

## LAB-0

→ ticket8 - ["SBI BANK.NS", "ICICI BANK.NS", "KOTAK BANK.NS"]  
 data = gf.download(tickets, start = "2021-01-01", end = "2024-12-30", group\_by='ticket')

import yfinance as gf import pandas as pd import matplotlib.pyplot as plt

→ SBI\_data = data['SBI BANK.NS']  
 SBI\_data['Daily Return'] = SBI\_data['close'].pct\_change()

Print("Daily returns for SBI Industries:")

Print(SBI\_data['Daily Return'].head())

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

plt.subplot(2, 1, 1)

SBI\_data['close'].plot(title="SBI Industries - closing price")

plt.subplot(2, 1, 2)

SBI\_data['Daily return'].plot(title="SBI Industries - Daily Return")

plt.tight\_layout()

plt.show()

→ Daily Returns for SBI Industries.

Date:

2024-01-01

NoN

2024-01-02

0.000589

2024-01-03

-0.015420

2024-01-04

0.016730

2024-01-05

-0.080516

(Date, closing, daily\_return)  
 ("Date, closing and daily\_return")

import pandas as pd

df = pd.DataFrame({'USN': [1, 2, 3, 4, 5], 'Name': ['A', 'B', 'C', 'D', 'E'], 'Marks': [10, 20, 30, 40, 50]})

df = pd.DataFrame(df)

print(df)

USN	Name	Marks
0	A	10
1	B	20
2	C	30
3	D	40
4	E	50

2. From sklearn.datasets import load\_diabetes

diabetes = load\_diabetes()

df = pd.DataFrame(diabetes.data, columns=diabetes.feature\_names)

df['target'] = diabetes.target

df.head()

3. file\_path = 'Sample\_Sales\_Data.csv'

df = pd.read\_csv(file\_path)

print(Sample\_data)

Product	Quarterly Price	Sales	Region
---------	-----------------	-------	--------

0 Laptop 5000 5000 North

1 Mouse 3000 3000 East

4. df.to\_csv('outhub.csv', index=False)

Print("Data saved to outhub.csv")

## LAB-1

① Data set considered : housing.csv

(i) To load .csv file into data frame.

import pandas as pd  
df = pd.read\_csv('housing.csv')

(ii) To display information of all columns

→ df.columns

→ Index ('longitude', 'latitude', 'housing\_median\_age',  
'total\_rooms', 'total\_bedrooms', 'population',  
'households', 'median\_house\_value')

(iii) To display statistical information of all numerical  
columns

→ df.describe()

longitude latitude housing\_median\_age total\_rooms population

	20640	20640	20640	20640	20640
cent	20640	20640	20640	20640	20640

	20433	20640	20640	20640
total_bedrooms	20433	20640	20640	20640
households	20640	20640	20640	20640

(iv) To display the counts of unique tickles for  
"ocean\_proximity" column.

c = df['ocean\_proximity'].unique()

print(c)

print(len(c))

[ 'NEAR BAY' 'H Ocean' INLAND 'NEAR OCEAN' ISLAND ]  
5

- (v) To display which columns in a data set having missing values count greater than 0

Print ( df. columns [ df. is null () . any () ] )

Index ( [ 'total\_bedrooms' ], dtypes = 'object' )

- II Data sets considered : diabetes.csv, income.csv  
import pandas as pd.

df1 = pd.read\_csv('diabetes.csv')

df2 = pd.read\_csv('income.csv')

- ① which columns in the data set had missing values.  
How did you handle them?

Print ( df. columns [ df. is null () . any () ] )

Index ( [ ], dtypes = 'object' )

- ② which categorical columns did you identify in data set?

How did you encode them?

categorical\_col = df.select\_dtypes(exclude= ['number']) columns.

diabetes: Gender / class.

adult: workclass, education, marital-status, occupation,  
relationship, race, gender, native-country  
income

## \* Linear Regression and Multiple Linear Regression

Find linear Regression of data of week and product

$x_i$  (week)  $y_i$  (Sales in thousands)

1	2
3	5
4	9

$$\text{Soln} \quad x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \quad y = \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$$

$$x^T x = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

$$(x^T x)^{-1} = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix}$$

$$(x^T x)^{-1} x^T = \begin{bmatrix} 1.5 & 0.5 \\ -0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.3 & -0.1 & 0.1 & 0.5 \end{bmatrix}$$

$$\hat{P} = ((x^T x)^{-1} x^T) y = \begin{bmatrix} 1 & 0.5 & 0 & 0.5 \\ -0.3 & -0.1 & 0.1 & 0.5 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$$

$$= \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$$

$$y = a_0 + a_1 x$$

$$y = -0.5 + 2.2x$$

for  $x = 5$

$$y = -0.5 + (2.2)5 = 10.5$$

\*

- ① import libraries
- ② import data - distribution
- ③ Analyze data - distribution
- ④ Distribution plot - visualization.
- ⑤ Relationship b/w variables
- ⑥ Shift the data
- ⑦ Train the model
- ⑧ Predict the result
- ⑨ Visualize prediction
- ⑩ check values of co-efficient & intercept

