

5/8/25

lab 0

→ Initialize Values directly into DF.

```
import pandas as pd  
df = pd.DataFrame({  
    "USN": [23, 21, 24, 25],  
    "NAME": ["Sam", "Ran", "DAM",  
             "rahul", "raksheth"],  
    "Marks": [45, 78, 24, 90, 60]})
```

3)

df

→ import datasets from sklearn.datasets
from sklearn.datasets import
load_diabetes
dia = load_diabetes()
df = pd.DataFrame(dia, data, columns=
dia.feature_names)

df

→ ~~Methods~~: Import datasets from a
~~specific~~.csv file

```
filepath = "/content/Sample-Sales-data.csv"  
df = pd.read_csv(filepath)  
df.head()
```

→ Method 4: Downloading datasets from existing dataset repositories like Kaggle.

```
df = pd.read_csv("diabetes.csv")  
df.head()
```

→ Using the code given in the above slides, do the exercise of the "Stock market Data Analysis", considering the following.

1. HDFC Bank Ltd, ICICI Bank Ltd, Kotak Mahindra Bank Ltd
stocks = ["HDFC Bank NS", "ICICI Bank NS", "Kotak Bank NS"]
2. Start date :- 2024-01-01, End Date :- 2024-12-31
3. Plot closing price and daily return for all three banks mentioned.

Lob - 1

Date _____
Page _____

- i) To load .csv into the data frame
- ii) To display info of all columns
- iii) To displays statistical info of all numerical
- iv) To displays the count of unique labels for "Ocean proximity" column
- v) To display which att in dataset have missing values count greater than 20%

import pandas as pd

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

print(df.info)

print(df.describe())

df[["Ocean proximity"]].count

missing_values = df.isnull.sum()

print(missing_values[missing_values > 0])

Import pandas as pd

Import numpy as np

diabetes_df = pd.read_csv('Diabetes.csv')
adult_df = pd.read_csv("path/adult.csv")

def preprocess_data(df):

missing_values = df.isnull.sum()

print(missing_values[missing_values])

numerical_cols = df.select(include_cols)

categorical_cols = df.select(include_cols)

import SimpleImputer(strategy='ffill',
 df[[categorical_cols]])

encoder = OneHotEncoder(drop='first',
 sparse_output=False)

encoded_categorical_df = pd.DataFrame(
 encoder.fit_transform(categorical_data, columns=
 encoder.get_feature_names_out(
 categorical_cols)))

12/03/2025

Lab-2

Date _____
Page _____

import

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

```
import pandas as pd  
from sklearn.model_selection import  
train_test_split.
```

```
from sklearn.preprocessing import  
LabelEncoder  
from sklearn.tree import DecisionTreeClassifier  
from sklearn import tree  
import matplotlib.pyplot as plt
```

```
data = pd.read_csv('seattle-weather.csv')
```

```
data_cleaned = data.drop(['date'], axis=1)  
label_encoder = LabelEncoder()  
data_cleaned['weather'] = label_encoder.  
fit_transform(data_cleaned['weather'])
```

```
X = data_cleaned.drop(['weather'], axis=1)  
y = data_cleaned['weather']
```

```
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.2,  
random_state=42)
```

```
Id3_classifier = DecisionTreeClassifier(  
criterion='entropy', max_depth=5,  
random_state=42)
```

id3-classifier.fit(x-train, y-train)

plt.figure(figsize=(15, 8))

tree.plot_tree(id3-classifier, feature_names=x.columns, classnames=list(label-encoder.classes_), filled=True)

plt.title("ID3 Decision Tree classifier
(Improved)"),

plt.show()

O/P

"Outlook":

"Sunny": {

"Humidity": {

"High": "No"

"Normal": "Yes"

}}

overcast: "Yes"

"Rain": {

"Wind": {

"Weak": "Yes"

"Strong": "No"

}}

Lab-3

Linear Regression

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

```
file_path = "Acontent/Book1.csv"  
df = pd.read_csv(file_path)
```

```
X = df[['weeks']].values  
y = df[['Sales in months']].values
```

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

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

```
week_to_predict = int(input("Enter the  
week number to predict Sales:"))
```

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

```
predicted_sales = future_week @ theta_best
```

```
print("Theta (Intercept and Slope):",  
theta_best.flatten())
```

```
print("Predicted Sales for week  
{} week-to-predict : {} predicted sales: {}")
```

```
plt.scatter(X, y, color='blue', label='ActualSales')
```

```
plt.plot(X, X_b @ theta_best, color='red',  
label='Regression Line')
```

plt.scatter([week_to_predict], predicted_sales,
color='green', marker='x', label='Prediction')

plt.xlabel('weeks')

plt.ylabel('Sales in Month')

plt.legend()

plt.title('Linear Regression using Matplotlib
Method')

