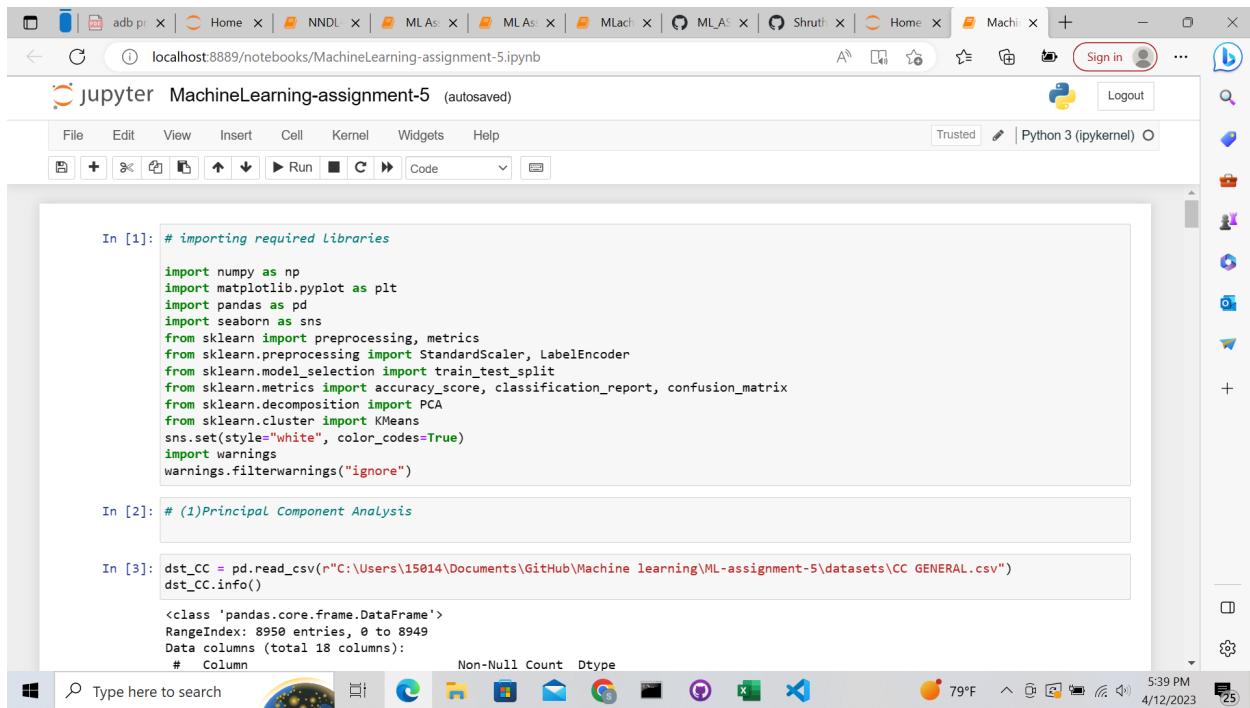


# Spring 2023 5710 Machine Learning: Assignment 5

## In class programming:

1. Principal Component Analysis
  - a. Apply PCA on CC dataset.
  - b. Apply k-means algorithm on the PCA result and report your observation if the silhouette score has improved or not?
  - c. Perform Scaling+PCA+K-Means and report performance.



The screenshot shows a Jupyter Notebook interface running on a Windows desktop. The browser tab is 'localhost:8889/notebooks/MachineLearning-assignment-5.ipynb'. The notebook title is 'MachineLearning-assignment-5' (autosaved). The toolbar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Trusted, Python 3 (ipykernel), and a sign-in button. The code cells are as follows:

```
In [1]: # importing required Libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
from sklearn import preprocessing, metrics
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
sns.set(style="white", color_codes=True)
import warnings
warnings.filterwarnings("ignore")

In [2]: # (1)Principal Component Analysis

In [3]: dst_CC = pd.read_csv(r"C:\Users\15014\Documents\GitHub\Machine learning\ML-assignment-5\datasets\CC GENERAL.csv")
dst_CC.info()
```

The status bar at the bottom shows the search bar, taskbar icons, battery level (79%), and the date/time (5:39 PM, 4/12/2023).

