

Name Shreya Mittal Std _____ Sec G E

Roll No. _____ Subject ML-LAB School/College _____

School/College Tel. No. _____ Parents Tel. No. _____

Sl. No.	Date	Title	Page No.	Teacher Sign / Remarks
1.	5/3/25	Lab 0 [todo list]	1-2	1/3/25
2.	5/3/25	Lab 1 [data processing]	3-4	1/3/25
3.	12/3/25	Lab 2 [ID3]	5-7	1/3/25
4.	19/3/25	Lab 3 [Linear Reg < Normal]	8-9	1/3/25
5.	2/4/25	Lab 4 [logistic Reg & Binary]	10-11	1/3/25
6.	2/4/25	Lab 5 [KNN]	12-13	1/3/25
7.	16/4/25	Lab 6 [Random forest]	14-15	1/3/25
8.	21/5/25	Lab 7 [Random forest]	15-15	1/3/25
9.	21/5/25	Lab 8 [Boosting - Method]	16-17	
10.	21/5/25	Lab 9 [k-means Algo]	18-19	
11.	21/5/25	Lab 10 [dimensionality red using PCA method]	20-21	

Lab = 0

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

To do:

Method 1: Initializing values directly into DataFrame.
Insert your known values, five rows of data with column headings as 'USN', 'Name', 'Marks'

→ import pandas as pd

```
data = {  
    'USN': ['IBM202266', 'IB-267', 'IB-268'],  
    'Name': ['Shreya', 'Shr', 'Shree', 'Shri'],  
    'Marks': [99, 98, 97, 96]  
}
```

df = pd.DataFrame(data)

print(df)

Method 2: Importing datasets from sklearn.datasets

Loading diabetes datasets sklearn.datasets.load_diabetes

→ from sklearn.datasets import load_diabetes

diabetes = load_diabetes()

df = pd.DataFrame(diabetes, columns = diabetes.
feature_names)

df['target'] = diabetes.target

print(df)

Method 3: Importing datasets from specific .csv file

~~file sample-sales-data.csv~~

→ path = r"/content/sample-sales-data.csv"

df = pd.read_csv(path)

print(df)

Method 4: Downloading datasets from existing dataset repositories like Kaggle, UCI, Moulody, KEEI, etc

→ path = r"/content/datasetofdiabetes.csv"

df = pd.read_csv(path)

print(df)

TODO:

1. HDFC Bank Ltd, ICICI Bank Ltd, Kotak Mahindra Bank Ltd
 tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTBANK.NS"]
 → import yfinance as yf
 import matplotlib.pyplot as plt
 tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTBANK.NS"]
2. start date: 2024-01-01, End date: 2024-12-31
 → data = yf.download(tickers, start='2024-01-01', end='2024-12-31', group_by=tickers)
 print(data)
3. Plot the closing price & daily returns for all the three banks mentioned
 → # HDFC BANK
 HDFC = data['HDFCBANK.NS']
 HDFC['Daily Return'] = HDFC['close'].pct_change()
 plt.figure(figsize=(12, 6))
 plt.subplot(2, 1, 1)
 HDFC['close'].plot(title='HDFCBANK - closing Price')
 plt.subplot(2, 1, 2)
 HDFC['Daily Return'].plot(title='HDFCBANK - Daily Return', color='orange')
 plt.tight_layout()
 plt.show()
- # KOTAK BANK
 KOTAK = data['KOTAKBANK.NS']
 KOTAK['Daily Return'] = KOTAK['close'].pct_change()
 plt.figure(figsize=(12, 6))
 plt.subplot(2, 1, 1)
 KOTAK['close'].plot(title='KOTAKBANK - closing Price')
 plt.subplot(2, 1, 2)
 KOTAK['Daily Return'].plot(title='KOTAKBANK - Daily Return', color='orange')
 plt.tight_layout()
 plt.show()

Lab = 1

Demonstrate various data pre-processing techniques for a given dataset

- i) To load .csv file into the data frame
→ import pandas as pd
import numpy as np

~~path = r'C:\housing.csv'~~
df = pd.read_csv(path)
print(df)

- ii) To display info of all the columns
→ print(df.info())

- iii) To display statistical info of all numerical
→ print(df.describe())

- iv) To display the count of unique labels for
"Ocean Proximity" column
→ print(df[["Ocean Proximity"]].value_counts())

- v) To display which attributes (columns) in a dataset
have missing values count greater than zero.
→ missing_values = df.isnull().sum()
print(missing_values[missing_values > 0])

1. Which columns in the dataset had missignvalues?
How did you handle them?

- None of colm had missignvalues, but If had, then
numerical column's NaN can be replaced with median
and categorical column's null can be replaced
with mode.

