

DAT565/DIT407 Assignment 5

Henning Anevik

Joar Forsberg

2024-10-09

Problem 1: Preprocessing the dataset

We tried both standard and min-max normalization. We saw min-max normalization made the clusters a bit more apparent when plotting pairs of features in second task. And so, in the end, we opted for min-max scaling B. This was easily done by calling on functions from the sklearn library.

Problem 2: Determining the appropriate number of clusters

From the inspecting the graph we see that three is an appropriate number of clusters, which can be seen in Figure 1. This is known as the "elbow method", which we were taught about in the lectures C.

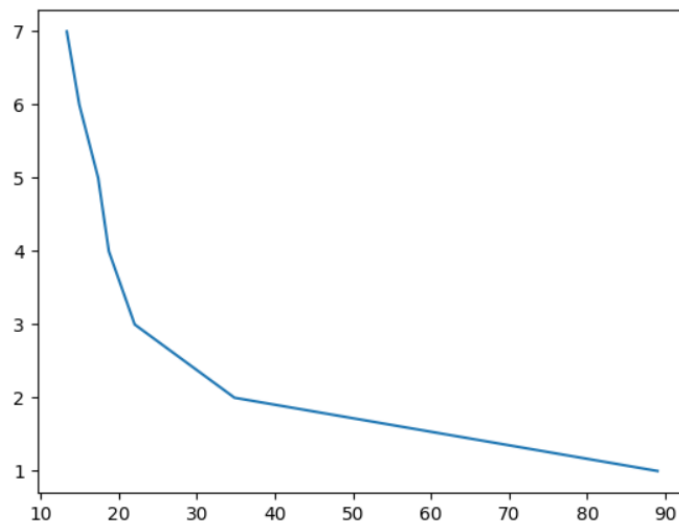


Figure 1: Elbow method for appropriate number of clusters

Problem 3: Visualizing the classes

a)

We see many plots of pairs of features result in straight lines. This makes a lot of sense since many features are highly correlated. For example the plot of kernel length vs. kernel groove length is almost a straight line, which makes a lot of sense seen in Figure 2.

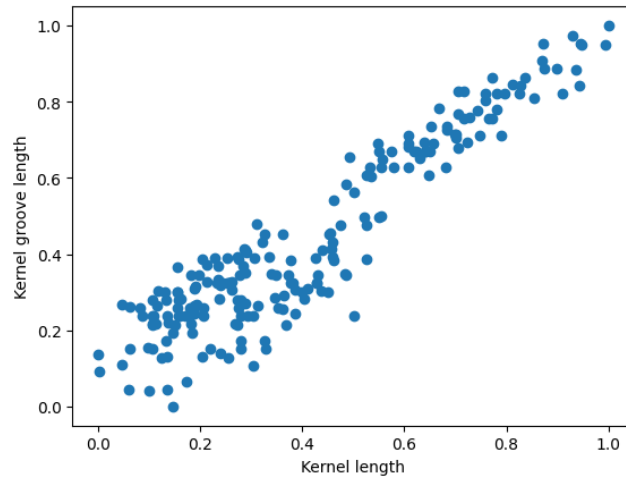


Figure 2: kernel length vs. kernel groove length

More interesting is are the graphs in which we can make out some semblance of clustering - where we are able to, somewhat, make out the different classes. The graph which displays the highest degree of clustering is the one of compactness and kernel groove length seen in Figure 3. We can see a pretty clear delineation between two of the classes and a less distinct one of the third class D.

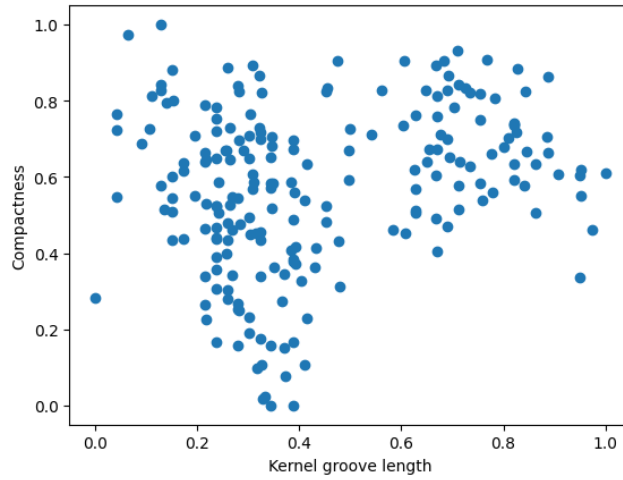


Figure 3: Compactness vs Kernel groove length

b)

The data seem to have a linear relationship in the two-dimensional plot. The Figure 4 suggests that the data might contain patterns in the random projection. But since Gaussian Random Projection involves randomness, repeating this can result in slightly different results due to the randomness, which was stated in the assignment E.

