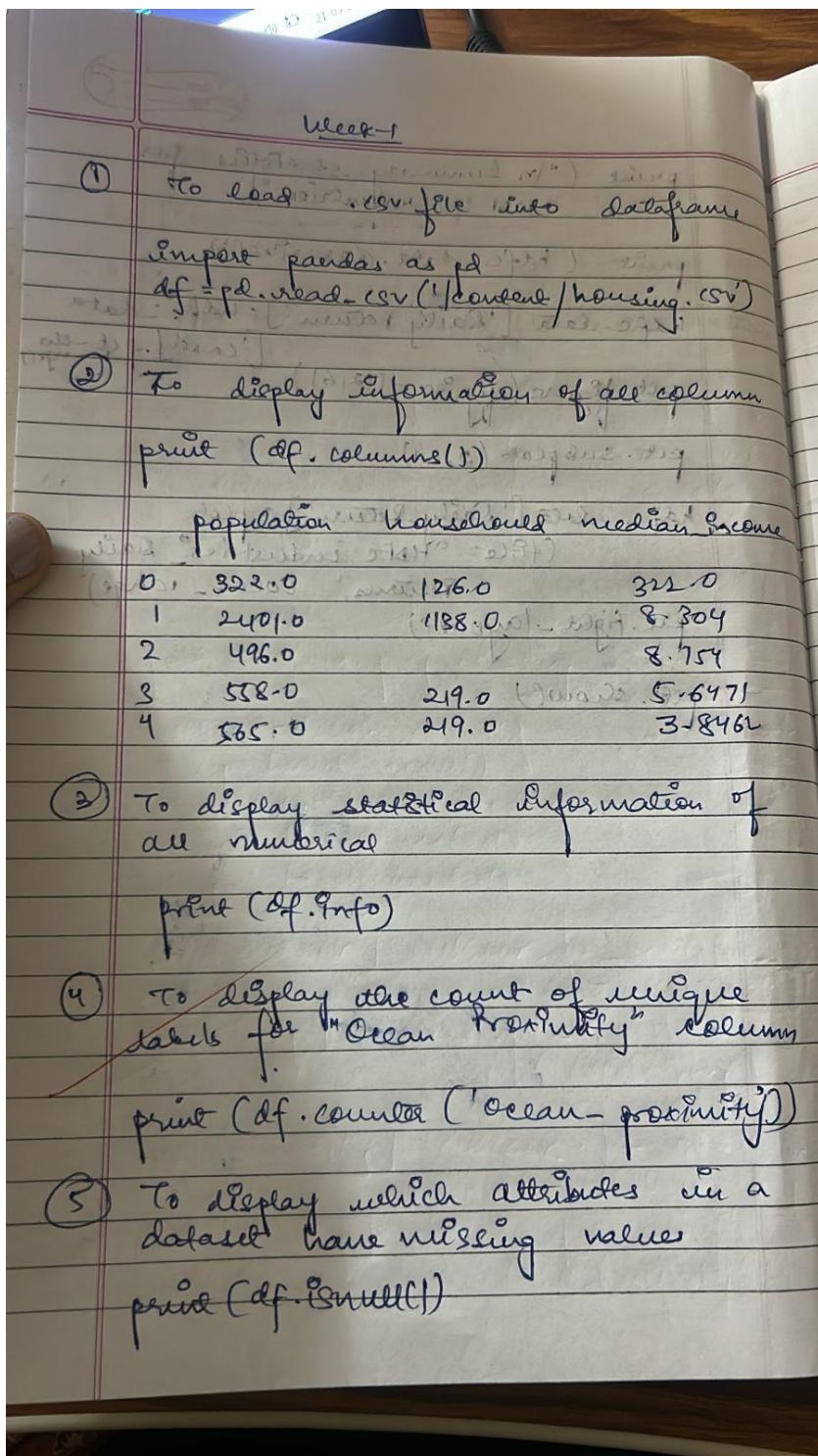
#ICICIBANK
ICICI=data['ICICIBANK.NS']
ICICI['Daily Return']=ICICI['Close'].pct_change()
print(ICICI)

plt.figure(figsize=(12,6))
plt.subplot(2,1,1)
ICICI['Close'].plot(title='ICICI BANK - Closing Price')
plt.subplot(2,1,2)
ICICI['Daily Return'].plot(title='ICICI BANK - Daily Return',color='orange')
plt.tight_layout()
plt.show()
```

## Program 2

Demonstrate various data pre-processing techniques for a given dataset

Screenshot



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

off	ocean proximity		
	NEAR BAY	NEAR BAY	INLAND
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			
26			
27			
28			
29			
30			
31			
32			
33			
34			
35			
36			
37			
38			
39			
40			
41			
42			
43			
44			
45			
46			
47			
48			
49			
50			
51			
52			
53			
54			
55			
56			
57			
58			
59			
60			
61			
62			
63			
64			
65			
66			
67			
68			
69			
70			
71			
72			
73			
74			
75			
76			
77			
78			
79			
80			
81			
82			
83			
84			
85			
86			
87			
88			
89			
90			
91			
92			
93			
94			
95			
96			
97			
98			
99			
100			
101			
102			
103			
104			
105			
106			
107			
108			
109			
110			
111			
112			
113			
114			
115			
116			
117			
118			
119			
120			
121			
122			
123			
124			
125			
126			
127			
128			
129			
130			
131			
132			
133			
134			
135			
136			
137			
138			
139			
140			
141			
142			
143			
144			
145			
146			
147			
148			
149			
150			
151			
152			
153			
154			
155			
156			
157			
158			
159			
160			
161			
162			
163			
164			
165			
166			
167			
168			
169			
170			
171			
172			
173			
174			
175			
176			
177			
178			
179			
180			
181			
182			
183			
184			
185			
186			
187			
188			
189			
190			
191			
192			
193			
194			
195			
196			
197			
198			
199			
200			
201			
202			
203			
204			
205			
206			
207			
208			
209			
210			
211			
212			
213			
214			
215			
216			
217			
218			
219			
220			
221			
222			
223			
224			
225			
226			
227			
228			
229			
230			
231			
232			
233			
234			
235			
236			
237			
238			
239			
240			
241			
242			
243			
244			
245			
246			
247			
248			
249			
250			
251			
252			
253			
254			
255			
256			
257			
258			
259			
260			
261			
262			
263			
264			
265			
266			
267			
268			
269			
270			
271			
272			
273			
274			
275			
276			
277			
278			
279			
280			
281			
282			
283			
284			
285			
286			
287			
288			
289			
290			
291			
292			
293			
294			
295			
296			
297			
298			
299			
300			
301			
302			
303			
304			
305			
306			
307			
308			
309			
310			
311			
312			
313			
314			
315			
316			
317			
318			
319			
320			
321			
322			
323			
324			
325			
326			
327			
328			
329			
330			
331			
332			
333			
334			
335			
336			
337			
338			
339			
340			
341			
342			
343			
344			
345			
346			
347			
348			
349			
350			
351			
352			
353			
354			
355			
356			
357			
358			
359			
360			
361			
362			
363			
364			
365			
366			
367			
368			
369			
370			
371			
372			
373			
374			
375			
376			
377			
378			
379			
380			
381			
382			
383			
384			
385			
386			
387			
388			
389			
390			
391			
392			
393			
394			
395			
396			
397			
398			
399			
400			
401			
402			
403			
404			
405			
406			
407			
408			
409			
410			
411			
412			
413			
414			
415			
416			
417			
418			
419			
420			
421			
422			
423			
424			
425			
426			
427			
428			
429			
430			
431			
432			
433			
434			
435			
436			
437			
438			
439			
440			
441			
442			
443			
444			
445			
446			
447			
448			
449			
450			
451			
452			
453			
454			
455			
456			
457			
458			
459			
460			
461			
462			
463			
464			

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(r"/content/Dataset of Diabetes .csv")
df.head(10)
df.shape
print(df.info())
print(df.describe())
missing_values = df.isnull().sum()

# Display columns with missing values
print(missing_values[missing_values > 0])

#Set the values to some value (zero, the mean, the median, etc.).
# Step 1: Create an instance of SimpleImputer with the median strategy for Age and mean stratergy for Salary
imputer1 = SimpleImputer(strategy="median")
imputer2 = SimpleImputer(strategy="mean")

df_copy=df

# Step 2: Fit the imputer on the "Age" and "Salary"column
# Note: SimpleImputer expects a 2D array, so we reshape the column
imputer1.fit(df_copy[["AGE"]])
imputer2.fit(df_copy[["BMI"]])

# Step 3: Transform (fill) the missing values in the "Age" and "Salary" column
df_copy["AGE"] = imputer1.transform(df[["AGE"]])
df_copy["BMI"] = imputer2.transform(df[["BMI"]])

# Verify that there are no missing values left
print(df_copy["AGE"].isnull().sum())
print(df_copy["BMI"].isnull().sum())
#Handling Categorical Attributes
#Using Ordinal Encoding for gender COlumn and One-Hot Encoding for City Column

# Initialize OrdinalEncoder
ordinal_encoder = OrdinalEncoder(categories=[["M", "F","f"]])
# Fit and transform the data
df_copy["Gender_Encoded"] = ordinal_encoder.fit_transform(df_copy[["Gender"]])
```

```

# Initialize OneHotEncoder
onehot_encoder = OneHotEncoder()

# Fit and transform the "City" column
encoded_data = onehot_encoder.fit_transform(df[["CLASS"]])

# Convert the sparse matrix to a dense array
encoded_array = encoded_data.toarray()

# Convert to DataFrame for better visualization
encoded_df = pd.DataFrame(encoded_array, columns=onehot_encoder.get_feature_names_out(["CLASS"]))
df_encoded = pd.concat([df_copy, encoded_df], axis=1)

df_encoded.drop("Gender", axis=1, inplace=True)
df_encoded.drop("CLASS", axis=1, inplace=True)

print(df_encoded.head())
normalizer = MinMaxScaler()
df_encoded[['BMI']] = normalizer.fit_transform(df_encoded[['BMI']])
df_encoded.head()
scaler = StandardScaler()
df_encoded[['AGE']] = scaler.fit_transform(df_encoded[['AGE']])
df_encoded.head()
df_encoded_copy1=df_encoded
df_encoded_copy2=df_encoded
df_encoded_copy3=df_encoded

Q1 = df_encoded_copy1['BMI'].quantile(0.25)
Q3 = df_encoded_copy1['BMI'].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
df_encoded_copy1['BMI'] = np.where(df_encoded_copy1['BMI'] > upper_bound, upper_bound,
                                    np.where(df_encoded_copy1['BMI'] < lower_bound, lower_bound, df_encoded_copy1['BMI']))

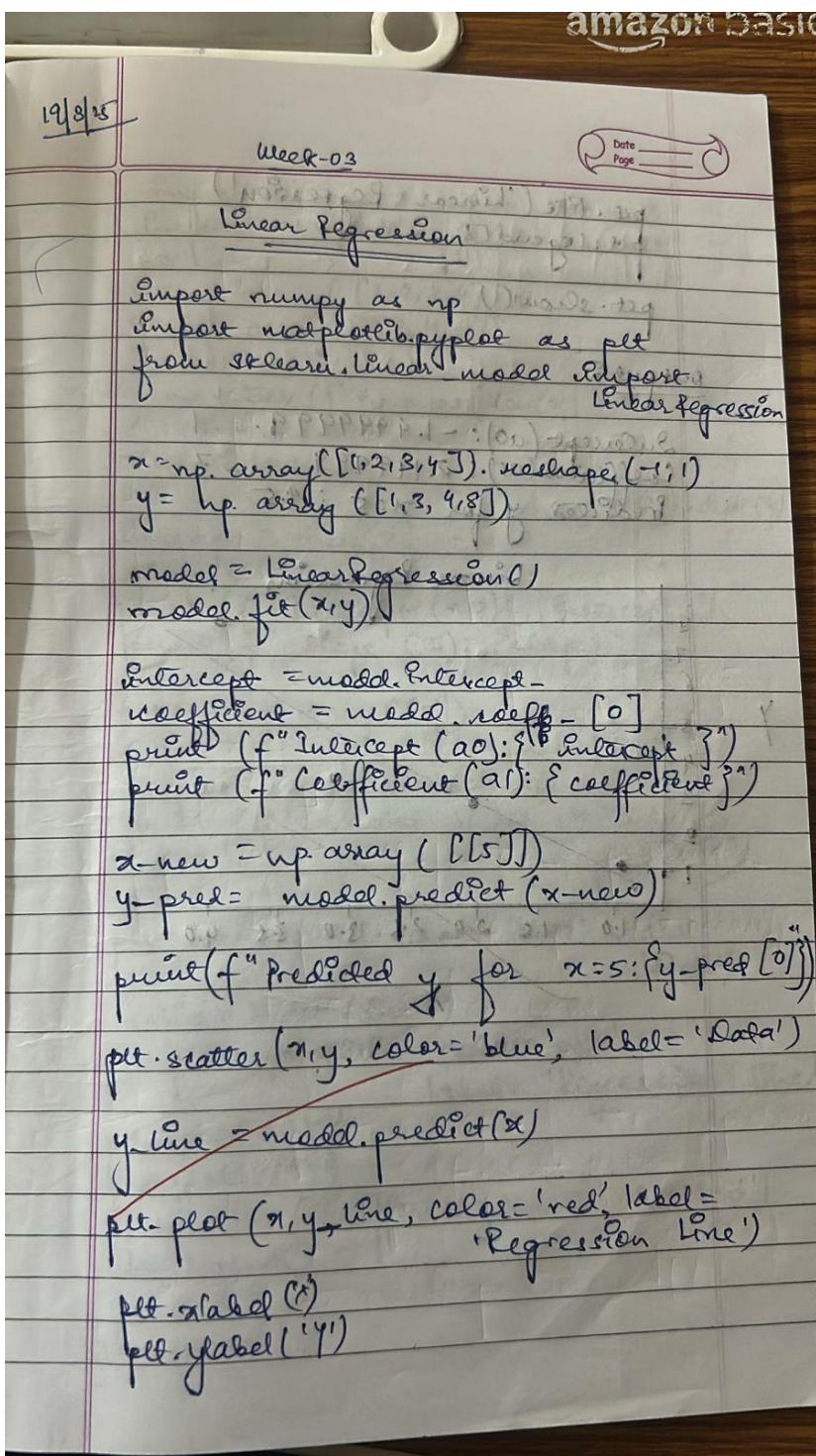
print(df_encoded_copy1.head())
df_encoded_copy2['BMI_zscore'] = stats.zscore(df_encoded_copy2['BMI'])
df_encoded_copy2['BMI'] = np.where(df_encoded_copy2['BMI_zscore'].abs() > 3, np.nan, df_encoded_copy2['BMI'])
# Replace outliers with NaN
print(df_encoded_copy2.head())
df_encoded_copy3['BMI_zscore'] = stats.zscore(df_encoded_copy3['BMI'])
median_salary = df_encoded_copy3['BMI'].median()
df_encoded_copy3['BMI'] = np.where(df_encoded_copy3['BMI_zscore'].abs() > 3, median_salary,
                                   df_encoded_copy3['BMI'])
print(df_encoded_copy3.head())