Q&A

2. Which categorical column did you identify in the dataset? How did you encode them?
- The diabetes dataset had gender, and class as categorical columns and adult dataset had workclass, education, marital-status, occup, reltⁿ, race, gender, native-country and income as categorical columns.
- Using Ordinal encoder, we can encode the categorical columns.
3. What is the diff b/w Min-Max scaling & standardization? When should you use either the
-
- Min-Max scaling, also called ~~normalization~~^{other}, transforms data to fit within a specific range (usually 0 to 1) by scaling based on the min-max values in the dataset.
 - Standardization scales data by subtracting mean and dividing by the standard deviation, resulting in distribution with a mean of 0 and standard deviation of 1.
 - Use min-max scaling when you need your data to be in specific range, while standardization, when your data is normally distributed and outliers might be present.

~~5/27/25~~
~~5/27/25~~

Lab = 2

ID 3

```
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
    self.value = value
    self.label = label
    self.children = {}
```

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

```

```

def print_tree(node, depth=0):
    if node.label is not None:
        print(f'{depth} {node.label}')
    else:
        print(f'{depth} {node.feature}')
        for value, child in node.children.items():
            print(f'{depth + 1} {value}')
            print_tree(child, depth + 1)

```

example dataset

```
data = pd.DataFrame( {  
    'outlook': ['S', 'S', 'O', 'Rain', 'Sun', ...],  
    'Temperature': ['Hot', 'Mild', 'Cold', ...],  
    'Humidity': ['High', 'Normal', ...],  
    'Wind': ['Weak', 'Strong', ...],  
    'PlayTennis': ['No', 'Yes', ...]  
})
```

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)

O/P

feature: Outlook

value: 0

feature: Humidity

value: 0

leaf: 0

value: 1

leaf: 1

value: 1

leaf: 1

value: 2

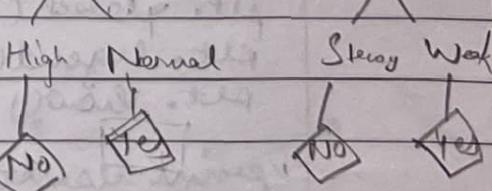
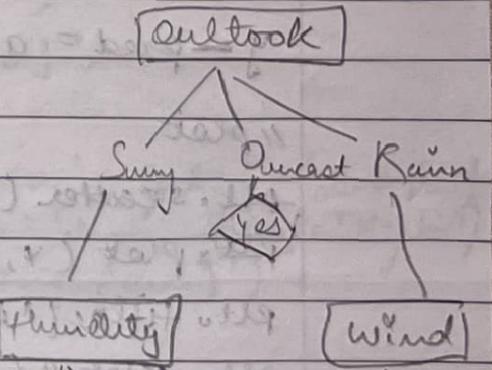
feature: Wind

value: 0

leaf: 0

value: 1

leaf: 0



Lab = 3Linear Regression— using eq^o (normal) form:

import numpy as np

import matplotlib.pyplot as plt

 $x = np.array([1, 2, 3, 4, 5])$ $y = np.array([1.2, 1.8, 2.6, 3.2, 3.8])$ $n = len(x)$ no. of data points $\text{sum-}x = np.sum(x)$ $\text{sum-}y = np.sum(y)$ $\text{sum-}xy = np.sum(x * y)$ $\text{sum-}x^2 = np.sum(x ** 2)$

