

CS 5710
MACHINE LEARNING

Name: Shiva Kandagatla

Student Id: 700745245

E-mail: SXK54250@ucmo.edu

Course: CS 5710

Assignment: Assignment #5

GitHub Link:

<https://github.com/Shiva-Kandagatla98/In-Class-Assignment-5.git>

Video Demo Link:

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.

Source Code:

```
question1.py U X
8
9 import pandas as pd
10 from sklearn.decomposition import PCA
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.model_selection import train_test_split
13 from sklearn.cluster import KMeans
14 from sklearn import metrics
15 import warnings
16 warnings.filterwarnings("ignore")
17
18 # Reading dataset CC.csv
19 import pandas as pd
20 from sklearn.decomposition import PCA
21 from sklearn.preprocessing import StandardScaler
22 from sklearn.model_selection import train_test_split
23 from sklearn.cluster import KMeans
24 from sklearn import metrics
25 from sklearn.linear_model import LogisticRegression
26
27 # Load dataset
28 df = pd.read_csv('CC GENERAL.csv')
29
30 # Replace null values with mean
31 df['MINIMUM_PAYMENTS'].fillna(df['MINIMUM_PAYMENTS'].mean(), inplace=True)
32 df['CREDIT_LIMIT'].fillna(df['CREDIT_LIMIT'].mean(), inplace=True)
33
34 # Apply PCA on dataset
35 pca = PCA(n_components=2)
36 X_pca = pca.fit_transform(df.iloc[:, 1:18])
37 pca_df = pd.DataFrame(data=X_pca, columns=['PC1', 'PC2'])
38 final_df = pd.concat([pca_df, df[['TENURE']]], axis=1)
39 print(final_df)
40
41 # Split into train and test sets
42 X = final_df.drop('TENURE', axis=1).values
43 y = final_df['TENURE'].values
44 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
45
```

```
question1.py X
41 # Split into train and test sets
42 X = final_df.drop('TENURE', axis=1).values
43 y = final_df['TENURE'].values
44 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
45
46 # Train logistic regression model on training set
47 lr = LogisticRegression()
48 lr.fit(X_train, y_train)
49 y_train_pred = lr.predict(X_train)
50 train_accuracy = metrics.accuracy_score(y_train, y_train_pred) * 100
51 print('Accuracy on training set with PCA: %.4f %%' % train_accuracy)
52
53 # Apply k-means clustering on original data and calculate Silhouette score
54 n_clusters = 2
55 kmeans = KMeans(n_clusters=n_clusters)
56 kmeans.fit(df.iloc[:, 1:18])
57 y_cluster_kmeans = kmeans.predict(df.iloc[:, 1:18])
58 score = metrics.silhouette_score(df.iloc[:, 1:18], y_cluster_kmeans)
59 print('Silhouette Score for original data with k-means: ', score)
60
61 # Apply scaling, PCA, and k-means clustering on data and calculate Silhouette score
62 scaler = StandardScaler()
63 X_scaled = scaler.fit_transform(df.iloc[:, 1:18])
64 pca2 = PCA(n_components=2)
65 X_pca2 = pca2.fit_transform(X_scaled)
66 pca_df2 = pd.DataFrame(data=X_pca2, columns=['PC1', 'PC2'])
67 final_df2 = pd.concat([pca_df2, df[['TENURE']]], axis=1)
68 n_clusters = 2
69 kmeans2 = KMeans(n_clusters=n_clusters)
70 kmeans2.fit(X_scaled)
71 y_cluster_kmeans2 = kmeans2.predict(X_scaled)
72 score2 = metrics.silhouette_score(X_scaled, y_cluster_kmeans2)
73 print('Silhouette Score for scaled data with PCA and k-means: ', score2)
74
```