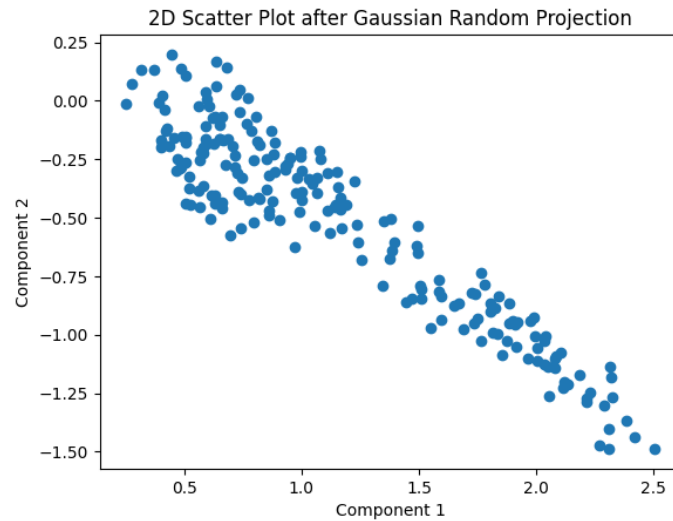


Figure 4: 2D Scatter Plot after Gaussian Random Projection

c)

The code for c) can be seen in F.

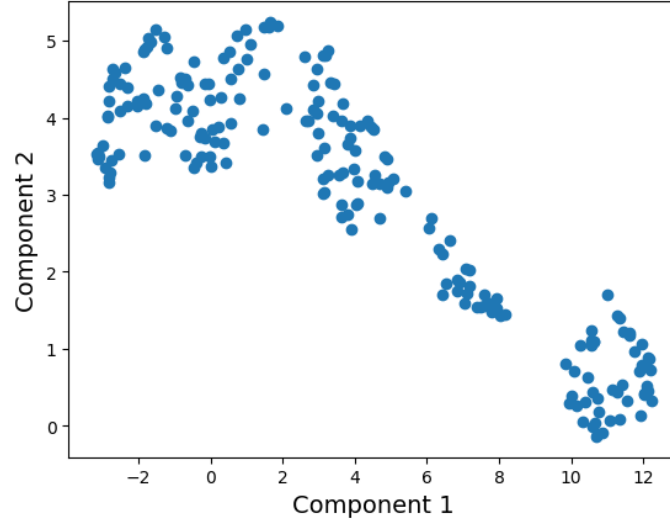


Figure 5: UMAP Riemannian manifold 2D

d)

It's not hard to imagine lines being drawn that separates the datapoints into three different sets. This is especially clear when looking at the graph we obtained by applying UMAP which displays in Figure 5 a high degree of delineation between all three clusters and in which the imagined lines become more apparent.

Problem 4: Evaluating clustering

Running G yield the following. A Rand index of 0.8641604010025062 and the highest accuracy and its associated mapping is (0.8857142857142857, 0: 2, 1: 3, 2: 1)

Problem 5: Agglomerative clustering

We see that the average linkage method works best H. This is a more "robust" linkage method since it considers all the points in a pair of clusters and not just the closest points, unlike single or complete linkage.

We also have well-behaved data; the classes are all roughly the same size and there are no particularly outlying datapoints that will have a disproportionate impact on the average. A plot a dendrogram of the hierarchical clustering can be seen in 6.

