

A decorative banner featuring five large, bold letters: 'I', 'N', 'D', 'E', and 'X'. Each letter is enclosed in a white square frame with a black border, and they are all set against a light blue background.

NAME: SANJANA SURESH STD.: _____ SEC.: _____ ROLL NO.: _____ SUB.: ML LAB

WEEK - 00

To - do EXERCISES

- ① initializing values directly into DataFrame - five rows of data with columns heading as USN, Name, Marks.

```
import pandas as pd
```

```
data = { "USN": [ "234", "235", "236", "237", "238" ],  
        "Name": [ "Asha", "Cathy", "Rosa", "Amy", "Gina" ],  
        "Marks": [ 85, 90, 78, 92, 88 ] }
```

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

mint (df)

	USN	Name	Mark
0	234	Ashra	85
1	235	Cathy	90
2	236	Rosa	78
3	237	Amy	92
4	238	Gina	88

- ### ② Importing datasets from sklearn.datasets

```
from sklearn.datasets import load_diabetes
```

```
import pandas as pd
```

diabetes = load-diabetes()

```
df_diabetes = pd.DataFrame (diabetes.data, columns = diabetes.feature_names)
```

print (df-diabetes)

③ Import datasets from a specific .csv file

import pandas as pd

df_csv = pd.read_csv (" / sample - sales - data . csv ")

print (df_csv. head ())

id	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 ())

ID	No-Patim	gender	AGE	Area	C7
0	502	F	50	4.7	46
1	735	M	26	4.5	62
2	420	F	50	4.7	

BMI CLASS

0	24.0	N
1	23.0	N
2	24.0	N
3	24.0	N
4	21.0	N

To - 00

Stock Market Data Analysis

1. ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]

2. Start date and end date.

3. Plot the closing price and daily return for all three banks.

import yfinance as yf

import pandas as pd

import matplotlib.pyplot as plt

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

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

print("First 5 rows of the dataset")

print(data.head(1))

print("\nShape of dataset")

print(data.shape)

print("\nColumn names")

print(data.columns)

hdfc_data = data['HDFCBANK.NS']

print("\nSummary statistics for HDFC Industries")

print(hdfc_data.describe())

hdfc_data['Daily Return'] = hdfc_data['Close'].pct_change()

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

plt.subplot(2, 1, 1)

hdfc_data['Daily_Return'] = hdfc_data.plot(title = "HDFC Industries - Daily Returns", color='orange')

plt.tight_layout()

plt.show()

WEEK 1

- ① To load .csv file into dataframe

```
import pandas as pd
df = pd.read_csv ('/content/housing.csv')
print (df.head())
print (df.columns)
print (df.info)
```

- ② To display information of all columns

```
print (df.columns())
```

	population	households	median-income	median-house-value	ocean-proximity
0	322.0	126.0	322.0		
1	2401.0	1131.0	8.904		
2	496.0		8.754		
3	558.0	219.0	5.6471		
4	565.0	219.0	3.8462		

- ③ To display statistical information of all numerical

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

- ④ To display the count of unique label for "Ocean Proximity" column

```
print (df.value_counts('ocean-proximity'))
```

- ⑤ To display which attributes in a dataset have missing values.

```
print (df.isnull())
```

off
0 NEAR RAY
1 NEAR RAY
2 INLAND

✓ Gen
05-03-2018

For Diabetes dataset, apply data preprocessing techniques -

import pandas as pd

import numpy as np

import os

import sklearn.preprocessing import LabelEncoder, MinMaxScaler, StandardScaler

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

Check for missing values

print("Missing values in Diabetes dataset : In", diabetes_df.isnull().sum())

Handle missing values with mean

numeric_cols = diabetes_df.select_dtypes(include = np.number).columns

diabetes_df[numeric_cols] = diabetes_df[numeric_cols].fillna(diabetes_df[numeric_cols].mean())

To fill missing categorical values with mode

for col in diabetes_df.select_dtypes(include = ['object']).columns :

diabetes_df[col].fillna(diabetes_df[col].mode()[0], inplace = True)

Identify categorical values and encode them.

