# **Assignment-5**

Name: Sammeta Jaya Sai Charan

**Student ID:** 700739775

Github link: https://github.com/saicharan255/Assignment5.git

Video link: https://drive.google.com/file/d/1q2rj-ejiSnqOe6WP17Gy-

EiiNlbS0OZL/view?usp=sharing

# 1. Principal Component Analysis

```
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")
dst_CC = pd.read_csv('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
1 BALANCE
2 BALANCE_FREQUENCY
3 PURCHASES
4 ONEOFF_PURCHASES
                                                8950 non-null object
                                                8950 non-null
                                                                   float64
                                               8950 non-null float64
                                               8950 non-null float64
8950 non-null float64
4 ONEOFF_PURCHASES
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
     CASH_ADVANCE_TRX
                                                8950 non-null
                                                                    int64
 12 PURCHASES_TRX
                                                8950 non-null
                                                                    int64
     CREDIT_LIMIT
                                                8949 non-null
 14 PAYMENTS
                                                8950 non-null
                                                                    float64
```

Firsty we have imported the required libraries and also imported the dataset CC GENERAL.csv, by using info() it shows the information of the dataset.

```
dst_CC.isnull().any()
     CUST_ID
                                            False
     BALANCE
                                            False
     BALANCE
             _FREQUENCY
                                            False
     PURCHASES
                                            False
    ONEOFF_PURCHASES
INSTALLMENTS_PURCHASES
                                           False
                                           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
CREDIT_LIMIT
                                            False
                                             True
     PAYMENTS
                                           False
     PRC_FULL_PAYMENT
                                            False
     TENURE
                                            False
     dtype: bool
[6] dst_CC.fillna(dst_CC.mean(), inplace=True)
     dst_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
```

We next checked the dataset for null values, if they have null values or true for null values then it replaces with NA because we are using fillna().

# a. Apply PCA on CC dataset

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

Now we are going to apply Principal component analysis on the CC dataset, we are using fit transform() and getting principal dataframe with principal component 1,2,3 and we get final dataframe is displayed by using head().

# b. Apply k-means algorithm on the PCA result

```
A = finalDf.iloc[:,0:-1]
     b = finalDf.iloc[:,-1]
     nclusters = 3 # this is the k in
     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)
     scr = metrics.silhouette_score(A, b_cluster_kmeans)
print("Sihouette Score: ",scr)
                       precision
                                        recall f1-score
C >
                                          1.00
                             0.00
                                          1.00
                                                       0.00
                             1.00
1.00
1.00
                                          0.00
                                                       0.00
                                          0.00
                                                       0.00
                   8
                                          0.00
                                                       0.00
                             1.00
                  10
                                          0.00
                  11
12
                                          0.00
                                                       0.00
         accuracy
                                          0.30
0.00
                                                       0.00
     weighted avg
                              0000
                                          0
0
0
                                                                      0]
                                                     9
9
9
          0
                                                                      øj
                              0
0
                                     0
0
                                                        0
0
                                                                      0]
0]
                                                         ø
                                                        ø
                                            ø
        188
                                                        ø
                                                                      0]
         284
     Accuracy of our training dataset with PCA: 0.0 Sihouette Score: 0.5109307274319468
```

In this step we apply k means algorithm on the obtained final data frame here we take k as 3 clusters and predicted the clusters for each data point. We then calculated the silhouette score using metrics.silhouette score().

# c. Perform Scaling+PCA+K-Means and report performance.

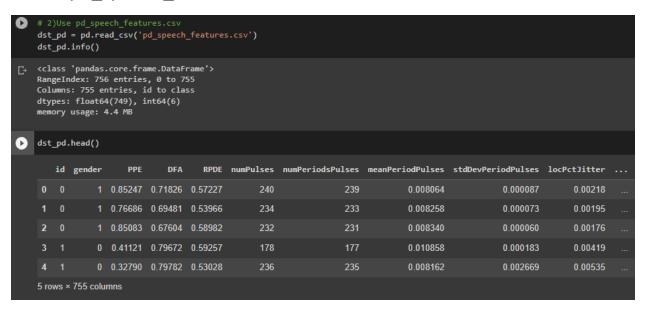
```
scaler = StandardScaler()
scaler.fit(a)
X_scaled_array = scaler.transform(a)
pca = PCA(3)
x_pca = pca.fit_transform(X_scaled_array)
prplDf = pd.DataFrame(data = x_pca, columns = ['principal cmpt 1', 'principal cmpt 2', 'principal cmpt 3'])
finalDf = pd.concat([prplDf, dst_CC.iloc[:,-1]], axis = 1)
finalDf.head()
                                                                    10.
   principal cmpt 1 principal cmpt 2 principal cmpt 3 TENURE
0
           -1.718893
                                                 0.535710
           -1.169308
                              2.509310
                                                 0.627778
            0.938415
3
           -0.907503
                              0.045856
                                                 1.521631
           -1 637829
                              -0 684972
                                                 0.425753
4
```