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2.0 & 2.1 \end{bmatrix} = f^{-1}(x)$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 5.0 & 5.1 \end{bmatrix} = f^{-1}(x)$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1.0 & 1.1 \end{bmatrix} = f^{-1}(x)$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 8.0 & 8.1 \end{bmatrix} = f^{-1}(x)$$

## \* Linear Regression.

```
import Pandas as pd
```

```
import numpy as np
```

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

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

```
from sklearn.impute import SimpleImputer
```

```
file_path = "canada_per_capita_income.csv"
```

```
data = pd.read_csv(file_path)
```

```
data.fillna(data.mean(numeric_only=True),  
inplace=True)
```

```
x = data[['year']]
```

```
y = data['per capita income(US$)']
```

```
model = LinearRegression()
```

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

```
year_2020 = np.array([[2020]])
```

```
income_2020 = model.predict(year_2020)
```

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

```
plt.plot(x, model.predict(x), color='red',  
linewidth=2, label='Regression Line')
```

```
plt.xlabel('year')
```

~~```
plt.ylabel('Per capita Income (US$)')
```~~~~```
plt.title('canada per capita income prediction')
```~~~~```
plt.legend()
```~~~~```
plt
```~~~~```
plt.show()
```~~

## OUTPUT

Predicted per capita income in  
in 2020 = \$41288.69

Estimated value of labor force participation  
in 2020 = 15.5%

(Actual value) was 15.5% in 2020.

Actual population growth rate = 0.5%

## \* Multiple Linear Regression.

From Sklearn. preprocessing import LabelEncoder

file\_path = "1000-companies.csv"

data = pd.read\_csv("content/1000-companies.csv")

data = fillna(data.mean(numeric\_only=True), inplace=True)

label\_encoder = LabelEncoder()

data['State'] = label\_encoder.fit\_transform(data['State'])

x = data.iloc[:, :-1]

y = data.iloc[:, -1]

model = LinearRegression()

model.fit(x, y)

test\_data = np.array([[91694.48, 15841.3, 11931.2, label\_encoder.transform(['Florida'])]]))

Predicted\_Profit = model.predict(test\_data)

Print(f"Predicted Profit: \${predicted\_Profit} ")

Output

Predicted Profit: \$ -2153.01

error $y_{\text{pred}} = \text{model. predict}(x)$  $\text{mae} = \text{mean-absolute-error}(y, y_{\text{pred}})$  $\text{mse} = \text{mean-squared-error}(y, y_{\text{pred}})$  $\text{rmse} = \text{np.sqrt(mse)}$  $R^2 = R^2 = \text{Score}(y, y_{\text{pred}})$ 

Printf("Salary prediction (Experience) - Error metrics")

Printf("MAE : { mae : 2f }")

Printf("mse : { mse : 2f }")

Printf("rmse : { rmse : 2f }")

Printf("R<sup>2</sup> Score : { r2 : 4f }")

Output

mae : 18202.74

mse : 571205771.52

rmse : 23899.91

R<sup>2</sup> Score : 0.9924

Printf (" Slope (Exponent): % model. coef-T03:3f")

Printf (" f " Intercept (C): % model. intercept-2f")

