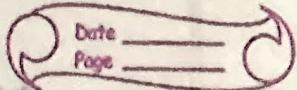


05/03

Week-0



TO-DO EXERCISES

- ① Initializing values directly into Dataframe - files means of data with column 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)

print(df)

	USN	Name	Marks
0	234	Asha	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")

	age	sex	bmi	bp	s1	s2	s3	s4	s5	s6
0	0.03807	0.50680	0.061696	0.02870.061	-0.034	-0.043	-0.017			
1	-0.00882	-0.04472	-0.052474							
2	0.085295	0.050680	0.061696	0.02870.061	-0.034	-0.043	-0.017			
3										
4										

- ③ Import datasets from a specific URL "http://...".csv files
- ```
import pandas as pd
df = pd.read_csv("sample-sales-data.csv")
print(df.head())
```

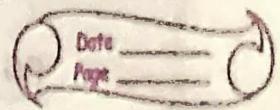
|   | Product  | Quantity | Price | Sales | Region |
|---|----------|----------|-------|-------|--------|
| 0 | Laptop   | 5        | 1000  | 5000  | North  |
| 1 | Mouse    | 15       | 20    | 300   | West   |
| 2 | Keyboard | 10       | 50    | 500   | East   |
| 3 | Monitor  | 8        | 200   | 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())
```

⑤ No\_Patients Gender Age urea or

|   |     |       |    |    |     |     |    |
|---|-----|-------|----|----|-----|-----|----|
| 0 | 502 | 179   | 75 | F  | 150 | 4.7 | 46 |
| 1 | 735 | 34221 | M  | 26 | 4.5 | 62  |    |
| 2 | 420 | 47975 | F  | 50 | 4.7 |     |    |



## BMI CLASS

0 24.0 N

1 23.0 N

2 24.0 N

3 24.0 N

4 24.0 N

## 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("Shape of dataset")
```

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

```
print("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()  
plt.figure(figsize=(12, 6))

plt.subplot(2, 1, 1)

hdfc\_data['Daily Return'] = plot  
(title="HDFC Industries' Daily  
Returns", color='orange')

plt.figure\_layout()

plt.show()

## Week-1

To load .csv file into dataframe

import pandas as pd

df = pd.read\_csv('content/housing.csv')

To display information of all column

print(df.columns())

population household median income

|   | population | household | median income |
|---|------------|-----------|---------------|
| 0 | 322.0      | 126.0     | 322.0         |
| 1 | 2401.0     | 1138.0    | 8.304         |
| 2 | 496.0      |           | 8.754         |
| 3 | 558.0      | 219.0     | 5.6471        |
| 4 | 565.0      | 219.0     | 318462        |

To display statistical information of all numerical

print(df.info)

To display the count of unique labels for "Ocean Proximity" column

print(df.value\_counts('ocean\_proximity'))

To display which attributes in a dataset have missing values

print(df.isnull())

off

0

ocean proximity

NEAR BAY

1

NEAR BAY

2

INLAND

## ID3 - Using weather dataset

Import pandas as pd

data = pd.read\_csv('Tennis.csv')

data

|   | outlook  | temp | humidity | windy | play |
|---|----------|------|----------|-------|------|
| 0 | Sunny    | hot  | high     | false | no   |
| 1 | Sunny    | hot  | high     | true  | no   |
| 2 | overcast | hot  | high     | false | yes  |
| 3 | raining  | mild | high     | false | yes  |
| : | :        | :    | :        | :     | :    |

from sklearn.preprocessing import LabelEncoder

outlook = LabelEncoder()

temp = LabelEncoder()

humidity = LabelEncoder()

windy = LabelEncoder()

play = LabelEncoder()

data['outlook'] = outlook.fit\_transform(data['outlook'])

data['temp'] = outlook.fit\_transform(data['temp'])

data['humidity'] = outlook.fit\_transform(data['humidity'])

data['windy'] = outlook.fit\_transform(data['windy'])

`data['play'] = outlook. fit_transform([data['play']]`

`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    | 1        | 1     | 0    |
| : | :       | :    | :        | :     | :    |

`features = cols = ('outlook', 'temp', 'humidity', 'windy')`

`x = data[features - cols]`

`y = data['play']`

`print(x)`

`print(y)`

`print(y)`

|  | outlook | temp | humidity | windy |
|--|---------|------|----------|-------|
|--|---------|------|----------|-------|

|   |   |   |   |   |
|---|---|---|---|---|
| 0 | 2 | 1 | 0 | 0 |
| 1 | 2 | 1 | 0 | 1 |
| 2 | 0 | 1 | 0 | 0 |
| 3 | 1 | 2 | 0 | 0 |
| 4 | 1 | 0 | 1 | 0 |
| 5 | 0 | 0 | 1 | 1 |
| : | : | : | : | : |

0 0 0 0 0  
1 1 0 0 0  
2 1 1 1 0  
3 1 1 1 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 |
|---------|------|----------|-------|
| 0       | 1    | 2        | 0     |
| 1       | 2    | 1        | 0     |
| 2       | 1    | 1        | 1     |

y-test

|          |      |     |
|----------|------|-----|
| 822      | play | no  |
| 935      | play | yes |
| 020      | play | yes |
| 01000000 | play | yes |

classifier score ( $x$ -test,  $y$ -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

feature\_importance\_df = feature\_importance\_df  
df.sort\_values

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

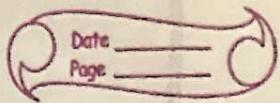
print(feature\_importance\_df)

## Feature Information Gain

|   |          |           |
|---|----------|-----------|
| 0 | outlook  | 0.574558  |
| 3 | windy    | 0.1279205 |
| 2 | humidity | 0.211237  |
| 1 | temp     | 0.000000  |

19/8/25

Week-03



## Linear Regression

Import numpy as np

Import matplotlib.pyplot as plt

from sklearn.linear\_model import

LinearRegression

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

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

