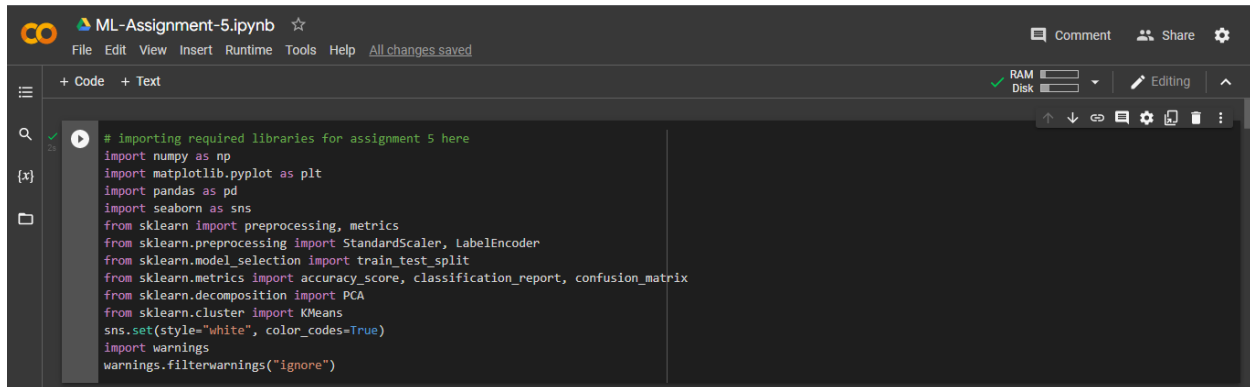


ML-Assignment-5

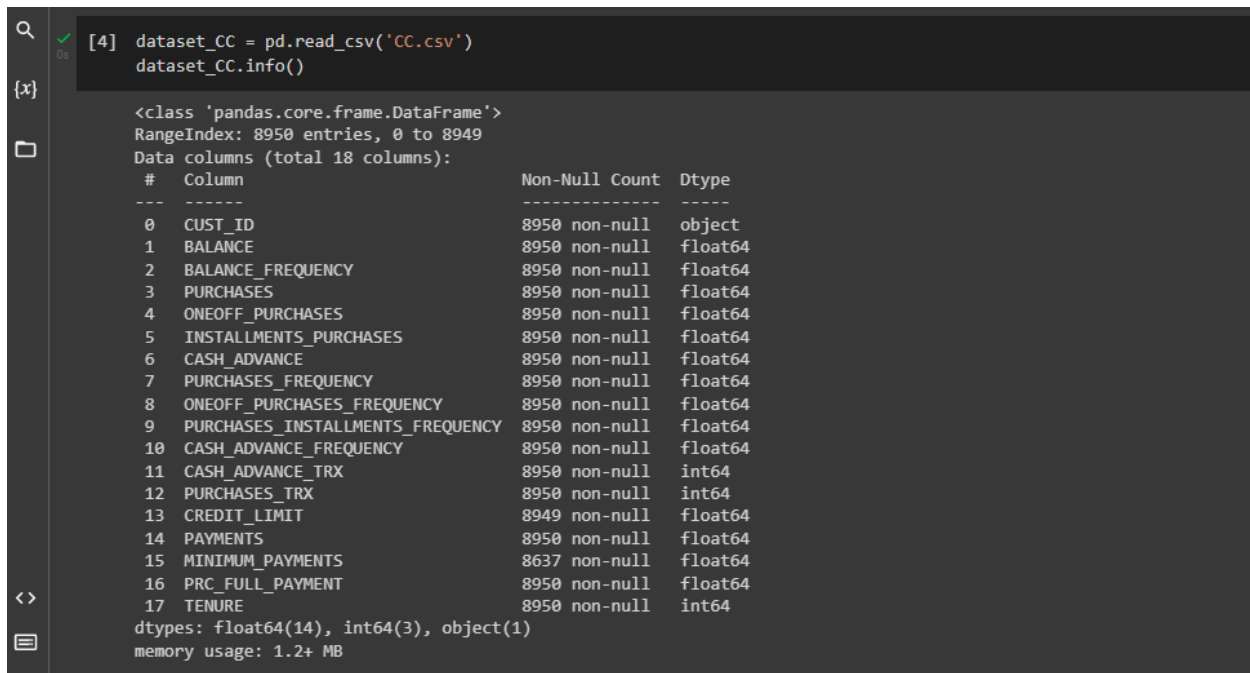
#Question 1:

Principal Component Analysis

- Applied PCA on CC dataset.
- Applied k-means algorithm on the PCA result.
- Perform Scaling+PCA+K-Means.



```
# importing required libraries for assignment 5 here
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")
```



```
[4] dataset_CC = pd.read_csv('CC.csv')
dataset_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
```

```
+ Code + Text
[5] dataset_CC.head()
```

