

```

1 import numpy as np
2 import pandas as pd
3 import math
4 import matplotlib.pyplot as plt
5 import matplotlib.cm as cm
6 import matplotlib.gridspec as gridspec
7 import seaborn as sns
8 from sklearn.cluster import KMeans
9 from yellowbrick.cluster import KElbowVisualizer
10 from sklearn.decomposition import PCA
11 from sklearn.metrics import silhouette_samples, silhouette_score
12 from sklearn.cluster import DBSCAN
13 import plotly.express as px
14 import umap.plot
15 from sklearn.manifold import TSNE
16
17 #####***** Function Hub *****#####
18
19 ##### Subpart 1 #####
20
21 def cat_plot(data, cat_features):
22     # Plot barplots
23     fig = plt.figure(figsize=(10,20))
24     gs = gridspec.GridSpec(5,2)
25     ax = {}
26
27     for ftr, i in zip(cat_features, range(len(cat_features))):
28         ax[i] = fig.add_subplot(gs[i])
29         ax[i] = sns.countplot(data, x=ftr)
30         ax[i].set_xticklabels(ax[i].get_xticklabels()) #, rotation=40
31         ax[i].set_xlabel(cat_features[i])
32     plt.tight_layout()
33     plt.show()
34
35
36 def num_plot(X):
37     # Plot histograms
38     X.hist(figsize=(20, 15))
39     plt.suptitle("Histograms of the Attributes", fontsize=20)
40     plt.show()
41
42     # Plot boxplots
43     X.boxplot(figsize=(6, 10))
44     plt.title("Boxplot of the Attributes")
45     plt.xticks(rotation=45)
46     plt.show()
47
48     # Plot pairwise scatterplots
49     sns.pairplot(X, corner=True)
50     plt.suptitle("Pairwise Scatterplots", fontsize=20)
51     plt.show()
52

```

```

53     # Heatmap of cross correlations
54     sns.heatmap(X.corr(numeric_only=False))
55     plt.title("Heatmap")
56     plt.show()
57
58
59 ##### Subpart 2 #####
60
61 def plot_outliers(outliers):
62     fig = plt.figure(figsize=(12, 6))
63     gs = gridspec.GridSpec(1,2)
64     ax = {}
65
66     ax[0] = fig.add_subplot(gs[0])
67     ax[0] = sns.scatterplot(x=range(0,440), y=outliers)
68     ax[0].set(xlabel = "Item", ylabel = "Soft-Min score", title = 'Soft
-Min scores (gamma = 1)')
69     ax[0].set_xticks(range(0,440,40))
70
71     ax[1] = fig.add_subplot(gs[1])
72     ax[1].set(ylabel = "Soft-Min score", title = 'Boxplot of Soft-Min
scores (gamma = 1)')
73     ax[1].set_xticks([])
74     ax[1] = plt.boxplot(outliers)
75     min, max = [item.get_ydata()[1] for item in ax[1]['whiskers']]
76
77     plt.show()
78     return min, max
79
80
81 ##### Subpart 3 #####
82
83 def distance_plots(mean_distance, var_distance, gamma_range, outliers
):
84     fig, axs = plt.subplots(2, 2, figsize=(14, 10))
85     axs = axs.flatten()
86
87     for i in range(len(mean_distance)):
88         if i in outliers:
89             axs[1].plot(gamma_range, mean_distance[i], linewidth=2, c
="blue")
90         else:
91             axs[1].plot(gamma_range, mean_distance[i], linewidth=0.3)
92
93     for i in range(len(var_distance)):
94         if i in outliers:
95             axs[0].plot(gamma_range, var_distance[i], linewidth=2, c=
"blue")
96         else:
97             axs[0].plot(gamma_range, var_distance[i], linewidth=0.3)
98
99     sns.scatterplot(x=gamma_range, y=np.var(mean_distance, axis=0),
ax=axs[3])

