

I N D E X

NAME: SNEHA SANTHOSH BHAT STD: 6th SEC: D ROLL NO: IBM21CS213

S. No.	Date	Title	Page No.	Teacher's Sign / Remarks
1.	21/3/24	Importing and Exporting Data using Pandas library functions	10	b
2.	28/3/24	End to End Machine Learning Project	10	c
3.	04/04/24	Simple Linear Regression and Multiple Linear Regression	10	d
4.	18/4/24	Decision Tree	10	e
5.	25/4/24	Logistic Regression	10	f
6.	09/05/24	Build KNN Classification model for a given dataset	10	g
7.	09/05/24	Build Support Vector machine model for a given dataset	10	h
8.	23/05/24	Build Artificial Neural Network model with back propagation for a given dataset	10	i
9.	23/05/24	Random Forest, Adaboost	10	j
10.	30/05/24	K-Means Algorithm to cluster data	10	k
11.	30/05/24	Dimensionality Reduction using PCA method	h	l

WEEK - 1

To import and export data using Pandas library functions

(i) importing data

```
import pandas as pd
```

```
airbnb_data = pd.read_csv("data/listings-austin.csv")
```

```
airbnb_data.head()
```

(ii) Reading data from URL

```
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
```

```
col_names = ["sepal_length-in-cm", "sepal_width-in-cm",  
            "petal_length-in-cm", "petal_width-in-cm",  
            "class"]
```

```
iris_data = pd.read_csv(url, names=col_names)
```

```
iris_data.head()
```

(iii) Exporting dataframe to csv file

```
iris_data.to_csv("cleaned_iris-data.csv")
```

OUTPUT :

	sepal-length-in-cm	sepal-width-in-cm	petal-length-in-cm	petal-width-in-cm	class
0	5.1	3.5	1.4	0.2	iris-setosa
1					
2					
3					
4					

28/3/24

WEEK - 2

End to End Machine Learning Project

1. Look at the Big Picture
2. Get the Data

- (i) Download the data

```
import os
```

```
import tarfile
```

```
import urllib
```

```
DOWNLOAD_ROOT = "https://raw.githubusercontent.com/"
```

```
housing_path = os.path.join("data", "01")
```

```
housing_url = download.root + "datasets/housing/  
housing.tgz"
```

```
def fetch_housing(housing_url=housing_url,  
                  housing_path=housing_path):  
    os.makedirs(name=housing_path, exist_ok=True)
```

```
    tgz_path = os.path.join(housing_path,  
                           "housing.tgz")
```

```
    fetch_housing_data()
```

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

```
import seaborn as sns
```

```
housing.hist(bins=50, figsize=(20, 15))
```

```
plt.show
```

```
housing['income_cat'] = pd.cut(x=housing['median_income'],  
                                bins=[0, 1.5, 3, 4.5, 6, np.inf],  
                                labels=[1, 2, 3, 4, 5])
```

```
housing['income_cat'].hist()
```

3. Discover and visualize the data to gain insights

~~```
strat_train_set.shape, strat_test_set.shape
```~~~~```
strat_test_set.reset_index().to_feather(fname=  
                                         'data/01/strat_test-  
                                         set.ft')
```~~

```

housing = strat_train_set.copy();
housing.shape
housing.plot(kind = 'scatter', x = 'longitude', y = 'latitude')
plt.show()
corr_matrix = housing.corr()
corr_matrix['median_house_value'].sort_values(
    ascending = False)

```

```

from pandas.plotting import scatter_matrix
attributes = ['median_house_value', 'median_income',
               'total_income', 'housing_median_age']
scatter_matrix(frame = housing[attributes], figsize=(12,8))
plt.show()

```

4. Prepare the Data for Machine Learning Algorithms

```

from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy = 'median')
housing_num = housing.drop(['ocean_proximity'], axis=1)
imputer.fit(housing_num)
imputer.statistics_
housing_num.median().values
X = imputer.transform(housing_num)
X.shape
housing_tr = pd.DataFrame(data = X, index = housing_num.index,
                           columns = housing_num.columns)