In [3]: `dst_CC = pd.read_csv(r"C:\Users\15014\Documents\GitHub\Machine learning\ML-assignment-5\datasets\CC GENERAL.csv")  
dst_CC.info()`

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 8950 entries, 0 to 8949  
Data columns (total 18 columns):  
 #   Column           Non-Null Count  Dtype     
 ---  --  
 0   CUST_ID          8950 non-null   object    
 1   BALANCE          8950 non-null   float64  
 2   BALANCE_FREQUENCY 8950 non-null   float64  
 3   PURCHASES         8950 non-null   float64  
 4   ONEOFF_PURCHASES 8950 non-null   float64  
 5   INSTALLMENTS_PURCHASES 8950 non-null   float64  
 6   CASH_ADVANCE      8950 non-null   float64  
 7   PURCHASES_FREQUENCY 8950 non-null   float64  
 8   ONEOFF_PURCHASES_FREQUENCY 8950 non-null   float64  
 9   PURCHASES_INSTALLMENTS_FREQUENCY 8950 non-null   float64  
 10  CASH_ADVANCE_FREQUENCY 8950 non-null   float64  
 11  CASH_ADVANCE_TRX 8950 non-null   int64     
 12  PURCHASES_TRX    8950 non-null   int64     
 13  CREDIT_LIMIT     8949 non-null   float64  
 14  PAYMENTS          8950 non-null   float64  
 15  MINIMUM_PAYMENTS 8637 non-null   float64  
 16  PRC_FULL_PAYMENT 8950 non-null   float64  
 17  TENURE            8950 non-null   int64     
dtypes: float64(14), int64(3), object(1)  
memory usage: 1.2+ MB
```

In [4]: `dst_CC.head()`

Out[4]:

	CUST_ID	BALANCE	BALANCE_FREQUENCY	PURCHASES	ONEOFF_PURCHASES	INSTALLMENTS_PURCHASES	CASH_ADVANCE	PURCHASES_FREQUENCY
0	C10001	40.900749	0.818182	95.40	0.00	95.4	0.000000	0.1661
1	C10002	3202.467416	0.909091	0.00	0.00	0.0	6442.945483	0.0001
2	C10003	2495.148862	1.000000	773.17	773.17	0.0	0.000000	1.0001
3	C10004	1666.670542	0.636364	1499.00	1499.00	0.0	205.788017	0.0831
4	C10005	817.714335	1.000000	16.00	16.00	0.0	0.000000	0.0831

In [5]: `dst_CC.isnull().any()`

Out[5]:

CUST_ID	False
BALANCE	False
BALANCE_FREQUENCY	False
PURCHASES	False
ONEOFF_PURCHASES	False
INSTALLMENTS_PURCHASES	False
CASH_ADVANCE	False
PURCHASES_FREQUENCY	False
ONEOFF_PURCHASES_FREQUENCY	False
PURCHASES_INSTALLMENTS_FREQUENCY	False
CASH_ADVANCE_FREQUENCY	False
CASH_ADVANCE_TRX	False

localhost:8889/notebooks/MachineLearning-assignment-5.ipynb

jupyter MachineLearning-assignment-5 (autosaved)

In [5]: `dst_CC.isnull().any()`

Out[5]:

CUST_ID	False
BALANCE	False
BALANCE_FREQUENCY	False
PURCHASES	False
ONEOFF_PURCHASES	False
INSTALLMENTS_PURCHASES	False
CASH_ADVANCE	False
PURCHASES_FREQUENCY	False
ONEOFF_PURCHASES_FREQUENCY	False
PURCHASES_INSTALLMENTS_FREQUENCY	False
CASH_ADVANCE_FREQUENCY	False
CASH_ADVANCE_TRX	False
PURCHASES_TRX	False
CREDIT_LIMIT	True
PAYMENTS	False
MINIMUM_PAYMENTS	True
PRC_FULL_PAYMENT	False
TENURE	False

dtype: bool

In [6]: `dst_CC.fillna(dst_CC.mean(), inplace=True)`  
`dst_CC.isnull().any()`

Out[6]:

CUST_ID	False
BALANCE	False
BALANCE_FREQUENCY	False
PURCHASES	False

localhost:8889/notebooks/MachineLearning-assignment-5.ipynb

jupyter MachineLearning-assignment-5 (autosaved)

In [6]: `dst_CC.fillna(dst_CC.mean(), inplace=True)`  
`dst_CC.isnull().any()`

Out[6]:

CUST_ID	False
BALANCE	False
BALANCE_FREQUENCY	False
PURCHASES	False
ONEOFF_PURCHASES	False
INSTALLMENTS_PURCHASES	False
CASH_ADVANCE	False
PURCHASES_FREQUENCY	False
ONEOFF_PURCHASES_FREQUENCY	False
PURCHASES_INSTALLMENTS_FREQUENCY	False
CASH_ADVANCE_FREQUENCY	False
CASH_ADVANCE_TRX	False
PURCHASES_TRX	False
CREDIT_LIMIT	False
PAYMENTS	False
MINIMUM_PAYMENTS	False
PRC_FULL_PAYMENT	False
TENURE	False

In [7]: `a = dst_CC.iloc[:,1:-1]`  
`b = dst_CC.iloc[:, -1]`  
`print(a.shape, b.shape)`

(8950, 16) (8950,)

In [8]:  `#(1a) Apply PCA on CC Dataset`