```
model = LinearRegression()
```

```
model.fit(x, y)
```

Intercept = model.intercept -

coefficient = model.coef\_[0]

```
print(f"Intercept (a0): {intercept}")
```

```
print(f"Coefficient (a1): {coefficient}")
```

```
x_new = np.array([5])
```

```
y_pred = model.predict(x_new)
```

```
print(f"Predicted y for x=5: {y_pred[0]}")
```

```
plt.scatter(x, y, color='blue', label='Data')
```

~~y\_line = model.predict(x)~~

```
plt.plot(x, y_line, color='red', label='Regression Line')
```

```
plt.xlabel('x')
```

```
plt.ylabel('y')
```

pre. fitRe ('Linear Regression')

pre. legend()

pre. show()

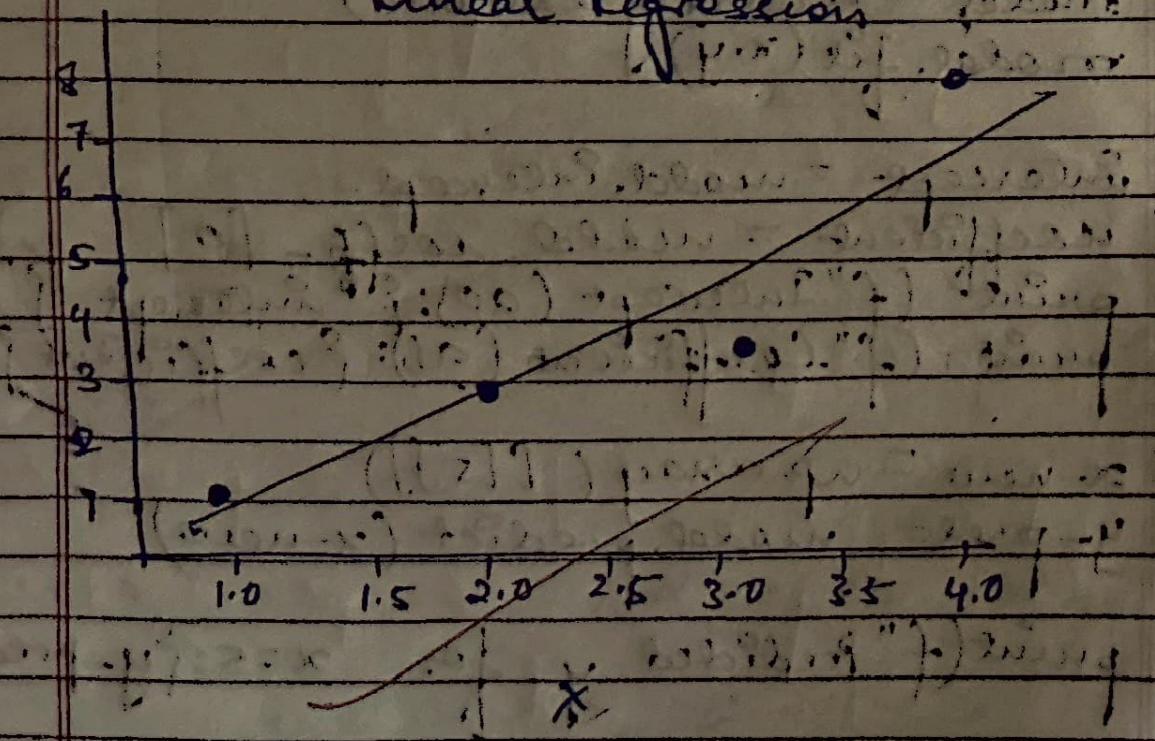
output:

Intercept (a0): -1.4999999...

Coefficient(a1): 2.199999...

Predicted y for x=5: 9.5

Linear Regression:



import numpy as np  
import matplotlib.pyplot as plt

$x = np.array([1, 2, 3, 4])$

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

$x = np.vstack([np.ones(len(x)), x]).T$

$\theta = np.linalg.inv(x.T \cdot \text{dot}(x)) \cdot \text{dot}(x.T) \cdot \text{dot}(y)$

Intercept =  $\theta[0]$

slope =  $\theta[1]$

print(f"Intercept (a0): {Intercept}")

print(f"slope (a1): {slope}")

$x\_new = np.array([5])$

$x\_new = np.array([1, x\_new[0]])$

$y\_pred = x\_new \cdot \text{dot}(\theta)$

print(f"Predicted y for x=5: {y\\_pred[0]}")

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

$y\_line = x \cdot \text{dot}(\theta)$

plt.plot(x, y\_line, color='red', label='Regression Line')

plt.xlabel('X')

plt.ylabel('Y')

plt.title('Different Matrix Method')

plt.legend()

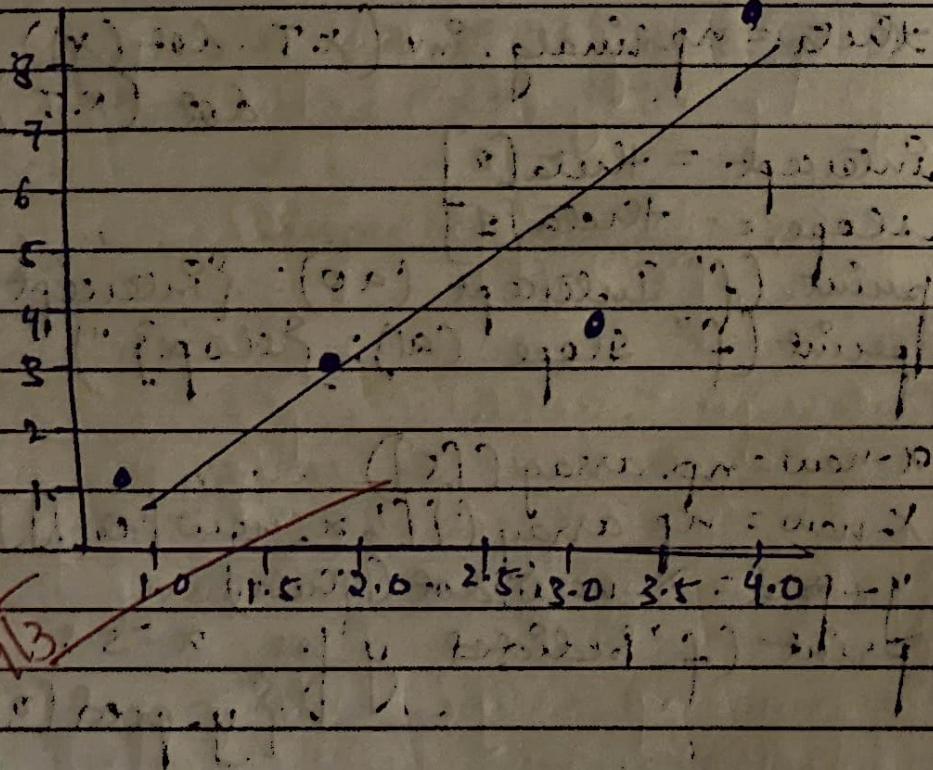
plt.show()

Intercept ( $a_0$ ): -1.5