```

~~housing_tr.head()~~

~~housing_cat = housing_cat[['ocean_proximity']]~~

~~housing_cat.head()~~

~~from sklearn.pipeline import Pipeline~~

~~from sklearn.preprocessing import StandardScaler~~

~~num_pipeline = Pipeline([~~

```

        ('imputer', SimpleImputer(strategy = 'median')),
        ('attribute_adder', CombinedAttributesAdder()),
        ('std_scalar', StandardScaler())
    ]
)
```

~~7)~~

```
housing_num_tr = num_pipeline.fit_transform(housing_num)
housing_num_tr.shape
```

5. Select and Train a model :

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

```
lin_reg = LinearRegression()
```

```
lin_reg.fit(X = housing_prepared, y = housing_labels)
```

```
housing_predictions = lin_reg.predict(housing_prepared)
```

```
lin_mse = mean_squared_error(housing_label, housing_predictions)
```

```
lin_rmse = np.sqrt(lin_mse)
```

```
lin_rmse
```

```
tree_reg = DecisionTreeRegressor()
```

```
tree_reg.fit(X = housing_prepared, y = housing_labels)
```

```
forest_reg = RandomForestRegressor()
```

```
forest_reg.fit(X = housing_prepared, y = housing_labels)
```

```
forest_reg.fi
```

6. Fine - Tune Your Model :

```
grid_search.best_params_
```

```
grid_search.best_estimator_
```

```
cat_encoder = full_pipeline.named_
```

```
transformers_['cat']
```

```
final_model = grid_search.best_estimator_
```

```
X_test_prepared = full_pipeline.transform(X=X_test)
```

```
final_mse = np.sqrt(final_mse)
```

```
final_rmse
```

```
squared_errors = (y_test - final_predictions)  
** 2
```

7. Launch , Monitor and Maintain your system

WEEK - 3

Simple Linear Regression and Multiple Linear Regression

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from pandas.core.common import random_state
from sklearn.linear_model import LinearRegression

```

```

df_sal = pd.read_csv(r'C:\Users\STUDENT\Downloads\Salary-Data.csv')
df_sal.head()
df_sal.describe()

```

```

plt.title('Salary Distribution Plot')
sns.distplot(df_sal['Salary'])
plt.show()

```

```

x = df_sal.iloc[:, :1]
y = df_sal.iloc[:, 1:]

```

```

X_train, X_test, y_train, y_test = train_test_split
(X, y, test_size=0.2,
random_state=0)

```

```

regressor = LinearRegression()
regressor.fit(X_train, y_train)

```

~~```

y_pred_test = regressor.predict(X_test)
y_pred_train = regressor.predict(X_train)
print(f'Coefficient: {regressor.coef_}')
print(f'Intercept: {regressor.intercept_}')

```~~

```
df_start = pd.read_csv(r'c:\Users\STUDENT\Downloads
df_start.head()
df_start.describe()
```

```
plt.title('Profit Distribution Plot')
sns.distplot(df_start['Profit'])
plt.show()
```

```
regressor = LinearRegression()
regressor.fit(X_train, y_train)
y_pred = regressor.predict(X_test)
```

```
np.set_printoptions(precision = 2)
```

```
result = np.concatenate((y_pred, reshape(len(y_pred),
y_test.reshape(len(y-test), 1)),
```

```
result
```

W  
Y

## WEEK - 4

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

```

import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

import sklearn.datasets as datasets
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, roc_auc_score
from sklearn.tree import plot_tree
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV,
from sklearn.metrics import confusion_matrix
RandomizedSearchCV
from sklearn.metrics import classification_report

```

url = " ... "

df = pd.read\_csv(url, names = ['sepal length(cm)', ...])  
df.head()

X = df.drop("Species", axis=1)

y = df["Species"]