```
In [9]: pca = PCA(3)
a_pca = pca.fit_transform(a)
prplDf = pd.DataFrame(data = a_pca, columns = ['principal cmpnt 1', 'principal cmpnt 2', 'principal cmpnt 3'])
finalDf = pd.concat([prplDf, dst_CC.iloc[:, -1]], axis = 1)
finalDf.head()
```

Out[9]:

	principal cmpnt 1	principal cmpnt 2	principal cmpnt 3	TENURE
0	-4326.383979	921.566882	183.708383	12
1	4118.916665	-2432.846346	2369.969289	12
2	1497.907641	-1997.578694	-2125.631328	12
3	1394.548536	-1488.743453	-2431.799649	12
4	-3743.351896	757.342657	512.476492	12

In [10]:  `#(1b) Apply K Means on PCA Result`

```
A = finalDf.iloc[:, 0:-1]
b = finalDf.iloc[:, -1]
```

In [11]:  `nclusters = 3 # this is the k in kmeans`

```
km = KMeans(n_clusters=nclusters)
km.fit(A)
```

In [10]:  `#(1b) Apply K Means on PCA Result`

```
A = finalDf.iloc[:, 0:-1]
b = finalDf.iloc[:, -1]
```

In [11]:  `nclusters = 3 # this is the k in kmeans`

```
km = KMeans(n_clusters=nclusters)
km.fit(A)

# predict the cluster for each data point
b_cluster_kmeans = km.predict(A)

# Summary of the predictions made by the classifier
print(classification_report(b, b_cluster_kmeans, zero_division=1))
print(confusion_matrix(b, b_cluster_kmeans))

train_acc = accuracy_score(b, b_cluster_kmeans)
print("\nAccuracy of our training dataset with PCA:", train_acc)

#Calculate silhouette Score
scr = metrics.silhouette_score(A, b_cluster_kmeans)
print("Silhouette Score: ", scr)
```

	precision	recall	f1-score	support
0	0.00	1.00	0.00	0.0
1	0.00	1.00	0.00	0.0

jupyter MachineLearning-assignment-5 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) ○

Run Cell Code

	precision	recall	f1 score	support
0	0.00	1.00	0.00	0.0
1	0.00	1.00	0.00	0.0
2	0.00	1.00	0.00	0.0
6	1.00	0.00	0.00	204.0
7	1.00	0.00	0.00	190.0
8	1.00	0.00	0.00	196.0
9	1.00	0.00	0.00	175.0
10	1.00	0.00	0.00	236.0
11	1.00	0.00	0.00	365.0
12	1.00	0.00	0.00	7584.0

accuracy 0.00 8950.0  
macro avg 0.70 0.30 0.00 8950.0  
weighted avg 1.00 0.00 0.00 8950.0

[	0	0	0	0	0	0	0	0	0]
[	0	0	0	0	0	0	0	0	0]
[	0	0	0	0	0	0	0	0	0]
[	175	28	1	0	0	0	0	0	0]
[	173	15	2	0	0	0	0	0	0]
[	169	27	0	0	0	0	0	0	0]
[	149	26	0	0	0	0	0	0	0]
[	188	47	1	0	0	0	0	0	0]
[	284	78	3	0	0	0	0	0	0]
[	5389	2069	126	0	0	0	0	0	0]

Accuracy of our training dataset with PCA: 0.0  
Sihouette Score: 0.5109307274319466

localhost:8889/notebooks/MachineLearning-assignment-5.ipynb

Sign in

## jupyter MachineLearning-assignment-5 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

In [12]:

```
#1c) Scaling +PCA + KMeans
a = dst_CC.iloc[:,1:1]
b = dst_CC.iloc[:,1]
print(a.shape,b.shape)
```

(8950, 16) (8950,)

In [13]:

```
#Scaling
scaler = StandardScaler()
scaler.fit(a)
X_scaled_array = scaler.transform(a)
#PCA
pca = PCA(3)
x_pca = pca.fit_transform(X_scaled_array)
prplDF = pd.DataFrame(data = x_pca, columns = ['principal cmpt 1', 'principal cmpt 2','principal cmpt 3'])
finalDF = pd.concat([prplDF, dst_CC.iloc[:,1:1]], axis = 1)
finalDF.head()
```

Out[13]:

	principal cmpt 1	principal cmpt 2	principal cmpt 3	TENURE
0	-1.718893	-1.072940	0.535721	12
1	-1.169307	2.509322	0.627928	12
2	0.938414	-0.382600	0.161225	12
3	-0.907503	0.045859	1.521694	12
4	-1.637830	-0.684976	0.425731	12

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) ○

```
In [14]: A = finalDF.iloc[:,0:-1]
b = finalDF["TENURE"]
print(A.shape,b.shape)

(8950, 3) (8950,)

In [15]: A_train, A_test, b_train, b_test = train_test_split(A,b, test_size=0.34,random_state=0)
nclusters = 3
# this is the k in kmeans
km = KMeans(n_clusters=nclusters)
km.fit(A_train,b_train)

# predict the cluster for each training data point
b_clust_train = km.predict(A_train)

# Summary of the predictions made by the classifier
print(classification_report(b_train, b_clust_train, zero_division=1))
print(confusion_matrix(b_train, b_clust_train))

train_acc = accuracy_score(b_train, b_clust_train)
print("Accuracy of our training dataset with PCA:", train_acc)

#Calculate silhouette Score
scr = metrics.silhouette_score(A_train, b_clust_train)
print("Silhouette Score: ",scr)
```

precision	recall	f1-score	support
-----------	--------	----------	---------

jupyter MachineLearning-assignment-5 (autosaved)

```
In [16]: # predict the cluster for each testing data point
b_clust_test = km.predict(A_test)

# Summary of the predictions made by the classifier
print(classification_report(b_test, b_clust_test, zero_division=1))
print(confusion_matrix(b_test, b_clust_test))

train_acc = accuracy_score(b_test, b_clust_test)
print("\nAccuracy of our training dataset with PCA:", train_acc)

#Calculate silhouette Score
scr = metrics.silhouette_score(A_test, b_clust_test)
print("Silhouette Score: ",scr)
```

	precision	recall	f1-score	support
0	0.00	1.00	0.00	0.0
1	0.00	1.00	0.00	0.0
2	0.00	1.00	0.00	0.0
6	1.00	0.00	0.00	65.0
7	1.00	0.00	0.00	55.0
8	1.00	0.00	0.00	68.0
9	1.00	0.00	0.00	57.0
10	1.00	0.00	0.00	85.0
11	1.00	0.00	0.00	103.0
12	1.00	0.00	0.00	2610.0
accuracy			0.00	3043.0
macro avg	0.70	0.30	0.00	3043.0

jupyter MachineLearning-assignment-5 (autosaved)

```
In [16]: # predict the cluster for each testing data point
b_clust_test = km.predict(A_test)

# Summary of the predictions made by the classifier
print(classification_report(b_test, b_clust_test, zero_division=1))
print(confusion_matrix(b_test, b_clust_test))