Slope (a1): 12, 20, 00, 26, 00.

Predicted  $y$  for  $n=5$ : 9.500...04

## Different Matrix Methods



1000 P. 2000

~~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

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

$y_{predicted}$

`model.score(X-test, y-test)`

`model.predict_proba(X-test)`

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

$y_{predicted}$

`model.coef_`

`model.intercept_`

`import math`

`def sigmoid(x):`

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

`def prediction_function(age):`

$$z = 0.127 * \text{age} - 4.943$$

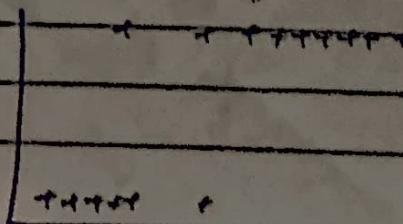
`y = sigmoid(z)`

`return y`

`age = 35`

`prediction_function(age)`

`0.3709834`



To do:

Q1.4 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 ( $a_0 = -5$  (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

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

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

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

by  $\alpha = F$

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

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

$\therefore$  probability = 0.646 (64.68%)

c) Predicted class : Pass

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

$$p_p = \frac{e^{z_1}}{\sum_j e^{z_j}}$$

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

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

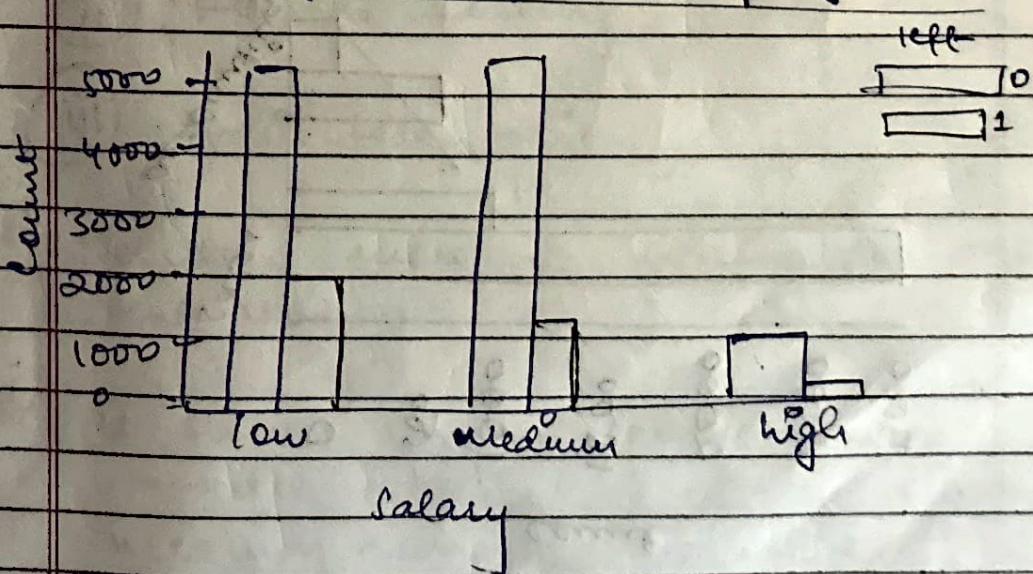
$$\text{Class 3 } (z=0) = 0.090 \text{ (9.00%)}$$

To do: Implementation - Logistic Regression  
(Binary classification)

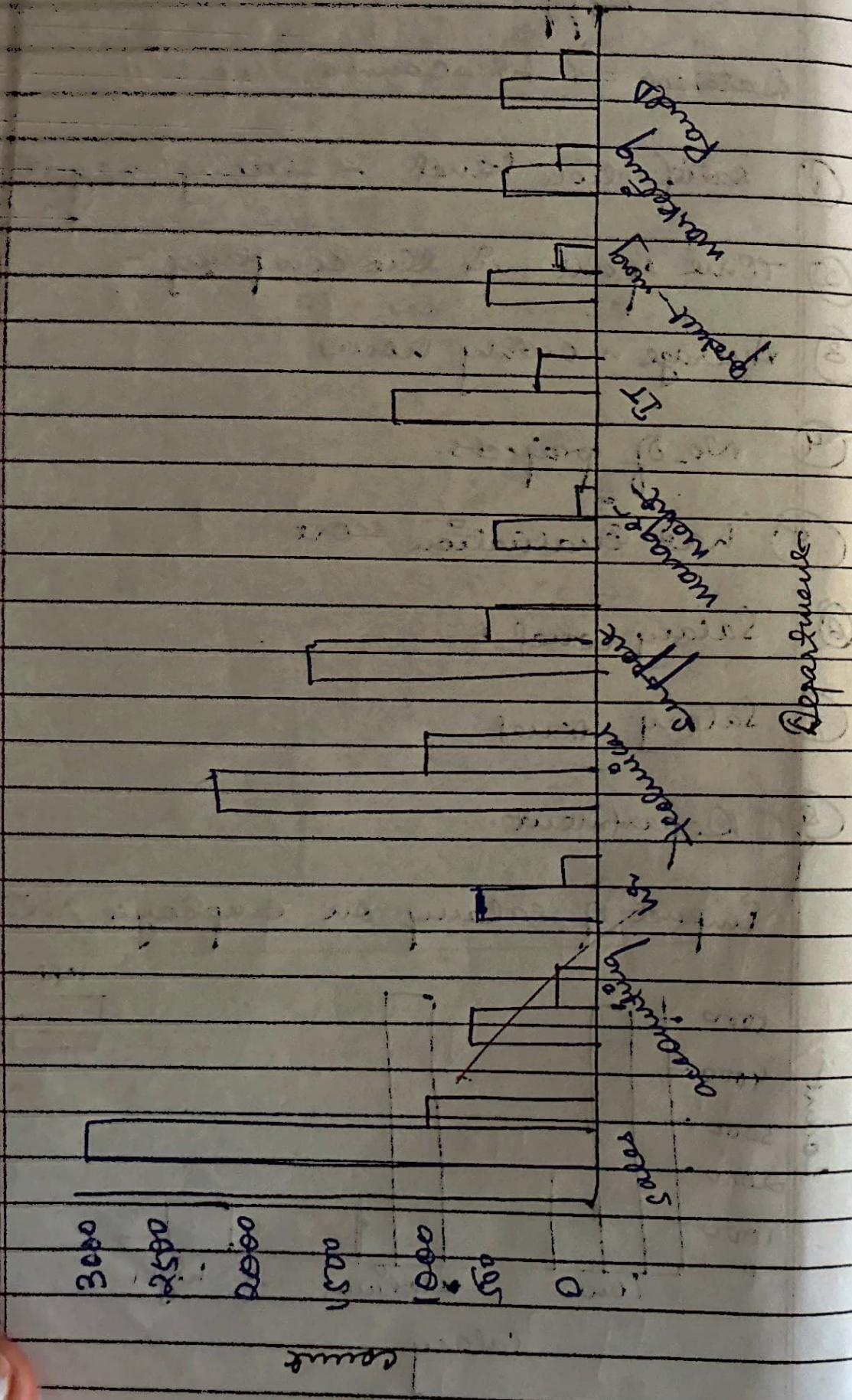
Dataset = ( HR - comma - sep. csv )

