

```
1 ###
2 import pandas as pd
3
4
5 import matplotlib.pyplot as plt
6
7 import seaborn as sns
8
9
10 import warnings
11 warnings.filterwarnings("ignore")
12
13 np.set_printoptions(suppress = True)
14
15
16 from sklearn.preprocessing import StandardScaler
17 from sklearn.metrics import silhouette_score
18 from sklearn.cluster import AgglomerativeClustering
19 from sklearn.metrics.pairwise import
    euclidean_distances
20 from sklearn.cluster import DBSCAN
21
22 from sklearn.decomposition import PCA, TruncatedSVD
23 from sklearn.metrics import silhouette_score
24 from sklearn.cluster import KMeans
25
26 # import functions from scipy to perform clustering
27 from scipy.cluster.hierarchy import linkage
28 from scipy.cluster.hierarchy import dendrogram
29 from scipy.cluster.hierarchy import cophenet
30
31 from sklearn.ensemble import RandomForestClassifier
32 from sklearn.metrics import accuracy_score
33 from sklearn.metrics import confusion_matrix
34 from sklearn.metrics import classification_report
35 from sklearn.tree import DecisionTreeClassifier
36
37 from sklearn.discriminant_analysis import
    LinearDiscriminantAnalysis
38
39 from surprise import KNNWithMeans, SVDpp
```

```
40 from surprise import accuracy
41 ### md
42 ### Section B 1
43 ###
44 df_data = pd.read_csv("diabetes.csv")
45 ###
46 df_data.head(10)
47 ###
48 df_data.shape
49 ###
50 df_data.dtypes
51 ###
52 df_data.describe()
53 ###
54 df_data.isnull().sum()
55 ###
56 df_data.Glucose = df_data.Glucose.replace(0,
    df_data['Glucose'].mean())
57 ###
58 def replace_q(s):
59     if s == '20c':
60         return 500.0
61     elif s=='!':
62         return 0.0
63     else:
64         return float(s)
65 ###
66 gmean = df_data.Glucose.mean()
67 df_data['Glucose'] = df_data['Glucose'].apply(
    lambda x : gmean if x == 0 else float(x))
68 ###
69 df_data.BloodPressure = df_data.BloodPressure.
    replace(0,df_data['BloodPressure'].mean())
70 df_data.SkinThickness = df_data.SkinThickness.
    replace(0, df_data['SkinThickness'].mean())
71 df_data.Insulin = df_data.Insulin.replace(0,df_data
    ['Insulin'].mean())
72 df_data.BMI = df_data.BMI.replace(0,df_data['BMI'].
    mean())
73 ###
74 df_data.describe()
```

```

75 #%%
76 plt.rcParams['figure.figsize'] = [15,8]
77 df_data.boxplot()
78
79 plt.show()
80 #%%
81 Q1 = df_data.quantile(0.25)
82 Q3 = df_data.quantile(0.75)
83
84
85 IQR = Q3 - Q1
86
87 df_data = df_data[~((df_data < (Q1 - 1.5*IQR)) | (
    df_data > (Q3 + 1.5*IQR))).any(axis=1)]
88
89 df_data.shape
90 #%%
91 plt.rcParams['figure.figsize'] = [15,8]
92 df_data.boxplot()
93
94 plt.show()
95 #%%
96 X = df_data.drop('Outcome' , axis=1)
97
98 y = df_data['Outcome']
99 #%%
100 X_scle = StandardScaler().fit_transform(X)
101 #%% md
102 90% variance
103
104
105 #%%
106 cov_matrix = np.cov(X_scle.T)
107 eig_vals,eig_vecs = np.linalg.eig(cov_matrix)
108
109 print("Eigen Vectors \n%s",eig_vecs)
110 print("Eigen Values \n%s",eig_vals)
111 #%%
112 tot = sum(eig_vals)
113 var_exp = [(i/tot)*100 for i in sorted(eig_vals,
    reverse=True)]

```

```

114 cum_var_exp = np.cumsum(var_exp)
115 print("Cumulative variance explained",cum_var_exp)
116 ###
117 plt.plot(cum_var_exp)
118 plt.axvline(x=5, ymax=90, ls='--', color='r')
119 plt.axhline(y=90, xmax=90, ls='--', color='r')
120 ###
121 print("Top 5 Eigen Vectors are :",sorted(eig_vals,
      reverse=True)[0:5])
122 ###
123 var_exp
124 ###
125 pca = PCA(0.90)
126 X_PCA = pca.fit(X_scle)
127 ###
128 X_PCA.explained_variance_
129 ###
130 X_PCA.components_
131 ###
132 X_PCA.explained_variance_[ :5]
133 ###
134 X_PCA.components_[ :5]
135 ### md
136 ### Section B 2
137 ###
138 # 95% variance
139 pca1 = PCA(0.95)
140 X_PCA_1 = pca1.fit_transform(X)
141 ###
142 from yellowbrick.cluster import
    SilhouetteVisualizer
143 n_clusters=[2,3,4,5,6,7,8,9,10]
144
145 for K in n_clusters:
146     model=KMeans(n_clusters=K,random_state=42)
147     viz = SilhouetteVisualizer(model).fit(X_PCA_1)
148     viz.show()
149 ###
150 km = KMeans(n_clusters=2,random_state=42)
151 model = km.fit(X_PCA_1)
152 df_data['Cluster'] = model.labels_