Output:

```
question1.py — In-Class-Assignment-5

question1.py ×
1
2 # Question - 1
3 #
4 # 1. Principal Component Analysis
5 #     a. Apply PCA on CC dataset.
6 #     b. Apply k-means algorithm on the PCA result and report your observation if the silhouette score has improved or not?
7 #     c. Perform Scaling+PCA+K-Means and report performance
8
9 import pandas as pd
10 from sklearn.decomposition import PCA
11 from sklearn.preprocessing import StandardScaler

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

/usr/bin/python3 "/Users/shiva/Downloads/Assignment 5/In-Class-Assignment-5/question1.py"

The default interactive shell is now zsh.
To update your account to use zsh, please run `chsh -s /bin/zsh`.
For more details, please visit https://support.apple.com/kb/HT208050.
● Shivas-Air-2:In-Class-Assignment-5 shiva$ /usr/bin/python3 "/Users/shiva/Downloads/Assignment 5/In-Class-Assignment-5/question1.py"
   PC1      PC2  TENURE
0 -4326.383956  921.566884    12
1  4118.916676 -2432.846347    12
2  1497.907660 -1997.578692    12
3  1394.548556 -1488.743450    12
4  -3743.351874   757.342659    12
...
8945 -4208.357938  1122.443274     6
8946 -4123.924001   951.683803     6
8947 -4379.444202   911.504566     6
8948 -4791.117744  1032.540944     6
8949 -3623.702749  1555.134769     6

[8950 rows x 3 columns]
Accuracy on training set with PCA: 84.0543 %
Silhouette Score for original data with k-means: 0.5117097284367592
Silhouette Score for scaled data with PCA and k-means: 0.20966187483954374
○ Shivas-Air-2:In-Class-Assignment-5 shiva$
```

-----*****-----

2. Use pd_speech_features.csv

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

Source Code:

```
question2.py U X
1 # 2. Use pd_speech_features.csv
2 # a. Perform Scaling
3 # b. Apply PCA (k=3)
4 # c. Use SVM to report performance
5
6 import pandas as pd
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.decomposition import PCA
9 from sklearn.model_selection import train_test_split
10 from sklearn.svm import SVC
11
12 # Load dataset
13 data = pd.read_csv("pd_speech_features.csv")
14 features = data.drop('class', axis=1).values
15 labels = data['class'].values
16
17 # Apply scaling
18 scaler = StandardScaler()
19 scaled_features = scaler.fit_transform(features)
20
21 # PCA
22 pca = PCA(n_components=3)
23 pca_components = pca.fit_transform(scaled_features)
24 pca_df = pd.DataFrame(data=pca_components,
25                       columns=['PC 1', 'PC 2', 'PC 3'])
26 final_df = pd.concat([pca_df, data[['class']]], axis=1)
27 print('Final dataset after PCA:')
28 print(final_df)
29
30 # Train and evaluate SVM model
31 X_train, X_test, y_train, y_test = train_test_split(scaled_features, labels, test_size=0.3, random_state=0)
32 svm_model = SVC(max_iter=1000)
33 svm_model.fit(X_train, y_train)
34 y_pred = svm_model.predict(X_test)
35 svm_acc = round(svm_model.score(X_train, y_train) * 100, 2)
36 print("SVM accuracy =", svm_acc)
37
```