```

### Program 3

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

Screenshot



pet. fitRe('Linear Regression')

pet. legend()

pet. show()

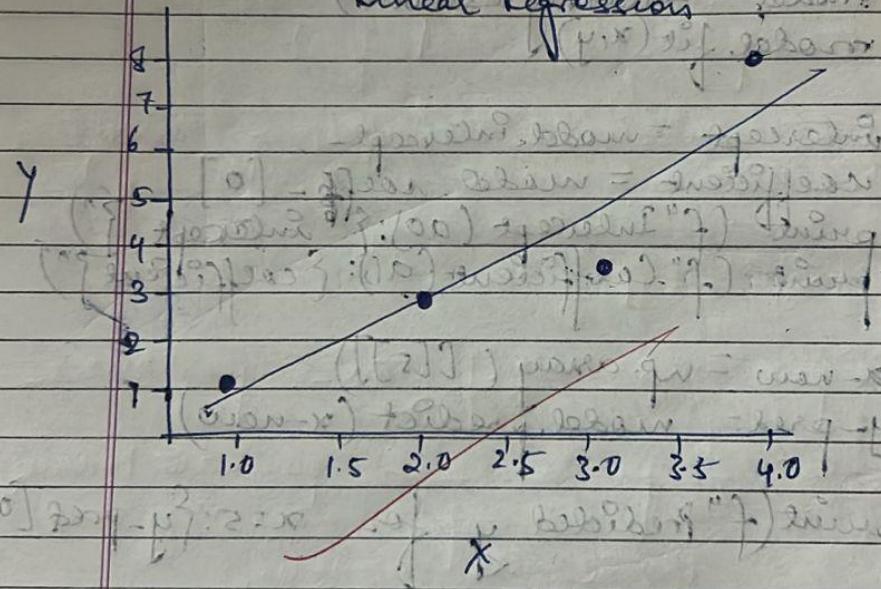
Output:

Intercept (a0): -1.4999999...1

Coefficient(a1): 2.199999...7

Predicted y for x=5: 9.5

Linear Regression = Relation



'slope' = 2.201, 'intercept' = 1.4999999999999998

(x) fibreg. blown = 9.5

'label' = 'y = 2.201x + 1.4999999999999998'

(1) fibreg. blown

(2) fibreg. blown

(3) fibreg. blown

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

```

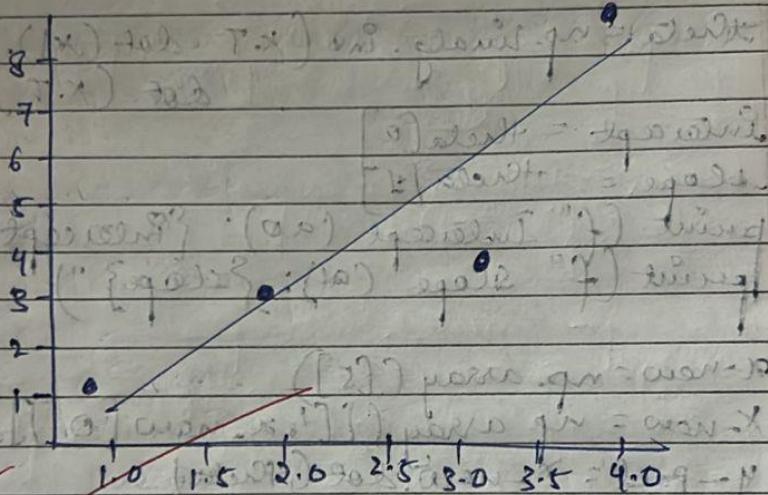
import numpy as np
import tensorflow as tf
x = np.array([1, 2, 3, 4])
y = np.array([1, 3, 4, 8])
x = np.vstack([np.ones(len(x)), x]).T
theta = np.linalg.inv(x.T.dot(x)).dot(y)
intercept = theta[0]
slope = theta[1]
print(f"Intercept (a0): {intercept}")
print(f"Slope (a1): {slope}")

x_new = np.array([5])
x_new = np.array([[1, x_new[0]]])
y_pred = x_new.dot(theta)
print(f"Predicted y for x=5: {y_pred[0]}")

plt.scatter(x, y, color='blue', label='Data')
y_line = x.dot(theta)
plt.plot(x, y_line, color='red', label='Regression Line')
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Different Matrix Method')
plt.legend()
plt.show()
  
```

Intercept ( $a_0$ ): -1.5  
Slope ( $a_1$ ): 2.2010000000000006  
Predicted  $y$  for  $x=5$ : 9.500000000000001

Different Matrix Method



Code:

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

data={"X":[1,2,3,4,5],
      "Y":[1.2,1.8,2.6,3.2,3.8]}
df=pd.DataFrame(data)
df

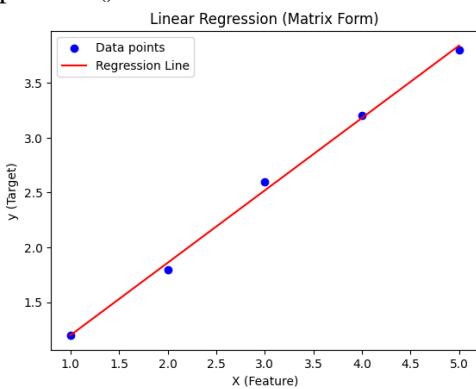
Xi=df["X"].mean()
Yi=df["Y"].mean()

df["Xi^2"]=[Xi**2 for Xi in df["X"]]
Xisq=df["Xi^2"].mean()

xiyi=[]
x=df["X"]
y=df["Y"]
for i in range(len(x)):
    xiyi.append(x[i]*y[i])
df["XiYi"]=xiyi
print(df["XiYi"])
XiYi2=df["XiYi"].mean()
print(XiYi2)
a1 = (df["XiYi"].sum() - len(df) * Xi * Yi) / (df["X"].apply(lambda x: x**2).sum() - len(df) * Xi**2)
a0 = Yi - a1 * Xi

x=9
Y=a0+a1*x
print(Y)

plt.scatter(df["X"], df["Y"], color='blue', label='Data points') # Scatter plot of original data
plt.plot(df["X"], a0 + a1 * df["X"], color='red', label='Regression Line') # Correct regression line
plt.title('Linear Regression (Matrix Form)')
plt.xlabel('X (Feature)')
plt.ylabel('y (Target)')
plt.legend()
plt.show()
```



```

import numpy as np
import matplotlib.pyplot as plt

X = np.array([1, 2, 3, 4])
y = np.array([1, 3, 4, 8])

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

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

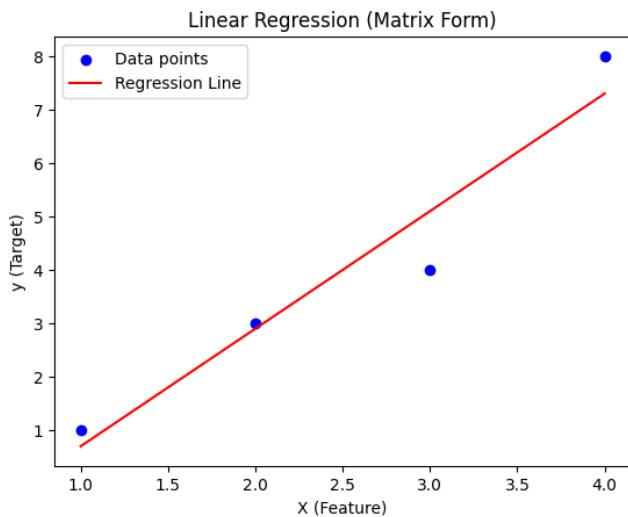
b, m = theta

y_pred = m * X + b
print(f"Slope (m): {m}")
print(f"Intercept (b): {b}")