categorical_cols = ['Gender', 'CLASS']

label_enc = LabelEncoder()

for col in categorical_cols :

diabetes_df[col] = label_enc.fit_transform(diabetes_df[col])

Apply min max Scaling

min_max_scaler = MinMaxScaler()

diabetes_df_scaled = pd.DataFrame(min_max_scaler.fit_transform(diabetes_df), columns = diabetes_df.columns)

Apply standardization

```

① std_scaler = StandardScaler()
diabetes_df_standardized = pd.DataFrame(standardizer.fit_transform(diabetes_df),
                                         columns=diabetes_df.columns)

i # Save the preprocessed datasets
d diabetes_df_scaled.to_csv("./content/diabetes-preprocessed.csv", index=False)
l vs. Makedir("./content", exist_ok=True)
l diabetes_df_standardized.to_csv("./content/diabetes-standardized.csv", index=False)
l print("Diabetes dataset processing completed. Preprocessed datasets saved.")

```

Q1. Which columns in the dataset had missing values? How did you handle them?

- Numeric columns such as glucose levels, blood pressure, insulin
- Categorical columns such as gender and CLASS

1 To handle numeric data, mean of the repetitive column was used -

2 categorical data, mode

Q2. Which categorical columns did you identify in the dataset? How did you encode?

Gender and CLASS are the two categorical columns.

Q3. Encoding method - Label encoding Male → 0
Female → 1

Q4. What is the difference between min-max scaling and standardization?

- Min Max Scaling - transforms data to fixed range [0,1]

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

Q5. Standardization - transforms data so that it has a mean of 0 and standard deviation of 1

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

19/03

KAR - 03

Build a linear regression model using

- 1) Simple linear regression
- 2) Linear regression in matrix

1) SIMPLE LINEAR REGRESSION

```
import numpy as np
```

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

```
from sklearn.linear_model import LinearRegression
```

```
xi = np.array([1, 2, 3, 4, 5]).reshape(-1, 1)
```

```
yi = np.array([1.2, 1.8, 2.6, 3.2, 3.8])
```

```
model = LinearRegression()
```

```
model.fit(xi, yi)
```

```
m = model.coef_[0]
```

```
c = model.intercept_
```

```
future_weeks = np.array([7, 9]).reshape(-1, 1)
```

```
predicted_sales = model.predict(future_weeks)
```

```
x_range = np.arange(1, 10, 0.1).reshape(-1, 1)
```

```
y_range = model.predict(x_range)
```

```
plt.scatter(xi, yi, color='blue', label='Actual Sales')
```

```
plt.plot(x_range, y_range, color='red', label='Regression Line: y = m * x + c')
```

```
plt.scatter(future_weeks, predicted_sales, color='green', marker='o', label='Predicted Sales (Weeks 7 and 9)')
```

```
plt.xlabel('Weeks')
```

```
plt.ylabel('Sales')
```

```
plt.title('Weekly sales prediction using Linear Regression')
```

```
plt.legend()
```

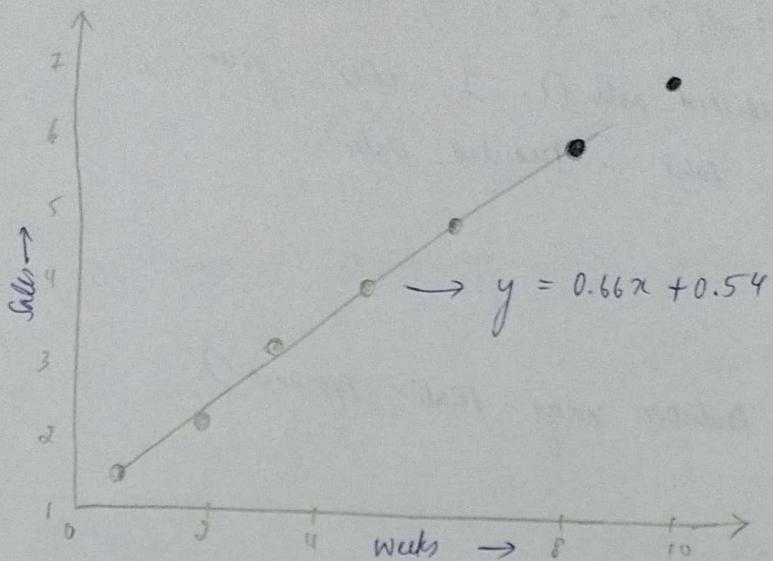
```
plt.grid(True)
```

plt. show()

