MACHINE LEARNING ASSIGNMENT - 5

Name: Thandarupalli Naveen Goud

UCM ID :nxt46830

700#:700734683

Video

Link: https://drive.google.com/file/d/12jD6GhkvIeV9HxbDZ40_DN5nNpiBidT4/view?usp=share_link

Github Link: https://github.com/nxt46830/ML-Assignment-1

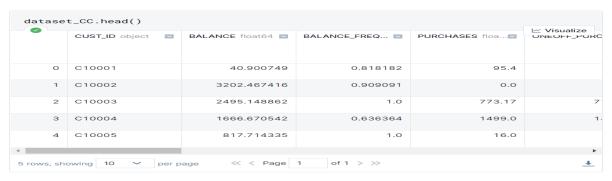
- 1. Principal Component Analysis
- a. Apply PCA on CC dataset.
- b. Apply k-means algorithm on the PCA result and report your observation if the silhouette score has improved or not?
- c. Perform Scaling+PCA+K-Means and report performance.

```
# 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_matr
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
sns.set(style="white", color_codes=True)
import warnings
warnings.filterwarnings("ignore")
```

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

```
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
    BALANCE
    BALANCE_FREQUENCY
                                                            float64
     PURCHASES
 4 ONEOFF_PURCHASES
                                                            float64
    INSTALLMENTS_PURCHASES
                                                            float64
 6 CASH_ADVANCE
                                                            float64
    PURCHASES FREQUENCY
                                                            float64
 8 ONEOFF_PURCHASES_FREQUENCY
                                                            float64
    PURCHASES_INSTALLMENTS_FREQUENCY 8950 non-null
                                                            float64
10 CASH_ADVANCE_FREQUENCY
                                          8950 non-null
                                                            float64
11 CASH_ADVANCE_TRX
                                          8950 non-null
                                                            int64
                                         8950 non-null
8950 non-null
8949 non-null
8950 non-null
12 PURCHASES TRX
                                                            int64
 13 CREDIT_LIMIT
                                                            float64
 14 PAYMENTS
                                                            float64
 15 MINIMUM_PAYMENTS
                                                            float64
 16 PRC_FULL_PAYMENT
                                          8950 non-null
                                                            float64
17 TENURE
                                          8950 non-null
                                                            int64
dtvpes: float64(14). int64(3). object(1)
```





```
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
dtype: bool
```

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

(8950, 16) (8950,)

#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', 'pri finalDf = pd.concat([principalDf, dataset_CC.iloc[:,-1]], axis = 1)
finalDf.head()
```

					∠ Visualize
	principal compo	principal compo	principal compo	TENURE int64	
0	-4326.383978558 351	921.5668815814 222	183.7083834738 5111	12	
1	4118.916664523 621	-2432.846345990 4533	2369.969289360 4197	12	
2	1497.907640740 3047	-1997.578694215 8522	-2125.631327723 396	12	
3	1394.548536133 8856	-1488.743452853 2237	-2431.799649021 8034	12	
4	-3743.351895614 359	757.3426565700 995	512.4764917625 604	12	
5 rows, she	owing 10 v per p	age « < Page	1 of 1 > >>		<u>+</u>

#1.b Apply K Means on PCA Result

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

	precision	recall	f1-score	su	pport
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	7:	584.0
accuracy			0.00	89	950.0
macro avg	0.70	0.30	0.00	89	950.0
weighted avg	1.00	0.00	0.00	89	950.0
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]
[28 175	1 0	0 0	0 0	0	0]
[15 173	2 0	0 0	0 0	0	0]
[27 169	0 0	0 0	0 0	0	0]
[26 149	0 0	0 0	0 0	0	0]
[47 188	1 0	0 0	0 0	0	0]
[78 284	3 0	0 0	0 0	0	0]
[2068 5390	126 0	0 0	0 0	0	0]]

Sihouette Score: 0.5109769750121257

^{&#}x27;\nSihouette Score- ranges from -1 to +1 , a high value indicates that the object is $v \triangleq$

```
#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', 'pri
finalDf = pd.concat([principalDf, dataset_CC.iloc[:,-1]], axis = 1)
finalDf.head()