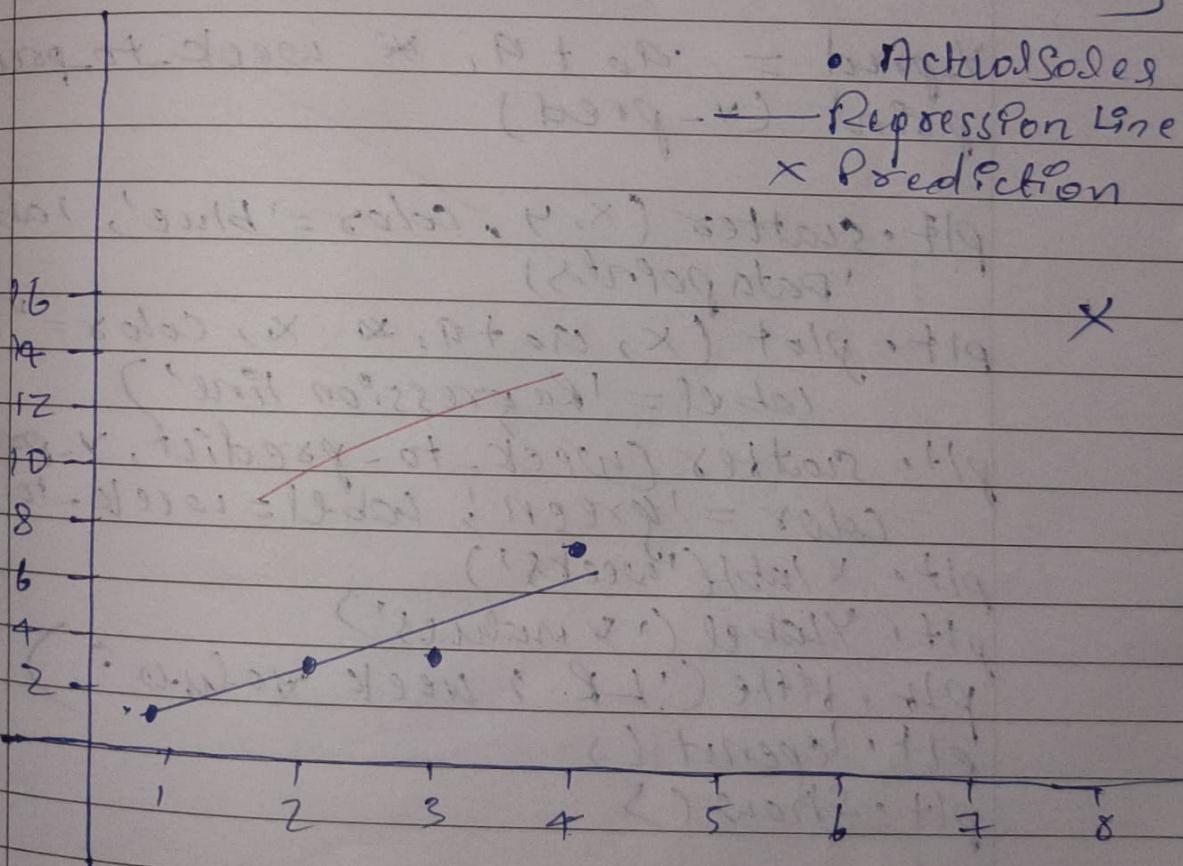
plt.show()

O/P

Enter the week number to predict

Sales : 8

theta (Intercept and slope) : [-1.05, 2.2]



2nd Approach.

```
import numpy as np
```

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

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

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

```
x_mean = np.mean(x)
```

```
y_mean = np.mean(y)
```

$$\alpha_1 = \frac{(np.mean(x * y) - x_mean * y_mean)}{(np.mean(x * x) - x_mean * x_mean)}$$

$$\alpha_0 = y_mean - \alpha_1 * x_mean$$

$$\text{week_to_predict} = 5$$

$$y_{\text{pred}} = \alpha_0 + \alpha_1 * \text{week_to_predict}$$

point (y-pred)

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

```
plt.plot(x, alpha_0 + alpha_1 * x, color='red', label='Regression line')
```

```
plt.scatter(week_to_predict, y_pred, color='green', label='week-to-predict')
```

```
plt.xlabel('weeks')
```

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

```
plt.title('LR : week values : )
```

```
plt.legend()
```

```
plt.show()
```

off

Intercept (a_0): -1.5

Slope (a_1): 2.2

Predicted value for week 5: 9.5

~~9.5~~
~~19.375~~

2/4/25

Lab- 4

Date _____
Page _____Logistic Regression

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 $a_1 = 0.8$.
 (Coefficient for study hours)

- Write a logistic regression equation for this problem.
- Calculate the probability that a student who studies for 7 hours will pass.
- Determine the predicted class (pass or fail) for this student based on a threshold of 0.5.