print(f"Equation of the regression line : $y = (m = .2f^3)x + (c = .2f^3)''$)

print(f"Predicted sales for week 7 : {predicted_sales[0]:.2f}''")

print(f"Predicted sales for week 9 : {predicted_sales[1]:.2f}''")



Equation of regression line : $y = 0.66x + 0.54$

Predicted sales for week 7 : 5.6

week 9 : 6.48

② LINEAR REGRESSION IN MATRIX FORM

```
import numpy as np
```

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

```
x_i = np.array([1, 2, 3, 4])
```

```
y_i = np.array([1, 3, 4, 8])
```

```
x = np.c_[np.ones(len(x_i)), x_i]
```

```
y = y_i.reshape(-1, 1)
```

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

```
c, m = theta.flatten()
```

```
future_week = np.array([1, 7, 9])
```

```
x_future = np.c_[np.ones(len(future_week)), future_week]
```

```
predicted_sales = x_future @ theta
```

$x\text{-range} = \text{np. linspace}(1, 10, 100)$

$y\text{-range} = c + m * x\text{-range}$

$\text{plt. scatter}(x_i, y_i, \text{color} = \text{'blue'}, \text{label} = \text{'Actual Sales'})$

$\text{plt. plot}(x\text{-range}, y\text{-range}, \text{color} = \text{'red'}, \text{label} = f\text{'Regression Line: } y = mx + c)$

$\text{plt. scatter}([7, 9], \text{predicted_sales}[1:], \text{color} = \text{'green'}, \text{marker} = \text{'o'}, \text{label} = \text{'Predicted Sales'})$

$\text{plt. xlabel}(\text{'Weeks'})$

$\text{plt. ylabel}(\text{'Sales'})$

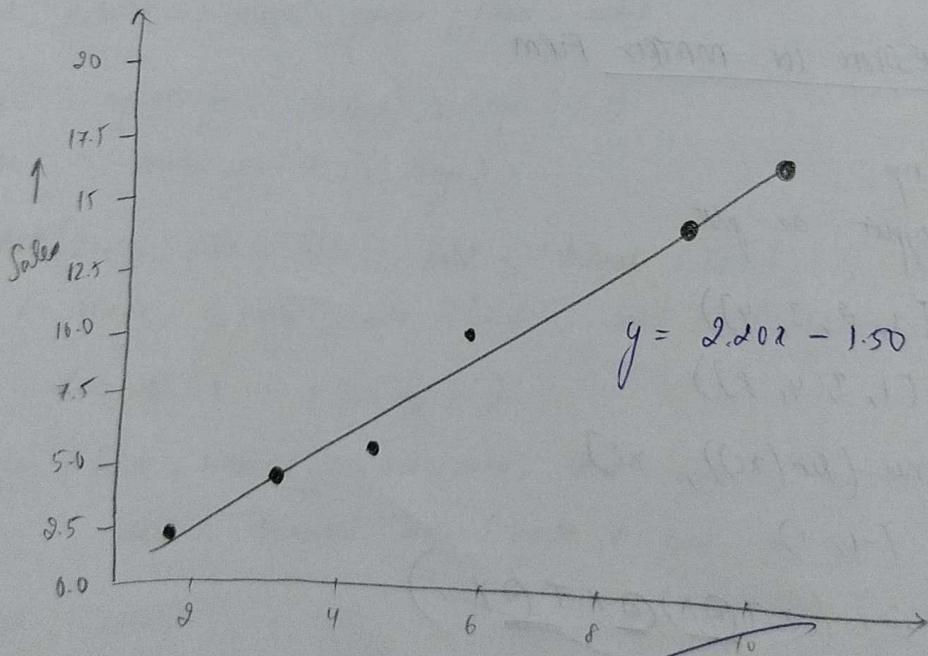
$\text{plt. title}(\text{'Weekly Sales Prediction using Matrix Approach'})$

$\text{plt. legend}()$

$\text{plt. grid}(\text{True})$

$\text{plt. show}()$

Q) Equation of the regression line $y = 2.20x - 1.50$.