✓ Visualize

 V
                                                                              Visualize
        principal compo...
                         principal compo...
                                           principal compo...
                                                            TENURE int64
        -1.718892260161
                         -1.072940187437
    0
                                            0.535622891588
                                                                          12
                    892
                                    9387
                                                      7592
                                            0.628180107453
        -1.169307579955
                          2.509321907735
                                                                          12
                   3465
                                    7274
                                                      5476
         0.938416272192
                          -0.382602616976
                                            0.160899824444
                                                                          12
     2
                   3274
                                    5072
                                                     46763
                          0.045859536408
     3
        -0.907502880168
                                            1.521747366473
                                                                          12
                   9199
                                  029006
                                                      7184
        -1.637828969277
                          -0.684976407724
                                            0.425558131372
                                                                          12
                   2323
                                    5449
                                                     08226
5 rows, showing 10 v per page
                               << < Page 1</p>
                                               of 1 > >>
                                                                                      <u>+</u>
```

S	principal compo	principal compo	principal compo	TENURE int64	∠ Visualiz
	,,				
0	-1.718892260161 892	-1.072940187437 9387	0.535622891588 7592	12	
1	-1.169307579955 3465	2.509321907735 7274	0.628180107453 5476	12	
2	0.938416272192 3274	-0.382602616976 5072	0.160899824444 46763	12	
3	-0.907502880168 9199	0.045859536408 029006	1.521747366473 7184	12	
4	-1.637828969277 2323	-0.684976407724 5449	0.425558131372 08226	12	
ows, sh	owing 10 V per p	age « < Page	1 of 1 > >>		

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

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

```
[14]
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.34,random_st
nclusters = 3
# this is the k in kmeans
km = KMeans(n_clusters=nclusters)
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
```

	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	Ø	0	0	0	0	0	0	0]
[0	0	0	0	0	0	0	0	0	0]
[1	LØ5	4	30	0	0	0	0	0	0	0]
[1	108	1	26	0	0	0	0	0	0	0]
[96	4	28	0	0	0	0	0	0	0]
[89	2	27	0	0	0	0	0	0	0]
[1	L07	6	38	0	0	0	0	0	0	0]
[1	184	11	67	0	0	0	0	0	0	0]
[34	108	723	843	0	0	0	0	0	0	0]]

Accuracy for our Training dataset with PCA: 0.0 Sihouette Score: 0.3816397521049892

^{&#}x27;\nSihouette Score- ranges from -1 to +1 , a high value indicates that the object is v $\mbox{$\updownarrow$}$

```
# 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
"""
```

precision	recall	f1-score	support
0.00	1.00	0.00	0.0
0.00	1.00	0.00	0.0
0.00	1.00	0.00	0.0
1.00	0.00	0.00	65.0
1.00	0.00	0.00	55.0
1.00	0.00	0.00	68.0
1.00	0.00	0.00	57.0
1.00	0.00	0.00	85.0
1.00	0.00	0.00	103.0
1.00	0.00	0.00	2610.0
		0.00	3043.0
0.70	0.30	0.00	3043.0
1.00	0.00	0.00	3043.0
	0.00 0.00 0.00 1.00 1.00 1.00 1.00 1.00	0.00 1.00 0.00 1.00 0.00 1.00 1.00 0.00 1.00 0.00 1.00 0.00 1.00 0.00 1.00 0.00 1.00 0.00	0.00 1.00 0.00 0.00 1.00 0.00 0.00 1.00 0.00 1.00 0.00 0.00

```
01
[[ 0
      0 0
             0
                0
                   0
                      0
                          0
                             0
[
  0
      0 0
             0
                0
                    0
                       0
                          0
                             0
                                0]
Γ
   0
      0
         0
            0
                0
                   0
                      0
                          0
                             0
                                0]
           0 0 0 0
                            0
                                0]
[ 41
      3 21
                          0
[ 43
      0 12 0 0 0 0 0 0 0]
57
      1 10
            0
               0 0 0
                         0 0
                                0]
                      0
                                0]
      0 22
                          0
                            0
35
            0
                0
                  0
[ 63
        17
             0
                0
                   0
                          0
                             0
                                0]
[ 68
            0
                            0
      4 31
                0
                                0]
                   0
                      0
                          0
[1775 382 453 0
                0
                          0 0 0]]
                    0
                      0
Accuracy for our Training dataset with PCA: 0.0
Sihouette Score: 0.38425040003352195
```

'\nSihouette Score- ranges from -1 to +1 , a high value indicates that the object is $v \diamondsuit$