Answer

a) $\frac{1}{1+e^{-z}}$ $z = a_0 + a_1 x$

$$\frac{1}{1+e^{-5+0.8 \times 7}}$$

b) $\frac{1}{1+e^{(-5+0.8 \times 7)}} = 0.85 > 0.5$

c) Threshold of 0.5

0.85 is less than 0.5 so the student passed.

Consider $z = [2, 1, 0]$ for three classes.

Apply Softmax function to find the probability values of three classes

$$\frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

$$e^2 = 7.389$$

$$e^1 = 2.718$$

$$e^0 = 1$$

Sum of exponentials.

$$\begin{aligned} S &= e^2 + e^1 + e^0 \\ &= 7.389 + 2.718 + 1 \\ &= 11.107 \end{aligned}$$

SoftMax Probabilities

$$\frac{e^2}{S} = \frac{7.389}{11.107} = 0.665$$

$$\frac{e^1}{S} = \frac{2.718}{11.107} = 0.245$$

$$\frac{e^0}{S} = \frac{1}{11.107} = 0.090$$

$$\text{Probabilities} = [0.665, 0.245, 0.090]$$

1. For dataset file "HR-comma-sep.csv"

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

Satisfaction level - Employees with lower satisfaction levels are more likely to leave

Salary - Low salaries could be a major reason for employees quitting.

Avg Monthly hours - Excessive working hours could lead to burnout.

2. For Zoo dataset

i) Did you perform any data preprocessor steps? If yes, what were they, and why were they necessary?

Yes,

Data cleaning : Checked for missing values and inconsistencies.

Encoding Categorical data : The column animal-name is categorical and was removed as it does not contribute to classification

Class Balance check : Ensured that dataset was balanced across class-type

i) were there any missing or inconsistent values in the dataset? How did you handle them?

- Non missing values
- Replaced duplicate entries

(ii) Interpretation of the confusion matrix

- Shows classification accuracy of predictors
- High diagonal indicates high performance
- e.g. No diagonal miss classification

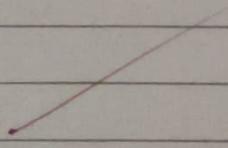
02/4/25

Lab. 5

Date _____
Page _____

Consider the following dataset, for $k=3$ and test data $(X, 35, 100)$ as $(\text{Person}, \text{Age}, \text{Salary})$ solve using kNN classifier model and predict the target.

Person	Age	Salary	Target	Distance
A	18	50	N	52.2
B	23	55	N	48.5
C	24	70	N	31.9
D	41	60	Y	40.4
E	43	70	Y	31.0
F	38	40	Y	60.0
X	35	100	?=Y	

$$\begin{aligned} k=1 &\rightarrow Y \\ k=2 &\rightarrow N \\ k=3 &\rightarrow Y \end{aligned} \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Majority Y}$$


For Iris dataset

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

→ Train KNN for different K values ($K=1 \text{ to } 10$)

Compute Accuracy Rate

$$\frac{\text{Correct Prediction}}{\text{Total Prediction}} \times 100$$



Compute Error rate:

$$\text{Error Rate} = 1 - \text{Accuracy rate}$$

For Diabetes dataset

Q What is the purpose of feature scaling?
How to perform it?