- Q 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 Salary Genus
- ⑧ Department

2) Impact of salary on employee retention



## Department vs. employee retention



Model Accuracy = 0.75866. 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

iv) There were no missing or incomplete values

v) 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.

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

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

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

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

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

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

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

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

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

$$[C, D, E] \rightarrow [N, Y, Y]$$

$$X(35, 100) \rightarrow Y$$

## Euis Dataset

- Here the K value has no impact on core accuracy for K value (1, 2, 3) but after that there were some fluctuations

K value can be chosen based on hit 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 algos like KNN
- improves convergence speed

### How to perform:

#### 1) Standardization (Z-score normalization)

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

use StandardScaler() in scikit-learn

## 2) Normalization (MinMaxScaling)

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

Use MinMaxScaler() in scikit-learn.

16/4/25

## Decision Tree Algorithm, etc.

- Input: dataset  $D$  with features  $F$  and target variable  $T$

④

## Random Forest Algorithm

- Start / Bootstrap sampling

Input: for each tree, randomly select  $N$  samples with replacement from the dataset.

- Build decision tree:

- for each tree, use the bootstrap sample
  - build a decision tree
- Randomly select  $k$  features for each split

- Make predictions:

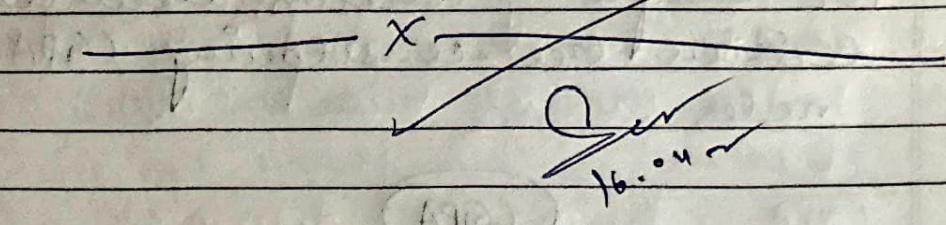
- Each tree makes a prediction for a new data point

- Combine predictions:

- Classification - Use majority voting to determine the final class
- Regression - Average the predictions

## (5) Out-of-Bag (OOB) Error (Optional):

- Use out-of-bag samples to estimate the error of the model.

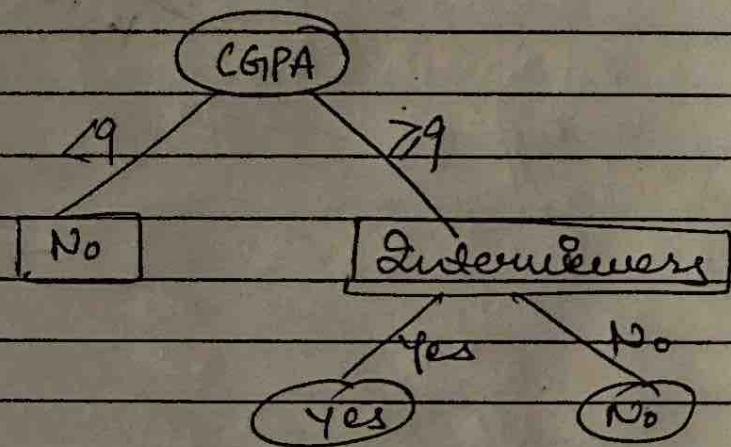


7/5/25

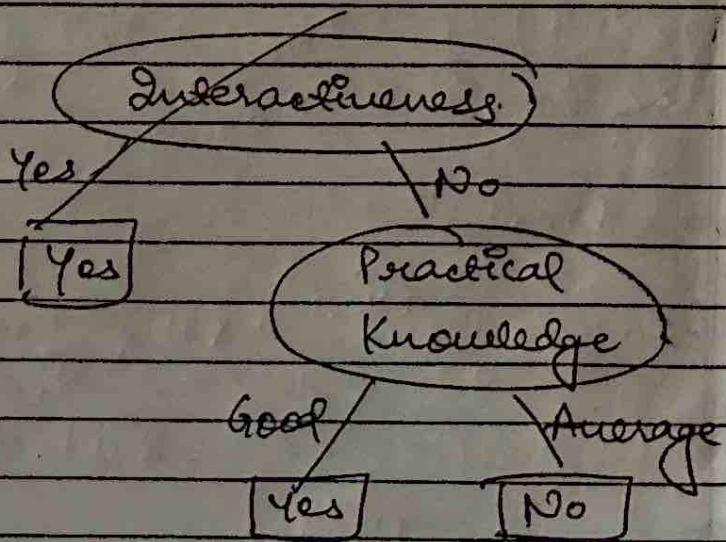
Week-7

- ③ Implement random forest ensemble method on a given dataset.

For sample S1 shown, draw the decision tree considering CGPA as root node.



For sample S2, shown, draw the decision tree considering Interactions as root node



Write the answer for the following question -

① For "iris.csv" dataset -

what is the best accuracy score and confusion matrix of the classifier you obtained and using how many trees?

→ Best accuracy score observed

$$\Rightarrow 1.0$$

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

→ Confusion Matrix

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

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

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

7/5/25

Lab-8

Q8

Implement Boosting ensemble method on a given dataset

Considering AdaBoost Algorithm for the following sample data, show the decision stamp calculation step for is the attribute, CGPA.