X\_train, X\_test, y\_train, y\_test = train\_test\_split

(X, y, test\_size = 0.3  
random\_state = 1)

dt = DecisionTreeClassifier(max\_depth = 3, min\_samples\_leaf = 10,  
random\_state = 1)  
dt.fit(X, y)

```
from IPython.display import Image
from sklearn.tree import export_graphviz
! pip install pydotplus
import pydotplus
features = X.columns
dot_data = export_graphviz(dt, out_file=None,
graph = pydotplus.graph_from_dot_data(dot_data)
Image(graph.create_png())
```

```
dt = DecisionTreeClassifier(random_state=1)
```

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

```
y_pred_train = dt.predict(X_train)
```

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

```
y_prob = dt.predict_proba(X_test)
```

```
print('Accuracy of Decision Tree-Train:', accuracy_score(y_pred_train, y_train))
```

```
print('Accuracy of Decision Tree-Test:', accuracy_score(y_pred, y_test))
```

```
print(classification_report(y_test, y_pred))
```

```
dt = DecisionTreeClassifier(random_state=1)
```

```
params = {'max_depth': [2, 3, 4, 5],
 'min_samples_split': [2, 3, 4, 5],
 'min_samples_leaf': [1, 2, 3, 4, 5]}
```

~~```
gsearch = GridSearchCV(dt, param_grid=params, cv=3)
```~~~~```
gsearch.fit(X, y)
```~~~~```
gsearch.best_params_
```~~

```

dt = DecisionTreeClassifier(** gsearch.best_params_,
                           random_state=1)
dt.fit(X_train, y_train)

y_pred_train = dt.predict(X_train)
y_prob_train = dt.predict_proba(X_train)[:, 1]
y_pred = dt.predict(X_test)
y_prob = dt.predict_proba(X_test)[:, 1]

```

Output:

petal width (cm) ≤ 0.3

gini = 0.667

samples = 150

value = [50, 50, 50]

True False

gini = 0.0

samples = 50

value = [50, 0, 0]

petal width (cm) ≤ 1.75

gini = 0.5

samples = 100

value = [0, 50, 50]

~~(1st X)~~
~~(2nd X)~~
~~(3rd X)~~

25/4/24

WEEK - 5

Predicting if a person would buy life insurance based on his age using logistic regression

```
import pandas as pd  
from matplotlib import pyplot as plt  
%matplotlib inline
```

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

```
plt.scatter (df.age, df.bought_insurance, marker = 'x',  
            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 = 0.8)
```

```
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()  
model.fit (X-train, y-train)
```

```
y-predicted = model.predict (X-test)  
model.predict_proba (X-test)  
model.score (X-test, y-test)
```

```
model.coef_
```

```
model.intercept_
```

```
import math
```

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

```
def prediction_function(age):
```

$$z = 0.042 * \text{age} - 1.53 \quad \# 0.04150133 \sim 0.042 \text{ and}$$

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

$$-1.52726963 \sim 1.53$$

```
return y
```

```
age = 35
```

Output :-

```
prediction_function(age)
```

0.4850044983805899

```
age = 43
```

Output :-

```
prediction_function(age)
```

0.568565299077705

Q2

Q5

Q8

Q11

Q14

Q17

Build KNN Classification model for a given dataset

```
import pandas as pd
import numpy as np
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics
import matplotlib.pyplot as plt
```

```
df = pd.read_csv('teleCust1000t.csv')
df.head()
```

```
df['custcat'].value_counts()
```

```
x = df.drop(['custcat'], axis=1)
x = x.values
x[0:5]
```

```
y = df['custcat'].values
y[0:5]
```

```
from sklearn.preprocessing import StandardScaler
```

```
x = StandardScaler().fit(x).transform(x.astype(float))
x[0:5]
```

```
from sklearn.model_selection import train_test_split
```

~~xc-train, xc-test, yc-train, yc-test = train_test_split
 $(x, y, test_size=0.2,$
 $random_state=4)$~~