Min-Max Scaling

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

Standardization

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

BB
24-11-2025

15/04/25

Date / /
Page / /

Lec - 6

Random forest Algorithm

Step 1: Select random k data points from the training set

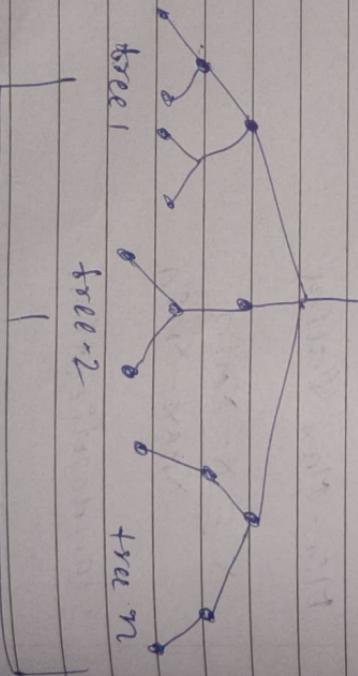
Step 2: Build the decision tree associated with selected data points

Step 3: Choose the number N for decision trees that you want to build

Step 4: Repeat Step 1 & 2

Step 5: For new data points, find the prediction of each different decision tree and assign the new data points to the category that wins the majority votes

Class: Apple, banana, trouser

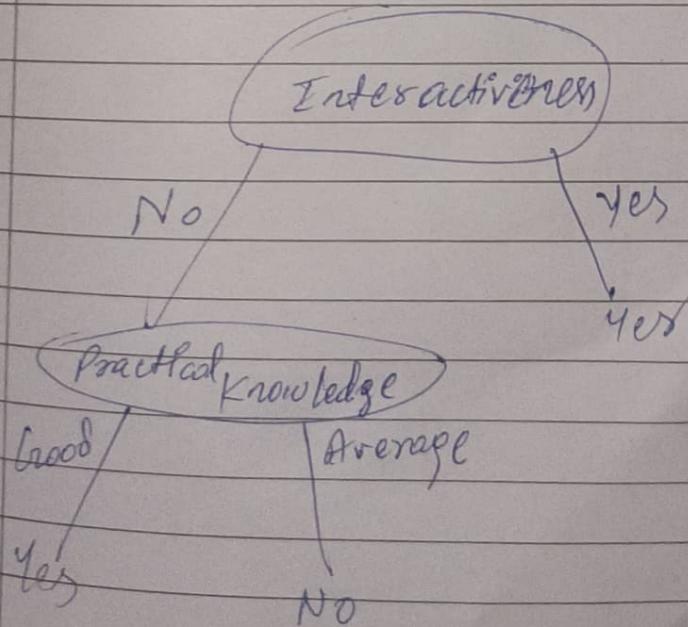
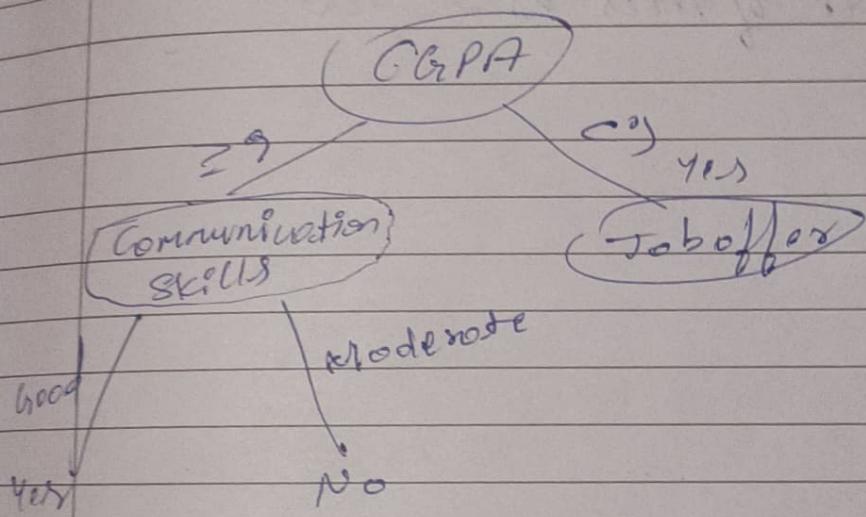


↓
Majority voting
final class

Random Forest

1) Sample SI

	CGPA	Interactivity	Communication Prof.	
1	≥ 9	Yes	G	G Y
2	< 9	No	M	G Y
3	≥ 9	No	M	A N
4	≥ 9	No	M	A N
5	≥ 9	Yes	M	G Y



1. For Iris.csv dataset

2. Best accuracy score: 100 %.

3. Confusion matrix:

	Predicted Setosa	Versicolor	Virginica
Setosa	50	0	0
Versicolor	0	50	0
Virginica	0	0	50

4. Number of trees used 100

AdaBoost

Practical

CGPA	Interactivity	Knowledge	Communication	Step 1 prob
≥ 9	Y	Good	Good	Y
< 9	N	Good	Moderate	Y
≥ 9	N	Avg	Moderate	N
< 9	N	Avg	Good	N
≥ 9	Y	Good	Moderate	Y
≥ 9	Y	Good	Moderate	Y

