

Lab 2:(I) Dataset: housing.csv

import pandas as pd

(i) Load csv file:

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

(ii) Display all columns:

df.columns

(iii) Statistical

index(['longitude', 'latitude', 'housing_median_age',
'total_rooms', ...], dtype='object')

(iv) Statistical information:

df.describe()

longitude latitude housing_median_age total_rooms total_bedrooms
count 20640 20640 20640 20640 20640population households median_income median_value
count 20640 20640 20640 20640

(v) Count of unique values for "ocean_proximity" column:

cnt = len(df['ocean_proximity'].unique())
print(cnt)

→ 5

(vi) display which ~~attribute~~ attributes showing missing
values count greater than 0.

`print(df1.columns[df1.isnull().any()])`

⇒ `Index(['total_bedroom'], dtype='object')`

(ii) Dataset: `diabetes.csv`

`df2 = pd.read_csv('diabetes.csv')`

(i) which columns have missing values:

`print(df2.columns[df2.isnull().any()])`

⇒ `Index([], dtype='object')`

(ii) Categorical columns:

`cat1 = df2.select_dtypes(exclude=['Number']).columns`
`print(cat1)`

⇒ `Index(['Gender', 'Class'], dtype='object')`

(iii)

Dataset: `adult.csv`

`df3 = pd.read_csv('adult.csv')`

(i) Missing values:

`print(df3.columns[df3.isnull().any()])`

⇒ `Index([], dtype='object')`

(ii) Categorical columns:

`cat2 = df3.select_dtypes(exclude=['Number']).columns`
`cat2`

⇒ `Index(['workclass', 'education', 'marital-status', 'occupation',
 'relationship', 'race', 'gender', 'native-country', income
 dtype='object')`

Lab-1:

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

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

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

hdfc_data = data['HDFCBANK.NS']

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

print("In Daily Returns for HDFC Industries: ")

print(hdfc_data['Daily Return'].head())

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

plt.subplot(2, 1, 1)

hdfc_data['close'].plot(title='HDFC Industries - (Closing Price)')

plt.subplot(2, 1, 2)

hdfc_data['Daily Return'].plot(title='HDFC ')

color="orange"

plt.title('Layout')

plt.show()

→ Daily Returns to HDFC Industries:

Data

2024-01-01

NaN

2024-01-02

0.000794

2024-01-03

-0.015420

2024-01-04

-0.018270

2024-01-05

-0.005116

import pandas as pd

1. data = {'USN': [1, 2, 3, 4, 5], 'Name': ['A', 'B', 'C', 'D', 'E'], 'Marks': [1020, 1040, 1030]}

df = pd.DataFrame(data)

print(df)

USN Name Marks

0	1	A	10
1	2	B	10
2	3	C	10
3	4	D	90
4	5	E	80

2. from sklearn.datasets import load_digits

digits = load_digits()

df = pd.DataFrame(digits.data, columns=digits.feature)

df['target'] = digits.target

df.head()

3. file_path = "sample_data.csv"

df = pd.read_csv(file_path)

print(sample_data)

	Product	Quantity	Price	Info	Region
0	Laptop	5	1000	\$200	North
1	Mouse	15	20	\$200	West

4. df.to_csv('output.csv', index=False)

print("Data saved to output.csv")

By using Python

Lab-4:

Linear Regression and Multiple Linear Regression:

Find linear regression of data of week and profit.

x_i (week)	y_i (Sale in thousands)
1	2.8
2	4.2
3	5.2
4	9.0

$$\beta = ((X^T X)^{-1} X^T) Y$$

$$X^T = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 9 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 4 & 9 & 16 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 9 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 20 \\ 20 & 120 \end{bmatrix} \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

$$(X^T X)^{-1} = \begin{bmatrix} 1.2115 & -0.192 \\ -0.192 & 0.0389 \end{bmatrix}$$

~~$$(X^T X)^{-1} X^T = \begin{bmatrix} 1.2115 & -0.192 \\ -0.192 & 0.0389 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 4 & 9 \end{bmatrix}$$~~

$$= \begin{bmatrix} 0.8295 & 0.4435 & 0.2515 & -0.5165 \\ -0.1152 & -0.0384 & 0 & 0.1576 \end{bmatrix}$$

~~$$(X^T X)^{-1} X^T Y = \text{Value}$$~~

~~$$= \begin{bmatrix} 0.8295 & 0.4435 & 0.2575 & -0.5165 \\ -0.1152 & -0.0384 & 0 & 0.1536 \end{bmatrix} \begin{bmatrix} R \\ Y \\ Z \\ Q \end{bmatrix}$$~~

~~$$= \begin{bmatrix} 0.038 \\ 0.9984 \end{bmatrix}$$~~

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

$$(X^T X)^{-1} X^T = \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ 0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix}$$

$$2Y = 2 \begin{bmatrix} 3 \\ 4 \\ 5 \\ 9 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$$

$$Y = \beta_0 + \beta_1 X$$

Q1. Predict Canada's per capita income in year 2020.

import pandas as pd
from sklearn import linear_model

```
df = pd.read_csv('canada-per-capita-income.csv')
```

```
df.columns = ['Year', 'per-capita income']
```

```
X = df[['Year']]
```

```
y = df['per-capita income']
```

```
reg = linear_model.LinearRegression()
```

```
reg.fit(X, y)
```

```
predicted = reg.predict([[2020]])
```

```
print(predicted)
```

$\Rightarrow 41288.694$

Q2. Predict salary of employee with 12 years of experience.