Week

Sales
19.03 no 23

02/04

LAB - 03LOGISTIC REGRESSION

- Q1. Consider a binary classification problem where we want to predict whether a student will pass or fail.
- Given $a_0 = -5$ (intercept)
 $a_1 = 0.8$ (coefficient)
- Q2. Write the logistic regression equation for this problem.

$$\text{sigmoid}(z) = p(x) = \frac{1}{1 + e^{-z}} \quad z = a_0 + a_1 x$$

$$z = -5 + 0.8x \quad p(x) = \frac{1}{1 + e^{-(5+0.8x)}}$$

- Q3. Calculate the probability that a student who studies for 7 hours will pass.

$$\text{Probability}(x/\text{pass}) = \frac{1}{1 + e^{5-0.8(7)}} = \frac{1}{1 + e^{-0.6}} = 0.6457$$

- Q4. Determine the predicted class (pass or fail) for this student based on a threshold.

If threshold = 0.5

$$p(x=7) < 0.5$$

Student will ~~fail~~ [pass]

Q2. Consider $z = [2, 1, 0]$ for three classes. Apply softmax function

$$\text{softmax}(z_k) = \frac{e^{z_k}}{\sum_{j=1}^k e^{z_j}}$$

$$\text{softmax}(z_1) = \frac{e^2}{e^2 + e^1 + e^0} = \frac{7.39}{7.39 + 2.72 + 1} \approx 0.665$$

$$\text{softmax}(z_2) = \frac{e^1}{e^2 + e^1 + e^0} = \frac{2.72}{7.39 + 2.72 + 1} \approx 0.244$$

$$\text{softmax}(z_3) = \frac{e^0}{e^2 + e^1 + e^0} = \frac{1}{7.39 + 2.72 + 1} \approx 0.095$$

Probabilities for three classes = ~~66.5%, 24.4%, 9.1%~~.

g. HR - comma - sep. .csv

(i) Which variables did you identify as having a direct and clear impact on employee retention? Why?

Variables such as satisfaction-level, average-monthly-hours, number-project, time-spend-company, salary.

satisfaction level : strong negative correlation

average-monthly-hours : strong positive correlation

number-project : positive correlation

time-spend-company : strong positive correlation

salary : strong negative correlation

(ii) What was the accuracy of your logistic regression model?

The accuracy was around (0.79) (79%).

f. Zoo.csv

(i) Did you perform any data preprocessing steps? If yes, what were they?

① Handling categorical data.

class-type → class-mapping

② Splitting dataset

80% training, 20% testing

③ Drop irrelevant columns

Drop animal_name, class as it is irrelevant.

④ Feature scaling

StandardScaler() to standardize numerical features

$$\text{mean} = 0$$

$$\text{variance} = 1$$

(ii) Were there any missing or inconsistent values in the dataset?

Ans No missing values.

(iii) What does the confusion matrix tell you about the performance of your

The confusion matrix represents how well the model classifies different classes.

Diagonal values : True classification

Other values : Incorrect classification

(iv) Which class types were most frequently misclassified?

Ans Class 1, 5, and 6 were misclassified.

Possible reasons

① Feature similarity

② Small sample size

12/03/25

LAB - 02

1D3 (Iterative Dichotomiser)

import pandas as pd

import math

from collections import Counter

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

df = df. dropna(inplace=True)

def entropy(data):

label = data['label'].tolist()

counts = Counter(label)

probabilities = [count / len(label) for count in counts.values()]

entropy_value = -sum(p * math.log2(p) for p in probabilities if p > 0)

return entropy_value

def gain(data, feature):

initial_entropy = entropy(data)

feature_values = data[feature].unique()

weighted_entropy = 0

for value in feature_values:

subset = data[data[feature] == value]

weighted_entropy += (len(subset) / len(data)) * entropy(subset)

return initial_entropy - weighted_entropy

def id3(data, features, target_attribute):