| CGPA | Predicted Job Offer | Actual Job Offer | Weight |
|------|---------------------|------------------|--------|
|------|---------------------|------------------|--------|

|          |     |     |       |
|----------|-----|-----|-------|
| $\geq 9$ | Yes | Yes | $4/6$ |
| $< 9$    | No  | Yes | $1/6$ |
| $\geq 9$ | Yes | No  | $1/6$ |
| $< 9$    | No  | No  | $1/6$ |
| $\geq 9$ | Yes | Yes | $1/6$ |
| $\geq 9$ | Yes | Yes | $1/6$ |

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

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

$$= 0.347$$

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

$$= 0.4488$$

$$\text{wt } (\alpha_j)_{j+1} = \frac{1}{6} \times e^{-0.347} = 0.1249$$

$$= 0.9428$$

$$w(\alpha_i) = \frac{1}{6} \times e^{0.347} = 0.2801$$

| CGPA | Predicted job offer | Actual job offer | Weight |
|------|---------------------|------------------|--------|
| 7.9  | Yes                 | Yes              | 0.1249 |
| <9   | No                  | Yes              | 0.1249 |
| 7.9  | Yes                 | No               | 0.2501 |
| <9   | No                  | No               | 0.1249 |
| 7.9  | Yes                 | Yes              | 0.1249 |
| >9   | Yes                 | Yes              | 0.1249 |

Re-write the answer for the following question -

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 estimates : 0.8277

Confusion matrix (10 estimates):

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

Best accuracy : 0.831 with n-estimators = 72.

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

$\rightarrow$  iteration 1:

| Record Number   | Close to c1<br>$(1.0, 1.0)$ | Close to c2<br>$(5.0, 7.0)$ | Assign 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 $(3.5, 4.5)$ | 4.30                        | 2.92                        | C2                |

Cluster 1  $\{R_1, R_2, R_3\}$  and Cluster 2  $\{R_4, R_5, R_6, R_7\}$

Their new centroids are -

$$C_1 = \frac{(1.0 + 1.5 + 3.0)}{3}, \frac{(1.0 + 2.0 + 4.0)}{3}, \\ = 1.83, \quad 2.33$$

$$C_2 = \frac{(5.0 + 3.5 + 4.5 + 3.5)}{4}, \frac{(7 + 5 + 5 + 4)}{4} \\ = 4.12, \quad 5.37$$

Iteration 2°

| Record No.                | Close to C1<br>(1.83, 2.33) | Close to C2<br>(4.12, 5.37) | Assign to<br>Cluster |
|---------------------------|-----------------------------|-----------------------------|----------------------|
| R <sub>1</sub> (1.0, 1.0) | 1.57                        | 5.64                        | C <sub>1</sub>       |
| R <sub>2</sub> (1.5, 2.0) | 0.47                        | 4.52                        | C <sub>1</sub>       |
| R <sub>3</sub> (3.0, 4.0) | 2.12                        | 1.63                        | C <sub>2</sub>       |
| R <sub>4</sub> (5.0, 7.0) | 5.57                        | 1.91                        | C <sub>2</sub>       |
| R <sub>5</sub> (3.5, 5.0) | 3.16                        | 0.72                        | C <sub>2</sub>       |
| R <sub>6</sub> (4.5, 5.0) | 3.58                        | 0.52                        | C <sub>2</sub>       |
| R <sub>7</sub> (3.5, 4.5) | 2.63                        | 1.05                        | C <sub>2</sub>       |

Cluster 1 {R<sub>1</sub>, R<sub>2</sub>} and Cluster 2 {R<sub>3</sub>, R<sub>4</sub>, R<sub>5</sub>, R<sub>6</sub>, R<sub>7</sub>}