```
print('Train set:', xc_train.shape, yc_train.shape)
print('Test set:', xc_test.shape, yc_test.shape)
```

$K_s = 10$

`mean_acc = []`

`for n in range(1, Ks):`

`neigh = K Neighbors Classifiers (n_neighbors=n). fit`
 (x_train, y_train)

`yhat = neigh.predict(x_test)`

`mean_acc.append(metrics.accuracy_score(y_test, yhat))`

`mean_acc`

`plt.plot(range(1, Ks), mean_acc)`

`plt.xlabel('No of k's')`

`plt.ylabel('Accuracy')`

`plt.show()`

`print("Best accuracy was {} with k = {}".format(max(mean_acc), max(i+1 for i, _ in enumerate(mean_acc))))`

`neigh = KNeighborsClassifier(n_neighbors=9).fit(x_train, y_train)`

`yhat = neigh.predict(x_test)`

`yhat[0:5]`

`print("Train set Accuracy:", metrics.accuracy_score(y_train, neigh.predict(x_train)))`

~~`print("Test set Accuracy:", metrics.accuracy_score(y_test, yhat))`~~

Output:-

Train set Accuracy: 0.5025

Test set Accuracy: 0.34

09/05/2024

WEEK - 7

Build Support vector machine model for a given dataset

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns  
%matplotlib inline
```

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

```
iris.head()
```

```
sns.pairplot(data=iris, hue='species', palette='Set2')
```

```
from sklearn.model_selection import train_test_split
```

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

```
y = iris.iloc[:, 4]
```

```
x_train, x_test, y_train, y_test = train_test_split  
(x, y, test_size=0.30)
```

```
from sklearn.svm import SVC
```

```
model = SVC()
```

```
model.fit(x_train, y_train)
```

```
pred = model.predict(x_test)
```

~~```
from sklearn.metrics import classification_report,
print(confusion_matrix(y-test, pred))
print(classification_report(y-test, pred))
```~~

```
model2 = SVC(C=5, kernel='rbf')
```

```
model2.fit(x-train, y-train)
```

```
pred = model2.predict(x-test)
```

```
from sklearn.metrics import classification_report, confusion_matrix
```

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

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

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.metrics import confusion_matrix, accuracy_score
```

```
from sklearn.model_selection import cross_val_score
```

```
classifier = KNeighborsClassifier(n_neighbors=3)
```

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

```
y-pred = classifier.predict(x-test)
```

```
cm = confusion_matrix(y-test, y-pred)
```

```
cm
```

```
accuracy = accuracy_score(y-test, y-pred) * 100
```

```
print('Accuracy of our model is equal'
```

```
+ str(round(accuracy, 2)) + '%.')
```

```
k-list = list(range(1, 50, 2))
```

```
cv-scores = []
```

```
for k in k-list:
```

```
knn = KNeighborsClassifier(n_neighbors=k)
```

```
scores = cross_val_score(knn, x-train, y-train,
cv=10, scoring='accuracy')
```

```
cv-scores.append(scores.mean())
```

A  
24  
34

$MSE = [1 - x \text{ for } x \text{ in cv-scores}]$

`plt.figure()`

`plt.figure(figsize=(15, 10))`

`plt.title('The optimal number of neighbors',  
 fontsize=20, fontweight='bold')`