~~if len(data[target_attribute].unique()) == 1 :~~

~~return data[target_attribute].iloc[0]~~

~~if len(features) == 0 :~~

~~return data[target_attribute].mode()[0]~~

best-feature = max (feature, key = lambda feature : gain (data, feature))

tree = {best-feature : {}}

feature = [f for f in feature if f != best-feature]

for value in data [best-feature].unique () :

subset = data [data [best-feature] == value].drop (columns = [best-feature])

if subset.empty :

tree [best-feature] [value] = data [target-attribute].mode () [0]

else :

tree [best-feature] [value] = id3 (subset, feature, target-attribute)

return tree

target-attribute = 'label'

features = [cn for cn in df.columns if cn != target-attribute]

decision-tree = id3 (df, feature, target-attribute)

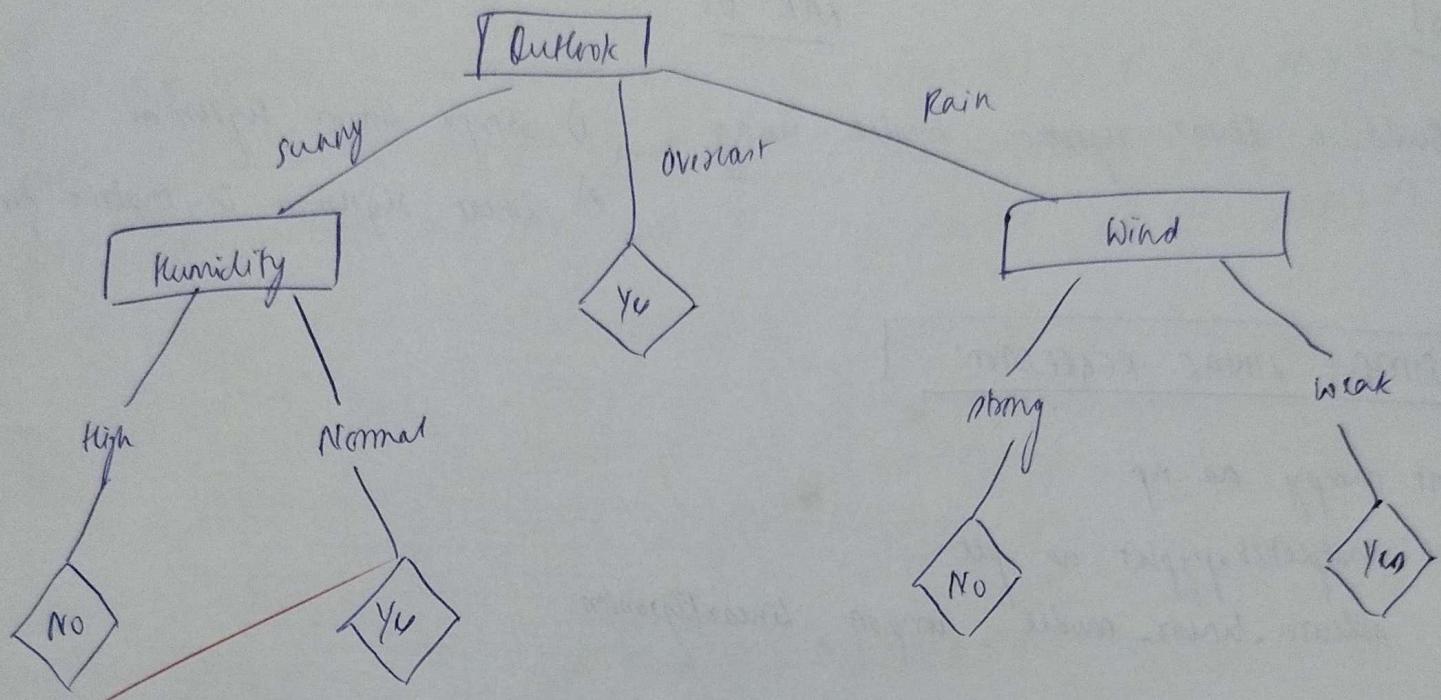
print (decision-tree)

off L'outlook': {'sunny': L'humidity': L'high': 'no', 'normal': 'yes'}

'overcast': 'yes',

'rainy': L'wind': L'weak': 'yes', 'strong': 'no'}

33



Arg(3) / 25

LAB - 04
KNN Classification

Q1. Consider the following dataset
Test data $(x, 35, 100)$ as $(\text{Person}, \text{Age}, \text{Salary})$