$$a_1 = \frac{(n * \text{sum-}xy - \text{sum-}x * \text{sum-}y)}{n * \text{sum-}x^2 - \text{sum-}x^2}$$

$$a_0 = (\text{sum-}y - a_1 * \text{sum-}x) / n$$

$$y = \text{pred} = a_1 * x + a_0$$

"plot"

plt.scatter(x, y, color='blue', label='data points')

plt.plot(x, ypred, color='red', label='Regre. line!')

plt.title('Linear Regression (using formula!)')

plt.xlabel('x')

plt.ylabel('y')

plt.legend()

plt.show()

"print result"

print(f" slope(a₁): {a1} ")" " (f" Intercept(a₀): {a0} ")" " (f" prediction(y = a₁ * x + a₀): ")

" " (f" x={x}, y-pred={y-pred} ")

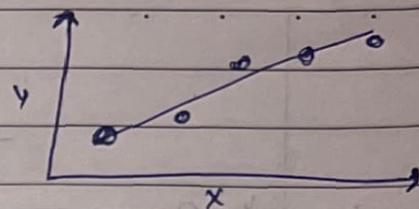
O/P

$$\text{slope}(a_0) = 0.66$$

$$\text{intercept}(a_1) = 0.54$$

$$\text{Predictor } (y = a_1 \times x + a_0)$$

$$x = 9, y_{\text{pred}} = 3.84$$



- using Matrix form:

import numpy as np.

import matplotlib.pyplot as plt

~~x = np.array([1 2 3 4 5 ...])~~

~~y = np.array([1.2, 1.8, 2.4, 2.8, ...])~~

"create matrix [1, x]"

x-matrix = np.c_[np.ones(len(x)), x]

"calculate theta (slope & intercept) using normal eqn"

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

"extract intercept & slope"

b, m = theta

"make prediction"

y-pred = m*x + b

"print result"

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

print(f" intercept (b): {b} ")

"plot"

~~plt.scatter(x, y, color='blue', label='Actual')~~

~~plt.plot(x, y-pred, color='red', label='Reg line')~~

~~plt.title('Linear Reg(Matrix form)')~~

~~plt.xlabel('X')~~

~~plt.ylabel('Y')~~

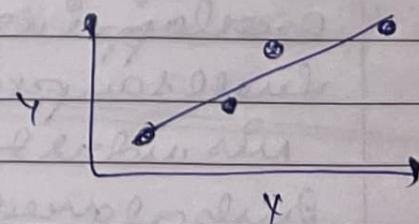
~~plt.legend()~~

~~plt.show()~~

O/P

slope (m): 0.66

intercept (b): 0.54



Lab-4 Logistic Regression

(1) "HR common - sep"

→ which variable did you identify as having a direct & clear impact on separation?

- Satisfaction level : low satisfy → high attrition
- No. of projects : too few or too many → high attrition
- Avg. monthly hours : extreme hours → high attrition
- Education status : * promotion in 5 years → high attrition
- Salary : low salary → higher attrition

→ what was the accuracy of log. reg model?
is it good accuracy?

• Accuracy : 76 - 80 %

→ it is decent but not perfect,
correct prediction 3 out of 4 cases.

e.g.: employee retention is influenced by
many complex factors that not in dataset

- missing external factors (job market, ..)
- class imbalance may effect prediction
- feature interactions (e.g. salary + work hours)
may not be captured well.

(2) Zoo dataset

(i) which class type most frequently misclassified?
similar species caused confusion.

overlapping features