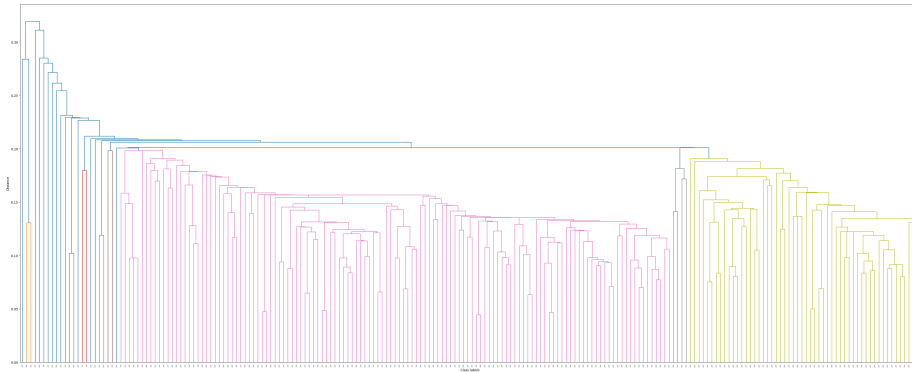


Figure 6: Plot of a dendrogram of the hierarchical clusterin

Its a bit hard to see in the Figure 6 but it showcase three different clusters one cluster to the right, one in the middle and lastly a cluster to the left.

A Libraries

```

1 !pip install umap-learn
2 import pandas as pd
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.metrics import mean_squared_error,
  ↳ rand_score
5 from sklearn.cluster import KMeans
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import itertools
9 from sklearn.random_projection import
  ↳ GaussianRandomProjection
10 import umap
11 from sklearn.preprocessing import StandardScaler,
  ↳ MinMaxScaler
12 from scipy.cluster.hierarchy import dendrogram,
  ↳ linkage
13 from sklearn.cluster import AgglomerativeClustering
14
15 num_clusters = 10

```

B Preprocessing the dataset

```

1 df = pd.read_csv("/content/seeds.tsv", delimiter='\t')
2 mapper = {'15.26' : "Area", # 1
3           '14.84' : "Perimeter", # 2
4           '0.871' : "Compactness", # 3

```

```

5         '5.763' : "Kernel_length", # 4
6         '3.312' : "Kernel_width", # 5
7         '2.221' : "Asymmetry_coefficient", # 6
8         '5.22'  : "Kernel_groove_length", # 7
9         '1'     : "Numerical_class_label"} # 8
10
11 first_row = pd.DataFrame({v: [k] for k, v in
    ↪ mapper.items()})
12
13 df = df.rename(columns=mapper)
14 df = pd.concat([first_row, df], ignore_index=True) #
    ↪ [209 rows x 8 columns]
15 df = df.astype("float")
16 print(df)
17
18 scaler = MinMaxScaler()
19 df_features = df.drop("Numerical_class_label", axis=1)
20 class_labels = df["Numerical_class_label"]
21 normalized_df =
    ↪ pd.DataFrame(scaler.fit_transform(df_features),
    ↪ columns=df_features.columns)
22 df = pd.concat([normalized_df, class_labels], axis=1)
23 print("\n")
24 print(df)

```

C Determining the appropriate number of clusters

```

1 features = df.iloc[:, :-1].to_numpy()
2 inertias = []
3 for n in range(7):
4     kmeans = KMeans(n_clusters=n + 1,
    ↪ random_state=0).fit(features)
5     inertias.append(kmeans.inertia_)
6
7 plt.plot(inertias, range(1, 8))
8 plt.show()

```

D Problem 3 a)

```

1 features = df.columns[:-1]
2 combinations = list(itertools.product(features,
    ↪ features))
3 cleaned_combinations = [tup for tup in combinations
    ↪ if tup[0] != tup[1]]
4

```

```

5 for pair in cleaned_combinations:
6     plt.scatter(df[pair[0]], df[pair[1]])
7     plt.xlabel(pair[0])
8     plt.ylabel(pair[1])
9     plt.show()

```

E Problem 3 b)

```

1 features = df.iloc[:, :-1].to_numpy()
2 reducer = GaussianRandomProjection(n_components=2,
    ↪ random_state=42)
3 embedding = reducer.fit_transform(features)
4 print(len(features[0])) # 7 x 210
5
6 plt.scatter(embedding[:, 0], embedding[:, 1])
7 plt.xlabel('Component_1')
8 plt.ylabel('Component_2')
9 plt.title('2D Scatter Plot after Gaussian Random
    ↪ Projection')
10 plt.show()

```

F Problem 3 c)

```

1 reducer = umap.UMAP(random_state=42)
2 embedding = reducer.fit_transform(features)
3
4 plt.scatter(embedding[:, 0], embedding[:, 1])
5 plt.xlabel('Component_1', fontsize=14)
6 plt.ylabel('Component_2', fontsize=14)
7 plt.show()