Slope (m): 2.2000000000000006
Intercept (b): -1.5

plt.scatter(X, y, color='blue', label='Data points')
plt.plot(X, y_pred, color='red', label='Regression Line')
plt.title('Linear Regression (Matrix Form)')
plt.xlabel('X (Feature)')
plt.ylabel('y (Target)')
plt.legend()
plt.show()

```



## Program 4

Build Logistic Regression Model for a given dataset

Screenshot

2/4/25

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

### Logistic Regression

(Insurance dataset)

```
import pandas as pd
from matplotlib import pyplot as plt
df = pd.read_csv('Insurance_data.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

$y_{predicted} = \text{model.predict}(x_{test})$

$y_{predicted}$

$\text{model.score}(x_{test}, y_{test})$

$\text{model.predict_proba}(x_{test})$

$y_{predicted} = \text{model.predict}([[60]])$

$y_{predicted}$

$\text{model.coef}$

$\text{model.intercept}$

Import math

def sigmoid(x):

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

def prediction\_function(age):

$$z = 0.127 * age - 4.973$$

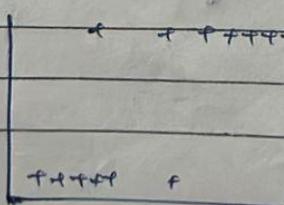
$y = \text{sigmoid}(z)$

return y

age = 35

$\text{prediction\_function}(age)$

0.3709834



To Do:

Q4.4 Consider a binary classification problem where we want to predict whether a student will pass or fail based on their study hours. The logistic regression model has been trained, and the learned parameters are ( $a_0 = -5$  (full intercept) and  $a_1 = 0.8$  (coeff for study hours))

a) Write the logistic regression eq<sup>n</sup> for this.  
b) Calculate the probability that a student who studies for 7 hours will pass.

c) Determine the predicted class (pass or fail) for this student based on the threshold of 0.5

~~2) Consider  $z = [2, 1, 0]$  for three classes. Apply softmax func<sup>n</sup> to find the probability values of three classes.~~

Sol<sup>n</sup> a) logistic regression equation

$$P(\text{pass}) = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$$

$$a_0 = -5, a_1 = 0.8$$

$$P(\text{pass}) = \frac{1}{1 + e^{-( -5 + 0.8x)}}$$

b)  $x=7$

$$P(\text{Pass}) = \frac{1}{1 + e^{-(540.8x)}} \quad (\text{Round } 4 \text{ s.f.})$$

$$\Rightarrow \frac{1}{1 + e^{-0.6}} = 0.646$$

$$\therefore \text{probability} = 0.646 (64.66\%)$$

c) Predicted class % Pass

$$2) z = [2, 1, 0]$$

$$P_p = \frac{e^{z^T}}{\sum_j e^{z_j}}$$

$$\text{Class 1 } (z=2) = 0.665 (66.52\%)$$

~~$$\text{Class 2 } (z=1) = 0.245 (24.47\%)$$~~

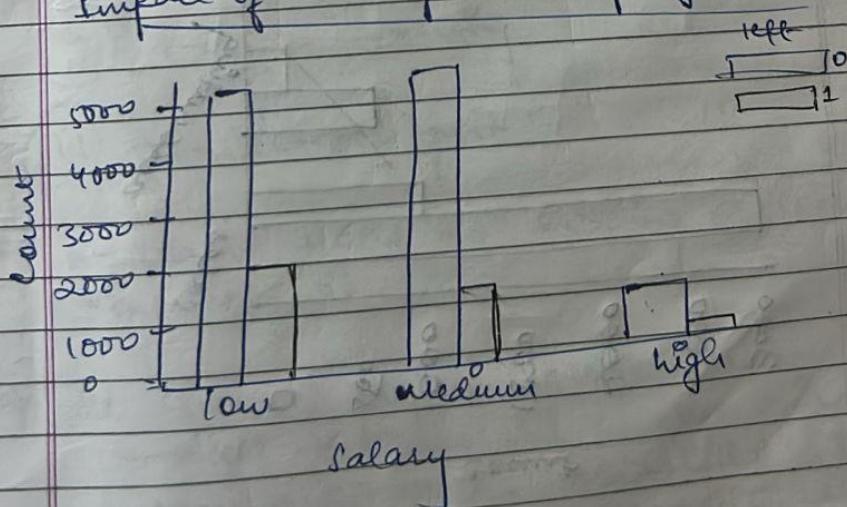
~~$$\text{Class 3 } (z=0) = 0.090 (9.00\%)$$~~

To do: Implementation - Logistic Regression  
(Binary classification)

Dataset = ( HR - comma - sep. csv )

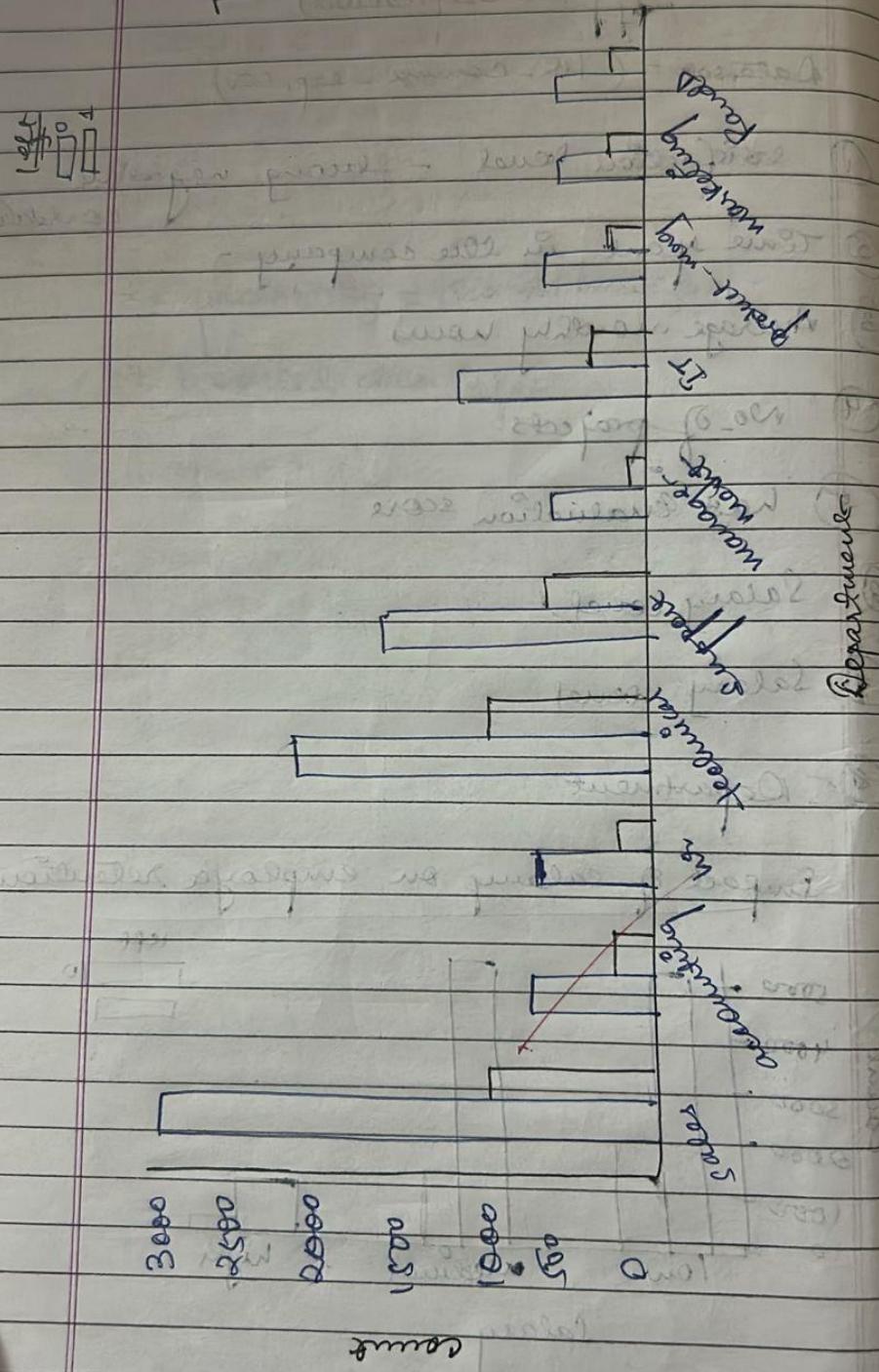
- P ① satisfaction level - strong negative correlation
- ② Time spent in the company -
- ③ Average monthly hours
- ④ No. of projects
- ⑤ last evaluation score
- ⑥ Salary level
- ⑦ Department

24 Impact of salary on employee retention



3P

## Department vs employee retention



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

$$\text{Model Accuracy} = 0.75866. \quad 67$$

	precision	recall	f1-score	support
0	0.80	0.92	0.85	2294
1	0.47	0.23	0.31	706

### Q14 HR - comma\_sep.csv :-

(i) employee retention has direct impact on user satisfaction level

(ii) Accuracy: 0.758

It is good because generally it ranges from 0.75 and 0.85

### Q15 Zoo dataset :

(i) Data cleaning: removed animal\_name column since it was unnecessary

(ii) Feature scaling

(iii) Train-test-split

2) No there were no missing or inconsistent values

3) It shows perfect performance  
 → all entries lie on the diagonal, meaning correctly classified into respective classes

4) None. The confusion matrix shows all classes were classified correctly.