In the scaling we used StandardScaler() and scaler.transform(A) for the dataframe a. In PCA we used pca.fit\_transform and got the principal dataframe with principal component 1,2,3 and we get final dataframe is displayed by using head().

```
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)
    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)
    scr = metrics.silhouette_score(A_train, b_clust_train)
    print("Sihouette Score: ",scr)
                                recall f1-score support
₽
                  precision
                        0.00
                                             0.00
                                                         0.0
                                  1.00
1.00
                                             0.00
                        0.00
                                             0.00
                                                        0.0
                        1.00
                                             0.00
                                  0.00
                                                      139.0
                        1.00
               8
9
                        1.00
                                   0.00
                                             0.00
                                                       128.0
                                   0.00
                                             0.00
                                                      118.0
               10
                        1.00
                                             0.00
                                   0.00
                                                       151.0
               12
                                   0.00
                                             0.00
                                                     4974.0
        accuracy
                                                      5907.0
5907.0
    weighted avg
                                             0.00
                                                         0]
0]
0]
              0
0
                   0
0
                              0
0
                                   0
0
                                              0
0
                                         0
                                                         øj
       96
89
                                                         0]
0]
                  28
27
                         0
0
                                                   0
0
                              0
0
                                              0
0
       185
                  66
                                         0
     F3393 739 842
                              0
                                   0
    Accuracy of our training dataset with PCA: 0.0
    Sihouette Score: 0.38120810623375717
```

We applied k means algorithm by taking k=3 clusters, predicted the each cluster from each training data point then the predictions were made by the classifier. We calculated the silhouette score using metrics.silhouette score().

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



Firstly we have imported the required pd speech features.csv dataset by using pd.read csv().

#### a. Perform Scaling

```
[ ] #a)Scaling Data
scaler = StandardScaler()
X_Scale = scaler.fit_transform(X)
```

We have done scaling using StandardScalar().

# b. Apply PCA (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()
    principal cmpnt 1 principal cmpnt 2 Principal cmpnt 3 class
0
            -10.047372
                                 1.471078
                                                    -6.846402
           -10.637725
                                 1.583750
                                                    -6.830976
            -13.516185
                                                    -6.818696
                                                    15.290899
3
             -9.155083
                                 8.833600
             -6.764469
                                 4.611465
```

we applied the PCA with k=3, used pca.fit\_transform() and got the principal dataframe with principal component 1,2,3 and we get final dataframe is displayed by using head().

# c. Use SVM to report performance

```
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))
glass_acc_svc = accuracy_score(y_pred,y_test)
print('accuracy is',glass_acc_svc )
scr= metrics.silhouette_score(X_test, y_pred)
print("Sihouette Score: ",scr)
              precision recall f1-score support
                 0.67 0.42 0.51
0.84 0.93 0.88
                                                   62
                                                   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
Sihouette Score: 0.2504463884705572
```

In the above we used support vector machine to obtain the performance svmClassifier.predict(). Predictions were made by the classifier and accuracy is calculated. Finally we calculated the silhouette score and got 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('Iris.csv')
    dst_iris.info()
_- <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 150 entries, 0 to 149
    Data columns (total 6 columns):
                      Non-Null Count Dtype
     # Column
                         150 non-null
                                            int64
     1 SepalLengthCm 150 non-null
2 SepalWidthCm 150 non-null
3 PetalLengthCm 150 non-null
                                           float64
                                            float64
                                          float64
     4 PetalWidthCm 150 non-null
                                           float64
    5 Species 150 non-null obj
dtypes: float64(4), int64(1), object(1)
                                            object
    memory usage: 7.2+ KB
```

We imported the required libraries linear discriminant analysis, loaded Iris.csv dataset from local drive.

```
dst_iris.isnull().any()
                     False
    SepalLengthCm
                   False
    SepalWidthCm
                   False
    PetalLengthCm False
    PetalWidthCm
                   False
    Species
                    False
    dtype: bool
[ ] x = dst_iris.iloc[:,1:-1]
    y = dst_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_state=0)
[ ] 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)
```

Firstly we have checked for null values, located the rows and columns, we used standardScalar() to perform the operation and trained, tested the dataset. Now we applied linear discriminant analysis LDA algorithm trained and tested the dataset and printed the result.

# 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