Their new centroids are:

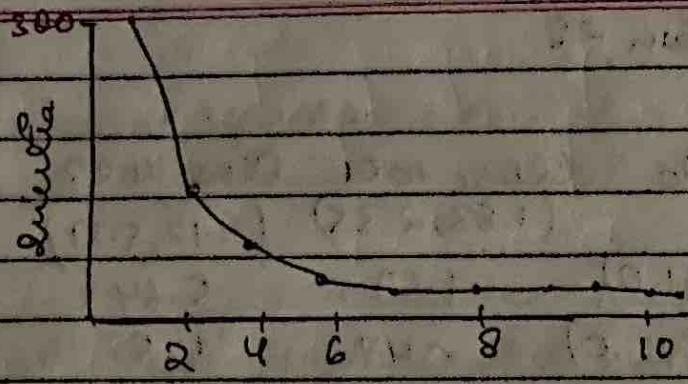
$$C_1 = (1.0 + 1.5)/2, (1.0 + 2.0)/2 \\ = 2.5/2 = 1.25, 3/2 = 1.5$$

$$C_2 = (3 + 5 + 3.5 + 4.5 + 3.5)/5 \rightarrow (4 + 7 + 5 + 5 + 4.5)/5 \\ = 18.5/5 = 3.9 \rightarrow 25.5/5 = 5.1$$

Q) Write the answers for the following questions:

① For "iris.csv" dataset,

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



No. of clusters

⑧ 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 :-

$$e_1 = \begin{bmatrix} 0.5594 \\ -0.8303 \end{bmatrix}, e_2 = \begin{bmatrix} 0.8303 \\ 0.5594 \end{bmatrix}$$

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

② Mean centre the data :-

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

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

③  $y_{\text{centered}} = \begin{bmatrix} 4.8 & 8.8 & 13.8 & 7.8 \\ 11.85 & 4.85 & 8.5 & 14.85 \end{bmatrix}$

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

② Using eq" projecting data onto first principle component -

$$z = e_1^T \cdot \text{centred}$$

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

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

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

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

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

Write the answer for the following question:

Q) For "heart.csv" dataset,

Compare the accuracy scores before and after applying PCA.

→ Model accuracy

Before

SVM: 0.8804

After

0.8424

Logistic Regression: 0.8533

0.8461

Random forest: 0.8859

0.8533

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