```

```

100
101     sns.scatterplot(x=gamma_range, y=np.mean(var_distance, axis=0),
102                     ax=axes[2])
103     axes[0].set_ylabel("MEAN of the relevance distances between
104                        components")
105     axes[1].set_ylabel("VARIANCE of the relevance distances between
106                        components")
107     #axes[2].set_ylabel("VARIANCE of the mean relevance distances
108                        between components")
109     axes[2].set_ylabel("Average MEAN of the relevance distances
110                        between components")
111     axes[3].set_ylabel("Average VARIANCE of the relevance distances
112                        between components")
113     for ax in axes:
114         ax.set_xlabel("Gamma")
115
116     plt.show()
117
118 def attribution_stat_plots(statistic, gamma_range, num_features,
119                           outliers, type):
120     fig = plt.figure(figsize=(14, 14))
121     gs = gridspec.GridSpec(3, 2)
122     ax = {}
123     for f_ind, feature in enumerate(num_features):
124         ax[f_ind] = fig.add_subplot(gs[f_ind])
125         ax[f_ind].set(xlabel="Gamma", ylabel="Bootstrap sample " +
126                      type + " relevance", title=feature)
127
128         for item in range(statistic.shape[0]):
129             if item in outliers:
130                 ax[f_ind] = plt.plot(gamma_range, statistic[item, :,
131                                f_ind], linewidth=2, c="blue")
132             else:
133                 ax[f_ind] = plt.plot(gamma_range, statistic[item, :,
134                                f_ind], linewidth=0.3)
135
136     plt.show()
137
138 def attribution_variance_means(mean, variance, gamma_range,
139                               num_features, outliers, inliers):
140     fig = plt.figure(figsize=(18, 18)) # (28,28)
141     gs = gridspec.GridSpec(2, 2)
142     ax = {}
143
144     ax[0] = fig.add_subplot(gs[0])
145     ax[0].set(xlabel="Gamma", ylabel="Average bootstrap sample MEAN
146          for INLIERS")
147     # alternative: "Variance of bootstrap sample mean for inliers"
148     for f_ind, feature in enumerate(num_features):
149         ax[0] = sns.scatterplot(x=gamma_range, y=np.var(mean[:, :,

```

```

140 f_ind][inliers], axis=0), legend=num_features)
141     ax[0].legend(labels=num_features)
142
143     ax[1] = fig.add_subplot(gs[1])
144     ax[1].set(xlabel="Gamma", ylabel="Average bootstrap sample MEAN
for OUTLIERS")
145     # alternative: "Variance of bootstrap sample mean for inliers"
146     for f_ind, feature in enumerate(num_features):
147         ax[1] = sns.scatterplot(x=gamma_range, y=np.var(mean[:, :,
f_ind][outliers], axis=0))
148     ax[1].legend(labels=num_features)
149
150     ax[2] = fig.add_subplot(gs[2])
151     ax[2].set(xlabel="Gamma", ylabel="Average bootstrap sample
VARIANCE for INLIERS")
152     for f_ind, feature in enumerate(num_features):
153         ax[2] = sns.scatterplot(x=gamma_range, y=np.mean(variance
[:, :, f_ind][inliers], axis=0), legend=num_features)
154     ax[2].legend(labels=num_features)
155
156     ax[3] = fig.add_subplot(gs[3])
157     ax[3].set(xlabel="Gamma", ylabel="Average bootstrap sample
VARIANCE for OUTLIERS")
158     for f_ind, feature in enumerate(num_features):
159         ax[3] = sns.scatterplot(x=gamma_range, y=np.mean(variance
[:, :, f_ind][outliers], axis=0))
160     ax[3].legend(labels=num_features)
161
162     plt.show()
163
164
165 def attribution_boxplots(statistic, stat_name):
166     fig = plt.figure(figsize=(14, 14))
167     gs = gridspec.GridSpec(2, 2)
168     ax = {}
169
170     ax[0] = fig.add_subplot(gs[0])
171     ax[0].set(ylabel= "Relevance " + stat_name + " over the Bootstrap
samples for OUTLIERS", title="Boxplot, gamma = 1")
172     ax[0] = statistic.boxplot()
173     ax[0].tick_params(axis='x', rotation=45)
174
175     ax[1] = fig.add_subplot(gs[1])
176     ax[1].set(title="Boxplot without outlier points")
177     ax[1].tick_params(axis='x', rotation=45)
178     ax[1] = sns.boxplot(data=statistic, showfliers=False)
179     plt.show()
180
181
182 ##### Subpart 4 #####
183
184 def silhouette_analysis(min_k, max_k, X, Umap = False):
185     """ Adapdet from:

```

```

186     https://scikit-learn.org/stable/auto\_examples/cluster/
187     plot\_kmeans\_silhouette\_analysis.html#sphx-glr-auto-examples-
188     cluster-plot-kmeans-silhouette-analysis-py
189     mapper_type="UMAP" or "t-SNE""
190     if Umap:
191         mapper = umap.UMAP(n_neighbors=4, min_dist=1, n_components=2
192 , metric='euclidean', random_state=42)
193     else:
194         mapper = TSNE(n_components=2, learning_rate='auto', init='pca
195 ', perplexity=30)
196     X_embedded = mapper.fit_transform(X)
197
198     range_n_clusters = np.arange(min_k, max_k+1)
199
200     for n_clusters in range_n_clusters:
201         # Create a subplot with 1 row and 2 columns
202         fig, (ax1, ax2) = plt.subplots(1, 2)
203         fig.set_size_inches(18, 7)
204
205         # The 1st subplot is the silhouette plot
206         # The silhouette coefficient can range from -1, 1 but in this
207         # example all
208         # lie within [-0.1, 1]
209         ax1.set_xlim([-0.1, 1])
210         # The (n_clusters+1)*10 is for inserting blank space between
211         # silhouette
212         # plots of individual clusters, to demarcate them clearly.
213         ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])
214
215         # Initialize the clusterer with n_clusters value and a random
216         # generator
217         # seed of 0 for reproducibility.
218
219         clusterer = KMeans(n_clusters=n_clusters, init='k-means++',
220 n_init = 100, random_state=42)
221         cluster_labels = clusterer.fit_predict(X)
222
223         # The silhouette_score gives the average value for all the
224         # samples.
225         # This gives a perspective into the density and separation of
226         # the formed
227         # clusters
228         silhouette_avg = silhouette_score(X, cluster_labels)
229         print("For n_clusters =", n_clusters,
230               "The average silhouette_score is :", silhouette_avg)
231
232         # Compute the silhouette scores for each sample
233         sample_silhouette_values = silhouette_samples(X,
234 cluster_labels)
235
236         y_lower = 10
237         for i in range(n_clusters):

```

```

229         # Aggregate the silhouette scores for samples belonging
    to
230         # cluster i, and sort them
231         ith_cluster_silhouette_values = sample_silhouette_values[
cluster_labels == i]
232
233         ith_cluster_silhouette_values.sort()
234
235         size_cluster_i = ith_cluster_silhouette_values.shape[0]
236         y_upper = y_lower + size_cluster_i
237
238         color = cm.nipy_spectral(float(i) / n_clusters)
239         ax1.fill_betweenx(np.arange(y_lower, y_upper), 0,
ith_cluster_silhouette_values,
240                             facecolor=color, edgecolor=color, alpha
=0.7,
241                             )
242
243         # Label the silhouette plots with their cluster numbers
    at the middle
244         ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
245
246         # Compute the new y_lower for next plot
247         y_lower = y_upper + 10 # 10 for the 0 samples
248
249         ax1.set_title("The silhouette plot for the various
clusters.")
250         ax1.set_xlabel("The silhouette coefficient values")
251         ax1.set_ylabel("Cluster label")
252
253         # The vertical line for average silhouette score of all
    the values
254         ax1.axvline(x=silhouette_avg, color="red", linestyle="--"
)
255
256         ax1.set_yticks([]) # Clear the yaxis labels / ticks
257         ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
258
259         # 2nd Plot showing the actual clusters formed on within a
    2D embedding of choice
260         colors = cm.nipy_spectral(cluster_labels.astype(float) /
n_clusters)
261         ax2.scatter(
262             X_embedded[:,1], X_embedded[:,0], marker=".", s=30,
lw=0, alpha=0.7, c=colors, edgecolor="k"
263         )
264         ax2.set_title("Data Embedding")
265         ax2.set_xlabel("Embedded dimension 2")
266         ax2.set_ylabel("Embedded dimension 1")
267
268         if Umap:
269             # Labeling the clusters
270             centers = clusterer.cluster_centers_