(Salary.csv)

$\Rightarrow 139574.041$

Q3. ~~Predict salaries for following candidates:~~

1) 2 years experience, 9 test score, 6 interview score

2) 11 years experience, 10 test score, 10 interview score

using dataset hiring.csv

```

from wordcloud import WordCloud
import pandas as pd
from sklearn import linear_model
df = pd.read_csv('hiring.csv')

```

```

df['experience'] = df['experience'].fillna(0)
df['experience'] = df['experience'].apply(lambda x: str(x).lower().strip())
df['test-score (out of 10)'] = df['test-score'].fillna(
    df['test-score'].median())

```

```

reg = linear_model.LinearRegression()
reg.fit(df[['experience', 'test-score', 'interview_score'],
           df['salary']])

```

```

print(reg.predict([[2, 9, 6]]))
print(reg.predict([[12, 10, 10]]))

```

$\Rightarrow [5305.967]$
 $[92002.183]$

d4. Predict profit using 1000 companies. Consider following.

91894.48 R&D spend, \$15846.3 Admin,
 31921.24 Marketing spend, Florida state

$\Rightarrow [511017.346]$

~~AB~~

Lab-5: Logistic Regression

(1) Binary Classification (Sigmoid Funcⁿ)

$$\theta_0 = -5, \theta_1 = 0.8$$

$$\text{Linear eq}^n \Rightarrow Z = \theta_0 + \theta_1 x$$

$p(x)$ for $x=7$

$$y = p(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}} = \frac{1}{1 + e^{-(5 + 0.8 \cdot 7)}} = 0.64$$

$$0.64 \geq 0.5 \Rightarrow 1 \text{ pass}$$

A student who studies for 7 hrs will pass.
The predicted class is pass.

(2) Multiclass Regression (Softmax funcⁿ)

$Z = [z_1, z_2, z_3]$ find probability values of the 3 classes.

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

$$\text{softmax}(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} = 0.665 = 66.5\%$$

$$\text{softmax}(z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} = 0.241 (= 24.1\%)$$

~~$$\text{softmax}(z_3) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}} = 0.091 = 9.1\%$$~~

1) insurance_data.csv:

```

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}

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

 $y_{\text{predicted}}$

$\text{model.score}(X_{\text{test}}, y_{\text{test}})$

$\text{model.predict_proba}(X_{\text{test}})$

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

 $y_{\text{predicted}}$ model.coef model.intercept import math

$\text{def sigmoid}(x):$

$\text{return } 1 / (1 + \text{math.exp}(-x))$

$\text{def prediction_function}(x):$

$z = 0.127 \times \text{age} - 4.973$

$y = \text{sigmoid}(z)$

$\text{return } y$

$\text{age} = 35$

Prediction function (age)

→ 0.3701839769532725

2)

iris.csv:

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

iris = pd.read_csv('iris.csv')

iris.head()

X = iris.drop('Species', axis='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(X_train, y_train)

y_pred = model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

print(accuracy)

⇒ 1.00

1. HR Comma Sep. SV:

- (i) Salary has a direct impact on employee retention because as per the data, employees who have less salary have a lower retention rate and vice versa.
- (ii) The accuracy of our ML Model is 79.63%. The accuracy is decent (neither good nor bad).

2. Zoo dataset:

- (i) I have dropped the columns 'animal name' and 'class-type' as they had categorical values.
- (ii) There weren't any missing values in dataset.
- (iii) The confusion matrix tell us that there are higher number of true positive values.
- (iv) Reptile class has been misclassified as Fresh class.

✓
✓
✓
✓

Lab - 03

Instance	α_2	α_3	Classification
----------	------------	------------	----------------

1	Hot	High	No
2	Hot	High	No
6	Cool	High	No
7	Hot	High	No
8	Hot	Normal	Yes

Attribute: α_2 values (α_2) = Hot, Cool

$$S_{\text{tot}} = [1_f, 3_g] \quad \text{Entropy}(\Sigma_{\text{tot}}) = -\frac{1}{2} \log(\frac{1}{2}) - \frac{3}{2} \log(\frac{3}{5})$$

$$S_{\text{hot}} = [1_f, 3_g] \quad \text{Entropy}(\Sigma_{\text{hot}}) = -\frac{1}{2} \log(\frac{1}{2}) - \frac{3}{4} \log(\frac{3}{4}) = 0.8112$$

$$S_{\text{cool}} = [0_f, 1_g] \quad \text{Entropy}(\Sigma_{\text{cool}}) = 0.0$$

$$\text{Gain}(\alpha_2) = \text{Entropy}(\Sigma_{\text{tot}}) - \frac{1}{2} \left[\text{Entropy}(\Sigma_{\text{hot}}) + \text{Entropy}(\Sigma_{\text{cool}}) \right]$$

$$= 0.9219 - \frac{1}{2} (0.8112) = 0.3219$$

Attribute: α_3

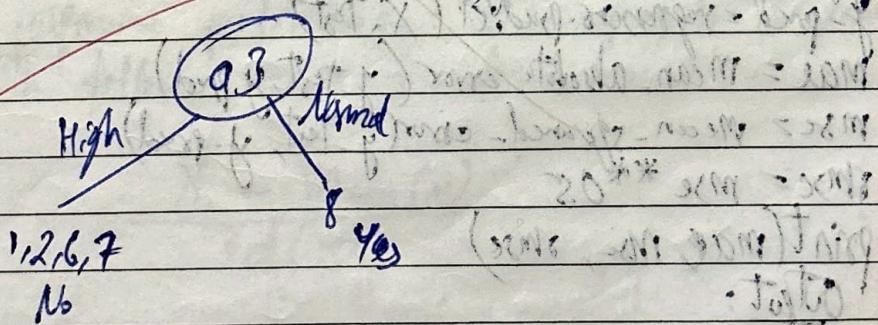
values (α_3) = High, Normal

$$S_{\text{tot}} = 0.7219$$

$$S_{\text{high}} = [0_f, 4_g] = 0$$

$$S_{\text{normal}} = 0$$

$$\text{Gain}(\alpha_3) = 0.9219 - (4/5)0 - (1/5)0 = 0.7219$$



2.29 : 34/11

2.18/21 : 32/11

2.22/21 : 32/11

1. The accuracy achieved is 1.0 (100%). Confusion matrix will show that all values along diagonal and zeros elsewhere. There are no misclassifications.
2. We need to visualize the tree, analyze the nodes and identify important features. In decision tree classifier discrete class labels are predicted for each instance. In regression tree continuous value for each instance. It aims to create split that minimize the variance of target variable.

Petrol consumption.csv:

```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error
df = pd.read_csv('petrol_consumption.csv')
X = df.drop('Petrol_Consumption', axis=1)
y = df['Petrol_Consumption']
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=42)
regressor = DecisionTreeRegressor()
regressor.fit(X_train, y_train)
y_pred = regressor.predict(X_test)
mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
rmse = mse ** 0.5
print(mae, mse, rmse)

```

Output:

MAE: 95.6

NSE: 1283.8

RMSE: 133.5357

Lab-6:

- Q. Consider the following dataset for $K=3$ and test data: $(X, 35, 10)$ as $(\text{Person}, \text{Age}, \text{Salary})$. Solve using KNN classifier model to predict the target.

Person	Age	Salary	Rank	Target dist.	Rank
A	18	150	N	82.81	5
B	22	55	N	96.57	9
C	24	70	N	31.95	2
D	41	60	Y	48.49	3
E	43	70	Y	81.89	1
F	38	90	Y	60.79	6
X	35	100			

$$K = 3$$

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

The target predicted is \hat{Y}

- Q1. To find value of K , perform iteration with range of K values. Perform evaluation for each K , and finally select the one with highest accuracy and lowest rate.

- Q2. Feature Scaling is used for standardization and to improve performance. To perform standard scalar on the following formula:

$$\frac{x_i - \text{mean}(x)}{\text{sd}(x)}$$

Original	Mean	Standard Deviation	Standardized
18	35	10	0.1
22	35	10	0.6
24	35	10	0.0