small sample size for rare animals leads to
misclassification

misclassified classes -

(2) Zoo dataset

- Performed data preprocessing steps?
• If yes then what & why necessary?
- Steps:

S (1)

- Converted categorical data into numerical.
— as ML model works with numerical data.

S (2)

- checked for missing values by handled them
— as missing values can cause errors or biased predictions

S (3)

- standardized features.

- as some features may have diff scales,
thus it ensures equal contributions to model learning.

→

- Any missing or inconsistent value? & how to handle
• No, there were no missing or inconsistent values in this dataset but if any, then
→ filling missing values using mode/mean

- Standardize levels to maintain consistency.

→

- confusion matrix tells?

- Show how well the model classifies each animal type

- High diagonal values → good prediction

- Off diagonal values → misclassification

Sub = 5
KNN

Q

Value for $K=3$, ($\times 35, 100$)

Person	Age	Salary K	Target	Distance	Rank
A	18	50	N	52.83	5
B	33	55	N	46.58	4
C	24	70	N	31.96	2
D	41	60	Y	40.44	3
E	43	78	Y	31.00	1
F	38	40	Y	80.08	6
X	35	100	?		

$$\text{distance} = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$\text{Rank 3} \rightarrow \begin{array}{l} E \rightarrow Y \\ C \rightarrow N \\ D \rightarrow Y \end{array} \quad \left. \begin{array}{l} x \rightarrow y \\ \hline \end{array} \right.$$

Assume: Rer $\times (35, 100)$ Target = Y (yes)

① Divide dataset

→ how to choose value? use Accuracy Rate & Error Rate

Choose K

start with diff values of K (eg

split the dataset into training & testing sets

train the KNN model with each K &
calculate the accuracy

use Acc Rate & Error Rate

- Accuracy Rate: measure how well the model classifies the data
- Error Rate: Error Rate = 1 - Acc. Rate
- Plot K vs Accuracy & K vs Error Rate to find optimal K.
- Small K may lead to overfitting, large K may lead to underfitting.

(2) Diabetes Dataset

→ What is purpose of feature scaling?
How to perform it?

- Purpose:
- ensures all features contribute equally to model
- improves the performance of distance-based Algs
- prevent dominance of features with large numerical values.

Perform:

- Standardization (Standard Scale)

$$\{ \mu = \text{mean} \quad \} \quad \{ \sigma = \text{std.} \quad \}$$

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

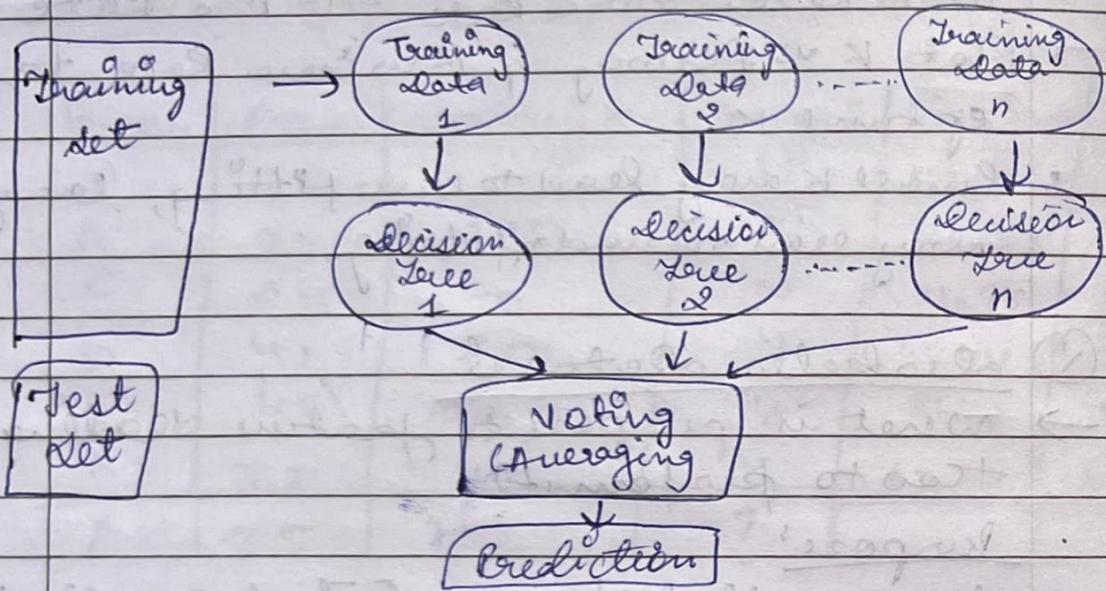
- Normalization (Min-Max Scaling) (Min-Max Scale)