Pred. 1 2 3 4 5 6 7 8 9  
1 0 12.0 15.0 18.0 11.0

MS =  $\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M \delta(y_i, \hat{y}_i)$

equation 10.100 and 10.101

10.100  $\hat{y}_i = \text{argmax}_{j \in \{1, 2, \dots, M\}} p_j(x_i)$

10.101  $p_j(x_i) = \frac{e^{w_j^T x_i + b_j}}{\sum_{k=1}^M e^{w_k^T x_i + b_k}}$

learning behavior - gradient descent (batch)  
iterations over  $\theta$ ,  $w_1, w_2, \dots, w_M$

gradient descent

gradient descent

gradient descent on weight vector  $\theta$   
(minimizing loss function)

gradient descent on weight vector  $\theta$

Code:

```
import pandas as pd
import matplotlib.pyplot as plt
import math
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

# Load dataset
df = pd.read_csv("/content/HR_comma_sep (2).csv")

# Scatter plot: Employee satisfaction vs Retention
plt.scatter(df.satisfaction_level, df.left, marker='+', color='red')
plt.xlabel("Satisfaction Level")
plt.ylabel("Left (1) / Stayed (0)")
plt.title("Impact of Satisfaction Level on Employee Retention")
plt.show()

# Define features (X) and target (y)
X = df[['satisfaction_level']]
y = df['left']

# Split dataset (90% train, 10% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.9, random_state=10)

# Train logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Predictions
y_predicted = model.predict(X_test)

# Model Accuracy
print(f"Model Accuracy: {model.score(X_test, y_test):.4f}")

# Probability predictions
print("Predicted Probabilities:")
print(model.predict_proba(X_test))

# Predict for a specific satisfaction level (e.g., 0.4)
predicted_status = model.predict([[0.4]])
print(f"Prediction for Satisfaction Level 0.4: {'Left' if predicted_status[0] == 1 else 'Stayed'}")

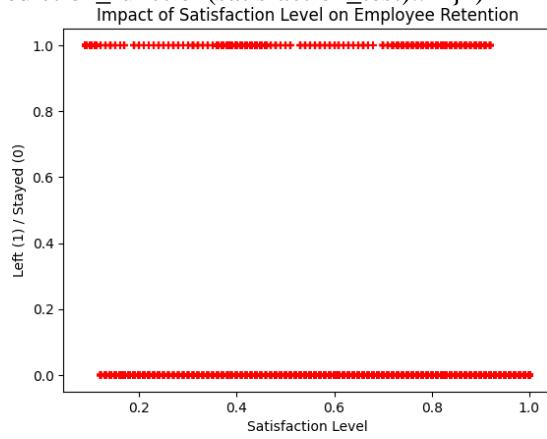
# Logistic function
def sigmoid(x):
    return 1 / (1 + math.exp(-x))
```

```

# Custom prediction function
m, b = model.coef_[0][0], model.intercept_[0]
def prediction_function(satisfaction):
    z = m * satisfaction + b
    y = sigmoid(z)
    return y

satisfaction_test = 0.4
print(f"Sigmoid Prediction for Satisfaction Level {satisfaction_test}:
{prediction_function(satisfaction_test):.4f}")

```



Model Accuracy: 0.7707

Predicted Probabilities:

```

[[0.81879598 0.18120402]
 [0.64435551 0.35564449]
 [0.67008191 0.32991809]
 ...
 [0.85026544 0.14973456]
 [0.93858587 0.06141413]
 [0.90306111 0.09693889]]

```

Prediction for Satisfaction Level 0.4: Stayed

Sigmoid Prediction for Satisfaction Level 0.4: 0.3644

```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, ConfusionMatrixDisplay
import matplotlib.pyplot as plt

```

# Load the Zoo dataset

```

file_path = "/content/zoo-data (1).csv"
zoo_data = pd.read_csv(file_path)

```

# Drop the 'animal\_name' column as it is not a relevant feature

```

X = zoo_data.drop(['animal_name', 'class_type'], axis=1) # Features

```

```

y = zoo_data['class_type'] # Target variable

```

```

# 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 Logistic Regression model for multi-class classification
model = LogisticRegression(multi_class='multinomial', solver='lbfgs', max_iter=200)
# Train the model
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)
# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy of the Multinomial Logistic Regression model: {accuracy:.2f}")

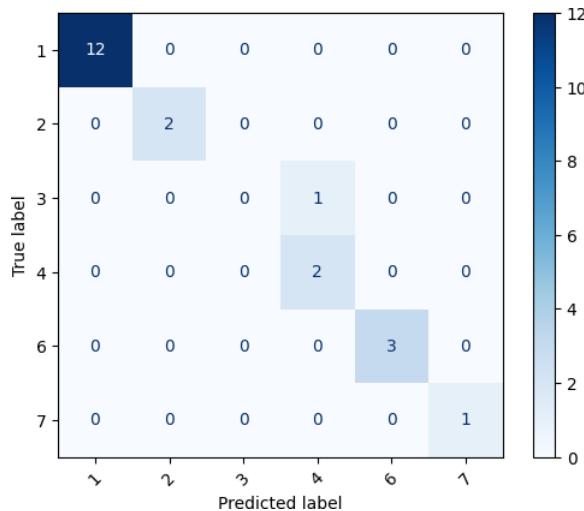
# Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Adjust display labels to match actual present labels in the test set
unique_classes_in_test = sorted(y_test.unique())

# Display confusion matrix
cm_display = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=unique_classes_in_test)
cm_display.plot(cmap='Blues', xticks_rotation=45)
plt.show()

```

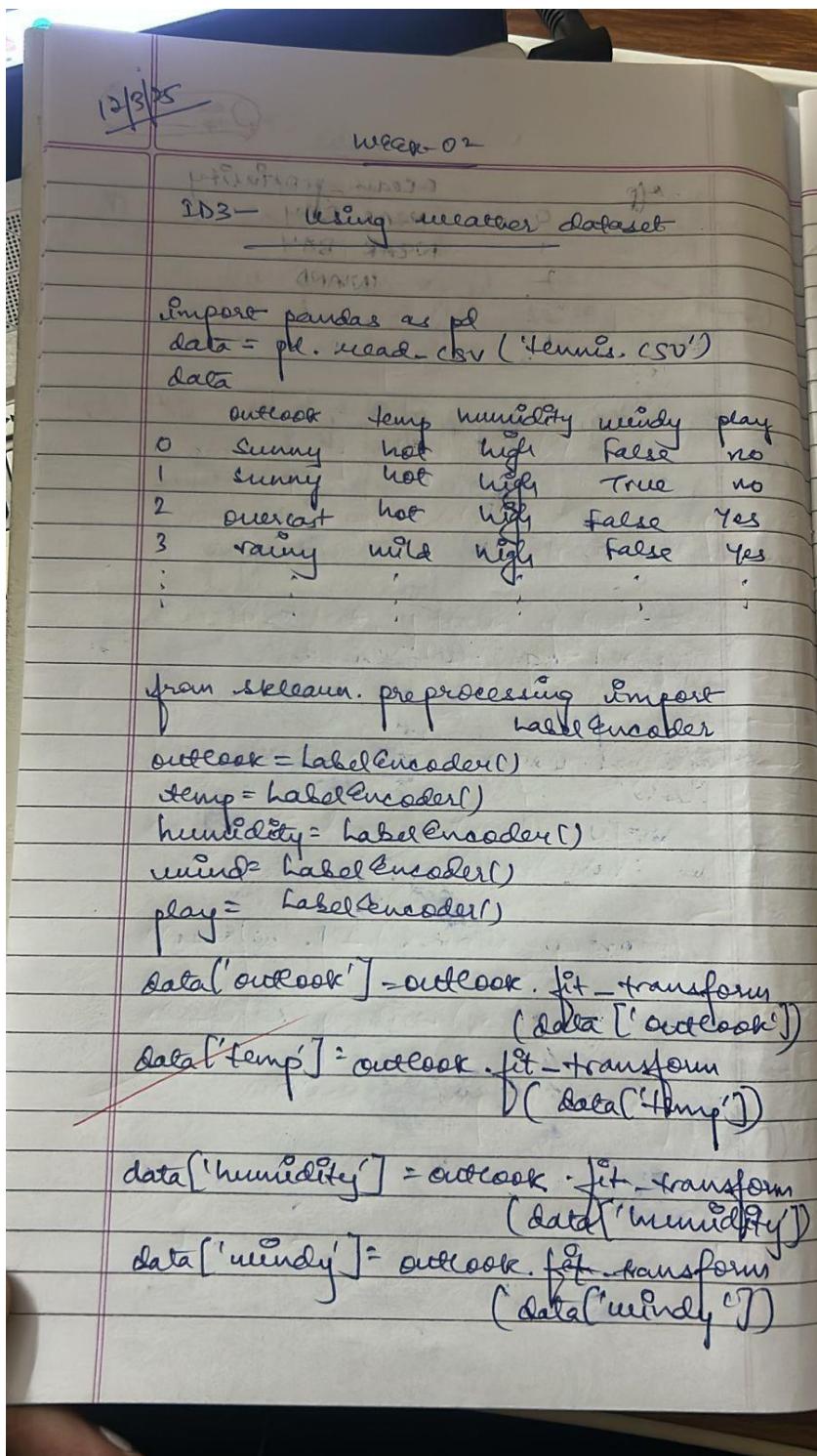
Accuracy of the Multinomial Logistic Regression model: 0.95