`plt.xlabel('Number of Neighbors K', fontsize=15)`

`plt.ylabel('Misclassification Error', fontsize=15)`

`sns.set_style("whitegrid")`

`plt.plot(k_list, MSE)`

`plt.show()`

`best_k = k_list[MSE.index(min(MSE))]`

`print("The optimal number of neighbors is %d." % best_k)`

Output :-

Accuracy of our model is equal 97.78%

The optimal number of neighbors is 1

✓  
16/24  
5/3  
✓  
2/3

## WEEK - 8

8. Build Artificial Neural Network model with back propagation on a given dataset

```
import numpy as np
```

```
X = np.array(([2, 9], [1, 5], [3, 6]), dtype = float)
```

```
y = np.array(([92], [86], [89]), dtype = float)
```

```
X = X / np.amax(X, axis = 0)
```

```
y = y / 100
```

```
epoch = 5000
```

```
lr = 0.1
```

```
inputlayer_neurons = 2
```

```
hiddenlayer_neurons = 3
```

```
output_neurons = 1
```

```
wh = np.random.uniform(size = (inputlayer_neurons,
 hiddenlayer_neurons))
bh = np.random.uniform(size = (1, hiddenlayer_neurons))
wout = np.random.uniform(size = (hiddenlayer_neurons,
 output_neurons))
bou = np.random.uniform(size = (1, output_neurons))
```

```
def sigmoid(x):
```

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

```
def derivatives_sigmoid(x):
```

```
 return x * (1 - x)
```

```
for i in range(epoch):
```

```
 hinpl = np.dot(X, wh)
```

```
 hinp = hinpl + bh
```

```
 hlayer_act = sigmoid(hinp)
```

```
 outinpl = np.dot(hlayer_act, wout)
```

```
 outinp = outinpl + bou
```

```
 output = sigmoid(outinp)
```

$$EO = y - \text{output}$$

outgrad = derivatives - sigmoid (output)

$$d\_output = EO * \text{outgrad}$$

$$EH = d\_output \cdot \text{dot} (w_{\text{out}}, T)$$

hiddengrad = derivatives - sigmoid (layer-act)

$$d\_hiddenlayer = EH * ddengrad$$

$$w_{\text{out}} += \text{layer-act}.T \cdot \text{dot} (d\_output) * lr$$

$$wh += x.T \cdot \text{dot} (d\_hiddenlayer) * lr$$

print ("Input: \n" + str(x))

print ("Actual Output: \n" + str(y))

print ("Predicted Output: \n", output)

Output :-

Input :

[[ 0.66666667 1. ]]

[ 0.3333333 0.55555556 ]

[ 1. 0.66666667 ] ]

Actual Output :

[ [ 0.92 ] ]

[ 0.86 ] ]

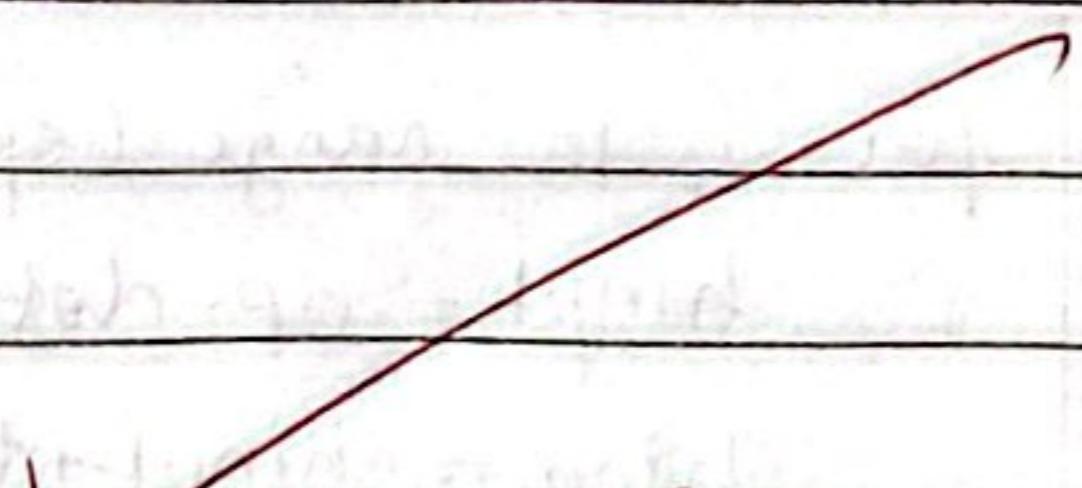
[ 0.89 ] ]