```

```

153 #%%
154 df_data.head()
155 #%%
156 df_data[df_data['Cluster']==0].describe()
157 #%%
158 df_data[df_data['Cluster']==1].describe()
159 #%% md
160 **Interpretation:** The above summary shows that
    the average Insulin in this cluster is 78 +/- 15
    dev. On average, the gl level is 113.86.
161 #%% md
162 # Optimal number of clusters is 2
163 #%% md
164 # Section B 3
165 #%%
166 link = linkage(X_PCA_1, method='single')
167 dendrogram(link)
168 plt.show()
169 #%%
170 link = linkage(X_PCA_1, method='centroid')
171 dendrogram(link)
172 plt.show()
173 #%%
174 link = linkage(X_PCA_1, method='ward')
175 dendrogram(link)
176 plt.show()
177 #%%
178 link = linkage(X_PCA_1, method='complete')
179 dendrogram(link)
180 plt.show()
181 #%% md
182 # Section B 4
183 #%%
184 model_Km = KMeans(n_clusters = 4, random_state=40)
185 model_Km.fit(X_PCA_1)
186
187 df_data['Cluster_KM'] = model_Km.labels_
188 #%% md
189 model = KMeans(n_clusters=4)
190 distances = model.fit_transform(X_PCA_1)
191 variance = 0

```

```
192 i = 0
193 for label in model.labels_:
194     variance = variance + distances[i][label]
195     i = i + 1
196 model.labels_
197 ###
198 #Bar plot
199 ### md
200 # Section C 1
201 ### md
202 # Normal
203 ###
204 X_train, X_test, y_train, y_test =
    train_test_split(X, y, test_size=0.20,
        random_state=40 )
205
206 rf = RandomForestClassifier(n_estimators=15)
207 model = rf.fit(X_train,y_train)
208 y_pred = model.predict(X_test)
209
210 print(accuracy_score(y_test, y_pred))
211 cm = confusion_matrix(y_test, y_pred)
212 cm
213 ### md
214 # PCA
215 ###
216 pca = PCA(0.90 )
217 X_PCA_2 = pca.fit_transform(X)
218 ###
219 X_train, X_test, y_train, y_test =
    train_test_split(X_PCA_2, y, test_size=0.20,
        random_state=40 )
220
221 rf = RandomForestClassifier(n_estimators=20)
222 model = rf.fit(X_train,y_train)
223 y_pred = model.predict(X_test)
224
225 print(accuracy_score(y_test, y_pred))
226 cm = confusion_matrix(y_test, y_pred)
227 cm
228 ### md
```

```
229 # SVD
230 #%%
231 svd = TruncatedSVD(8)
232 X_SVD = svd.fit_transform(X)
233 #%%
234 np.cumsum(svd.explained_variance_/sum(svd.
    explained_variance_)*100)
235 #%%
236 svd = TruncatedSVD(3)
237 X_SVD = svd.fit_transform(X)
238 #%%
239 X_train, X_test, y_train, y_test =
    train_test_split(X_SVD, y, test_size=0.20,
        random_state=40 )
240
241 rf = RandomForestClassifier(n_estimators=20)
242 model = rf.fit(X_train,y_train)
243 y_pred = model.predict(X_test)
244
245 print(accuracy_score(y_test, y_pred))
246 cm = confusion_matrix(y_test, y_pred)
247 cm
248 #%% md
249 # LDA
250 #%%
251 #num of classes (2) ==> n_classes - 1 = 1
252 lda = LinearDiscriminantAnalysis(n_components = 1)
253 X_lda = lda.fit_transform(X, y)
254 #%%
255 X_train, X_test, y_train, y_test =
    train_test_split(X_lda, y, test_size=0.20,
        random_state=40 )
256
257 rf = RandomForestClassifier(n_estimators=20)
258 model = rf.fit(X_train,y_train)
259 y_pred = model.predict(X_test)
260
261 print(accuracy_score(y_test, y_pred))
262 cm = confusion_matrix(y_test, y_pred)
263 cm
264 #%%
```

```

265
266 ### md
267 # Scetion C 2
268 ###
269 df_rating = pd.read_csv('ratings.csv')
270 df_rating.head()
271 ###
272 df_rating.drop('timestamp', axis=1, inplace=True)
273 df_rating.head()
274 ###
275 df_rating[['movieId', 'rating']].groupby('movieId')
    .mean().sort_values(by='rating', ascending=False)
276 ###
277 popularity_tabel = df_rating.groupby('movieId').
    agg({'rating':'mean', 'userId':'count'})
278 popularity_tabel.rename({'userId':'TotalCount'},
    axis=1, inplace=True)
279 popularity_tabel = popularity_tabel.sort_values('
    rating', ascending= False)
280 popularity_tabel.head()
281 ###
282 popularity_tabel[popularity_tabel['TotalCount'
    ] > 250].head(5)
283 ### md
284 # Scetion C 3
285 ###
286 df_rating
287 ###
288 from surprise import Dataset,Reader
289 from surprise.model_selection import
    train_test_split,cross_validate
290 from surprise import KNNWithMeans,SVDpp
291 from surprise import accuracy
292 ###
293 reader = Reader(rating_scale=(1,5))
294 rating_data = Dataset.load_from_df(df_rating[['
    userId','movieId','rating']],reader)
295
296 [trainset,testset] = train_test_split(rating_data,
    test_size=0.15,shuffle=True)
297 trainsetfull=rating_data.build_full_trainset()

```



```
298 #%%
299 print("Number of users :",trainsetfull.n_users,'\n
      ')
300 print("Number of movies :",trainsetfull.n_items,'\n
      n')
301 #%%
302 alg=SVDpp()
303 alg.fit(trainsetfull)
304 #%%
305 alg.predict(uid=710,iid=166534)
306 #%%
307 results=cross_validate(algo=alg,data=rating_data,
      measures=['rmse'],cv=5,return_train_measures=True)
308 #%%
309
```