## Program 5

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

Screenshot



data['play'] = outlook.fit\_transform  
data

	outlook	temp	humidity	windy	play
0	2	1	0	0	0
1	2	1	0	1	0
2	0	1	0	0	1
3	1	2	0	0	1
4	1	0	1	0	1
5	0	0	0	0	0
6	0	0	0	0	0
7	0	0	0	0	0
8	0	0	0	0	0
9	0	0	0	0	0
10	0	0	0	0	0
11	0	0	0	0	0
12	0	0	0	0	0
13	0	0	0	0	0
14	0	0	0	0	0
15	0	0	0	0	0
16	0	0	0	0	0
17	0	0	0	0	0
18	0	0	0	0	0
19	0	0	0	0	0
20	0	0	0	0	0
21	0	0	0	0	0
22	0	0	0	0	0
23	0	0	0	0	0
24	0	0	0	0	0
25	0	0	0	0	0
26	0	0	0	0	0
27	0	0	0	0	0
28	0	0	0	0	0
29	0	0	0	0	0
30	0	0	0	0	0
31	0	0	0	0	0
32	0	0	0	0	0
33	0	0	0	0	0
34	0	0	0	0	0
35	0	0	0	0	0
36	0	0	0	0	0
37	0	0	0	0	0
38	0	0	0	0	0
39	0	0	0	0	0
40	0	0	0	0	0
41	0	0	0	0	0
42	0	0	0	0	0
43	0	0	0	0	0
44	0	0	0	0	0
45	0	0	0	0	0
46	0	0	0	0	0
47	0	0	0	0	0
48	0	0	0	0	0
49	0	0	0	0	0
50	0	0	0	0	0
51	0	0	0	0	0
52	0	0	0	0	0
53	0	0	0	0	0
54	0	0	0	0	0
55	0	0	0	0	0
56	0	0	0	0	0
57	0	0	0	0	0
58	0	0	0	0	0
59	0	0	0	0	0
60	0	0	0	0	0
61	0	0	0	0	0
62	0	0	0	0	0
63	0	0	0	0	0
64	0	0	0	0	0
65	0	0	0	0	0
66	0	0	0	0	0
67	0	0	0	0	0
68	0	0	0	0	0
69	0	0	0	0	0
70	0	0	0	0	0
71	0	0	0	0	0
72	0	0	0	0	0
73	0	0	0	0	0
74	0	0	0	0	0
75	0	0	0	0	0
76	0	0	0	0	0
77	0	0	0	0	0
78	0	0	0	0	0
79	0	0	0	0	0
80	0	0	0	0	0
81	0	0	0	0	0
82	0	0	0	0	0
83	0	0	0	0	0
84	0	0	0	0	0
85	0	0	0	0	0
86	0	0	0	0	0
87	0	0	0	0	0
88	0	0	0	0	0
89	0	0	0	0	0
90	0	0	0	0	0
91	0	0	0	0	0
92	0	0	0	0	0
93	0	0	0	0	0
94	0	0	0	0	0
95	0	0	0	0	0
96	0	0	0	0	0
97	0	0	0	0	0
98	0	0	0	0	0
99	0	0	0	0	0
100	0	0	0	0	0
101	0	0	0	0	0
102	0	0	0	0	0
103	0	0	0	0	0
104	0	0	0	0	0
105	0	0	0	0	0
106	0	0	0	0	0
107	0	0	0	0	0
108	0	0	0	0	0
109	0	0	0	0	0
110	0	0	0	0	0
111	0	0	0	0	0
112	0	0	0	0	0
113	0	0	0	0	0
114	0	0	0	0	0
115	0	0	0	0	0
116	0	0	0	0	0
117	0	0	0	0	0
118	0	0	0	0	0
119	0	0	0	0	0
120	0	0	0	0	0
121	0	0	0	0	0
122	0	0	0	0	0
123	0	0	0	0	0
124	0	0	0	0	0
125	0	0	0	0	0
126	0	0	0	0	0
127	0	0	0	0	0
128	0	0	0	0	0
129	0	0	0	0	0
130	0	0	0	0	0
131	0	0	0	0	0
132	0	0	0	0	0
133	0	0	0	0	0
134	0	0	0	0	0
135	0	0	0	0	0
136	0	0	0	0	0
137	0	0	0	0	0
138	0	0	0	0	0
139	0	0	0	0	0
140	0	0	0	0	0
141	0	0	0	0	0
142	0	0	0	0	0
143	0	0	0	0	0
144	0	0	0	0	0
145	0	0	0	0	0
146	0	0	0	0	0
147	0	0	0	0	0
148	0	0	0	0	0
149	0	0	0	0	0
150	0	0	0	0	0
151	0	0	0	0	0
152	0	0	0	0	0
153	0	0	0	0	0
154	0	0	0	0	0
155	0	0	0	0	0
156	0	0	0	0	0
157	0	0	0	0	0
158	0	0	0	0	0
159	0	0	0	0	0
160	0	0	0	0	0
161	0	0	0	0	0
162	0	0	0	0	0
163	0	0	0	0	0
164	0	0	0	0	0
165	0	0	0	0	0
166	0	0	0	0	0
167	0	0	0	0	0
168	0	0	0	0	0
169	0	0	0	0	0
170	0	0	0	0	0
171	0	0	0	0	0
172	0	0	0	0	0
173	0	0	0	0	0
174	0	0	0	0	0
175	0	0	0	0	0
176	0	0	0	0	0
177	0	0	0	0	0
178	0	0	0	0	0
179	0	0	0	0	0
180	0	0	0	0	0
181	0	0	0	0	0
182	0	0	0	0	0
183	0	0	0	0	0
184	0	0	0	0	0
185	0	0	0	0	0
186	0	0	0	0	0
187	0	0	0	0	0
188	0	0	0	0	0
189	0	0	0	0	0
190	0	0	0	0	0
191	0	0	0	0	0
192	0	0	0	0	0
193	0	0	0	0	0
194	0	0	0	0	0
195	0	0	0	0	0
196	0	0	0	0	0
197	0	0	0	0	0
198	0	0	0	0	0
199	0	0	0	0	0
200	0	0	0	0	0
201	0	0	0	0	0
202	0	0	0	0	0
203	0	0	0	0	0
204	0	0	0	0	0
205	0	0	0	0	0
206	0	0	0	0	0
207	0	0	0	0	0
208	0	0	0	0	0
209	0	0	0	0	0
210	0	0	0	0	0
211	0	0	0	0	0
212	0	0	0	0	0
213	0	0	0	0	0
214	0	0	0	0	0
215	0	0	0	0	0
216	0	0	0	0	0
217	0	0	0	0	0
218	0	0	0	0	0
219	0	0	0	0	0
220	0	0	0	0	0
221	0	0	0	0	0
222	0	0	0	0	0
223	0	0	0	0	0
224	0	0	0	0	0
225	0	0	0	0	0
226	0	0	0	0	0
227	0	0	0	0	0
228	0	0	0	0	0
229	0	0	0	0	0
230	0	0	0	0	0
231	0	0	0	0	0
232	0	0	0	0	0
233	0	0	0	0	0
234	0	0	0	0	0
235	0	0	0	0	0
236	0	0	0	0	0
237	0	0	0	0	0
238	0	0	0	0	0
239	0	0	0	0	0
240	0	0	0	0	0
241	0	0	0	0	0
242	0	0	0	0	0
243	0	0	0	0	0
244	0	0	0	0	0
245	0	0	0	0	0
246	0	0	0	0	0
247	0	0	0	0	0
248	0	0	0	0	0
249	0	0	0	0	0
250	0	0	0	0	0
251	0	0	0	0	0
252	0	0	0	0	0
253	0	0	0	0	0
254	0	0	0	0	0
255	0	0	0	0	0
256	0	0	0	0	0
257	0	0	0	0	0
258	0	0	0	0	0
259	0	0	0	0	0
260	0	0	0	0	0
261	0	0	0	0	0
262	0	0	0	0	0
263	0	0	0	0	0
264	0	0	0	0	0
265	0	0	0	0</td	

and or if. 0.333333 = 1/3rd of data

1 0 0 1 0 0

2 1 0 0 0 0

3 1 0 0 0 0

play (rainy) 33.333333 33.333333 33.333333

4 0 1 1 1 0 0

5 0 1 1 1 0 0

6 0 1 1 1 0 0

Name: play, type float64

from sklearn.model\_selection import train\_test\_split

x-train, x-test, y-train, y-test =

train\_test\_split(x, y, test\_size=0.2)

from sklearn.tree import DecisionTreeClassifier

classifier = DecisionTreeClassifier(criterion='entropy')

classifier.fit(x-train, y-train)

classifier.predict(x-test)

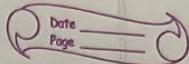
outlook temp humidity windy

9	1	2	0	1	0
0	2	1	0	0	1

$y_{\text{test}}$

play  
no  
yes  
0  
1  
0  
0

no play  
yes play



classifier score ( $x_{\text{test}}, y_{\text{test}}$ )

1.0

from sklearn import tree  
tree.plot\_tree(classifier)

clf = DecisionTreeClassifier(criterion='entropy')

clf.fit(x, y)

importances = clf.feature\_importances\_

feature\_importance\_df = pd.DataFrame()

'Feature': x.columns

'Information Gain': importances

~~3) feature\_importance\_df = feature\_importance\_df = feature\_importance\_df~~

df.set\_index

(by='Information Gain',  
ascending=False)

print(feature\_importance\_df)

## Feature Information Gain

0	outlook	0.514558
3	windy	0.1279205
2	humidity	0.211237
1	temp	0.00000

At 70% - (0.514558) = 0.1829205

0.18

outlook = wind (negative)

= humidity (negative) = 0.1279205  
(negative)

(0.1279205) \* 0.18 = 0.0229856

humidity = 0.0229856

Temp = 0.00000

outlook = 0.514558

windy = 0.1279205

humidity = 0.211237

temp = 0.00000

(0.211237 \* 0.18) + (0.1279205 \* 0.18) + (0.514558 \* 0.18) = 0.1829205

Code:

```
import numpy as np
import pandas as pd
from collections import Counter
class Node:
    def __init__(self, feature=None, value=None, label=None):
        self.feature = feature # Attribute to split on
        self.value = value    # Value of the attribute
        self.label = label    # Label if it's a leaf node
        self.children = {}    # Dictionary of child nodes

    def entropy(y):
        counts = np.bincount(y)
        probabilities = counts / len(y)
        return -np.sum([p * np.log2(p) for p in probabilities if p > 0])

    def information_gain(X, y, feature):
        total_entropy = entropy(y)
        values, counts = np.unique(X[:, feature], return_counts=True)
        weighted_entropy = sum((counts[i] / sum(counts)) * entropy(y[X[:, feature] == v]) for i, v in enumerate(values))
        return total_entropy - weighted_entropy

    def best_feature_to_split(X, y):
        gains = [information_gain(X, y, i) for i in range(X.shape[1])]
        return np.argmax(gains)

def id3(X, y, features):
    if len(set(y)) == 1:
        return Node(label=y[0])
    if len(features) == 0:
        return Node(label=Counter(y).most_common(1)[0][0])
    best_feature = best_feature_to_split(X, y)
    node = Node(feature=features[best_feature])
    feature_values = np.unique(X[:, best_feature])
    for value in feature_values:
        sub_X = X[X[:, best_feature] == value]
        sub_y = y[X[:, best_feature] == value]
        if len(sub_y) == 0:
            node.children[value] = Node(label=Counter(y).most_common(1)[0][0])
        else:
            node.children[value] = id3(np.delete(sub_X, best_feature, axis=1), sub_y, features[:best_feature] + features[best_feature+1:])
    return node

    if node.label is not None:
        print(f"{' ' * depth}Leaf: {node.label}")
        return
    print(f"{' ' * depth}Feature: {node.feature}")


```

```

for value, child in node.children.items():
    print(f"{' ' * depth}Value: {value}")
    print_tree(child, depth + 1)
# Example dataset
data = pd.DataFrame({
    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain', 'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny',
    'Overcast', 'Overcast', 'Rain'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'Normal', 'Normal',
    'High', 'Normal', 'High'],
    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong',
    'Weak', 'Strong'],
    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
})
X = data.iloc[:, :-1].apply(lambda col: pd.factorize(col)[0]).to_numpy()
y = pd.factorize(data['PlayTennis'])[0]
features = list(data.columns[:-1])
decision_tree = id3(X, y, features)
print_tree(decision_tree)

```

Feature: Outlook

Value: 0

  Feature: Humidity

  Value: 0

    Leaf: 0

  Value: 1

    Leaf: 1

  Value: 1

    Leaf: 1

  Value: 2

    Feature: Wind

      Value: 0

        Leaf: 1

      Value: 1

        Leaf: 0

## Program 6

Build KNN Classification model for a given dataset

Screenshot

KNN

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

Consider the following dataset for  $K=3$  and test data  $(X, 35, 100)$ . Use (Person, Age, Salary) to solve using KNN classifier model and predict the target!

Person	Age	Salary	K	Target
A	18	50		N
B	23	55		N
C	24	70		N
D	41	60		Y
E	43	70		Y
F	38	40		Y
X	35	100		?

$D = \sqrt{(Age - A_{age})^2 + (Salary - A_{salary})^2}$

$A \rightarrow \sqrt{(35-18)^2 + (100-50)^2} = 52.81$

$B \rightarrow \sqrt{(35-23)^2 + (100-55)^2} = 46.56$

$C \rightarrow \sqrt{(35-24)^2 + (100-70)^2} = 31.96$

$D \rightarrow \sqrt{(35-41)^2 + (100-60)^2} = 40.45$

$E \rightarrow \sqrt{(35-43)^2 + (100-70)^2} = 31.05$

$F \rightarrow \sqrt{(35-38)^2 + (100-40)^2} = 60.08$

$\boxed{[C, D, E] \rightarrow [N, Y, Y]}$

$\boxed{| X(35, 100) \rightarrow ?|}$

### Euis Dataset

→ Here the K value has no impact on core accuracy for K value (1, 2, 6) but after that there were some fluctuations

K value can be chosen based on hit and miss and trial and based on the size of the dataset.