	CUST_ID	BALANCE	BALANCE_FREQUENCY	PURCHASES	ONEOFF_PURCHASES	INSTALLMENTS_PURCHASES	CASH_ADVANCE	PURCHASES_FREQUENCY	ONEOFF_PURCHASES_FREQUENCY	PUR
0	C10001	40.900749	0.818182	95.40	0.00	95.4	0.000000	0.166667	0.000000	
1	C10002	3202.467416	0.909091	0.00	0.00	0.0	6442.945483	0.000000	0.000000	
2	C10003	2495.148862	1.000000	773.17	773.17	0.0	0.000000	1.000000	1.000000	
3	C10004	1666.670542	0.636364	1499.00	1499.00	0.0	205.788017	0.083333	0.083333	
4	C10005	817.714335	1.000000	16.00	16.00	0.0	0.000000	0.083333	0.083333	

```
[6] dataset_CC.isnull().any()
```

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

```
dataset_CC.fillna(dataset_CC.mean(), inplace=True)
dataset_CC.isnull().any()
```

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

```
[8] x = dataset_CC.iloc[:,1:-1]
y = dataset_CC.iloc[:, -1]
print(x.shape,y.shape)
```

(8950, 16) (8950,)

Q

{x}

✓ [9] #1.a Apply PCA on CC Dataset

pca = PCA(3)
x_pca = pca.fit_transform(x)
principalDf = pd.DataFrame(data = x_pca, columns = ['principal component 1', 'principal component 2', 'principal component 3'])
finalDf = pd.concat([principalDf, dataset_CC.iloc[:, -1]], axis = 1)
finalDf.head()

principal component 1 principal component 2 principal component 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

✓ #1.b Apply K Means on PCA Result

x = finalDf.iloc[:, 0:-1]
y = finalDf.iloc[:, -1]

RAM
Disk

+

Code

+

Text

✓ [34]

nclusters = 3 # this is the k in kmeans
km = KMeans(n_clusters=nclusters)
km.fit(X)

predict the cluster for each data point
y_cluster_kmeans = km.predict(X)

Summary of the predictions made by the classifier
print(classification_report(y, y_cluster_kmeans, zero_division=1))
print(confusion_matrix(y, y_cluster_kmeans))

train_accuracy = accuracy_score(y, y_cluster_kmeans)
print("\nAccuracy for our Training dataset with PCA:", train_accuracy)

#Calculate sihouette Score
score = metrics.silhouette_score(X, y_cluster_kmeans)
print("Sihouette Score: ", score)

=====
Sihouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.
=====

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 0.00 225.0

```
+ Code + Text
RAM
Disk

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

Accuracy for our Training dataset with PCA: 0.0
Silhouette Score: 0.5109307274319468
'\nSilhouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring cluster s.\n'

[12] #1.c Scaling +PCA + KMeans
x = dataset_CC.iloc[:,1:-1]
y = dataset_CC.iloc[:, -1]
print(x.shape,y.shape)

(8950, 16) (8950,)
```

```
#Scaling
scaler = StandardScaler()
scaler.fit(x)
X_scaled_array = scaler.transform(x)
#PCA
pca = PCA(3)
x_pca = pca.fit_transform(X_scaled_array)
principalDf = pd.DataFrame(data = x_pca, columns = ['principal component 1', 'principal component 2', 'principal component 3'])
finalDf = pd.concat([principalDf, dataset_CC.iloc[:, -1]], axis = 1)
finalDf.head()

principal component 1 principal component 2 principal component 3 TENURE
0 -1.718891 -1.072943 0.535835 12
1 -1.169311 2.509327 0.627689 12
2 0.938417 -0.382605 0.161430 12
3 -0.907502 0.045858 1.521756 12
4 -1.637827 -0.684979 0.425877 12

[14] X = finalDf.iloc[:,0:-1]
y = finalDf["TENURE"]
print(X.shape,y.shape)

(8950, 3) (8950,)
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.34, random_state=0)
n_clusters = 3
# this is the k in kmeans
km = KMeans(n_clusters=n_clusters)
km.fit(X_train, y_train)

# predict the cluster for each training data point
y_clus_train = km.predict(X_train)

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

train_accuracy = accuracy_score(y_train, y_clus_train)
print("Accuracy for our Training dataset with PCA:", train_accuracy)

#Calculate sihouette Score
score = metrics.silhouette_score(X_train, y_clus_train)
print("Sihouette Score: ", score)

'''
Sihouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.
'''
```

	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	139.0
7	1.00	0.00	0.00	135.0
8	1.00	0.00	0.00	128.0
9	1.00	0.00	0.00	118.0
10	1.00	0.00	0.00	151.0