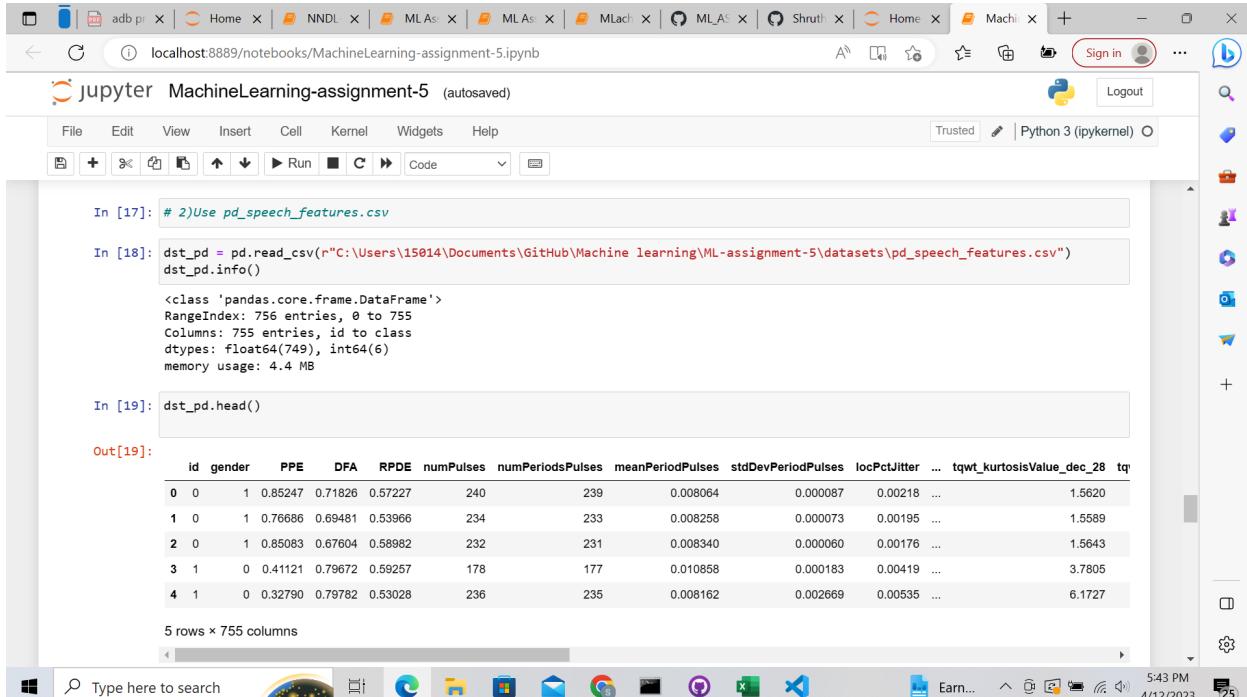
train_acc = accuracy_score(b_test, b_clust_test)
print("\nAccuracy of our training dataset with PCA: 0.0
Sihouette Score: 0.38364322686865926
```

	precision	recall	f1-score	support
0	0.00	1.00	0.00	0.0
1	0.00	1.00	0.00	0.0
2	0.00	1.00	0.00	0.0
6	1.00	0.00	0.00	65.0
7	1.00	0.00	0.00	55.0
8	1.00	0.00	0.00	68.0
9	1.00	0.00	0.00	57.0
10	1.00	0.00	0.00	85.0
11	1.00	0.00	0.00	103.0
12	1.00	0.00	0.00	2610.0
accuracy			0.00	3043.0
macro avg	0.70	0.30	0.00	3043.0
weighted avg	1.00	0.00	0.00	3043.0