Output:

```
question2.py x
1 # 2. Use pd_speech_features.csv
2 # a. Perform Scaling
3 # b. Apply PCA (k=3)
4 # c. Use SVM to report performance
5
6 import pandas as pd
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.decomposition import PCA
9 from sklearn.model_selection import train_test_split
10 from sklearn.svm import SVC
11

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
The default interactive shell is now zsh.
To update your account to use zsh, please run `chsh -s /bin/zsh`.
For more details, please visit https://support.apple.com/kb/HT208050.
● Shivas-Air-2:In-Class-Assignment-5 shiva$ /usr/bin/python3 "/Users/shiva/Downloads/Assignment 5/In-Class-Assignment-5/question1.py"
      PC1      PC2  TENURE
0 -4326.383956  921.566884    12
1  4118.916676 -2432.846347    12
2  1497.907660 -1997.578692    12
3  1394.548556 -1488.743450    12
4 -3743.351874  757.342659    12
...      ...      ...      ...
8945 -4208.357938 1122.443274     6
8946 -4123.924001  951.683803     6
8947 -4379.444202  911.504566     6
8948 -4791.117744 1032.540944     6
8949 -3623.702749 1555.134769     6

[8950 rows x 3 columns]
Accuracy on training set with PCA: 84.0543 %
Silhouette Score for original data with k-means: 0.5117097284367592
Silhouette Score for scaled data with PCA and k-means: 0.20966187483954374
● Shivas-Air-2:In-Class-Assignment-5 shiva$ /usr/bin/python3 "/Users/shiva/Downloads/Assignment 5/In-Class-Assignment-5/question2.py"
Final dataset after PCA:
      PC 1      PC 2      PC 3  class
0 -10.047372  1.471076 -6.846408     1
1 -10.637725  1.583749 -6.830979     1
2 -13.516185 -1.253542 -6.818700     1
3 -9.155084  8.833601 15.290900     1
4 -6.764471  4.611467 15.637122     1
...      ...      ...      ...
751 22.322682  6.481912  1.458757     0
752 13.442877  1.449413  9.352299     0
753  8.270263  2.391285 -0.908676     0
754  4.011760  5.412258 -0.847130     0
755  3.993112  6.072418 -2.020722     0

[756 rows x 4 columns]
SVM accuracy = 91.68
○ Shivas-Air-2:In-Class-Assignment-5 shiva$
```

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

Source Code:

```
question3.py X
1 # 3. Apply Linear Discriminant Analysis (LDA) on Iris.csv dataset to reduce dimensionality of data tok=2.
2
3 import pandas as pd
4 from sklearn.preprocessing import StandardScaler, LabelEncoder
5 from matplotlib import pyplot as plt
6 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
7 import warnings
8 warnings.filterwarnings("ignore")
9 import seaborn as sns
10 sns.set(style="white", color_codes=True)
11
12
13 import pandas as pd
14 from sklearn.preprocessing import StandardScaler, LabelEncoder
15 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
16 import seaborn as sns
17 import matplotlib.pyplot as plt
18
19 # Load iris dataset
20 df = pd.read_csv("iris.csv")
21
22 # Standardize features
23 scaler = StandardScaler()
24 X = scaler.fit_transform(df.iloc[:, :4].values)
25
26 # Encode class labels
27 label_encoder = LabelEncoder()
28 y = label_encoder.fit_transform(df['Species'].values)
29
30 # Apply LDA
31 lda = LinearDiscriminantAnalysis(n_components=2)
32 X_lda = lda.fit_transform(X, y)
33
34 # Create new dataframe with LDA components and class labels
35 data = pd.DataFrame(X_lda, columns=['LD1', 'LD2'])
36 data['class'] = y
37
38 print(data)
```

Output:

```
● Shivas-Air-2:In-Class-Assignment-5 shiva$ /usr/bin/python3 "/Users/shiva/Downloads/Assignment 5/In-Class-Assignment-5/question3.py"
LD1      LD2      class
0  9.423452 -0.513976  0
1  8.751900 -1.591678  0
2  8.973004 -1.068204  0
3  8.170186 -1.435135  0
4  9.249789 -0.136869  0
..      ...      ...
145 -6.167264 2.006912  2
146 -6.462099 1.062334  2
147 -6.447135 2.055824  2
148 -7.018075 2.727998  2
149 -6.838386 2.108245  2
[150 rows x 3 columns]
○ Shivas-Air-2:In-Class-Assignment-5 shiva$
```

-----*****-----

4. Briefly identify the difference between PCA and LDA

PCA is an unsupervised technique that seeks to reduce the dimensionality of a dataset while preserving most of its variance. It does so by transforming the data into a new coordinate system where the new axes (principal components) are orthogonal and ordered by the amount of variance they capture. PCA is often used for data visualization, data compression, and feature extraction.

LDA, on the other hand, is a supervised technique that seeks to find a linear combination of features that maximizes the separation between different classes of data. It does so by finding a projection that maximizes the between-class scatter and minimizes the within-class scatter. LDA is often used for classification and feature extraction.

In summary, PCA is an unsupervised technique that preserves the maximum amount of variance in the data while reducing dimensionality, while LDA is a supervised technique that seeks to maximize the separation between classes by finding a projection that preserves the class structure of the data.

-----*****-----