### Diabetes dataset

#### Purpose -

- ensures all features contribute equally by bringing them to a similar scale
- prevents large scale features from dominating distance based algo like KNN.
- Improves convergence speed

#### How to perform:

#### 17 Standardization (z-score normalization)

$$Z = \frac{x - M}{S}$$

$$20.00 = (40 - 0.5) + (28 - 22) \rightarrow Z = 3$$

use StandardScalar() in select - learn

$$Z = \frac{X - M}{S}$$

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

## 2) Normalization (MinMaxScaling)

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

$$\text{Max} - \text{Min}$$

Use MinMaxScalar() in scikit-learn.

8

Code:

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

# Function to train and evaluate KNN model
def knn_classification(data_path, target_column, dataset_name, k=5):
    # Load dataset
    df = pd.read_csv(data_path)

    # Split features and target
    X = df.drop(columns=[target_column])
    y = df[target_column]

    # Split 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)

    # Feature scaling for better performance
    scaler = StandardScaler()
    X_train = scaler.fit_transform(X_train)
    X_test = scaler.transform(X_test)

    # Train KNN model
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)

    # Make predictions
    y_pred = model.predict(X_test)

    # Evaluate model
    accuracy = accuracy_score(y_test, y_pred)
    print(f'Accuracy of KNN on {dataset_name} dataset: {accuracy:.4f}')
    print("Classification Report:")
    print(classification_report(y_test, y_pred))

    # Confusion matrix
    cm = confusion_matrix(y_test, y_pred)
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
    plt.xlabel('Predicted')
    plt.ylabel('Actual')
    plt.title(f'Confusion Matrix - {dataset_name}')
    plt.show()
```

```
# Run KNN classification on both datasets
knn_classification('/content/iris (3).csv', 'species', 'Iris', k=5)
knn_classification('/content/diabetes.csv', 'Outcome', 'Diabetes', k=5)
```

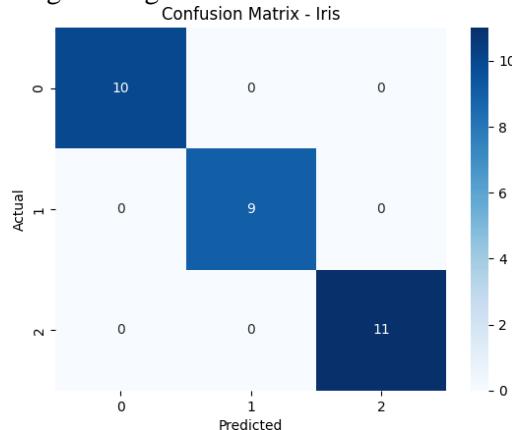
Accuracy of KNN on Iris dataset: 1.0000

Classification Report:

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

setosa	1.00	1.00	1.00	10
versicolor	1.00	1.00	1.00	9
virginica	1.00	1.00	1.00	11

accuracy		1.00	30	
macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30



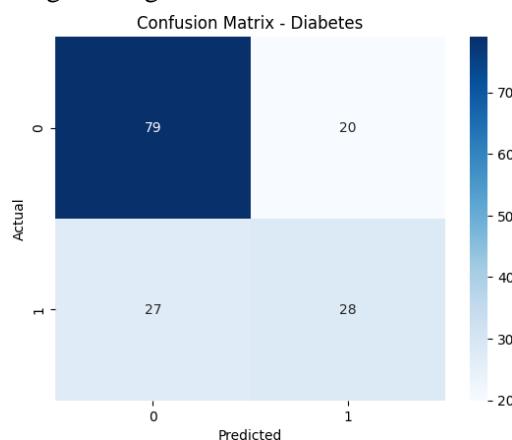
Accuracy of KNN on Diabetes dataset: 0.6948

Classification Report:

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.75	0.80	0.77	99
1	0.58	0.51	0.54	55

accuracy		0.69	154	
macro avg	0.66	0.65	0.66	154
weighted avg	0.69	0.69	0.69	154



```

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

# Load dataset
df = pd.read_csv('/content/heart.csv')

# Define features and target
X = df.drop(columns=['target']) # Assuming 'target' is the classification column
y = df['target']

# Split data
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 = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Find the best K value
k_values = range(1, 21)
accuracy_scores = []
for k in k_values:
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    accuracy_scores.append(accuracy_score(y_test, y_pred))

best_k = k_values[np.argmax(accuracy_scores)]
print(f'Best K value: {best_k}')

# Train model with best K
best_model = KNeighborsClassifier(n_neighbors=best_k)
best_model.fit(X_train, y_train)
y_pred = best_model.predict(X_test)

# Evaluate model
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy with best K ({best_k}): {accuracy:.4f}')
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Confusion matrix

```

```

cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title(f'Confusion Matrix - KNN (K={best_k})')
plt.show()

```

```

# Plot K values vs. Accuracy
plt.plot(k_values, accuracy_scores, marker='o')
plt.xlabel('K Value')
plt.ylabel('Accuracy')
plt.title('K Value vs Accuracy')
plt.show()

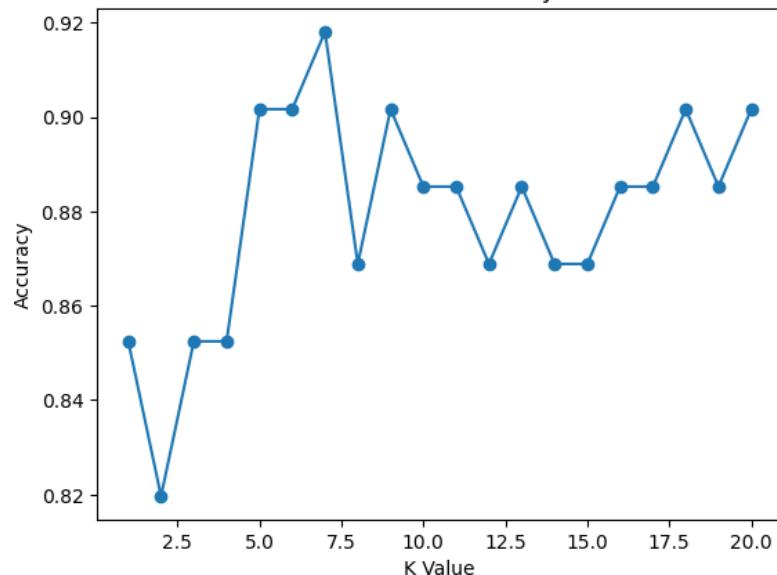
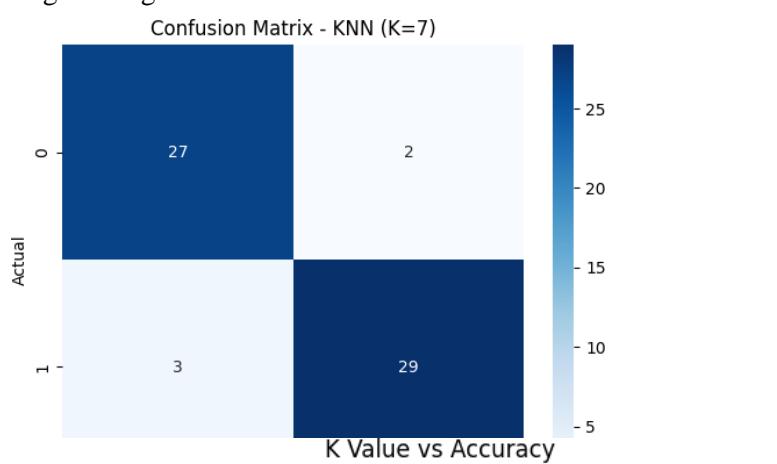
```

Best K value: 7

Accuracy with best K (7): 0.9180

Classification Report:

	precision	recall	f1-score	support
0	0.90	0.93	0.92	29
1	0.94	0.91	0.92	32
accuracy		0.92	0.92	61
macro avg	0.92	0.92	0.92	61
weighted avg	0.92	0.92	0.92	61



## Program 7

Build Support vector machine model for a given dataset

Code:

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

# Define the Linear SVM class
class LinearSVM:
    def __init__(self, learning_rate=0.001, reg_strength=0.1, num_iterations=1000):
        self.learning_rate = learning_rate
        self.reg_strength = reg_strength
        self.num_iterations = num_iterations

    def fit(self, X, y):
        # Initialize weights and bias
        num_samples, num_features = X.shape
        self.W = np.zeros(num_features) # Weights
        self.b = 0 # Bias

        # Gradient Descent
        for _ in range(self.num_iterations):
            # Compute the margin (decision function)
            margins = 1 - y * (np.dot(X, self.W) + self.b)
            # Compute gradient
            dw = -2 * np.sum(y * margins) / num_samples + 2 * self.reg_strength * self.W
            db = -2 * np.sum(y * margins) / num_samples

            # Update weights and bias
            self.W -= self.learning_rate * dw
            self.b -= self.learning_rate * db

    def predict(self, X):
        # Make predictions
        return np.sign(np.dot(X, self.W) + self.b)

# Generate toy data (binary classification)
np.random.seed(42)
num_samples = 100
X = np.random.randn(num_samples, 2)
y = np.ones(num_samples)
y[X[:, 0] < X[:, 1]] = -1 # Assign different class based on condition

# Train the Linear SVM
svm = LinearSVM(learning_rate=0.001, reg_strength=0.1, num_iterations=1000)
svm.fit(X, y)
```

```

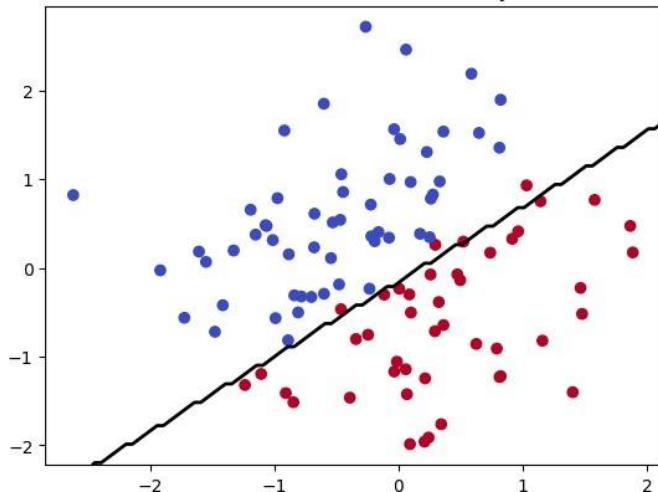
# Predict
y_pred = svm.predict(X)

# Visualize the decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='coolwarm')
ax = plt.gca()
xlim = ax.get_xlim()
ylim = ax.get_ylim()
xx, yy = np.meshgrid(np.linspace(xlim[0], xlim[1], 100), np.linspace(ylim[0], ylim[1], 100))
Z = svm.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contour(xx, yy, Z, levels=[0], linewidths=2, colors='black')
plt.title("Linear SVM Decision Boundary")
plt.show()

# Print accuracy (simple comparison)
accuracy = np.mean(y_pred == y)
print(f"Accuracy: {accuracy * 100:.2f}%")