```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.datasets import load_iris
accuracy_score
iris = load_iris()
X = iris.data
y = iris.target
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=42)

```

~~KNN~~

KNN = KNeighborsClassifier(n_neighbors=3)

KNN.fit(X_train, y_train)

y_pred = KNN.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

conf_matrix = confusion_matrix(y_test, y_pred)

class_report = classification_report(y_test, y_pred)

print(accuracy)

print(conf_matrix)

print(class_report)

Output:

Accuracy : 1.0

Confusion Matrix :

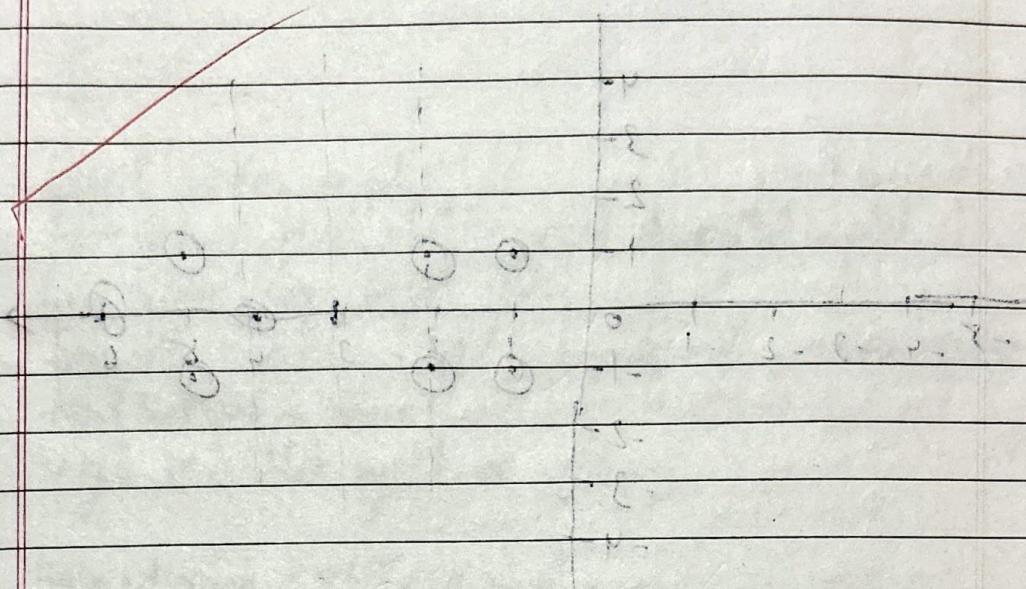
$$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

Classification Report:

	Precision	recall	F1-score	Support
0	1.00	1	1	10
1	1.00	1	1	9
2	1.00	1	1	11

accuracy = 1.0 Precision recall f1-score: Support
 Macro avg + +
 weighted avg + +

(0.0) (-2) (1.0) (0.8) (1.0) (1.0) (1.0) (1.0)



$$\hat{y} = \{x\} \cdot 2 + \{z\} = 1.2$$

$$\hat{y} = \{x\} \cdot 2 + \{z\} = 1.2$$

$$y = 2x + b + 2z$$

$$y = 2x + 2z + b + 2z$$

$$y = 2x + 2z + b + 2z$$

$$y = 2x + b + 2z$$

$$y = 2x + b + 2z$$

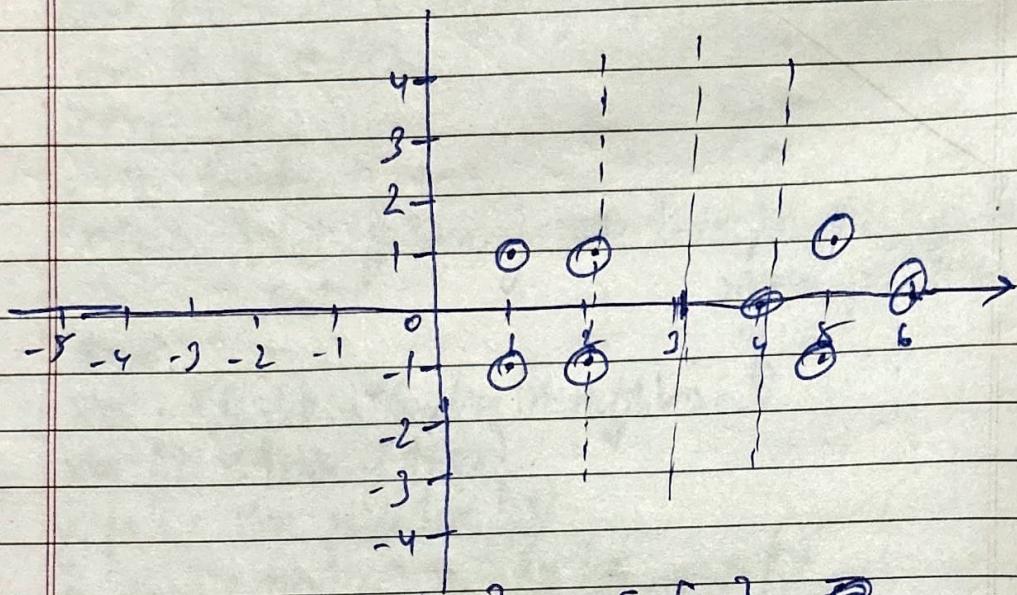
$$y = 2x + b + 2z$$

$$(2) \cdot (1) \cdot (1) x^2 + (2) \cdot (1) x + (2) x^2$$

Lab 7:

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

$$(1,1) (2,1) (1,-1) (2,-1) (4,0) (5,1) (5,-1) (6,0)$$



$$\bar{S}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad \bar{S}_2 = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \quad \bar{S}_3 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

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

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

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

$$\alpha_1 \bar{s}_1 \bar{s}_3 + \alpha_2 \bar{s}_2 \bar{s}_3 + \alpha_3 \bar{s}_3 \bar{s}_3 = -1$$

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

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

$$\alpha_1(9) + \alpha_2(9) + \alpha_3(17) = -1 \quad \alpha_3 = -7/2$$

$$\begin{aligned} w &= \alpha_1 \bar{s}_1 + \alpha_2 \bar{s}_2 + \alpha_3 \bar{s}_3 \\ &= \frac{13}{4} \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \frac{13}{4} \begin{bmatrix} 2 \\ -1 \end{bmatrix} + -\frac{7}{2} \begin{bmatrix} 4 \\ 0 \end{bmatrix} = \begin{bmatrix} -17 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \end{aligned}$$

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

intercept on x -axis = 3

$z(1) \Rightarrow$ line parallel to y -axis

iris.csv:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, confusion_matrix
import random as rn
import matplotlib.pyplot as plt
```

iris = pd.read_csv('iris.csv')