```

G Evaluating clustering

```

1 features = df.iloc[:, :-1].to_numpy()
2 kmeans = KMeans(n_clusters=3,
    ↪ random_state=0).fit(features)
3 print("Rand_index:", rand_score(df["Numerical_class_
    ↪ label"], kmeans.labels_))
4
5 unique_mappings = [dict(zip([0, 1, 2], perm)) for
    ↪ perm in itertools.permutations([1, 2, 3])] #
    ↪ clusters are keys
6 #print(unique_mappings)
7
8 accuracy_list = []

```

```

9  for mapping in unique_mappings: # {0: 2, 1: 3, 2: 1}
10     sum = 0
11     c1_class = mapping[0]          # {2 : 0, 3 : 1, 1 : 2}
12     c2_class = mapping[1]
13     c3_class = mapping[2]
14     for i in range(210):
15         if kmeans.labels_[i] == 0:
16             if int(df.iloc[i, -1]) == c1_class:
17                 sum = sum + 1
18         if kmeans.labels_[i] == 1:
19             if int(df.iloc[i, -1]) == c2_class:
20                 sum = sum + 1
21         if kmeans.labels_[i] == 2:
22             if int(df.iloc[i, -1]) == c3_class:
23                 sum = sum + 1
24     accuracy_list.append((float(sum/210), mapping))
25
26 max_accuracy = max(accuracy_list, key=lambda x: x[0])
27 print("The highest accuracy and its associated
      ↪ mapping:", max_accuracy)

```

H Agglomerative clustering

```

1  def accuracy(model):
2      unique_mappings = [dict(zip([0, 1, 2], perm)) for
      ↪ perm in itertools.permutations([1, 2, 3])] #
      ↪ clusters are keys
3
4      accuracy_list = []
5      for mapping in unique_mappings:
6          sum = 0
7          c1_class = mapping[0]
8          c2_class = mapping[1]
9          c3_class = mapping[2]
10         for i in range(210):
11             if model.labels_[i] == 0:
12                 if int(df.iloc[i, -1]) == c1_class:
13                     sum = sum + 1
14             if model.labels_[i] == 1:
15                 if int(df.iloc[i, -1]) == c2_class:
16                     sum = sum + 1
17             if model.labels_[i] == 2:
18                 if int(df.iloc[i, -1]) == c3_class:
19                     sum = sum + 1
20
21         accuracy_list.append((float(sum/210), mapping))
22     return max(accuracy_list, key=lambda x: x[0])
23 features = df.iloc[:, :-1].to_numpy()

```



```

24 clustering_ward =
    ↪ AgglomerativeClustering(n_clusters=3,
    ↪ linkage="ward").fit(features)
25 clustering_complete =
    ↪ AgglomerativeClustering(n_clusters=3,
    ↪ linkage="complete").fit(features)
26 clustering_average =
    ↪ AgglomerativeClustering(n_clusters=3,
    ↪ linkage="average").fit(features)
27 clustering_single =
    ↪ AgglomerativeClustering(n_clusters=3,
    ↪ linkage="single").fit(features)
28
29 clustering_ward_acc = accuracy(clustering_ward)
30 clustering_complete_acc =
    ↪ accuracy(clustering_complete)
31 clustering_average_acc = accuracy(clustering_average)
32 clustering_single_acc = accuracy(clustering_single)
33
34 print("Ward_acc:", clustering_ward_acc)
35 print("Complete_acc:", clustering_complete_acc)
36 print("Average_acc:", clustering_average_acc) # {0:
    ↪ 1, 1: 2, 2: 3})
37 print("Single_acc:", clustering_single_acc) # The
    ↪ highest accuracy and its associated mapping:
    ↪ (0.8857142857142857, {0: 2, 1: 3, 2: 1})
38
39 # Plotting dendrogram for task 5
40
41 linkage_matrix = linkage(features)
42 fig = plt.figure(figsize=(50, 20))
43 labels = df["Numerical_class_
    ↪ label"].astype(str).tolist()
44 dn = dendrogram(linkage_matrix, labels=labels,
    ↪ color_threshold=0.2)
45 plt.xlabel("Class_labels")
46 plt.ylabel("Distance")
47 plt.show()

```