```

Linear SVM Decision Boundary

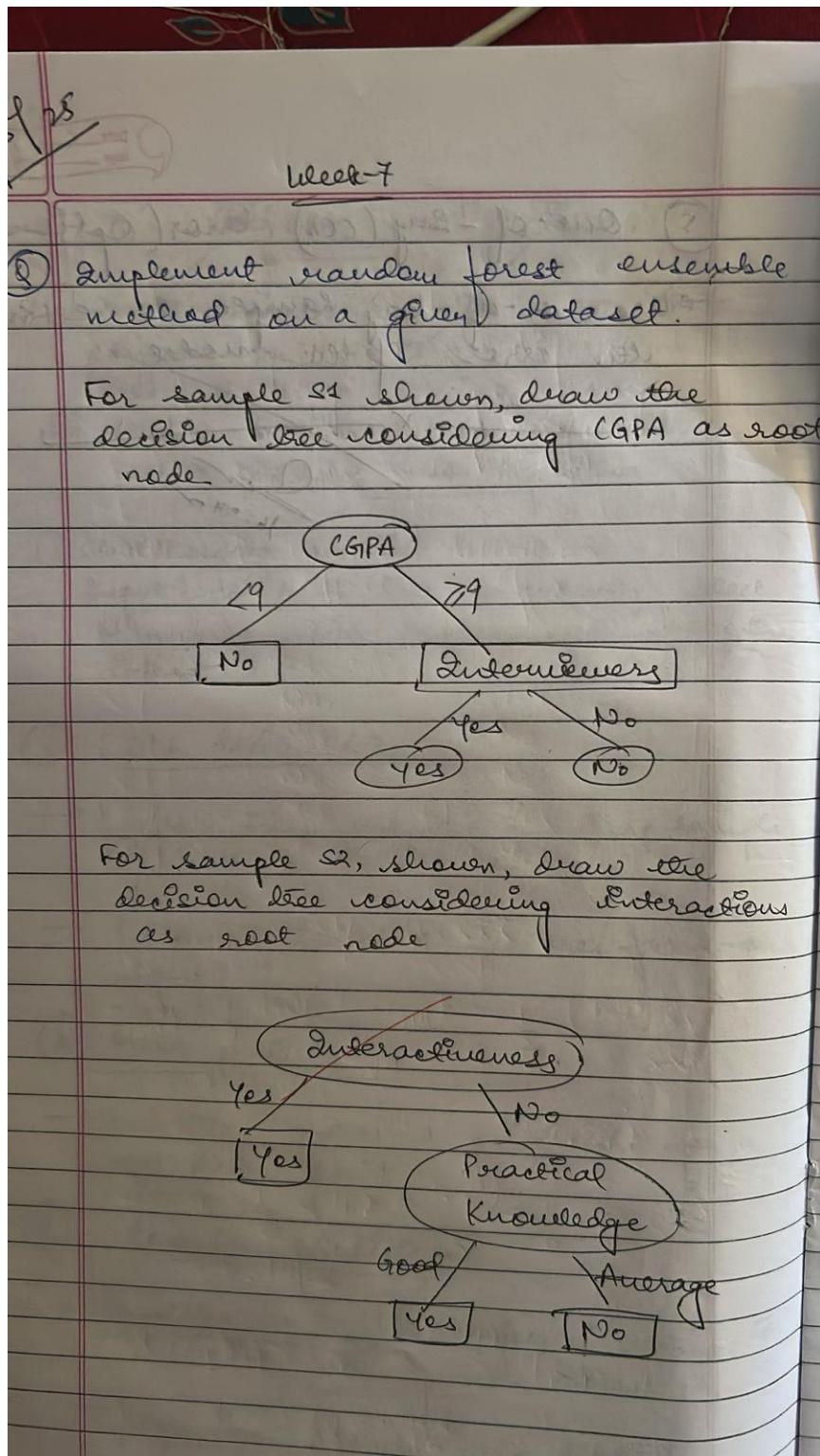


Accuracy: 96.00%

## Program 8

Implement Random forest ensemble method on a given dataset

Screenshot



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

Write the answer for the following question -

① for "iris.csv" dataset -  
what is the best accuracy score  
and confusion matrix of the  
classification you observed and  
using how many trees?

⇒ Best accuracy score Observed  
⇒ 1.0

⇒ No. of trees (n-estimators)  
= best\_n

⇒ Confusion Matrix

$$\begin{bmatrix} [16 & 0 & 0] \\ [0 & 14 & 0] \\ [0 & 0 & 15] \end{bmatrix}$$

Code:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt

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

# Assuming last column is the label
X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values

# Split into training and test sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
# 1. Train RF Classifier with default n_estimators=10
rf_default = RandomForestClassifier(n_estimators=10, random_state=42)
rf_default.fit(X_train, y_train)
y_pred_default = rf_default.predict(X_test)
accuracy_default = accuracy_score(y_test, y_pred_default)

print(f"Default RF Accuracy (n_estimators=10): {accuracy_default:.4f}")

best_accuracy = 0
best_n = 0
accuracies = []

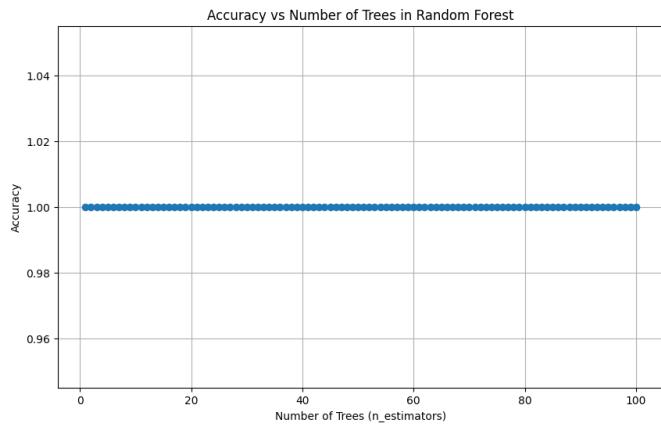
for n in range(1, 101):
    rf = RandomForestClassifier(n_estimators=n, random_state=42)
    rf.fit(X_train, y_train)
    y_pred = rf.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    accuracies.append(acc)

    if acc > best_accuracy:
        best_accuracy = acc
        best_n = n

print(f"Best RF Accuracy: {best_accuracy:.4f} with n_estimators = {best_n}")
# Plot accuracy vs. number of trees
plt.figure(figsize=(10, 6))
plt.plot(range(1, 101), accuracies, marker='o')
plt.title("Accuracy vs Number of Trees in Random Forest")
plt.xlabel("Number of Trees (n_estimators)")
```

```
plt.ylabel("Accuracy")
plt.grid(True)
plt.show()
```

Default RF Accuracy (n\_estimators=10): 1.0000  
Best RF Accuracy: 1.0000 with n\_estimators = 1



## Program 9

Implement Boosting ensemble method on a given dataset

Screensho

7/5/25      Lab-8

⑧ Implement Boosting ensemble method on a given dataset  
 Considering AdaBoost Algorithm for the following sample data, show the decision stamp calculation steps for the attribute CGPA.

CGPA	Predicted Job offer	Actual Job offer	Weight
≥ 9	Yes	Yes	1/6
< 9	No	Yes	1/6
≥ 9	Yes	No	1/6
< 9	No	No	1/6
≥ 9	Yes	Yes	1/6
≥ 9	Yes	Yes	1/6

$$E_{CGPA} = 2 \times 1/6 = 0.333$$

$$\alpha_{CGPA} = \frac{1}{2} \ln \left( \frac{1 - E_{CGPA}}{E_{CGPA}} \right) = \frac{1}{2} \frac{(1 - 0.333)}{0.333} = 0.347$$

$$Z_{CGPA} = \frac{1}{6} \times 4 \times e^{-0.347} + \frac{1}{6} \times 2 \times e^{0.347} = 0.9428$$

$$wt(\delta_i)_{i=1} = \frac{\frac{1}{6} \times e^{-0.347}}{0.9428} = 0.1249$$

$$out(d_j^0)_{i+1} = \frac{\frac{1}{6} \times e^{0.347}}{0.9428} = 0.2801$$

CGPA Predicted job actual weight  
offer offer job offer

7.9	Yes	Yes	0.1249
<9	No	Yes	0.1249
≥9	Yes	No	0.2501
<9	No	No	0.1249
≥9	Yes	Yes	0.1249
≥9	Yes	Yes	0.1249

Re-write the answer for the following question-

Q7 For "income.csv" dataset,  
what is the best accuracy score  
and confusion matrix of the classifier  
you obtained and using how many  
trees?

⇒ Accuracy with 10 estimators : 0.8277

Confusion matrix (10 estimators):

$$\begin{bmatrix} 1072 & 387 \\ 2138 & 1408 \end{bmatrix}$$

Best accuracy % 0.831 with  
n-estimators = 72.

t

Code:

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

# Step 1: Load the dataset
df = pd.read_csv("/content/income.csv")

# Step 2: Split into features and target
X = df.drop(columns=['income_level'])
y = df['income_level']

# Step 3: Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Step 4: AdaBoost with 10 estimators
model_10 = AdaBoostClassifier(n_estimators=10, random_state=42)
model_10.fit(X_train, y_train)
y_pred_10 = model_10.predict(X_test)
accuracy_10 = accuracy_score(y_test, y_pred_10)
conf_matrix_10 = confusion_matrix(y_test, y_pred_10)

print("Accuracy with 10 estimators:", round(accuracy_10, 4))
print("Confusion Matrix (10 estimators):\n", conf_matrix_10)

# Step 5: Fine-tune number of trees (1 to 50)
best_accuracy = 0
best_n = 0
accuracies = []

for n in range(1, 51):
    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)
    accuracies.append(acc)

    if acc > best_accuracy:
        best_accuracy = acc
        best_n = n
```

```

print(f"\nBest Accuracy: {round(best_accuracy, 4)} with n_estimators = {best_n}")

# Step 6: Plot accuracy vs. number of estimators
plt.figure(figsize=(10, 6))
plt.plot(range(1, 51), accuracies, marker='o', linestyle='-', color='blue')
plt.title('Accuracy vs Number of Trees (n_estimators)')
plt.xlabel('Number of Trees')
plt.ylabel('Accuracy')
plt.grid(True)
plt.tight_layout()
plt.show()

```

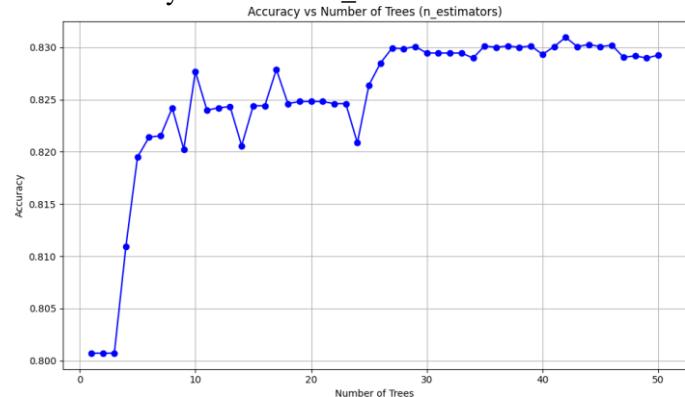
Accuracy with 10 estimators: 0.8277

Confusion Matrix (10 estimators):

`[[10722 387]`

`[ 2138 1406]]`

Best Accuracy: 0.831 with n\_estimators = 42



## Program 10

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

Screenshot

7/5/25  
Lab-9

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

Compute two clusters using K-means algorithm for clustering where initial cluster centers are  $(1.0, 1.0)$  and  $(5.0, 7.0)$ . Execute for 2 iterations