Predicted Output :

[ [ 0.76630321 ] ]

[ 0.75772239 ] ]

[ 0.76682594 ] ]



W/  
3-V

## WEEK - 9

9a, 9b. Random Forest, AdaBoost

```
import os
import sys
from tempfile import NamedTemporaryFile
from urllib.request import urlopen
from urllib.parse import unquote, urlparse
from urllib.error import HTTPError
from zipfile import ZipFile
import tarfile
import shutil
```

CHUNK\_SIZE = 40960

DATA\_SOURCE\_MAPPING = 'https://kaggle.com/...'

KAGGLE\_INPUT\_PATH = '/kaggle/input'

KAGGLE\_WORKING\_PATH = '/kaggle/working'

KAGGLE\_SYMLINK = 'kaggle'

```
for data_source_mapping in DATA_SOURCE_MAPPING.split(','):
 directory, download_url_encoded = data_source_mapping.split(';')
 download_url = unquote(download_url_encoded)
 filename = urlparse(download_url).path
 destination_path = os.path.join(KAGGLE_INPUT_PATH, directory)
 print('Data source import complete')
```

import numpy as np

~~import pandas as pd~~

~~import numpy as np~~

from sklearn.svm import SVC, SVR

from sklearn import preprocessing

import matplotlib.pyplot as plt

```
%matplotlib inline
from sklearn.model_selection import train_test_split
from sklearn.ensemble import (RandomForestClassifier,
 ExtraTreesClassifier,
 AdaBoostClassifier)
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics

import os
print(os.listdir("../input"))

dataset = pd.read_csv("../input/faults.csv")
dataset[0:6]

corr_matrix = input_x.corr().abs()
upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape),
 k=1).astype(np.bool))
to_drop = [column for column in upper.columns if any(upper[column] > 0.95)]
print(to_drop)

X_train, X_test, y_train, y_test = train_test_split(
 input_x_norm, y,
 test_size=0.3)

n_classes = 7
n_estimators = 100
RANDOM_SEED = 13

names = ['DecisionTreeClassifier', 'RandomForestClassifier',
 'ExtraTreesClassifier', 'AdaBoostClassifier']

models = [DecisionTreeClassifier(max_depth=None),
 RandomForestClassifier(n_estimators=n_estimators),
 ExtraTreesClassifier(n_estimators=n_estimators),
 AdaBoostClassifier(DecisionTreeClassifier(max_depth=None),
 n_estimators=n_estimators)]
```

for counter, model in enumerate(models):

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

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

print("Accuracy " + names[counter] + ":")

metrics.accuracy\_score(y\_test, y\_pred))

score = model.evaluate(x\_test, y\_test, batch\_size=128)

print(score)

print(model.metrics\_names)

Y\_predicted = model.predict(x\_test, batch\_size=32, verbose=0)

percent = 0

for i in range(0, len(Y\_predicted)):

class\_id\_predicted = np.argmax(Y\_predicted[i])

class\_id\_real = np.argmax(y\_test[i])

if class\_id\_predicted == class\_id\_real:

percent += 1

print('val accuracy:', percent / len(Y\_predicted))

Output:-

val accuracy: 0.7409948542024014

W  
Jish

30/5/24

## WEEK - 10

Build k-Means algorithm to cluster a set of data stored in a .csv file

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

```
data = pd.read_csv('Iris.csv')
```

```
X = data[['SepalLengthCm', 'SepalWidthCm',
 'PetalLengthCm', 'PetalWidthCm']]
```

```
model = KMeans(n_clusters=3)
```

```
model.fit(X)
```

```
plt.figure(figsize=(14, 7))
```

```
plt.subplot(1, 2, 2)
```