Person	Age	Salary	Target	Distance	Rank
A	18	50	N	52.81	5
B	23	55	N	46.52	4
C	24	70	N	31.96	2
D	41	60	Y	40.45	3
E	41	70	Y	31.06	1
F	38	40	Y	60.08	6
X	35	100	?		

For person A, Distance = $\sqrt{(35-18)^2 + (100-50)^2}$
 $= 52.81$

For person B, Distance = $\sqrt{(35-23)^2 + (100-55)^2} = 46.52$

For person C, Distance = $\sqrt{(35-24)^2 + (100-70)^2} = 31.96$

For person D, Distance = $\sqrt{(35-41)^2 + (100-60)^2} = 40.45$

For $k = 1$, target = Y
 $k = 2$ target = N
 $k = 3$ target = Y

\therefore For test data $(x, 35, 100)$ target = Y

Q2. For Iris dataset,

How to choose K value? Demonstrate using accuracy rate and error rate?

Ans To determine best value of K, check ① accuracy rate
② error rate

Best K value : 6

Highest accuracy rate b/w $k=1$ to 20 was 96.67% .

lowest error rate : 3: 33%

Diabetes dataset

What is the purpose of feature scaling?

- Since it is distance based algorithm, if dataset contains features with different scales, the model may be biased toward features with larger numerical values.
- StandardScaler() was applied

$$\text{mean} = 0$$

$$S.D = 1$$

Gen
or. 24. msl

LABRANDOM FOREST CLASSIFIER

Implement the Random Forest classifier using scikit-learn -

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

```
data = load_iris()
```

```
X = data.data
```

```
y = data.target
```

$x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train-test-split}(X, y, \text{test-size} = 0.3,$
 $\text{random-state} = 42)$

rf_classifier = RandomForestClassifier(n_estimators = 10, random_state = 42).

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

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

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

print(f"Accuracy: {accuracy:.4f} %")

✓
16-11-2023

07/05

For "iris. cov" dataset -

Q1. What is the best accuracy score and confusion matrix of the classifier you observed and using how many trees?

Ans

Best observed accuracy score = 1.00 (100%)
Perfect accuracy was achieved using number of trees
n_estimators : 1.

Confusion matrix.

Actual \ Predicted	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	9	0
virginica	0	0	11

07/05/25

LAB - 08Boosting ensemble methodAdaboost algorithm

CGPA	Interactivity	Practical Knowledge	Communication Skill	Job Profile
≥ 9	Yes	Good	Good	Yes
< 9	No	Good	Moderate	Yes
≥ 9	No	Average	Moderate	No
< 9	No	Average	Good	No
≥ 9	Yes	Good	Moderate	Yes
≥ 9	Yes	Good	Moderate	Yes

Step 1 : Consider 4 decision stumps -

D_1 (CGPA), D_2 (Interactivity), D_3 (Practical knowledge), D_4 (Communication skill)

Step 2 : Assign initial weight for each item = $1/6$

Step 3 : Iterate for each weak classifier $(102.0 - 1) / 102.0 \approx 1/102$

Decision stump for 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	$1/6$

$$\epsilon_{CGPA} = 2 \times \frac{1}{6} = 0.333$$

$$Z_{CGPA} = \frac{1}{6} \times 4 \times e^{-0.333}$$

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

$$\frac{1}{6} \times 2 \times e^{0.333}$$

$$= 0.347$$

$$Z_{CGPA} = 0.9428$$

$$d) w_t(d_j)_{i+1} = \frac{\frac{1}{t} \times e^{-0.547}}{0.9428} = 0.1249$$

$$w_t(d_j)_{i+1} = \frac{\frac{1}{t} \times e^{0.547}}{0.9428} = 0.2501$$

I D.R for Interaction

Interaction		Predicted	Actual	Weight
y_u	y_v	y_u	y_v	0.1249
No	No	y_u	y_v	0.2501
No	No	No	No	0.2501
No	No	No	No	0.1249
y_u	y_v	y_u	y_v	0.1249
y_u	y_v	y_u	y_v	0.1249