Output =

Slope = 0.55

Intercept = -T0 214.44

$$222.0 = -2 \times 0.55x + 214.44$$

AV  
11/3/28

$$2.05 \text{ (d) } b \text{ (b) } 2 = 0$$

$$\text{model} = [b + ax] = - (1)$$

$$222.0 = -2 \times 0.55x + 214.44$$

$$222.0 = -2 \times 0.55x + 214.44$$

$$(1) a_0 = -5$$

$$a_1x = 0.8$$

$$\text{logistic regression } p(x) = \frac{1}{1 + e^{-(a_0 + a_1x)}}$$

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

$$(2) x = 7$$

$$h(x) = \frac{1}{1 + e^{(-5 + 0.8 \times 7)}} = 0.6957$$

3. Threshold of 0.5

$$h(x) = 0.6957$$

$$h(x) \geq 0.5$$

Thus  $y = 1$  (Pass)

$$y = \begin{cases} 1 & \text{if } h(x) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$(4) z = [2, 1, 0] \rightarrow \text{classes}$$

$$\text{Softmax } (z_i) = \frac{e^{z_i}}{\sum_{j=1}^3 e^{z_j}}$$

$$\text{Softmax } (z_1) = \frac{e^2}{e^2 + e^1 + e^0} = 0.665$$

$$\text{Softmax}(z_2) = \frac{e^1}{e^2 + e^1 + e^0} = 0.244$$

$$\text{Softmax}(z_3) = \frac{e^0}{e^2 + e^1 + e^0} = 0.092$$

∴ Probability of 3 classes are  
66.8% 24.4% 9.1%

\* Logistic = Regression Binary.

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 = "+",  
columns = '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)

y\_predict = model.predict(X\_test)

model.score(X\_test, y\_test)

y\_predicted = model.predict([E60])

y\_predicted.

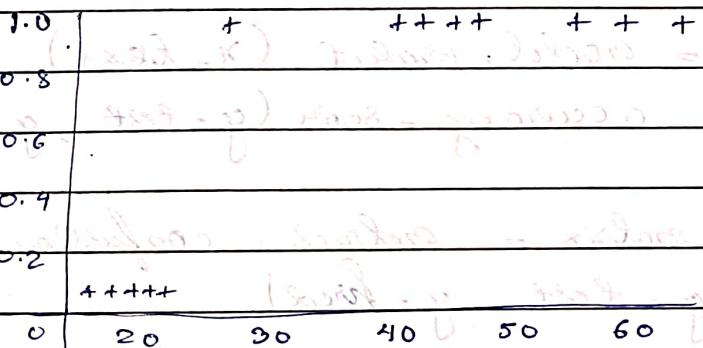
import math  
def sigmoid(x):  
 return 1 / (1 + math.exp(-x))

def prediction = sigmoid(age)  
 $\hat{y} = 0.127 * \text{age} - 1.973$

$y = \text{sigmoid}(\hat{y})$   
return y

age = 35  
prediction = sigmoid(age)

output is - 0.37098341



\* Logistic Regression - Multiclass

import pandas as pd  
from sklearn.datasets import load\_iris.

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

from sklearn.linear\_model import  
LogisticRegression

from sklearn metrics import accuracy\_score  
from sklearn import matrices

iris = pd.read\_csv("iris.csv")

iris.head(5)

X = iris.drop('Species', axis=1).values # columns  
y = iris.Species

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

model = LogisticRegression(multi\_class="multinomial")

model.fit = model.fit(X\_train, y\_train)

accuracy = accuracy\_score(y\_test, y\_pred)

confusion\_matrix = metrics.confusion\_matrix(y\_test, y\_pred)

cm\_display = metrics.ConfusionMatrixDisplay(confusion\_matrix=confusion\_matrix, display\_labels=[0, 1, 2])  
(cm\_display = cm\_display.plot(display\_labels=[0, 1, 2]))

cm\_display.plot()

plt.show()

OUTPUT

Accuracy of multinomial logistic regression model  
on the test set: 1.00

|            | Setosa          | versicolor | virginica |
|------------|-----------------|------------|-----------|
| Setosa     | 50              | 0          | 0         |
| versicolor | 0               | 9          | 0         |
| virginica  | 0               | 0          | 11        |
