Assignment 5 Documentation Yelamati Pavan Kalyan Rao 700729168

1. CC Dataset---PCA---K Means Clustering

- Read cc dataset using pandas.
- Check head to see columns and type of data.
- Drop CUST_ID column.
- Check for null values using isnull().sum().
- Fill missing values with column means using fillna() method.
- Use the k=5 to fit, predict KMeans on the data.
- Calculate silhouette score. In our case its 0.379 for k=5.
- Now use n_components=5 to fit, transform cc data with PCA.
- Use the k=5 to fit, predict KMeans on the transformed data.
- Calculate silhouette_score. In our case its 0.416 for k=5.
- Now do MinMaxScaling on cc data to bring uniformity to various columns. Then do PCA on scaled data with same configuration above.
- Use the k=5 to fit, predict KMeans on the scale transformed data. Calculate silhouette_score. In our case its 0.383 for k=5.

Initial KMeans fit on original data:

```
In [9]: ##elbow is at k=5
km = KMeans(n_clusters=5, random_state=0)
km.fit_predict(cc)
score = silhouette_score(cc, km.labels_, metric='euclidean')
print('Initial Silhouetter Score: %.3f' % score)
Initial Silhouetter Score: 0.379
```

KMeans fit on PCA transformed data:

Using PCA has improved clustering silhouette score.

KMeans fit on MinMax scaled PCA transformed data:

2. Speech Dataset---PCA---SVM

- Read speech dataset using pandas.
- Check head to see columns and type of data.
- Drop "id" column as it is not useful for modeling.
- Check for null values using isnull().sum(). No missing values found.
- Separate X and Y variables with Y as "class".
- Fit, predict SVM using this X and Y. Calculate accuracy. Its 0.756.
- Do MinMax scaling and Fit, predict SVM using this scaled_X and Y. Calculate accuracy. Its 0.87.
- Now Perform PCA on scaled data. Use this to fit SVM again and calculate accuracy.
- With n_components=3,5,10,30,50,100 test for accuracies. We observe an increase in accuracy as n_components increase. As small number of components doesn't hold(represent) all the required information of variables.

Initial SVM accuracy:

Accuracy of SVM on scaled data:

Accuracy of SVM on scaled PCA data: (n_components=3,5,10,30,50,100)

```
In [39]:  pca = PCA(n components=3)
              pca.fit(x speech transformed mms)
              speech pca transformed = pd.DataFrame(pca.transform(x speech transformed mms))#
 In [40]: H clf = svm.SVC()
              clf.fit(speech pca transformed, y speech)
             y_pred=clf.predict(speech_pca_transformed)
             print(accuracy score(y speech, y pred))
              0.8042328042328042
In [41]:  pca = PCA(n components=5)
             pca.fit(x speech transformed mms)
             speech pca transformed = pd.DataFrame(pca.transform(x speech transformed mms))
In [42]: H clf = svm.SVC()
             clf.fit(speech_pca_transformed, y_speech)
             y pred=clf.predict(speech pca transformed)
             print(accuracy_score(y_speech, y_pred))
             0.8359788359788359
In [43]:  pca = PCA(n_components=10)
             pca.fit(x speech transformed mms)
             speech pca transformed = pd.DataFrame(pca.transform(x speech transformed mms))
In [44]: H clf = svm.SVC()
             clf.fit(speech_pca_transformed, y_speech)
            y pred=clf.predict(speech pca transformed)
             print(accuracy score(y speech, y pred))
             0.8664021164021164
```

```
In [45]: pca = PCA(n components=30)
             pca.fit(x speech transformed mms)
             speech pca transformed = pd.DataFrame(pca.transform(x speech transformed mms))
In [46]: H clf = svm.SVC()
             clf.fit(speech_pca_transformed, y_speech)
             y_pred=clf.predict(speech_pca_transformed)
             print(accuracy_score(y_speech, y_pred))
             0.9074074074074074
In [47]:  Pca = PCA(n_components=50)
             pca.fit(x speech transformed mms)
             speech pca_transformed = pd.DataFrame(pca.transform(x_speech_transformed_mms));
In [48]: H clf = svm.SVC()
             clf.fit(speech_pca_transformed, y_speech)
             y_pred=clf.predict(speech_pca_transformed)
             print(accuracy_score(y_speech, y_pred))
             0.9206349206349206
In [49]:  ▶ pca = PCA(n_components=100)
             pca.fit(x speech transformed mms)
             speech pca transformed = pd.DataFrame(pca.transform(x speech transformed mms))
In [50]: M clf = svm.SVC()
             clf.fit(speech pca transformed, y speech)
             y pred=clf.predict(speech pca transformed)
             print(accuracy score(y speech, y pred))
             0.9298941798941799
```

3. IRIS Dataset---LDA

- Read iris dataset using pandas.
- Check head to see columns and type of data.
- Separate X and Y variables with Y as "Species".
- Use LDA from sklearn to transform the data with n_components=2.
- Check for the transformed variables.

```
In [29]: ► Ida = LDA(n components=2)
            iris x transformed = lda.fit(iris x, iris y).transform(iris x)
In [30]:  print(iris_x_transformed)
               1.04013304CT01 Z.Z100/131
             [-9.96642412e+00 1.39824256e+00]
             [-9.19297836e+00 3.48117706e-01]
             [-9.36227503e+00 8.31991272e-01]
             [-9.32391418e+00 9.16315048e-01]
             [-8.76248233e+00 -5.62718432e-02]
             [-8.88178087e+00 8.74105504e-01]
             [-9.87055894e+00 1.00209591e+00]
             [-7.60804498e+00 -8.20642233e-03]
             [-7.77791864e+00 -3.20242776e-01]
             [-7.96141242e+00 -6.57043051e-01]
             [-8.00979719e+00 3.34031714e-01]
             [-8.92701582e+00 6.16194234e-01]
             [-8.96319773e+00 5.77438977e-01]
              -7.88178440e+00 -1.83196764e-01]
             [-7.77863981e+00 -3.08300446e-01]
             [-8.32075648e+00 7.97946692e-01]
              [-9.80491422e+00 1.96064097e+00]
              [-1.00940635e+01 2.49438691e+00]
              [-8.14751552e+00 -5.83887141e-02]
```

3. PCA vs LDA

- Linear discriminant analysis is very similar to PCA both look for linear combinations of the features which best explain the data.
- The main difference is that the Linear discriminant analysis is a supervised dimensionality reduction technique that also achieves classification of the data simultaneously.
- While Principal component analysis is an unsupervised Dimensionality reduction technique, it ignores the class label.