$$\xi_{\text{Interaction}} = 1 * 0.2501 = \underline{0.2501}$$

$$\alpha_{\text{Interaction}} = \frac{1}{2} \ln \frac{(1 - 0.2501)}{0.2501}$$

$$= \underline{0.5490}$$

$$\begin{aligned} R_{\text{Interaction}} &= 0.1249 * 4 * e^{-0.547} + 0.2501 * 1 * e^{-0.547} \\ &\quad + 0.2501 * 1 * e^{0.547} \\ &= 0.5490 \end{aligned}$$

$$w_t(d_j)_{i+1} = \frac{0.1249 * e^{-0.547}}{0.866} = 0.0892$$

$$w_t(d_j)_{i+1} = \frac{0.2501 * e^{-0.547}}{0.866} = 0.1667$$

$$w_t(d_j)_{i+1} = \frac{0.2501 * e^{0.547}}{0.866} = 0.5001$$

III AS for [practical knowledge]

- No misclassification

IV AS for [communication skill]

$\alpha_{CGRM} = 0.347$	$\alpha_{minor \ art} = 0.5470$	$\alpha_{CMM} = -0.5465$	Final prediction
y_1	y_1	y_1	y_1
No	No	No	No
y_1	No	No	y_1
No	No	y_1	No
y_1	y_1	No	y_1
y_1	y_1	0.0	y_1

Income. cov

Q1. What is the best accuracy score and confusion matrix of the classifier?

Ans But accuracy score = 0.8140

No. of estimators = 73

Confusion matrix (n - estimators = 10)

		0	1	
0	6782	632	$(0.1 + 0.8 + 0.1) = 1.0$	
1	1144	1211	$(0.1 + 0.1 + 0.1) = 0.3$	

S

07/05/2025

KAB-07

Kmean clustering

- Q1. For the given data, compute 2 cluster using K mean algorithm for clustering when initial cluster centers are $(1.0, 1.0)$ and $(5.0, 7.0)$

$$\text{Adm} \quad \begin{array}{l} \text{No. of clusters} \\ \text{centrid } C_1 = (1.0, 1.0) \\ \text{centrid } C_2 = (5.0, 7.0) \end{array}$$

Record Number	close to C_1	close to C_2	Assign to cluster
$R_1 (1.0, 1.0)$	0.0	7.21	cluster
$R_2 (1.5, 2.0)$	1.12	6.12	cluster 1
$R_3 (3.0, 4.0)$	3.61	3.61	cluster 1
$R_4 (5.0, 7.0)$	7.21	0.0	cluster 2
$R_5 (3.5, 5.0)$	4.12	2.5	cluster
$R_6 (4.5, 5.0)$	5.71	2.06	cluster 2
$R_7 (3.5, 4.5)$	4.30	2.92	cluster 2