|            | Setosa          | versicolor | virginica |
|            | Predicted label |            |           |

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}^{-1}$$

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}^{-1}$$

$$0 = 0 + 0 + 0$$

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}^{-1}$$

$$0 = 0 + 0 + 0$$

$$0 = 0 + 0 + 0$$

$$0 = 0 + 0 + 0$$

## LAB-4 Decision trees

| instance | a <sub>2</sub> | a <sub>3</sub> | classification |
|----------|----------------|----------------|----------------|
| 1        | Hot            | High           | No             |
| 2        | Not            | High           | No             |
| 6        | Cool           | High           | No             |
| 7        | Not            | High           | No             |
| 8        | Not            | Normal         | Yes            |

$$\text{Entropy}(S) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right)$$

$$= 0.7219$$

for a<sub>2</sub>:

$$S_{hot}[1+, 3-] = \frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \left(\frac{3}{4}\right) \log_2 \left(\frac{3}{4}\right)$$

$$= 0.8113$$

$$S_{cool}[0+, 1-] = 0$$

$$Gain(S, a_2) = 0.7219 - \frac{4}{5} \times 0.8113 - 0 = 0.7286$$

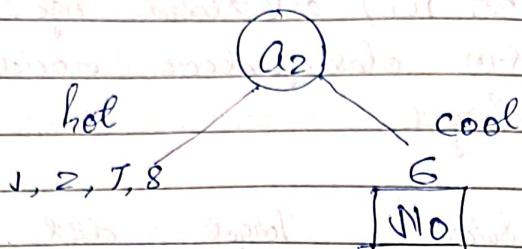
for a<sub>3</sub>:

$$S_{high}[0+, 4-] = 0$$

$$S_{normal}[1+, 0-] = 0$$

$$Gain(S, a_3) = 0.7219 - 0 - 0 = 0.7219$$

$a_2$  is the root node since Gain( $S, a_2$ ) is high



| instance | $a_2$ | $a_3$  | classification |
|----------|-------|--------|----------------|
| 1        | Hot   | High   | No             |
| 2        | Hot   | High   | No             |
| 7        | Hot   | High   | No             |
| 8        | Hot   | Normal | Yes            |

$$\text{Entropy}(S) = -\frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \frac{3}{4} \log_2 \left(\frac{3}{4}\right)$$

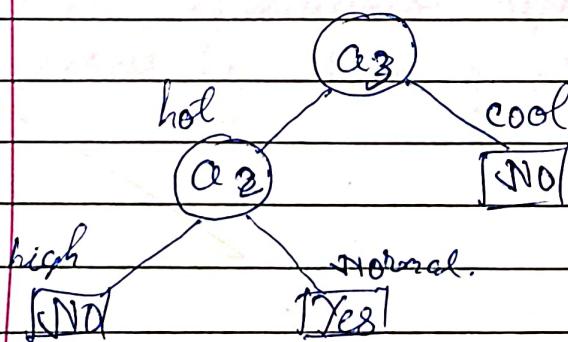
$$= 0.8113$$

for  $a_3$ :

$$\text{S}_{\text{high}}[0+, 3-] = 0$$

$$\text{S}_{\text{normal}}[1+, 0-] = 0$$

~~$$\text{Gain}(S, a_3) = 0.8113 - 0 - 0 = 0.8113$$~~



## LAB - 6

Date J 4 / 25  
Page \_\_\_\_\_

- \* consider the following dataset for  $K=3$  and test data  $(x, 35, 100)$  of (Reason, Age, Salary) & solving using KNN classification model & predict the target.

| Person | Age | Salary | Target | dist  | Rank |
|--------|-----|--------|--------|-------|------|
| A      | 18  | 50     | N      | 52.81 | 5    |
| B      | 23  | 55     | N      | 46.57 | 9    |
| C      | 24  | 70     | N      | 31.95 | 2    |
| D      | 41  | 60     | Y      | 10.49 | 3    |
| E      | 43  | 70     | Y      | 31.64 | 1    |
| F      | 38  | 40     | Y      | 60.07 | 6    |
| X      | 35  | 100    | ?      |       | 8    |

$$K = 3$$

$$\text{Rank } 1 = g_1 \quad \text{Rank } 2 = N \quad \text{Rank } 3 = Y$$

$$\text{Rank } 2 = N$$

$$\text{Rank } 3 = Y$$

The target predicted is Y