X = iris.drop(['species'], axis=1)

y = iris['species']

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

svm_rbf = OneVsRestClassifier(SVC(kernel='rbf', random_state=42,

probability=True))

svm_rbf.fit(X_train, y_train)

y_pred = svm_rbf.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

print(accuracy)

conf_matrix = confusion_matrix(y_test, y_pred)

plt.sns.heatmap(conf_matrix, annot=True)

fpr = {}

tpr = {}

roc_auc = {}

for i in range(26):

fpr[i], tpr[i], _ = roc_curve(y_test, letterf[:, i], y_pred_prob[:, i])

~~roc_auc(i) = roc(Avr[i], tpr[i])~~

~~for i in range(26):~~

~~plt.plot(Avr[i], tpr[i])~~

~~for i~~

Kernels = ['linear', 'rbf']

for K in Kernels:

SVM_classifier = SVC(K=K)

SVM_classifier.fit(X2_train, y_train)

y_pred = SVM_classifier.predict(X2_test)

accuracy = accuracy_score(y_test, y_pred)

print(K, ":", accuracy)

cm = confusion_matrix(y_test, y_pred)

sns.heatmap(cm, annot=True)

plt.show()

~~Accuracy of SVM classifier with linear kernel: 1.0
Confusion Matrix:~~

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

~~Accuracy of SVM classifier with rbf kernel: 1.0
Confusion matrix:~~

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

- Q1. Accuracy of SVM classifier with Linear Kernel: 1.0
with rbf Kernel: 1.0

Since both kernels perform equally well, using the one with simplicity and computational efficiency could be better choice.
In this case linear kernel would be a more practical choice due to its simplicity and efficiency.

- Q2. The confusion matrix shows the model's prediction breakdown highlighting true positives on the diagonal and misclassifications.
The AUC score measures model performance in binary classification. For multi-class, AUC is calculated for each class against the others. AUC of 1 indicates perfect performance. The combination of accuracy and AUC provides insight into the model's ability to differentiate b/w classes. The letter-recognition dataset is more complex than the iris dataset, so SVM performance is expected to be lower.

Lab 8 - Random Forest:

Decision Tree

→ A decision tree is a tree like model of decisions along with possible outcomes in a diagram.

→ There is always a scope for overfitting.

→ Results are not discrete.

→ Requires less computation.

→ Easy to visualize.

Parameters of Random Forest

- Max features: Max. no. of features RF is allowed to be in an individual tree. Features are num, sqrt, log2.
- n_estimators: Number of trees you want to build before taking max. voting or averaging of prediction.
- min_sample_leaf:
- n_jobs: No. of processors allowed.
- random_state:
- oob_score: cross-validation method

Algorithm:

- 1) Select random samples from a given data or training set.
- 2) The algorithm will construct a decision tree for every training data.
- 3) Voting will take place by averaging the decision tree.

4) Finally, select the most voted prediction result as the final prediction result.

Accuracy with 10 n-estimators = $10 = 1.0$

" with 50 " = $50 = 0.10$

" 100 " = $100 = 0.0$

~~accuracy decreases as number of estimators increases~~

~~estimation time, otherwise not feasible~~

~~accuracy becomes stable after 100 estimators~~

~~changes begin but no further~~

~~accuracy remains stable after 100 estimators~~

P _{f1}	P _{f2}	P _{f3}	P _{f4}	P _{f5}	P _{f6}
F	G	H	I	J	K
P1	Z	R	S	T	U

$$\hat{y} = \frac{(f_1 + f_2 + f_3 + \dots + f_n)}{n} L = \bar{x}$$

$$2\hat{y} = \frac{(P_1 + P_2 + P_3 + \dots + P_n)}{n} L = \bar{x}$$

~~existing output~~

$$(x - \bar{x})^2 = \sum (x_i - \bar{x})^2$$

$$[(P_1 - \bar{P})^2 + (P_2 - \bar{P})^2 + (P_3 - \bar{P})^2] L$$

~~P1~~

$$(x - \bar{x})^2 = \sum (x_i - \bar{x})^2$$

$$[(P_1 - \bar{P})(P_2 - \bar{P}) + (P_2 - \bar{P})(P_3 - \bar{P})] L$$

~~L~~

Lab-11:PCA:

1. Calculate mean
2. Calculation of covariance matrix
3. Eigenvalues of covariance matrix
4. Computation of the eigen vectors - unit eigenvectors
5. Computation of first principal components
6. Geometrical meaning of first principal components.
- Q. Given the table, reduce dimension from 2 to 1 using PCA:

Feature	$Eg\ 1$	$Eg\ 2$	$Eg\ 3$	$Eg\ 4$
x_1	9	8	13	7
x_2	11	4	5	14

Mean:

$$1) \bar{x}_1 = \frac{1}{4}(9+8+13+7) = 8$$

$$\bar{x}_2 = \frac{1}{4}(11+4+5+14) = 8.5$$

2) Covariance matrix:

$$\begin{aligned} \text{Cov}(x_1, x_1) &= \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)^2 \\ &= \frac{1}{3} [(9-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2] \\ &= 14 \end{aligned}$$

$$\begin{aligned} \text{Cov}(x_1, x_2) &= \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2) \\ &= \frac{1}{3} [(9-8)(11-8.5) + (8-8)(8-8.5) + (13-8)(5-8.5) + (7-8)(14-8.5)] \\ &= -11 \end{aligned}$$

$$\text{cov}(X_2, X_1) = \text{cov}(X_1, X_2)$$

$$= -11$$

$$\text{cov}(X_2, X_2) = \frac{1}{3} \sum_{k=1}^n (x_{2k} - \bar{x}_2)^2$$

$$= \frac{1}{3} [(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2]$$

$$= 23$$

$$\textcircled{2} \quad S = \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) \\ \text{cov}(X_2, X_1) & \text{cov}(X_2, X_2) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\textcircled{3} \quad \det(S - \lambda I) = \begin{vmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{vmatrix}$$

$$= (14-\lambda)(23-\lambda) - 11^2$$

$$\lambda = \frac{1}{2} (37 \pm \sqrt{565}) = 30.3849, 6.6151$$

$$= \lambda_1, \lambda_2$$

$$\textcircled{4} \quad U = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = (S - \lambda I) X$$

$$= \begin{bmatrix} 14-\lambda u_1 & -11 \\ -11 & 23-\lambda \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$= \begin{bmatrix} (14-\lambda_1)u_1 - 11u_2 \\ -11u_1 + (23-\lambda_1)u_2 \end{bmatrix} = 0$$

$$(14-\lambda_1)u_1 - 11u_2 = 0$$

$$-11u_1 + (23-\lambda_1)u_2 = 0$$

$$\frac{u_1}{11} = \frac{u_2}{14-\lambda_1} = t_1$$

$$u_1 = 11t$$

$$u_2 = (14-\lambda_1)t$$

$$U_1 = \begin{bmatrix} 11 \\ 14-\lambda \end{bmatrix}$$

$$\|U_1\| = \sqrt{11^2 + (14-\lambda)^2} = \sqrt{11^2 + (14 - 30.7849)^2} = 19.7348$$

$$E_1 = \begin{bmatrix} 11/\|U_1\| \\ (14-\lambda)/\|U_1\| \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$E_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

5)

$$\begin{bmatrix} X_{1K} \\ X_{2K} \end{bmatrix}$$

$$e_1^T \begin{bmatrix} X_{1K} - \bar{x}_1 \\ X_{2K} - \bar{x}_2 \end{bmatrix} = (0.5574 - 0.8303) \begin{bmatrix} X_{1K} - \bar{x}_1 \\ X_{2K} - \bar{x}_2 \end{bmatrix} \\ = 0.5574(X_{1K} - \bar{x}_1) - 0.8303(X_{2K} - \bar{x}_2)$$

$$X_{(2-2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 11 - (\bar{x}_1 - \bar{x}_2) \\ 14 - (\bar{x}_1 - \bar{x}_2) \end{bmatrix}$$

$$= \begin{bmatrix} 0.1111 - 0.1111 \\ 0.1111 - 0.1111 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$= 0.1111 - 0.1111$$