	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	139.0
7	1.00	0.00	0.00	135.0
8	1.00	0.00	0.00	128.0
9	1.00	0.00	0.00	118.0
10	1.00	0.00	0.00	151.0
11	1.00	0.00	0.00	262.0
12	1.00	0.00	0.00	4974.0
accuracy			0.00	5907.0
macro avg	0.70	0.30	0.00	5907.0
weighted avg	1.00	0.00	0.00	5907.0

```
[[ 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  0  0  0]
 [105 30  4  0  0  0  0  0  0  0  0]
 [108 26  1  0  0  0  0  0  0  0  0]
 [ 96 28  4  0  0  0  0  0  0  0  0]
 [ 89 27  2  0  0  0  0  0  0  0  0]
 [107 38  6  0  0  0  0  0  0  0  0]
 [185 66 11  0  0  0  0  0  0  0  0]
 [3398 844 732  0  0  0  0  0  0  0  0]]

Accuracy for our Training dataset with PCA: 0.0
Sihouette Score: 0.381863045581129
'\nSihouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring cluster
s.\n'
```

```
# predict the cluster for each testing data point
y_clus_test = km.predict(X_test)

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

train_accuracy = accuracy_score(y_test, y_clus_test)
print("\nAccuracy for our Training dataset with PCA:", train_accuracy)

#Calculate sihouette Score
score = metrics.silhouette_score(X_test, y_clus_test)
print("Sihouette Score: ", score)

'''
Sihouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.
'''
```

```

{X}  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      0.00      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 0 0 0]
 [ 41 21 3 0 0 0 0 0 0 0 0]
 [ 43 12 0 0 0 0 0 0 0 0 0]
 [ 57 10 1 0 0 0 0 0 0 0 0]
 [ 36 21 0 0 0 0 0 0 0 0 0]
 [ 63 17 5 0 0 0 0 0 0 0 0]
 [ 69 30 4 0 0 0 0 0 0 0 0]
 [1769 448 393 0 0 0 0 0 0 0 0]]

Accuracy for our Training dataset with PCA: 0.0
Silhouette Score: 0.3844803724503762
'\nSilhouette Score- ranges from -1 to +1 , a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring cluster s.\n'

```

#Question-2

Used pd_speech_features.csv

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

```

[18] dataset_pd = pd.read_csv('pd_speech_features.csv')
dataset_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

[19] dataset_pd.head()

   id  gender  PPE   DFA  RPDE  numPulses  numPeriodsPulses  meanPeriodPulses  stdDevPeriodPulses  locPctJitter  ...  tqwt_kurtosisValue_dec_28  tqwt_kurti
0  0      0      1  0.85247  0.71826  0.57227      240           239           0.008064           0.000087           0.00218  ...              1.5620
1  0      0      1  0.76686  0.69481  0.53966      234           233           0.008258           0.000073           0.00195  ...              1.5589
2  0      0      1  0.85083  0.67604  0.58982      232           231           0.008340           0.000060           0.00176  ...              1.5643
3  1      0      0  0.41121  0.79672  0.59257      178           177           0.010858           0.000183           0.00419  ...              3.7805
4  1      0      0  0.32790  0.79782  0.53028      236           235           0.008162           0.002669           0.00535  ...              6.1727

5 rows x 755 columns

```

```
principal component 1 principal component 2 Principal Component 3 class
```

principal component 1	principal component 2	Principal Component 3	class
-10.047372	1.471076	-6.846401	1
-10.637725	1.583748	-6.830976	1
-13.516185	-1.253543	-6.818696	1
-9.155083	8.833598	15.290908	1
-6.764470	4.611466	15.637125	1

```
[24] 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)
```

```

#2.c Support Vector Machine's

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
score = metrics.silhouette_score(X_test, y_pred)
print("Silhouette Score: ", score)

```

	precision	recall	f1-score	support
0	0.67	0.42	0.51	62
1	0.84	0.93	0.88	196
accuracy			0.81	258
macro avg	0.75	0.68	0.70	258
weighted avg	0.80	0.81	0.79	258

```

[[ 26 36]
 [ 13 183]]
accuracy is 0.810077519379845
Silhouette Score: 0.25044639806389946

```

#Question-3

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

```

#3. Apply Linear Discriminant Analysis (LDA) on Iris.csv dataset to reduce dimensionality of data to k=2.
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
dataset_iris = pd.read_csv('Iris.csv')
dataset_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

```

```

[28] dataset_iris.isnull().any()

Id                False
SepalLengthCm     False
SepalWidthCm      False
PetalLengthCm     False
PetalWidthCm      False
Species           False
dtype: bool

```

```

[29] x = dataset_iris.iloc[:,1:-1]
     y = dataset_iris.iloc[:, -1]
     print(x.shape, y.shape)

(150, 4) (150,)

```



```
[33] X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=0)
plt.show()

[34] sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
le = LabelEncoder()
y = le.fit_transform(y)

[35] 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)
```

#Question-4

Difference between PCA & LDA:

1. Both LDA and PCA rely on linear transformations and aim to maximize the variance in a lower dimension. PCA is an unsupervised learning algorithm while LDA is a supervised learning algorithm. This means that PCA finds directions of maximum variance regardless of class labels while LDA finds directions of maximum class separability.

PCA: It reduces the features into a smaller subset of orthogonal variables, called principal components – linear combinations of the original variables. The first component captures the largest variability of the data, while the second captures the second largest, and so on.

LDA: LDA finds the linear discriminants in order to maximize the variance between the different categories while minimizing the variance within the class.

GitHub Link:

<https://github.com/Sanjana9791/MachineLearningAssignment5.git>

Video Link:

https://drive.google.com/file/d/1S_CnaMZ5ojzbQjfu92k5QIA_3nJ7JX88/view?usp=sharing