```

```

271         # U-Map transformation of the cluster centers
272         centers_mapped = mapper.transform(centers)
273         # Draw white circles at cluster centers
274         ax2.scatter(centers_mapped[:,1], centers_mapped[:,0]
275 ], marker="o", c="white", alpha=1, s=200,
276                     edgecolor="k",
277 )
278         for i, c in enumerate(centers_mapped):
279             ax2.scatter(c[1], c[0], marker="$d$" % i, alpha=
280 1, s=50, edgecolor="k")
281
282     plt.suptitle(
283         "Silhouette analysis for KMeans clustering on sample
284 data with n_clusters = %d"
285         % n_clusters, fontsize=14, fontweight="bold",
286     )
287
288     plt.show()
289
290 def visualise_kmeans(X, clusterer):
291     # Compute UMAP for the visualisation
292     umapper = umap.UMAP(n_neighbors=4, min_dist=1, n_components=2,
293 metric='euclidean', random_state=42)
294     X_embedded_umap = umapper.fit_transform(X)
295     # Compute t-SNE for the visualisation
296     mapper = TSNE(n_components=2, learning_rate='auto', init='pca',
297 perplexity=30)
298     X_embedded = mapper.fit_transform(X)
299
300     fig = plt.figure(figsize=(18, 9))
301     gs = gridspec.GridSpec(1, 2)
302     ax = {}
303
304     # UMAP
305     ax[0] = fig.add_subplot(gs[0])
306     ax[0] = sns.scatterplot(x=X_embedded_umap[:, 1], y=
307 X_embedded_umap[:, 0], hue=clusterer.labels_, palette="tab10")
308     ax[0].set_xlabel("UMAP dimension 2")
309     ax[0].set_ylabel("UMAP dimension 1")
310     ax[0].set_title('UMAP')
311
312     # Labeling the clusters
313     centers = clusterer.cluster_centers_
314     # U-Map transformation of the cluster centers
315     centers_mapped = umapper.transform(centers)
316     # Draw white circles at cluster centers
317     ax[0].scatter(centers_mapped[:, 1], centers_mapped[:, 0], marker=
318 "o", c="white", s=200, edgecolor="k")
319     for i, c in enumerate(centers_mapped):
320         ax[0].scatter(c[1], c[0], marker="$d$" % i, c="white", s=50
321 , edgecolor="k")

```

```

316
317     # t-SNE
318     ax[1] = fig.add_subplot(gs[1])
319     ax[1] = sns.scatterplot(x=X_embedded[:, 1], y=X_embedded[:, 0],
hue=clusterer.labels_, palette="tab10")
320     ax[1].set_xlabel("t-SNE dimension 2")
321     ax[1].set_ylabel("t-SNE dimension 1")
322     ax[1].set_title('t_SNE')
323
324     plt.show()
325
326
327 def clusters_stats(df, clustering):
328     def means(df, clustering):
329         plots_per_row = 2
330         plots_per_column = math.ceil(len(set(clustering.labels_)) /
plots_per_row)
331         fig = plt.figure(figsize=(14 * plots_per_row, 4 *
plots_per_column)) # (28,28)
332         gs = gridspec.GridSpec(plots_per_column, plots_per_row)
333         ax = {}
334
335         for cluster in set(clustering.labels_):
336             filter = clustering.labels_ == cluster
337
338             n = sum(clustering.labels_ == cluster)
339             mean = df[filter].mean()
340             sd = df[filter].std()
341
342             ax[cluster] = fig.add_subplot(gs[cluster])
343             ax[cluster] = plt.errorbar(df.columns.values.tolist(),
mean, sd, marker='o', linestyle='None')
344             plt.xlabel('Attribute')
345             plt.ylabel('Mean +/- SD')
346             plt.title('Cluster ' + str(cluster) + ", size: " + str(n
))
347             plt.xticks(rotation=45)
348
349             fig.suptitle("Statistics of individual features for the
clusters: Mean and Standard Deviation", fontsize=20)
350             plt.show()
351
352     def boxplots(df, clustering):
353         K = len(set(clustering.labels_))
354         fig, ax = plt.subplots(nrows=1, ncols=6, sharex=True, sharey=
True, figsize=(28, K))
355         cluster_ind = []
356
357         for i, feature in enumerate(df.columns):
358             sns.boxplot(ax=ax[i], data=df, x=feature, y=clustering.
labels_, orient="h", palette="tab10")
359             cluster_ind.append("Cluster " + str(i))
360             plt.yticks(range(0, K), cluster_ind)

```



```
361         fig.suptitle("Statistics of individual features for the  
clusters: Boxplots", fontsize=20)  
362  
363         plt.show()  
364  
365     means(df, clustering)  
366     boxplots(df, clustering)  
367
```