Verginica

0

0

5.6

Step ①

CGPA	Predicted	Actual	Weight
≥ 9	Y	Y	$1/6$
< 9	N	Y	$1/6$
≥ 9	Y	N	$1/6$
< 9	N	N	$1/6$
≥ 9	Y	Y	$1/6$
≥ 9	Y	Y	$1/6$

Step ②

$f(x_i) \times \text{no of wrong predictions}$
 $\text{error} = 2 \times \frac{1}{6} = \frac{1}{3}$

Step 3

$$\alpha = \frac{1}{2} \log \left(\frac{1 - \text{error}}{\text{error}} \right)$$

$$\alpha = \frac{1}{2} \log \left(\frac{1 - \frac{1}{3}}{\frac{1}{3}} \right) = 0.347$$

Step 4:

$$z = \text{wt (correct)} * \text{No. of correct} * e^{-a} + \\ \text{wt (incorrect)} * \text{No. of Incorrect} * e^a$$

$$= \frac{1}{6} * 4 * e^{-0.347} + \frac{1}{6} * 2 * e^{0.347}$$

$$= 0.9428$$

Step 5:

Update the weight

$$\underline{w(\text{adj})_{\text{curr}}} \text{ of correct} * e^{-a}$$

$$2$$

$$\frac{\frac{1}{6} * e^{-0.347}}{0.9428} = 0.1249$$

$$\underline{w(\text{adj})_{\text{curr}}} \text{ of incorrect} * e^a$$

$$0.9428$$

$$\frac{\frac{1}{6} * e^{0.347}}{0.9428} = 0.2501$$

CGPA	Predicted	Actual	Weight
>=9	Y	Y	0.1249
<9	N	Y	0.2501
>=9	Y	N	0.2501
<9	N	N	0.7248
>=9	Y	Y	0.1249
>=9	Y	Y	0.1249

LAB 9

Date _____
Page _____

K-means, 2nd iteration

Record Number	A	B
R1	1.0	1.0
R2	1.5	2.0
R3	3.0	4.0
R4	5.0	7.0
R5	3.5	5.0
R6	4.5	5.0
R7	3.5	4.5

$$\text{Formula } d = \sqrt{(x^2 - x_1)^2 + (y^2 - y_1)^2}$$

Record	Point(A,B)	d(0)	d(1)	Clusters
R1	(1.0, 1.0)	0.00	7.2	C1
R2	(1.5, 2.0)	0.12	6.20	C1
R3	(3.0, 4.0)	5.61	3.61	C1/C2 (C1)
R4	(5.0, 7.0)	7.21	0.00	C2
R5	(3.5, 5.0)	5.00	2.50	C2
R6	(4.5, 5.0)	5.32	2.24	C2
R7	(3.5, 4.5)	1.30	3.20	C2

C1 cluster :- R1, R2, R3

C2 cluster :- R4, R5, R6, R7

C1 Mean

$$\frac{1+1.5+3}{3} = 1.83 \quad \frac{1+2+4}{3} = 2.3$$

C2 Mean

$$\frac{5+3.5+4.5+3.5}{4} = 4.125, \quad \frac{7+5+5+4.5}{4} = 5.375$$

$$C_1 = (1.83, 2.33)$$

$$C_2 = (4.125, 5.375)$$

Record	Point(A,B)	d_1	$d_1 C_2$	cluster
R1	(1, 1)	1.57	5.47	C1
R2	(1.5, 2)	0.47	4.32	C1
R3	(3, 4)	7.98	1.53	C2
R4	(5, 7)	2.91	1.76	C2
R5	(3.5, 5)	2.68	0.73	C2
R6	(4.5, 5.0)	3.14	0.52	C2
R7	(3.5, 2.5)	2.28	1.00	C2

Lab - 10

PCA

Feature	Ex1	Ex2	Ex3	Ex4
x_1	4	8	13	7
x_2	11	4	5	14

$$d_1 = 30.3849$$

$$d_2 = 6.6151$$

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

$$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

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

② mean centre the data

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

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

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

③ $Z = X^T e_1^T x_{\text{centered}}$ -1.301

$$x_1 = (0.5574)(-4) + (-0.8303)(2.5) = 6$$

$$x_2 = (0.5574)(0) + (-0.8303)(-4.5) = 3.7365$$

$$x_3 = (0.5574)(5) + (-0.8303)(-3.5) = 5.6952$$

$$x_4 = (0.5574)(-1) + (-0.8303)(5.5) = -5.12475$$