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

scale data range 0 to 1

Lab 6

Random Forest



→ The following steps are the working of Random Forest Algorithm:

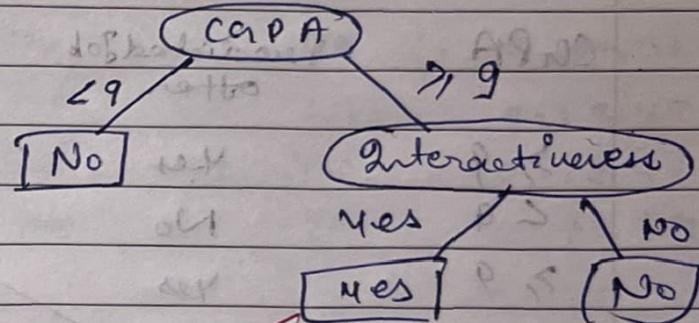
- Step 1: Select random samples from a given data or training set.
- Step 2: This Algo will construct a decision tree for every training data.
- Step 3: Voting will take place by averaging the decision tree.
- Step 4: Finally, select the ~~most in selected~~ most voted prediction result as the final prediction result.

Ge
16.04

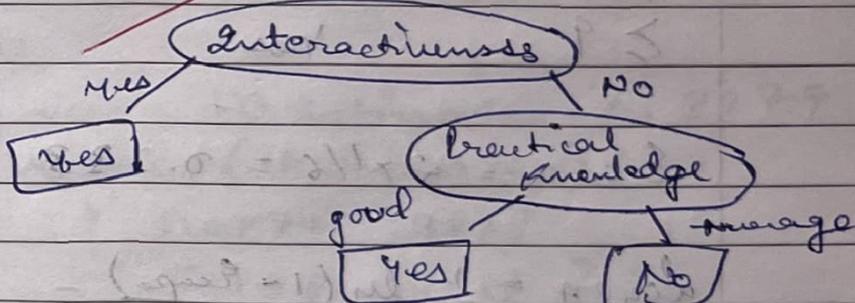
Lab = 7

Q Implement random forest useful method

For $S_1 \rightarrow$



For $S_2 \rightarrow$



Q Write answers.

(1) what is the best accuracy score & confusion matrix of the classification you observed & using how many trees?

→ Best Accuracy Score Observed = 1.0

→ No. of Trees (n = estimators) = best-n.

→ Confusion matrix = $\begin{bmatrix} 16 & 0 & 0 \\ 0 & 14 & 0 \\ 0 & 0 & 15 \end{bmatrix}$

Lab - 8

Q Implement Boosting ensemble method

Ch PA	Predicted Job Offer	Actual Job Offer	Weight
≥ 9	Yes	Yes	1/6
$04 < 9 \leq 11$	No	Yes	1/6
$(a) \geq 9$	Yes	No	1/6
< 9	No	No	1/6
≥ 9	Yes	Yes	1/6
≥ 9	Yes	Yes	1/6

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

$$\alpha_{\text{GPA}} = \frac{1}{2} \ln \frac{(1 - E_{\text{GPA}})}{E_{\text{GPA}}} = \frac{1}{2} \cdot \frac{(1 - 0.13)}{0.33} = 0.347$$

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

$$w_{t+1}(d_j) = \frac{\frac{1}{6} \times e^{-0.347}}{0.9428} = 0.1249$$

$$w_t(d_j) = \frac{\frac{1}{6} \times e^{0.347}}{0.9428} = 0.4501$$

CGPA	Predicted Job offer	Actual job offer	Weight
≥ 9	Yes	Yes	0.1249
< 9	No	Yes	0.1249
7.9	Yes	No	0.12501
< 9	No	No	0.1249
≥ 9	Yes	Yes	0.1249
≥ 9	Yes	Yes	0.1249

Q: Write answer.

1) Best Acc Score & Confusion Matrix

→ Accuracy with 10 estimates = 0.8277
confusion matrix (10 estimates)

$$\begin{bmatrix} [10722 & 387] \\ [2138 & 1408] \end{bmatrix}$$

Best Accuracy : 0.831 with n-estimates = 42

∴ we obtained our result

$$E[(0.4 + 0.3 + 0.1)] = E(0.8 + 2.1 + 0.1) = 0$$

$$88.8 - E(0.8) = 88.1 - 85.2 =$$

$$\sqrt{(0.2 + 2 + 0)} = \sqrt{(2.4 + 2.4 + 2.8 + 0.2)} = \sqrt{8.8} = 2.955$$

$$\sqrt{1.4} = \sqrt{1^2 + 0.4} =$$

Lab = 9

Q) Build K-Means algorithm to cluster

Iteration 1:

Record NO.	closest to C1 (1.0, 1.0)	closest to C2 (5.0, 2.0)	Assigned Cluster
R ₁ (1.0, 1.0)	0.0	7.21	C ₁
R ₂ (1.5, 2.0)	1.12	6.12	C ₂
R ₃ (3.0, 4.0)	3.61	3.61	C ₁
R ₄ (5.0, 2.0)	7.21	0.0	C ₂