```
[[ 0  0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0]
 [ 21 41  3  0  0  0  0  0  0  0]
 [ 12 43  0  0  0  0  0  0  0  0]
 [ 10 57  1  0  0  0  0  0  0  0]
 [ 22 35  0  0  0  0  0  0  0  0]
 [ 17 63  5  0  0  0  0  0  0  0]
 [ 30 69  4  0  0  0  0  0  0  0]
 [ 450 1765 395  0  0  0  0  0  0  0]]
```

## 2. Use pd\_speech\_features.csv

- Perform Scaling
- Apply PCA (k=3)
- Use SVM to report performance



In [17]: # 2) Use pd\_speech\_features.csv

In [18]: dst\_pd = pd.read\_csv(r"C:\Users\15014\Documents\GitHub\Machine learning\ML-assignment-5\datasets\pd\_speech\_features.csv")  
dst\_pd.info()

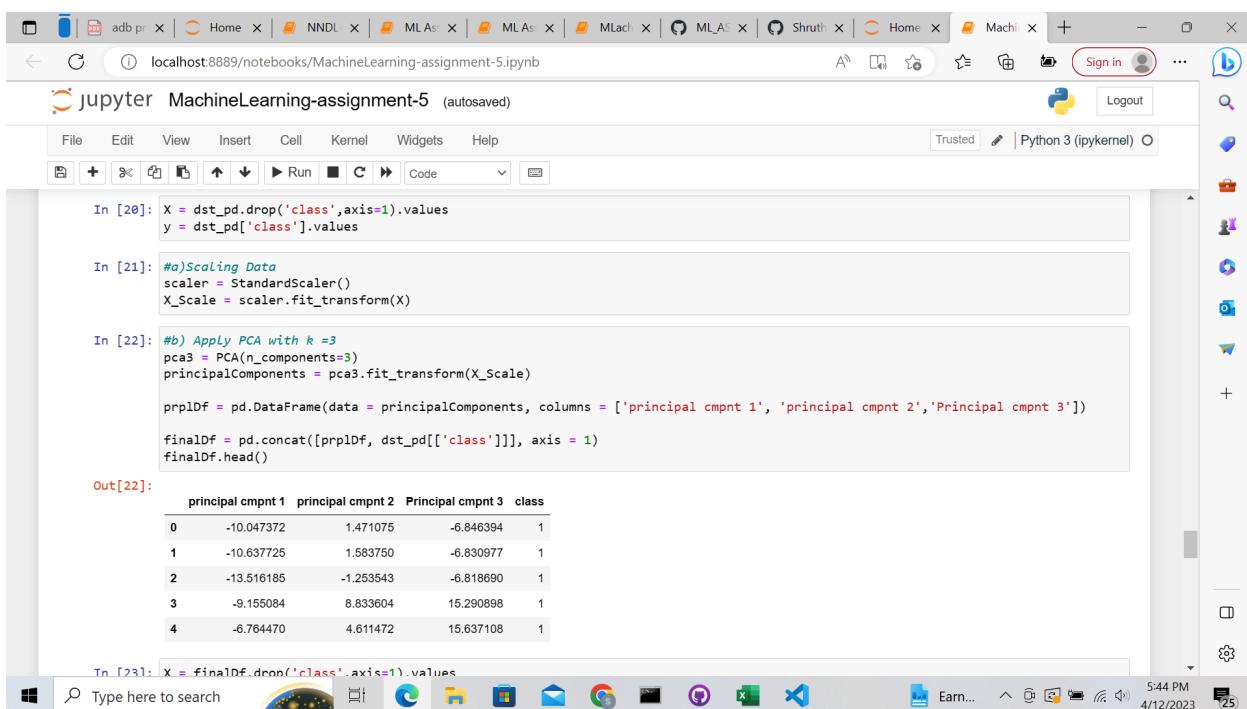
```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 756 entries, 0 to 755  
Columns: 755 entries, id to class  
dtypes: float64(749), int64(6)  
memory usage: 4.4 MB
```

In [19]: dst\_pd.head()

Out[19]:

	id	gender	PPE	DFA	RPDE	numPulses	numPeriodsPulses	meanPeriodPulses	stdDevPeriodPulses	locPctJitter	...	twqt_kurtosisValue_dec_28	tq
0	0	1	0.85247	0.71826	0.57227	240	239	0.008064	0.000087	0.00218	...	1.5620	
1	0	1	0.76686	0.69481	0.53966	234	233	0.008258	0.000073	0.00195	...	1.5589	
2	0	1	0.85083	0.67604	0.58982	232	231	0.008340	0.000060	0.00176	...	1.5643	
3	1	0	0.41121	0.79672	0.59257	178	177	0.010858	0.000183	0.00419	...	3.7805	
4	1	0	0.32790	0.79782	0.53028	236	235	0.008162	0.002669	0.00535	...	6.1727	

5 rows x 755 columns



In [20]: X = dst\_pd.drop('class', axis=1).values  
y = dst\_pd['class'].values

In [21]: #a) Scaling Data  
scaler = StandardScaler()  
X\_Scale = scaler.fit\_transform(X)

In [22]: #b) Apply PCA with k = 3  
pca3 = PCA(n\_components=3)  
principalComponents = pca3.fit\_transform(X\_Scale)

prplDf = pd.DataFrame(data = principalComponents, columns = ['principal cmpnt 1', 'principal cmpnt 2', 'Principal cmpnt 3'])

finalDf = pd.concat([prplDf, dst\_pd[['class']]], axis = 1)

finalDf.head()

Out[22]:

	principal cmpnt 1	principal cmpnt 2	Principal cmpnt 3	class
0	-10.047372	1.471075	-6.846394	1
1	-10.637725	1.583750	-6.830977	1
2	-13.516185	-1.253543	-6.818690	1
3	-9.155084	8.833604	15.290898	1
4	-6.764470	4.611472	15.637108	1

In [23]: X = finalDf.drop('class', axis=1).values

jupyter MachineLearning-assignment-5 (autosaved)

In [23]:

```
X = finalDf.drop('class',axis=1).values
y = finalDf['class'].values
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.34,random_state=0)
```

In [24]: #c) Support Vector Machine(SVM)

```
from sklearn.svm import SVC

svmClassifier = SVC()
svmClassifier.fit(X_train, y_train)

y_pred = svmClassifier.predict(X_test)

# Summary of the predictions made by the classifier
print(classification_report(y_test, y_pred, zero_division=1))
print(confusion_matrix(y_test, y_pred))
# Accuracy score
glass_acc_svc = accuracy_score(y_pred,y_test)
print('accuracy is',glass_acc_svc )