$$= 0.1111 - 0.1111$$

$$= 0.1111 - 0.1111$$

$$= 0.1111 - 0.1111$$

Code for Randomforest:

```

from sklearn.datasets import load_iris
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

X,y = iris.data, iris.target
y_binarized = LabelBinarizer().fit_transform(y, classes=[0,1,2])
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=42)

n_estimators_list = [50, 100, 200]
all_scores = []
rf_model_default = RandomForestClassifier(n_estimators=10, random_state=42)

# n_estimators_default = RandomForestClassifier()
# rf_model_default.fit(X_train, y_train)
# y_pred_default = rf_model_default(X_test)
# accuracy_default = accuracy_score(y_test, y_pred_default)
# print(accuracy_default)
# best_accuracy = 0
# best_n_estimators = 10

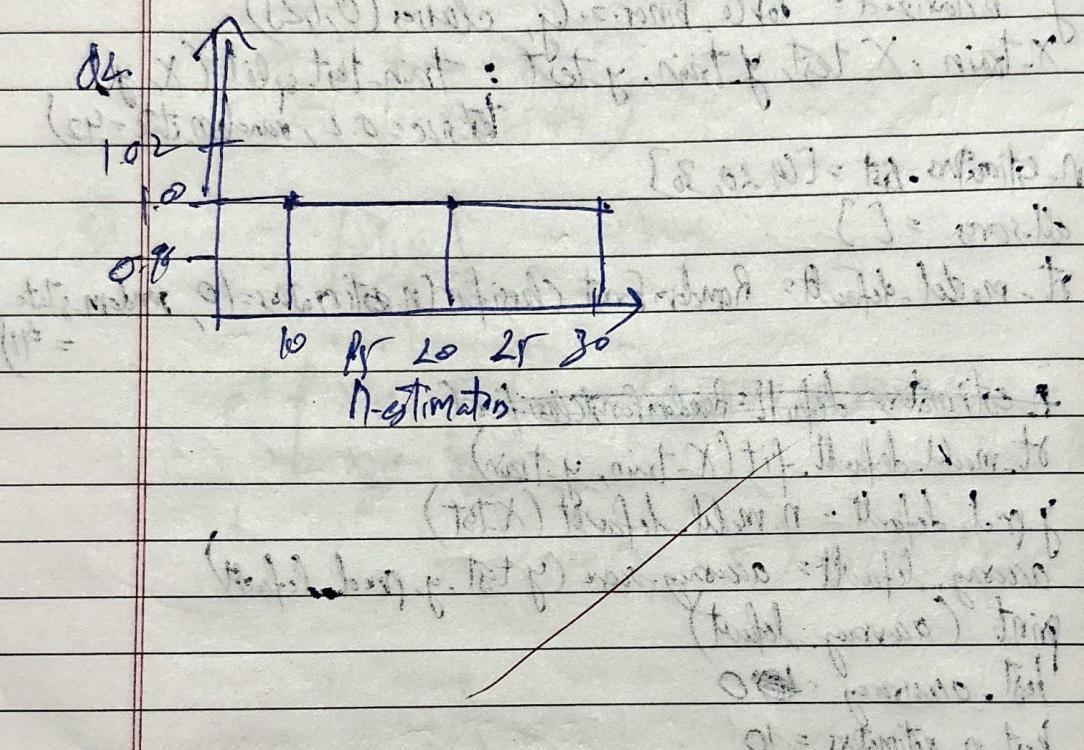
for n_estimators in [10, 50, 100, 200, 500]:
    rf_model = RandomForestClassifier(n_estimators=n_estimators, random_state=42)
    rf_model.fit(X_train, y_train)
    y_pred = rf_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    if accuracy > best_accuracy:
        best_accuracy = accuracy
        best_n_estimators = n_estimators
    print(f'best_accuracy) {best_accuracy}

```

Output:

Accuracy with default $H = 100\%$,
Best accuracy $\approx 100\%$.

Q1. Accuracy remains same for $N = 10, 20, 80$
accuracy = 1.00



($\text{cr} = \text{state variable, estimated by direct function}$) $\text{fwd} = \text{forward - library}$

Lab-9 : Boosting:

(Q1) Boosting is an ensemble learning technique that combines multiple weak learners to create a strong learner. It works by sequentially training models, where each new model focuses on mistakes made by the previous ones. The models' predictions are weighted based on their performance.

Q2. Parameters of AdaBoost classifiers:

- `n_estimators (int)`
- `base_estimator (object, default=None)` = base model to be used for boosting
- `n_estimators (int)` - no. of boosting rounds
- `learning_rate - scaling factor`
- `algorithm ('SAMME', 'SAMME.R', default='SAMME.R')`
- `random_state - controls randomness`
- `loss ('linear', 'square', 'exponential')` - loss function

Q3. Algorithm:

- 1) Initialize: set the weights of all training samples equally.
- 2) For each initialization:
 - Train a base learner on weighted training set
 - Calculate error of the model on training data.
 - Compute model's weight based on error.
 - Update sample weights: increase weight of misclassified samples
- 3) Combine: The final model is the weighted sum of the base learners, where learners with lower errors have higher weights.

Code :

```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
data = pd.read_csv('income.csv')
X = data.drop(columns=['income-level'])
y = data['income-level']
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=42)
ada_model = AdaBoostClassifier()
ada_model.fit(X_train, y_train)
y_pred = ada_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
print(accuracy * 100)
print(conf_matrix)

```

Output :

Accuracy : ~~83.29%~~

Confusion matrix:

7003	911
1223	1132

Lab-10: K-Means:

1. Algorithm:

- 1) Initialize: Randomly select K points as initial cluster centroids.
- 2) Assign: Assign each data point to the closest centroid.
- 3) Update: Compute new centroids as the ~~mean~~ mean of the points in each cluster.
- 4) Repeat steps 2 and 3 until centroids do not change.

2. How to determine no. of clusters (K)

- 1) Elbow Method
- 2) Silhouette Analysis
- 3) Gap Statistic
- 4) Domain Knowledge or business requirements

3. $SSE =$

$$\sum_{i=1}^K \sum_{x \in C_i} \|x - \mu_i\|^2$$

(i = cluster)

μ_i = centroid of cluster i

x = data point in cluster (i)

$\|x - \mu_i\|^2$ = squared Euclidean distance

4. We run K-means for a range of K : (1 to 10)

5.

Plot SSE vs (K)

Observe the point where SSE dramatically drops and levels off

It suggests diminishing returns by adding more clusters.