R ₅ (3.0, 5.0)	4.12	2.5	C ₂
R ₆ (4.5, 5.0)	5.31	2.06	C ₂
R ₇ (3.5, 4.5)	4.30	2.92	C ₂

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

Their new centroids are:

$$C_1 = (1.0 + 1.5 + 3.0)/3, \quad (1.0 + 2.0 + 4.0)/3 \\ = 5.5/3 = 1.83, \quad 7.0/3 = 2.33$$

~~$$C_2 = (5.0 + 3.5 + 4.5 + 3.5)/4, \quad (7 + 5 + 5 + 4)/4 \\ = 16.5/4 = 4.12, \quad 21.0/4 = 5.25$$~~

Iteration 2:

Record No	Close to C ₁ (1.0, 8.3, 2.33)	Close to C ₂ (4.12, 5.37)	Assigned to cluster
R ₁ (1.0, 1.0)	1.57	5.64	C ₁
R ₂ (1.5, 2.0)	0.47	4.52	C ₁
R ₃ (3.0, 4.0)	2.12	1.63	C ₂
R ₄ (5.0, 7.0)	5.57	1.91	C ₂
R ₅ (3.5, 5.0)	3.16	0.73	C ₂
R ₆ (4.5, 5.0)	3.58	0.52	C ₂
R ₇ (3.5, 4.5)	2.63	1.05	C ₂

Cluster 1 {R₁, R₂} & cluster 2 {R₃, R₄, R₅, R₆, R₇}.
their new centroids are -

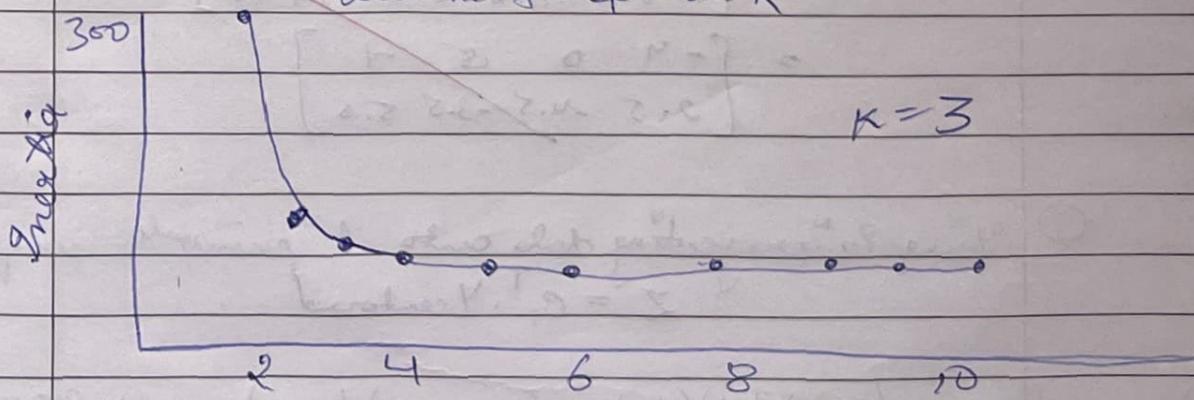
$$\begin{aligned}C_1 &= (1.0 + 1.5)/2, (1.0 + 2.0)/2 \\&= 2.5/2 = 1.25, 3/2 = 1.5\end{aligned}$$

$$\begin{aligned}C_2 &= (3+5+3.5+4.5+3.5)/5, (4+7+5+5+4.5)/5 \\&= 19.5/5 = 3.9, 25.5/5 = 5.1\end{aligned}$$

Q) Write answer.

- 1) draw the elbow plot. & K value?

elbow Plot for optimal K



Lab = 10

Q Implement dimensionality reduction using PCA method

feature	Example 1	Example 2	Example 3	Example 4
x_1	4	8	13	7
x_2	11	4	5	14

Eigen values :

$$\lambda_1 = 30.3849$$

$$\lambda_2 = 6.6151$$

Eigen Vectors -

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

$$e_2 = \begin{bmatrix} 0.8303 \\ 0.5514 \end{bmatrix}$$

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

② Mean center the data -

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

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

$$\begin{aligned} \text{Y centered} &= \begin{bmatrix} 4-8 & 8-8 & 13-8 & 7-8 \\ 11-8.5 & 4-8.5 & 5-8.5 & 14-8.5 \end{bmatrix} \\ &= \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix} \end{aligned}$$

③ Using Eqⁿ projecting data onto 1st principle component -

$$Z = e_1^T \cdot Y_{\text{centered}}$$

$$Z_1 = (0.5574)(-4) + (0.8303)(-4.5) = -4.38535$$

$$Z_2 = (0.5574)(0) + (0.8303)(-4.5) = 3.73635$$

$$Z_3 = (0.5574)(5) + (-0.8303)(-3.7) = 5.69305$$

$$Z_4 = (0.5574)(-1) + (-0.8303)(5.3) = -5.12405$$

$$Z = [-4.80535 \quad 3.73635 \quad 5.69305 \quad -5.12405]$$

Q Write answers -

1) Report the accuracy of scores before & after PCA

⇒ Model accuracy :

Before :

SVM : 0.8804

Logistic Regression : 0.8533

Random Forest : 0.8859

After :

SVM : 0.8424

Logistic Regression : 0.8641

Random Forest : 0.8533

~~DATA
17/5/2025~~