#Calculate silhouette Score
scr= metrics.silhouette_score(X_test, y_pred)
print("Silhouette Score: ",scr)
```

	precision	recall	f1-score	support
0	0.67	0.42	0.51	62
1	0.84	0.93	0.88	196

5:44 PM  
4/12/2023

jupyter MachineLearning-assignment-5 (autosaved)

In [23]:

```
svmClassifier = SVC()
svmClassifier.fit(X_train, y_train)

y_pred = svmClassifier.predict(X_test)

# Summary of the predictions made by the classifier
print(classification_report(y_test, y_pred, zero_division=1))
print(confusion_matrix(y_test, y_pred))
# Accuracy score
glass_acc_svc = accuracy_score(y_pred,y_test)
print('accuracy is',glass_acc_svc )

#Calculate silhouette Score
scr= metrics.silhouette_score(X_test, y_pred)
print("Silhouette Score: ",scr)
```

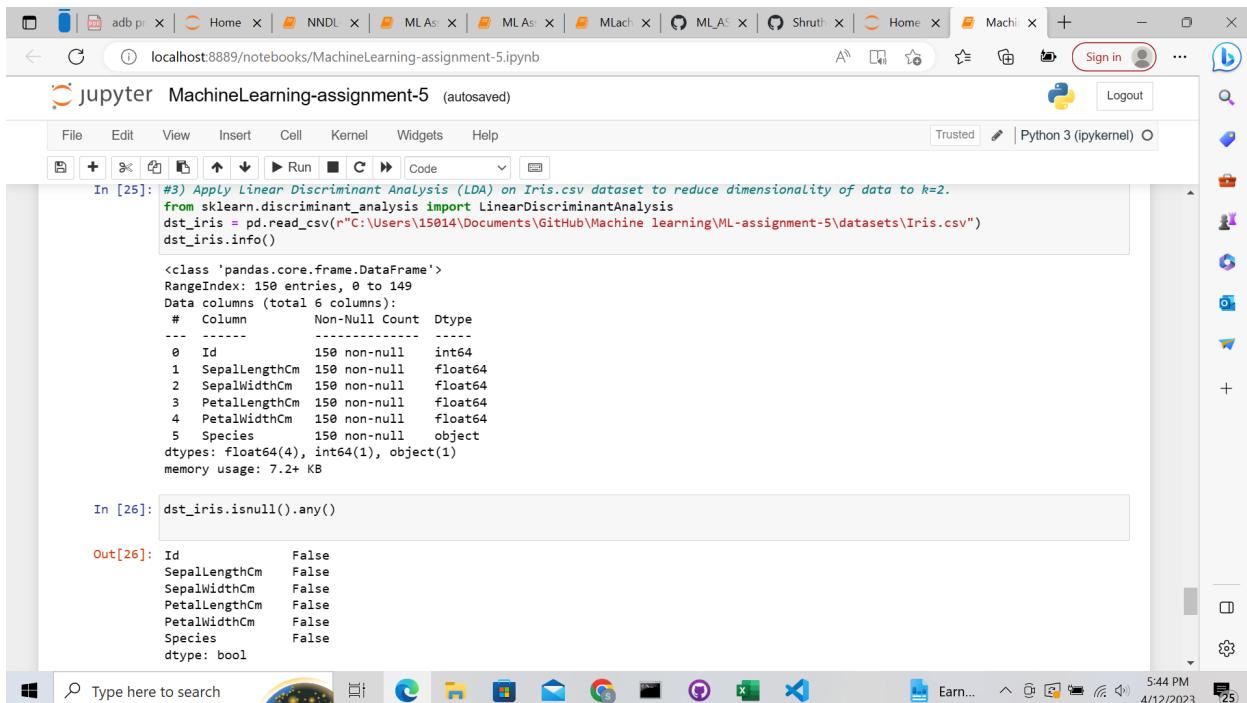
	precision	recall	f1-score	support
0	0.67	0.42	0.51	62
1	0.84	0.93	0.88	196

	accuracy	macro avg	weighted avg	
0	0.81	0.75	0.80	258
1	0.25	0.68	0.81	258