$$\text{cluster 1} = \{R_1, R_2, R_3\}$$

$$\text{cluster 2} = \{R_4, R_5, R_6, R_7\}$$

Now centrid,

$$C_1 = \frac{(1.0 + 1.5 + 3.0)}{3}, \quad \frac{(1.0 + 2.0 + 4.0)}{3}$$

$$= [1.83, 2.33]$$

$$C_2 = [4.12, 5.37]$$

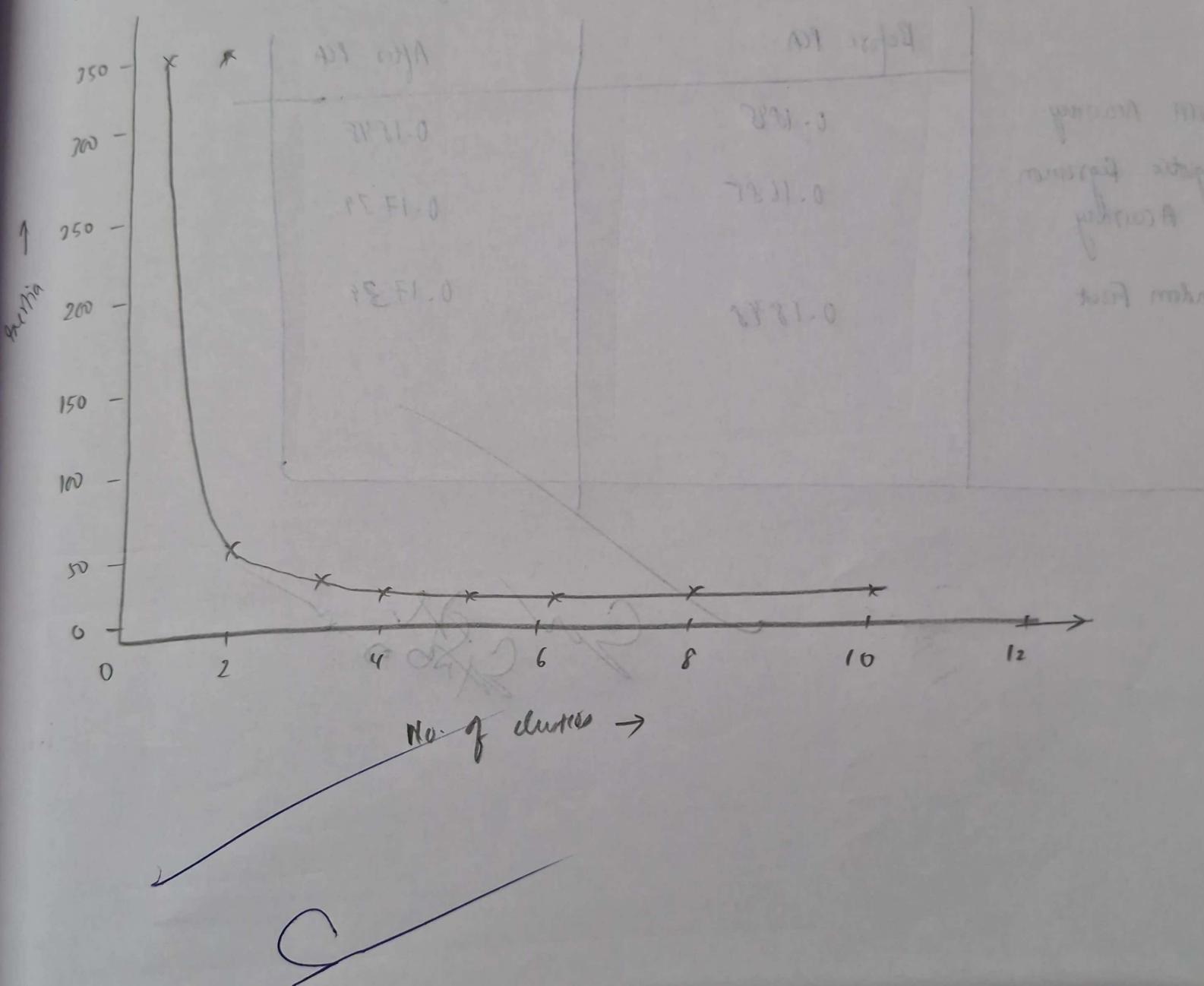
Ques. UV

Q1 - 831

Draw the elbow plot. What was the optimal k value obtained.

Q2 - 1000

Optimal value of $k = 3$



07/05/25

$$\frac{\text{LAB} - 10}{\text{PCA}}$$

[heart.csv]

Report the accuracy now before and after applying PCA

	Before PCA	After PCA
SVM Accuracy	0.1048	0.1848
Logistic Regression Accuracy	0.1685	0.1739
Random Forest	0.1848	0.1739

Surprise! The accuracy increased after applying PCA.