6. Parameters:

- n-clusters : no of clusters (k)
 - init - method to initialize centroids
 - n-init - times algo to run
 - max_iter - max iteration per run
 - tol - convergence tolerance
 - algorithm - auto, full
 - random_state - seed for reproducibility.

Code:

```
import pandas as pd  
import numpy as np  
from sklearn.cluster import KMeans  
import random  
import matplotlib.pyplot as plt
```

np.random.seed(42)

Names = [f' Person - {i}' for i in range(1, 51)]

`ages = np.random.randint(20, 60, size=50)`

Incomes = np.random.randint(10000, 12000, size=50)

income - df = pd. Data frame

Name : Name,

~~Age : 20 yrs
Income : in Gms~~

~~diff. pd. read-^{er} - cor (income. cor.)~~

`df-numeric = df.drop('Name', axis=1)`

z-scores = Standard Scores

Cdf -> cdf scaled fit transform (Cdf - numeric)

X_{train} , X_{test} , y_{train} , y_{test} pdff-scaled, test size = 0.2,
random_state = 42

$SSE = []$

$K_range = range(1, 11)$

for K in K_range :

$Km = KMeans(n_clusters = K, random_state = 42)$

$Km.fit(Xtrain)$

$de.append(Km.inertia_)$

$plt.plot(Krange, SSE, marker = 'o')$

$plt.show()$

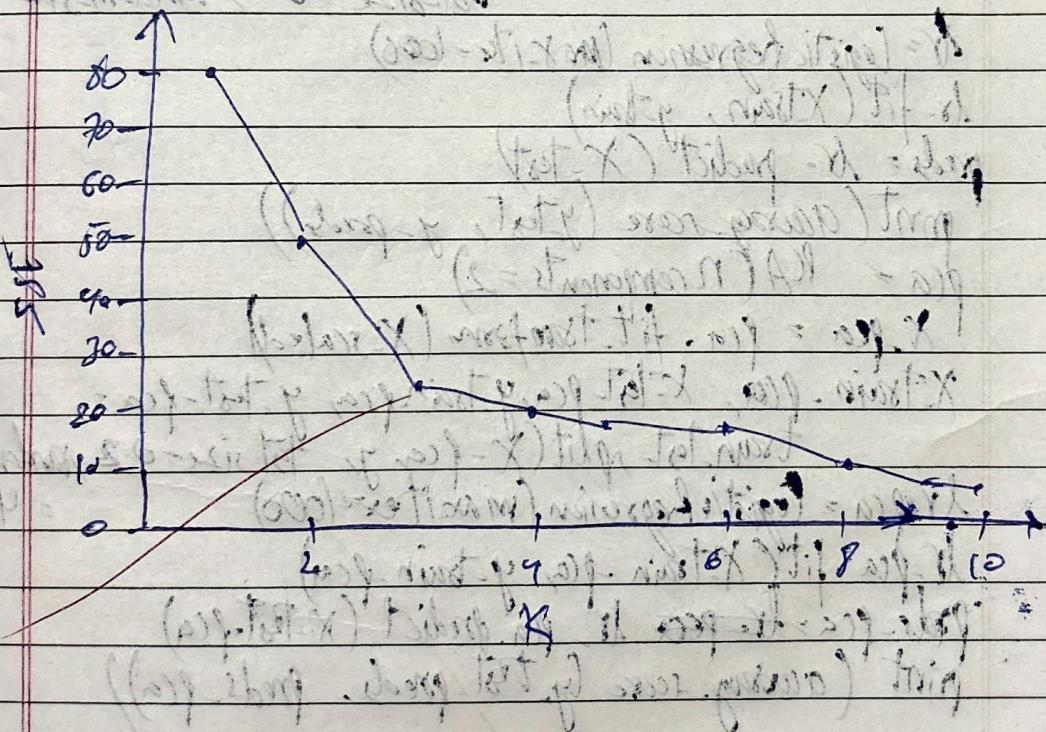
optimal $K = 3$

$Kmeans = KMeans(n_clusters = optimal_K, random_state = 42)$

$Kmeans.fit(Xtrain)$

$y_pred = Kmeans.predict(Xtest)$

$print(y_pred)$



$PRED = (train, test)$ of X for prediction

Code for PCA:

```
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.decomposition import PCA
```

digit = load_digits()

X, y = digit.data, digit.target

scaler = StandardScaler()

X_scaled = scaler.fit_transform(X)

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

lr = LogisticRegression(max_iter=1000)

lr.fit(X_train, y_train)

preds = lr.predict(X_test)

print(accuracy_score(y_test, preds))

pca = PCA(n_components=2)

X_pca = pca.fit_transform(X_scaled)

X_train_pca, X_test_pca, y_train_pca, y_test_pca =

train_test_split(X_pca, y, test_size=0.2, random_state=42)

lr_pca = LogisticRegression(max_iter=1000)

lr_pca.fit(X_train_pca, y_train_pca)

preds_pca = lr_pca.predict(X_test_pca)

print(accuracy_score(y_test, preds_pca))

Accuracy without pca = 0.9722

Accuracy with pca (2 components) = 0.5389