→ Iteration 1:

Record Number	Closest to C1 $(1.0, 1.0)$	Closest to C2 $(5.0, 7.0)$	Assigned to cluster
R1 $(1.0, 1.0)$	0.0	7.21	C1
R2 $(1.5, 2.0)$	1.12	6.12	C2
R3 $(3.0, 4.0)$	3.61	3.61	C1
R4 $(5.0, 7.0)$	7.21	0.0	C2
R5 $(3.5, 5.0)$	4.12	2.5	C2
R6 $(4.5, 5.0)$	5.31	2.06	C2
R7 $(3.5, 4.5)$	4.36	2.92	C2

Clusters  $\{R_1, R_2, R_3\}$  and Cluster 2  $\{R_4, R_5, R_6, R_7\}$

Their new centroids are-

$$C1 = \frac{(1.0 + 1.5 + 3.0)}{3}, \frac{(1.0 + 2.0 + 4.0)}{3}, \\ = 5.5/3 = 1.83, \quad 7.0/3 = 2.33$$
$$C2 = \frac{(5.0 + 3.5 + 4.5 + 3.5)}{4}, \frac{(7 + 5 + 5 + 4)}{4} \\ = 16.5/4 = 4.12, \quad 21.5/4 = 5.37$$

Iteration 2:

Record No.	Close to C1 (1.83, 2.33)	Close to C2 (4.12, 5.37)	Assign to Cluster
R1 (1.0, 1.0)	1.57	5.64	C1
R2 (1.5, 2.0)	0.47	4.52	C1
R3 (3.0, 4.0)	2.12	1.63	C2
R4 (5.0, 7.0)	5.57	1.91	C2
R5 (3.5, 5.0)	3.16	0.72	C2
R6 (4.5, 5.0)	3.58	0.52	C2
R7 (3.5, 4.5)	2.63	1.05	C2

Cluster 1 {R1, R2} and Cluster 2 {R3, R4, R5, R6, R7}

Their new centroids are:

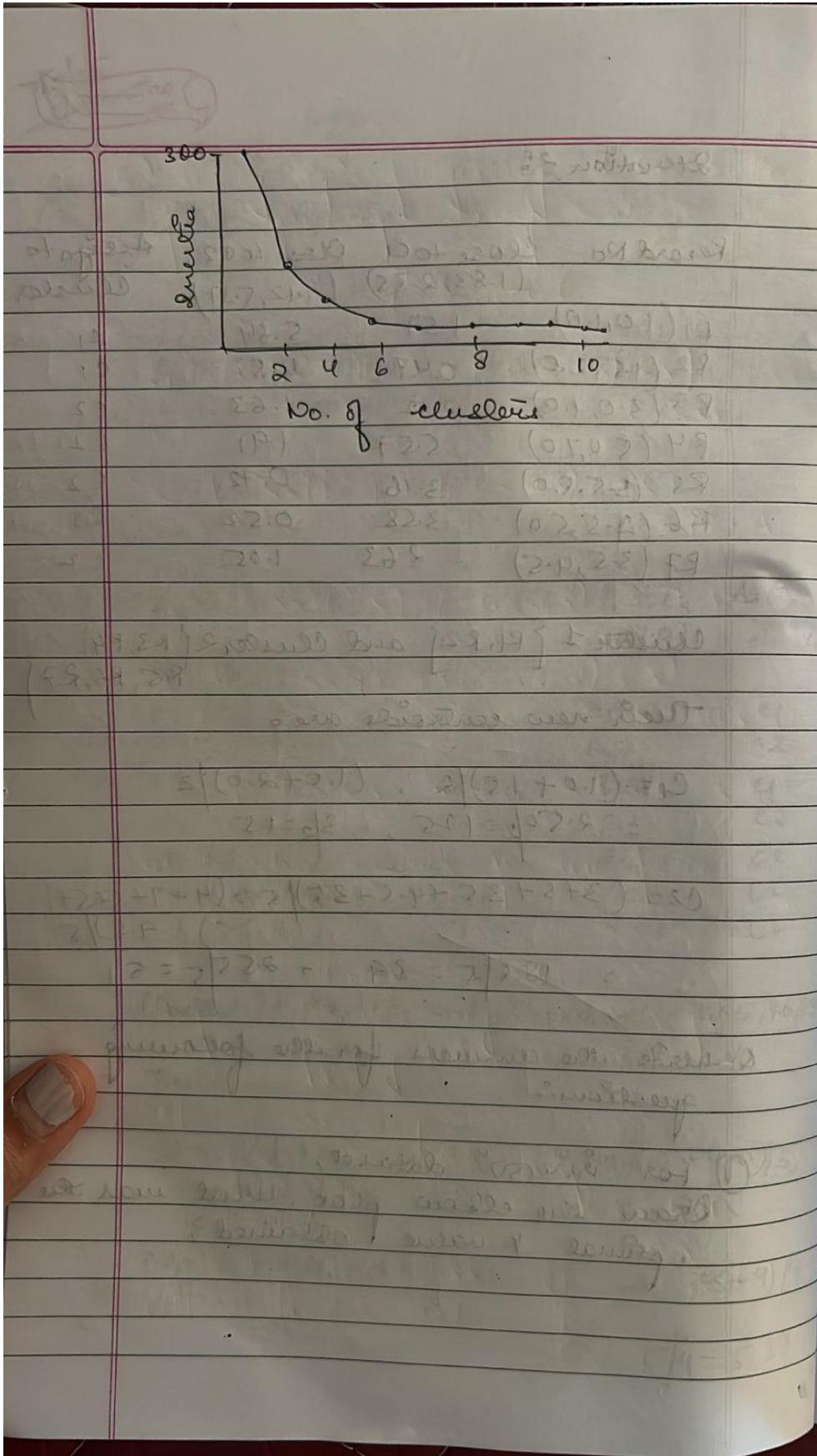
$$C_1 = (1.0 + 1.5)/2, (1.0 + 2.0)/2 \\ = 2.5/2 = 1.25, 3/2 = 1.5$$

$$C_2 = (3 + 5 + 3.5 + 4.5 + 3.5)/5 \rightarrow (18.5)/5 = 3.9 \rightarrow 25.5/5 = 5.1$$

Q. Write the answers for the following questions:

① For "iris.csv" dataset,

Draw the elbow plot. What was the optimal K value obtained?



Code:

```
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

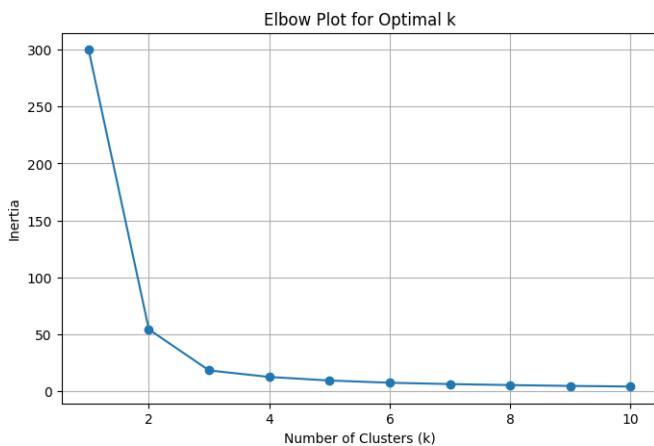
# Load the dataset
df = pd.read_csv("/content/iris (2).csv")

# Select only petal length and petal width
X = df[['petal_length', 'petal_width']]

# Optional: Standardize the data
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:
    model = KMeans(n_clusters=k, random_state=42, n_init=10)
    model.fit(X_scaled)
    inertia.append(model.inertia_)

# Plot the elbow graph
plt.figure(figsize=(8, 5))
plt.plot(k_range, inertia, marker='o')
plt.title('Elbow Plot for Optimal k')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia')
plt.grid(True)
plt.show()
```



## Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

Screenshot

lab-10

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

① Implement Dimensionality reduction using Principal Component Analysis (PCA) method %

Given the data in table, reduce the dimensions from 2 to 1 using PCA.  
Compute the 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 :-  
 $\lambda_1 = 30.3844$   
 $\lambda_2 = 6.6151$

Eigen vectors :-  
 $e_1 = \begin{bmatrix} 0.5594 \\ -0.8303 \end{bmatrix}$     $e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$

① data matrix =  $\begin{bmatrix} 4 & 8 & 13 & 7 \\ 11 & 4 & 5 & 14 \end{bmatrix}$

② Mean centre the data %

Mean  $x_1 = \frac{4+8+13+7}{4} = 8$

Mean  $x_2 = \frac{11+4+5+14}{4} = 8.5$

③  $Y_{\text{centered}} = \begin{bmatrix} 4-8 & 8-8 & 13-8 & 7-8 \\ 11-8.5 & 4-8.5 & 5-8.5 & 14-8.5 \end{bmatrix}$

$= \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix}$

- ② Using eq" projecting data onto first principle component

$$z = c_1^T \cdot Y_{\text{centered}}$$

$$z_1 = (0.5574)(-4) + (-0.8303)(25) = 4.30535$$

$$z_2 = (0.5574)(0) + (-0.8303)(-4.5) = 3.73635$$

$$z_3 = (0.5574)(5) + (-0.8303)(-3.7) = 5.69305$$

$$z_4 = (0.5574)(-1) + (-0.8303)(5.3) = -5.12405$$

$$z = [-4.30535 \quad 3.73635 \quad 5.69305 \quad -5.12405]$$

Write the answer for the following question:

- ① For "heart.csv" dataset, repeat the accuracy scores before and after applying PEA.

→ Model accuracy %

Before

SVM: 0.8804

Logistic Regression: 0.8533

Random Forest: 0.8859

After

0.8424

0.8461

0.8533

Code:

```
import pandas as pd
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.decomposition import PCA
from sklearn.metrics import accuracy_score

# Load dataset
df = pd.read_csv("/content/heart (1).csv") # Update to match your file path if needed

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

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

# Encode categorical columns
for col in categorical_cols:
    if X[col].nunique() == 2:
        X[col] = LabelEncoder().fit_transform(X[col])
    else:
        X = pd.get_dummies(X, columns=[col])

# Scale 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)

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

# Train and evaluate models (without PCA)
```

```

print("⌚ Accuracy without PCA:")
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    print(f"{name}: {accuracy_score(y_test, y_pred):.4f}")

# Apply PCA (reduce to 5 components)
pca = PCA(n_components=5)
X_pca = pca.fit_transform(X_scaled)
X_train_pca, X_test_pca, y_train_pca, y_test_pca = train_test_split(X_pca, y, test_size=0.2, random_state=42)

# Train and evaluate models (with PCA)
print("\n"# Accuracy with PCA:")
for name, model in models.items():
    model.fit(X_train_pca, y_train_pca)
    y_pred_pca = model.predict(X_test_pca)
    print(f"{name}: {accuracy_score(y_test_pca, y_pred_pca):.4f}")

```

⌚ Accuracy without PCA:

SVM: 0.8804

Logistic Regression: 0.8533

Random Forest: 0.8859

# Accuracy with PCA:

SVM: 0.8424

Logistic Regression: 0.8641

Random Forest: 0.8533