```
[[ 26  36]
 [ 13 183]]
accuracy is 0.810077519379845
Silhouette Score:  0.2504464167386183
```

5:44 PM  
4/12/2023

### 3. Apply Linear Discriminant Analysis (LDA) on Iris.csv dataset to reduce dimensionality of data to k=2.



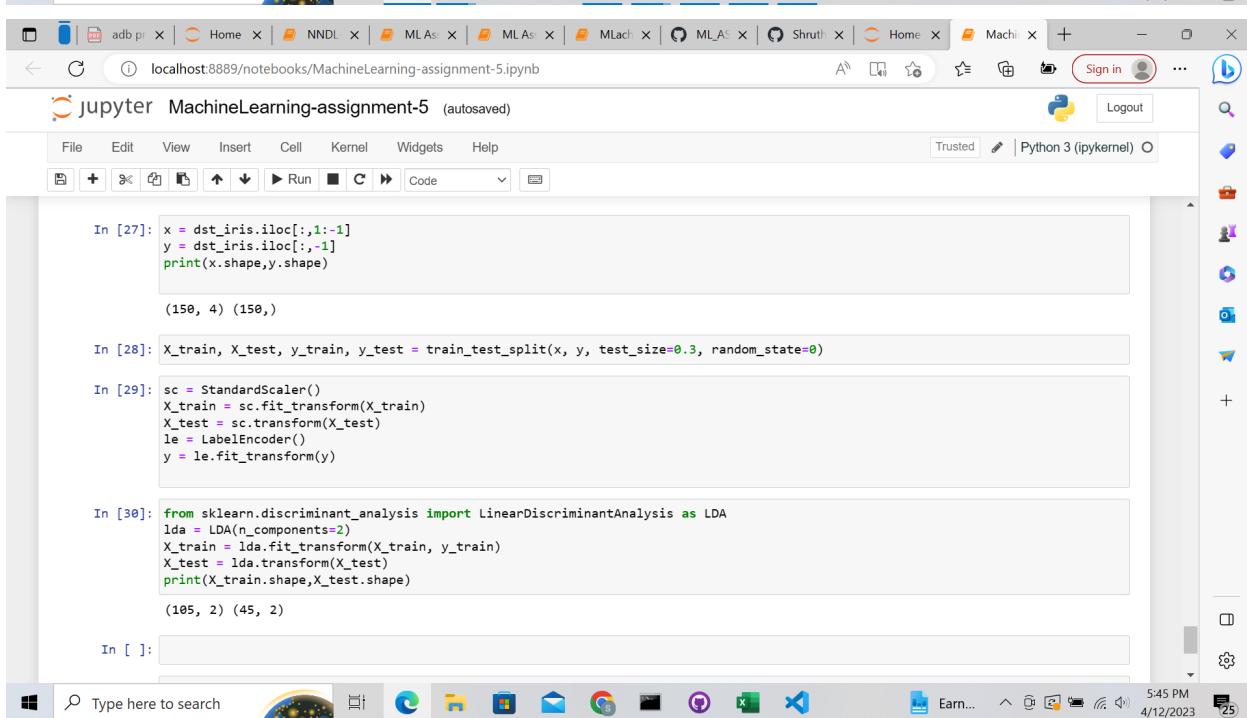
```
In [25]: #3) Apply Linear Discriminant Analysis (LDA) on Iris.csv dataset to reduce dimensionality of data to k=2.
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
dst_iris = pd.read_csv(r"C:\Users\15014\Documents\GitHub\Machine learning\ML-assignment-5\datasets\Iris.csv")
dst_iris.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Id          150 non-null    int64  
 1   SepalLengthCm 150 non-null    float64 
 2   SepalWidthCm  150 non-null    float64 
 3   PetalLengthCm 150 non-null    float64 
 4   PetalWidthCm  150 non-null    float64 
 5   Species      150 non-null    object  
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB

In [26]: dst_iris.isnull().any()

Out[26]:
```

Column	Non-Null Count	Dtype
Id	150	int64
SepalLengthCm	150	float64
SepalWidthCm	150	float64
PetalLengthCm	150	float64
PetalWidthCm	150	float64
Species	150	object



```
In [27]: x = dst_iris.iloc[:,1:-1]
y = dst_iris.iloc[:,1]
print(x.shape,y.shape)

(150, 4) (150,)

In [28]: X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=0)

In [29]: sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
le = LabelEncoder()
y = le.fit_transform(y)

In [30]: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
lda = LDA(n_components=2)
X_train = lda.fit_transform(X_train, y_train)
X_test = lda.transform(X_test)
print(X_train.shape,X_test.shape)

(105, 2) (45, 2)

In [ ]:
```

4. Briefly identify the difference between PCA and LDA

Both LDA and PCA are based on linear transformations and aim to maximize the variance of the lower dimension. PCA is an unsupervised learning algorithm while LDA is a supervised learning algorithm. That is, PCA finds indications of maximum variance regardless of class labels, while LDA finds indications of maximum class separation

GitHub Repo link : <https://github.com/ShruthiVallapReddy/ML-assignment-5.git>