```

import pandas as pd
import numpy as np
import mathlib.lib.phylo as phl
from sklearn.model_selection import train_test_split

```

```
iris = pd.read_csv("iris.csv")
```

```
diabetes = pd.read_csv('diabetes')
```

```
heart = pd.read_csv("heart")
```

```
X_iris = iris.iloc[:, :-1]
```

```
y_iris = iris.iloc[:, -1]
```

```
X_train, X_test, y_train, y_test = train_test
```

```
-split(X=iris, Y=iris, test_size=0.2,
```

```
random_state=42)
```

Knn = K Neighbour classifier (n\_neighbour=3)

Knn.fit(X\_train, y\_train)

y\_pred = Knn.predict(X\_test)

```
Print("KNN classification for IRIS dataset: ")
```

```
Print("Accuracy: " accuracy_score(y_test, ypred))
```

```
Print("confusion matrix: ", confusion_matrix(y_test, ypred))
```

```
Print("classification Report: ", classification_report(y_test, ypred))
```

```
(y-test, y-pred)
```

$X_{\text{diabetes}} = \text{diabetes}[:, :-1]$

$y_{\text{diabetes}} = \text{diabetes}[:, -1]$

$X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train}, \text{test}$

- Split ( $X_{\text{diabetes}}$ ,  $y_{\text{diabetes}}$ , test size = 0.2  
random\_state = 42)

Scaler = StandardScaler()

$X_{\text{train}} = \text{Scaler}.fit(X_{\text{train}})$

$X_{\text{test}} = \text{Scaler}.transform(X_{\text{test}})$

## Output

KNN classification of this Data Set

Accuracy : 1.0

Confusion matrix :

$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

Classification Report

|            | Precision | Recall | F1 Score | Support |
|------------|-----------|--------|----------|---------|
| Sub38a     | 1.00      | 1.00   | 1.00     | 10      |
| Mavigcolor | 1.00      | 1.00   | 1.00     | 9       |
| Maviginv   | 1.00      | 1.00   | 1.00     | 11      |

Iris.csv

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

iris = pd.read\_csv("iris1.csv")

X\_iris = iris.iloc[:, :-1]

y\_iris = iris.iloc[:, -1]

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X\_iris, y\_iris, test\_size=0.2, random\_state=42)

dtc = DecisionTreeClassifier()

dtc.fit(X\_train, y\_train)

y\_pred = dtc.predict(X\_test)

print(accuracy\_score(y\_test, y\_pred))

print(confusion\_matrix(y\_test, y\_pred))

print(classification\_report(y\_test, y\_pred))

## Output

Accuracy 1.0

confusion matrix

|    |   |    |
|----|---|----|
| 10 | 0 | 0  |
| 0  | 9 | 0  |
| 0  | 0 | 11 |

# classificador Report:

|           | Precision | recall | f1 score | Support |
|-----------|-----------|--------|----------|---------|
| Setas     | 1.00      | 1.00   | 1.00     | 10      |
| verbos    | 1.00      | 1.00   | 1.00     | 9       |
| virgíñica | 1.00      | 1.00   | 1.00     | 11      |

accuracy =  $\frac{10 + 9 + 11}{30} = 1.00$

macro avg. 1.00 1.00 1.00 30

weighted avg 1.00 1.00 1.00 30

~~precision = recall = f1 score = support = accuracy = 1.00~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

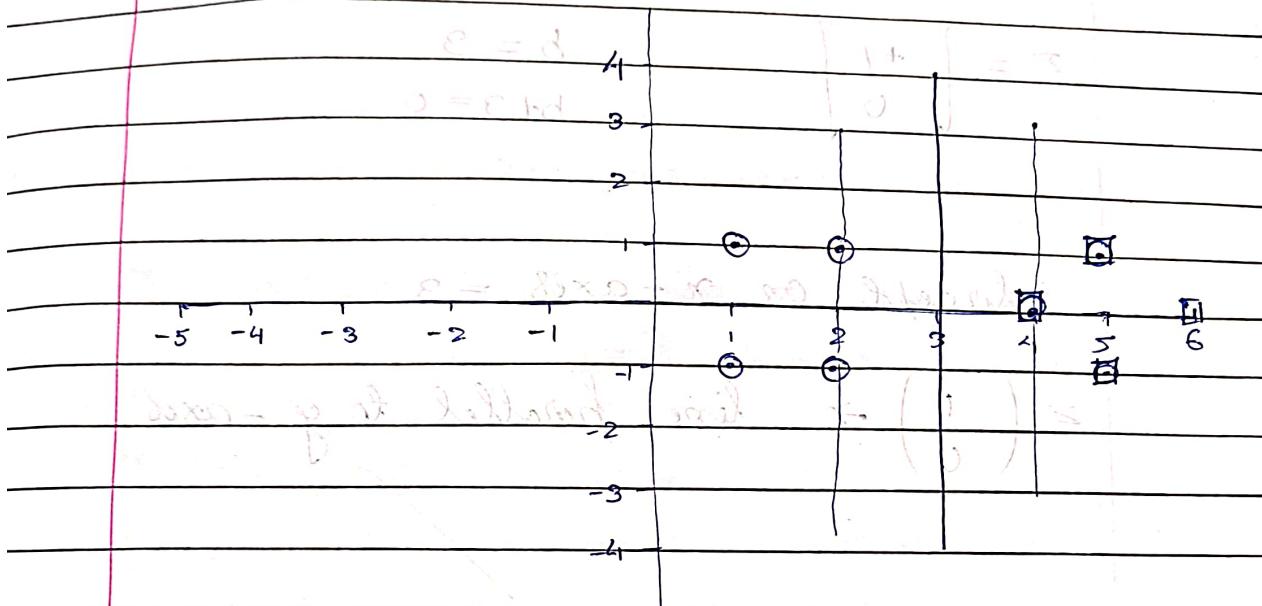
~~(precision = recall = f1 score = support = accuracy = 1.00)~~

~~(precision = recall = f1 score = support = accuracy = 1.00)~~

\* SVM

Draw an optimal hyperplane using linear SVM to classify the following points.

$\{(1, 1), (2, 1), (1, -1), (2, -1)\}$  - +ve labelled  
 $\{(4, 0), (5, 1), (5, -1), (6, 0)\}$  - -ve labelled.



$$\vec{s}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad \vec{s}_2 = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \vec{s}_3 = \begin{bmatrix} 4 \\ 0 \end{bmatrix}$$