```
[16]
    #2. Use pd_speech_features.csv
    # a. Perform Scaling
    # b. Apply PCA (k=3)
    # c. Use SVM to report performance
.
                                                                                         [17]
    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
                                                                                        [18]
   dataset_pd.isnull().any()
   id
                                False
   gender
                                False
   PPE
                                False
   DFA
                                False
   RPDE
                                False
   tqwt_kurtosisValue_dec_33
                                False
   tqwt_kurtosisValue_dec_34
                                False
    tqwt_kurtosisValue_dec_35
                                False
    tqwt_kurtosisValue_dec_36
                                False
   class
                                False
   Length: 755, dtype: bool
                                                                                        [19]
   X = dataset_pd.drop('class',axis=1).values
   y = dataset_pd['class'].values
                                                                                         [20]
     #Scaling Data
     scaler = StandardScaler()
    X_Scale = scaler.fit_transform(X)
```

```
[21]
# Apply PCA with k = 3
pca3 = PCA(n_components=3)
principalComponents = pca3.fit_transform(X_Scale)
principalDf = pd.DataFrame(data = principalComponents, columns = ['principal comp
finalDf = pd.concat([principalDf, dataset_pd[['class']]], axis = 1)
finalDf.head()
                                                                           Visualize
       principal compo... v principal compo... v
                                        Principal Compo...
                                                         class int64
                                                                       V
      -10.04737208020
                         1.471075534307
                                         -6.846406626217
                                                                       1
                                  3976
                   64
                                                   6795
       -10.63772513419
                         1.583748747034
                                         -6.830979740294
                                                                       1
                 8256
                                   128
                                                    065
       -13.51618523377
                        -1.253542561594
                                         -6.818698677650
                                                                       1
                  7032
                                  8577
                                                    746
    3 -9.155083698658
                        8.833598702496
                                         15.29089452140
                                                                       1
                                                   9589
                  137
                                   97
      -6.764469799613
                         4.611466268338
                                         15.63711379027
                                                                       1
```

253

 \equiv

```
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_st
```

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

```
precision recall f1-score support
          0
                 0.67
                        0.42
                                   0.51
                                               62
                 0.84
                         0.93
                                   0.88
                                              196
                                              258
   accuracy
                                   0.81
  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
Sihouette Score: 0.25044636702010714
```

3. Apply 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 dimensi
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
```

```
[25]
dataset_iris.isnull().any()
Ιd
               False
               False
SepalLengthCm
SepalWidthCm
              False
PetalLengthCm
             False
PetalWidthCm
              False
Species
              False
dtype: bool
                                                                                 [26]
x = dataset_iris.iloc[:,1:-1]
y = dataset_iris.iloc[:,-1]
print(x.shape,y.shape)
(150, 4) (150,)
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_s
```

```
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
le = LabelEncoder()
y = le.fit_transform(y)

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)
```

4. Briefly identify the difference between PCA and LDA

```
#4. Briefly identify the difference between PCA and LDA

"""Both LDA and PCA rely on linear transformations and aim to maximize the varian

'Both LDA and PCA rely on linear transformations and aim to maximize the variance in 

captures the largest variability of the data, while the second captures the secon

'It reduces the features into a smaller subset of orthogonal variables, called princif.

[32]

maximize the variance between the different categories while minimizing the varia

'LDA finds the linear discriminants in order to maximize the variance between the diff.
```

Principal Component Analysis:

- 1.PCA is a technique in unsupervised machine learning that is used to minimize dimensionality. The key idea of the vital component analysis (PCA) is to minimize the dimensionality of a data set consisting of several variables, either firmly or lightly, associated with each other while preserving to the maximum degree the variance present in the dataset.
- 2. This is achieved by translating the variables into a new collection of variables that are a mixture of our original dataset's variables or attributes so that maximum variance is preserved.

Linear Discriminant Analysis (LDA):

- 1.LDA is a technique of supervised machine learning. The critical principle of linear discriminant analysis (LDA) is to optimize the separability between the two classes to identify them in the best way we can determine. LDA is similar to PCA, which helps minimize dimensionality. Still, by constructing a new linear axis and projecting the data points on that axis, it optimizes the separability between established categories.
- 2.LDA does not function on finding the primary variable; it merely looks at what kind of point/features/subspace to distinguish the data offers further discrimination.

Conclusion:

The PCA and LDA are implemented in dimensionality reduction, which means a linear relationship between input and output variables. But it is possible to apply the PCA and LDA together and see the difference in their outcome. While PCA and LDA work on linear issues, they do have differences.