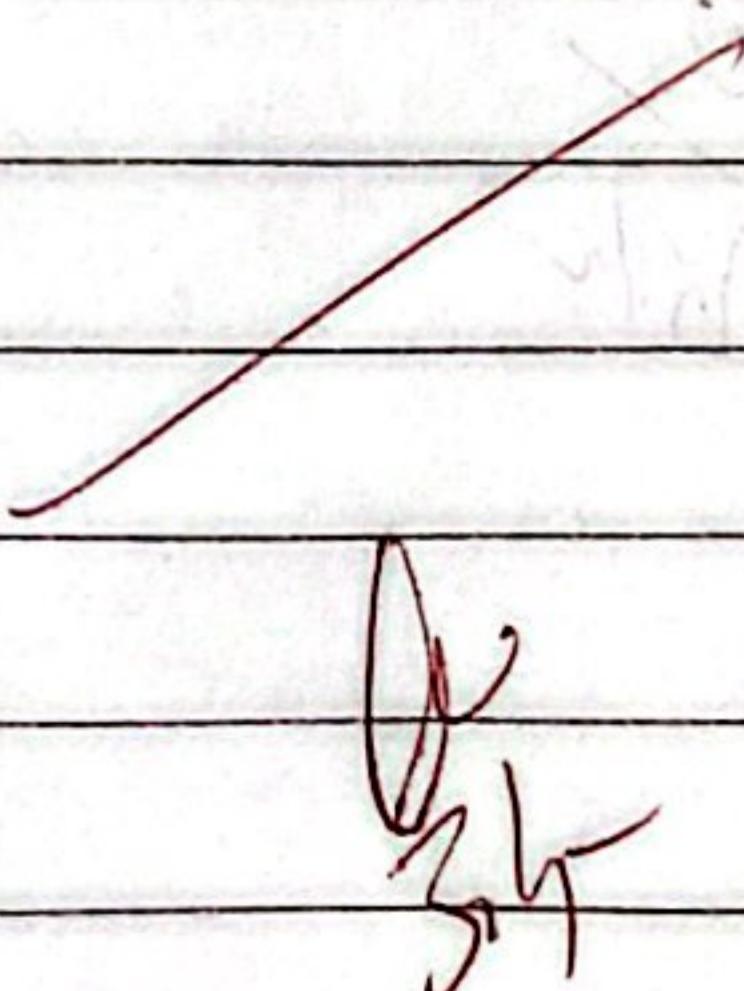
```
plt.scatter(X['PetalLengthCm'], X['PetalWidthCm'],
 c=model.labels_, cmap
 ='viridis', s=40)
```

```
plt.title('K-Means Clustering')
```

```
plt.xlabel('Petal Length')
```

```
plt.ylabel('Petal Width')
```

```
plt.show()
```



## WEEK - 11

Implement Dimensionality reduction using Principle Component Analysis (PCA) method

```
import seaborn as sns
```

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

```
import pandas as pd
```

```
import numpy as np
```

```
%matplotlib inline
```

```
from sklearn.datasets import load_breast_cancer
```

```
cancer = load_breast_cancer()
```

```
cancer.keys()
```

```
print(cancer['DESCR'])
```

```
df = pd.DataFrame(cancer['data'], columns=cancer['feature_names'])
```

```
df.head()
```

```
cancer['target_names']
```

```
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
```

```
scaler.fit(df)
```

```
scaled_data = scaler.transform(df)
```

```
from sklearn.decomposition import PCA
```

```
pca = PCA(n_components=2)
```

pca.fit(scaled\_data)

x\_pca = pca.transform(scaled\_data)

scaled\_data.shape

x\_pca.shape

plt.figure(figsize=(8, 6))

plt.scatter(x\_pca[:, 0], x\_pca[:, 1])

plt.xlabel('First Principal Component')

plt.ylabel('Second Principal Component')

pca.components\_

df\_comp = pd.DataFrame(pca.components\_,

columns=cancer

df\_comp

['feature-names'])

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

sns.heatmap(df\_comp, cmap='plasma')

✓  
W  
3011✓