$$\tilde{\vec{s}}_1 = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}, \quad \tilde{\vec{s}}_2 = \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix}, \quad \tilde{\vec{s}}_3 = \begin{bmatrix} 4 \\ 0 \\ 1 \end{bmatrix}$$

$$\alpha_1 \tilde{\vec{s}}_1 \tilde{\vec{s}}_1 + \alpha_2 \tilde{\vec{s}}_2 \tilde{\vec{s}}_2 + \alpha_3 \tilde{\vec{s}}_3 \tilde{\vec{s}}_1 = +1$$

$$\alpha_1 \tilde{\vec{s}}_1 \tilde{\vec{s}}_2 + \alpha_2 \tilde{\vec{s}}_2 \tilde{\vec{s}}_2 + \alpha_3 \tilde{\vec{s}}_3 \tilde{\vec{s}}_2 = +1$$

~~$$\alpha_1 \tilde{\vec{s}}_1 \tilde{\vec{s}}_3 + \alpha_2 \tilde{\vec{s}}_2 \tilde{\vec{s}}_3 + \alpha_3 \tilde{\vec{s}}_3 \tilde{\vec{s}}_3 = -1$$~~

$$\alpha_1(6) + \alpha_2(4) + \alpha_3(9) = +1 \quad \alpha_1 = 13/4$$

$$\alpha_1(4) + \alpha_2(6) + \alpha_3(4) = +1 \quad \alpha_2 = 13/4$$

$$\alpha_1(9) + \alpha_2(9) + \alpha_3(17) = -1 \quad \alpha_3 = -1/2$$

$$w = \alpha_1 s_1 + \alpha_2 s_2 + \alpha_3 s_3$$

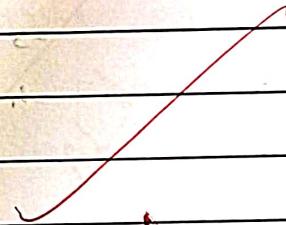
$$= B/4 \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} + 13/4 \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} + -7/2 \begin{bmatrix} 4 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 3 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} +1 \\ 0 \end{bmatrix} \quad b = 3$$

$$\therefore b+3=0$$

intercept on  $x$ -axis = 3

$\Sigma = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow$  line parallel to  $y$ -axis



$$\begin{bmatrix} p \\ q \\ r \end{bmatrix} = 28 \quad \begin{bmatrix} s \\ t \\ u \end{bmatrix} = 18 \quad \begin{bmatrix} v \\ w \\ x \end{bmatrix} = 18$$

~~gut~~

$$\begin{bmatrix} p \\ q \\ r \end{bmatrix} = 28 \quad \begin{bmatrix} s \\ t \\ u \end{bmatrix} = 38 \quad \begin{bmatrix} v \\ w \\ x \end{bmatrix} = 18$$

$$1+ = 2e^{2x} + 2e^{3x} + 2e^{2x}$$

$$1+ = 2e^{2x} + 2e^{3x} + 2e^{2x}$$

$$1+ = 2e^{2x} + 2e^{3x} + 2e^{2x}$$

$$10 = 10 \quad 1+ = (p)e^{2x} + (q)e^{3x} + (r)x$$

$$10 = 10 \quad 1+ = (p)e^{2x} + (q)e^{3x} + (r)x$$

$$10 = 10 \quad 1+ = (p)e^{2x} + (q)e^{3x} + (r)x$$

import hand as pd  
from sklearn.model\_selection import train-test-split  
from sklearn.svm import SVC  
from sklearn.preprocessing import LabelEncoder  
from sklearn.metrics import accuracy\_score,  
confusion\_matrix.

iris\_df = pd.read\_csv("content/iris(1)(2).csv")

X = iris\_df.drop(["Species"], axis=1)  
y = LabelEncoder().fit\_transform(iris\_df["Species"])

X\_train, X\_test, y\_train, y\_test = train-test-split  
(X, y, test\_size=0.5, random\_state=42)

Svm\_linear = SVC(kernel='linear')

Svm\_linear.fit(X\_train, y\_train)

y\_pred\_linear = svm\_linear.predict(X\_test)

Print("Linear Kernel: ")

Print("Accuracy:", accuracy\_score(y\_test, y\_pred\_linear))

Print("Confusion Matrix: \n", confusion\_matrix(y\_test, y\_pred\_linear))

Svm\_rbf = SVC(kernel='rbf')

Svm\_rbf = svm\_rbf.predict(X\_test)

Print ("\\nRBF Kernel")  
Print ("Accuracy : " accuracy - &ou (y-test,  
y-pred - rbf))

Print ("confusion matrix : ") confusion-matrix  
(y-test, y-pred - rbf))

From OUTPUT (matrix) we find that - 10 0 0

linear Kernel

Accuracy: 1.0

confusion matrix:

$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \end{bmatrix}$$

((0.00) = 0.2222222222222222)

RBF Kernel:

Accuracy: 1.0

confusion matrix:

$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

W  
M

Q. what is the accuracy score of the classifier using the linear kernel on RBF kernel.

Ans Linear Kernel Accuracy =  $\approx 0.98$

RBF Kernel Accuracy =  $\approx 1.00$

Q. which kernel gave better performance on the IRIS dataset

Ans RBF Kernel

## 1. Difference B/w Decision Tree And Random Forest

### Feature Decision Tree vs Random Forest

**Definition** A Single tree-based model that splits data based on features (usually trained with random samples).

|               |                                           |                                        |
|---------------|-------------------------------------------|----------------------------------------|
| Bias-variance | Low bias, high-variance (overfits easily) | low variance and better generalization |
|---------------|-------------------------------------------|----------------------------------------|

|          |                                          |                                                  |
|----------|------------------------------------------|--------------------------------------------------|
| Accuracy | Less accurate than RF for large datasets | Generally more accurate due to ensemble learning |
|----------|------------------------------------------|--------------------------------------------------|

|             |                                    |                                                    |
|-------------|------------------------------------|----------------------------------------------------|
| overfitting | Prone to overfitting on noisy data | less prone due to averaging across multiple trees. |
|-------------|------------------------------------|----------------------------------------------------|

|                  |                   |                     |
|------------------|-------------------|---------------------|
| Interpretability | Easy to interpret | Harder to interpret |
|------------------|-------------------|---------------------|

|                |                       |                           |
|----------------|-----------------------|---------------------------|
| Training Speed | Faster (single model) | Slower. (multiple models) |
|----------------|-----------------------|---------------------------|

|                  |                   |                                                |
|------------------|-------------------|------------------------------------------------|
| Prediction Speed | Faster (one tree) | Slower (aggregates prediction from many trees) |
|------------------|-------------------|------------------------------------------------|

## \* Parameters of Random Forest classifier ()

- (1) n-estimators
- (2) criterion
- (3) max\_depth
- (4) min\_sample\_split
- (5) min\_weight\_fraction\_leaf
- (6) min\_sample\_leaf
- (7) max\_features
- (8) max\_leaf\_nodes
- (9) min\_impurity\_decrease
- (10) bootstrap
- (11) oob\_score
- (12) n\_jobs
- (13) random\_state
- (14) verbose
- (15) class\_weight
- (16) warm\_start

## \* Algorithm of Random Forest.

(1) Input: Dataset with features and target

2. For  $i=1$  to  $n$ -estimators:-

- o Draw a bootstrap sample from training set
- o Train a Decision Tree on this sample

- At each node, select a random subset of features.

- use the best split among the selected features.

- Grow the tree to maximum depth.

### 3. Aggregate the predictions:

- For classification : use majority voting from all trees.

- For regression : use average prediction from all trees.

### 4. Output : Final Prediction

Code

```
import pandas as pd  
from sklearn.model_selection import train_test_split  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy_score, confusion_matrix  
from sklearn.preprocessing import LabelEncoder
```

```
clf = pd.read_csv("/content/train-train.csv")
```

```
Print("columns: \n"), df.columns
```

```
label_encoders = {}
```

```
for column in df.columns:
```

```
    if df[column].dtypes == 'object':
```

```
        le = LabelEncoder()
```

```
        df[column] = le.fit_transform(df[column].astype('str'))
```

```
label_encoders[column] = le
```

```
X = df.iloc[:, :-1]
```

```
y = df.iloc[:, -1]
```

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

```
rf_model = RandomForestClassifier(n_estimators = 100,  
random_state = 42)
```

```
rf_model.fit(X_train, y_train)
```

```
y_pred = rf_model.predict(X_test)
```

accuracy = accuracy - score (y-test, y-pred)

conf. matrix = confusion matrix (y-test, y-pred)

Print(f"\nAccuracy Score: {accuracy : 4f}")

Print(" ") in confusion matrix")

Print(conf-matrix)

### Output Data As ("array") Is:

columns =

Index(['Passenger Id', 'Survived', 'Pclass',  
'Name', 'Sex', 'Age', 'Sib Sp',  
'Parch', 'Ticket', 'Fare', 'cabin',  
'Embarked'], 0, type = object)

Accuracy Score = 0.8659

Confusion Matrix =

[ [ 26 0 17 ] ]

[ [ 0 12 5 ] ]

[ [ 82 0 11 ] ]

- ① write an algorithm of K means?
  - (i) choose K initial centroids
  - (ii) Assign each point to the nearest centroid.
  - (iii) Recalculate centroids.
  - (iv) Repeat until convergence.

- ② How to determine no. of clusters?
  - o elbow method:
  - o silhouette method:

Sum of Squared errors (SSE):

$$SSE = \sum_{i=1}^k \sum_{j \in C_i} (x_j - \mu_i)^2$$

where  $c_i$  = cluster

$\mu_i$  = centroid.

- ③ what is the formula of sum of squared errors? Plot SSE vs no. of cluster

- ④ Desirable elbows technique.

= Plot SSE against K. The elbow is where SSE start decreasing at a slower rate.  
This is the optimal K.

- ⑤ Discuss all the parameters used in Kmeans()

$\Rightarrow$  n-cluster, n-point, a+01, i+01

import pandas as pd.

import matplotlib.pyplot as plt.

From sklearn . data set import load\_iris.

from sklearn . preprocessing import StandardScaler

from sklearn . cluster import KMeans.

iris = load\_iris()

data = pd. DataFrame (iris . data, columns = iris . feature

x = data [ [ "Petal length (cm)", "Petal width (cm)" ] ]

scaler = StandardScaler().

x\_scaled = scaler. fit\_transform(x)

inertia = [ ]

K-range = range (1, 11)

For K is in K-range :

KMeans = KMeans (n\_clusters = K, random\_state = 42)

KMeans . Fit (x\_scaled)

inertia . append (KMeans . inertia - )

Plt. figure (fig size (8,5))

Plt. plot (K-range, inertia, marker = 'o')

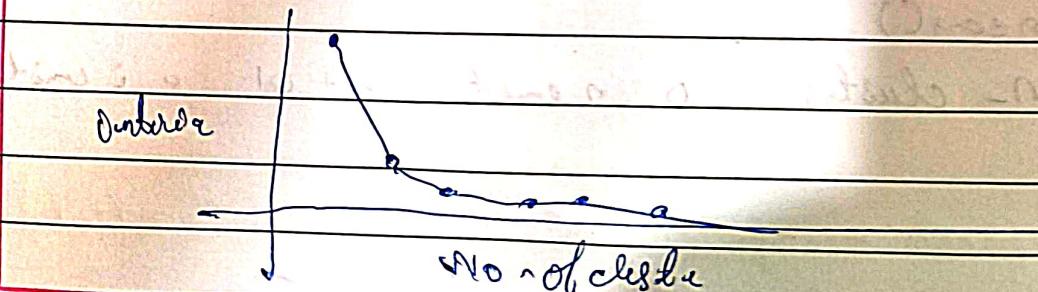
Plt. title ("elbow method for optimal K")

Plt. xlabel ("Number of clusters (K)")

Plt. ylabel ("inertia")

Plt. grid (True)

Plt. show()



## LAB-11 :- PCA

from sklearn.datasets import load\_digits

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

from sklearn.metrics import accuracy\_score.

digits = load\_digits()

X = digits.data

y = digits.target

Scalar = StandardScaler()

X\_scaled = scalar.fit\_transform(X)

Pca = PCA(n\_components=2)

X\_Pca = Pca.fit\_transform(X\_scaled)

X\_train, X\_test, y\_train, y\_test = train\_test\_split(skl)

(X\_Pca, y, test\_size=0.2, random\_state=42)

model = LogisticRegression()

model.fit(X\_train, y\_train)

y\_pred = model.predict(X\_test)

accuracy = accuracy\_score(y\_test, y\_pred)

Score = model.score(X\_test, y\_test)

Print("Model Score (Accuracy using score) :  
", round(score, 2))

Print("Accuracy using PCA with 2 component, noen,

DotNet  
Bokeh

A37

⇒ OUTPUT

Model Score :- 0.5389

Accuracy :- 0.5319

(Model - 3)

Model - 3

(Model - 3) Model - 3

(Model - 3) Model - 3

(Model - 3) Model - 3

(Model - 3) Model - 3

Model - 3 Model - 3

Model - 3 Model - 3

Model - 3 Model - 3

Model - 3