

DAT565/DIT407 Assignment 5

Saif Sayed
gussayedfa@student.gu.se

Gona Ibrahim Abdulrahman
gusibrigo@student.gu.se

2024-05-02

This is a report for assignment 5 for the course *Introduction to Data Science & AI* from Chalmers and Gothenburg University.

Problem 1: Loading and Visualizing the classes in the data

Firstly, we read the dataset into Pandas, providing the appropriate delimiter (`\t`) and specifying that there is no header in the file. We then add the column names as per the table in the problem description.

(A) In task A, we create scatter plots between each pair of features using nested loops. For each pair of features, we iterate over the unique class labels and scatter the points with the corresponding class label color. We label the x-axis and y-axis with the feature names, add a legend, and provide a title for each plot.

Some scatter plots exhibit a discernible pattern where the data points appear to align in a structured manner, while others display a more scattered distribution. For example, when examining the scatter plot between perimeter and area, the points form a linear pattern, suggesting a strong correlation between the two variables. This indicates that, regardless of the species, the area of the seed is dependent on its perimeter.

The scatter plots incorporate color coding, allowing us to distinguish between different species based on the legend. Notably, the points corresponding to each species tend to cluster together in distinct regions of the graph, effectively separating them from points belonging to other species.

One particularly intriguing scatter plot would be either the scatter plot between Perimeter and Area or Length of Kernel and Perimeter. These plots demonstrate a clear distinction between species, as the points representing each species remain separate and do not overlap with those of other species. This distinct separation of data points by species makes these plots particularly compelling for visualizing the species classification (see **Figure 1**).

(B) We use the Gaussian random projection class from `scikit-learn` to per-

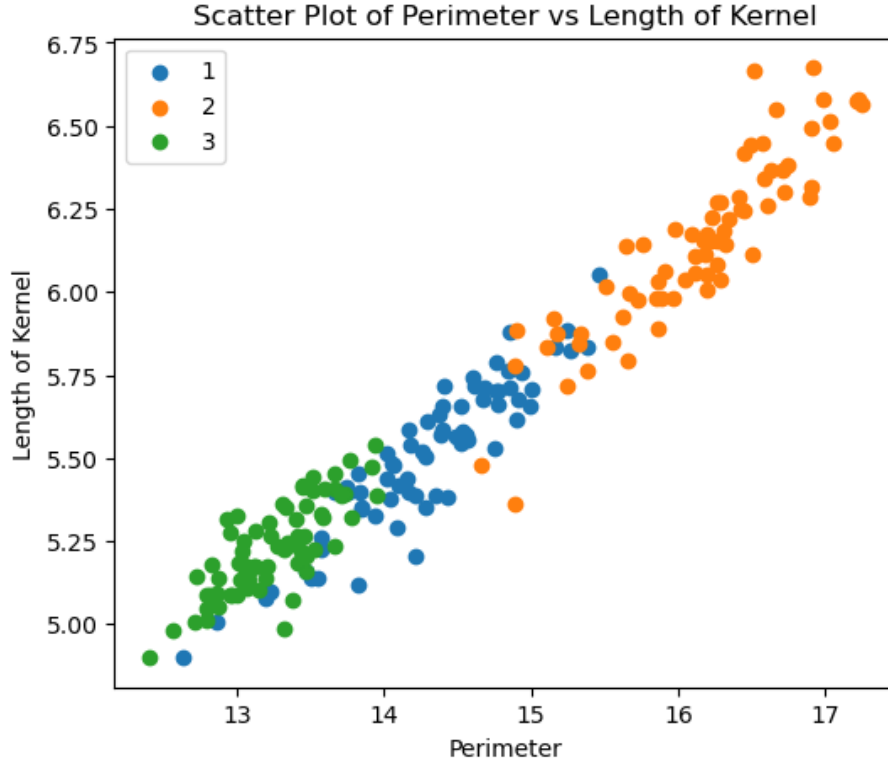


Figure 1: Length of Kernel Vs. Perimeter

form dimensionality reduction and visualize the data. After projecting the data, we create a scatter plot of the projected data. Each point in the scatter plot represents a wheat kernel, and the color of each point represents the class label of the wheat kernel. The scatter plot remains discernible (see **Figure 2**).

(C) We apply UMAP to project the data onto a two-dimensional space using the UMAP class. Then, we create a scatter plot of the projected data, where each point represents a wheat kernel. The color of each point corresponds to the class label of the wheat kernel (see **Figure 3**).

(D) We have observed that not all scatter plots demonstrate a linear relationship. The plots that do exhibit linearity involve pairs of features that are correlated, such as Length of Kernel, Width of Kernel, Perimeter, and Area. These features are directly correlated with each other and display a linear relationship. As each species possesses distinct values for these properties, the plots reveal linearity and distinguishability among seed species. Consequently, clustering algorithms can effectively classify the species based on these specific features. However, for other features that lack separability and exhibit overlapping clusters, the task of classifying the seeds becomes more challenging.

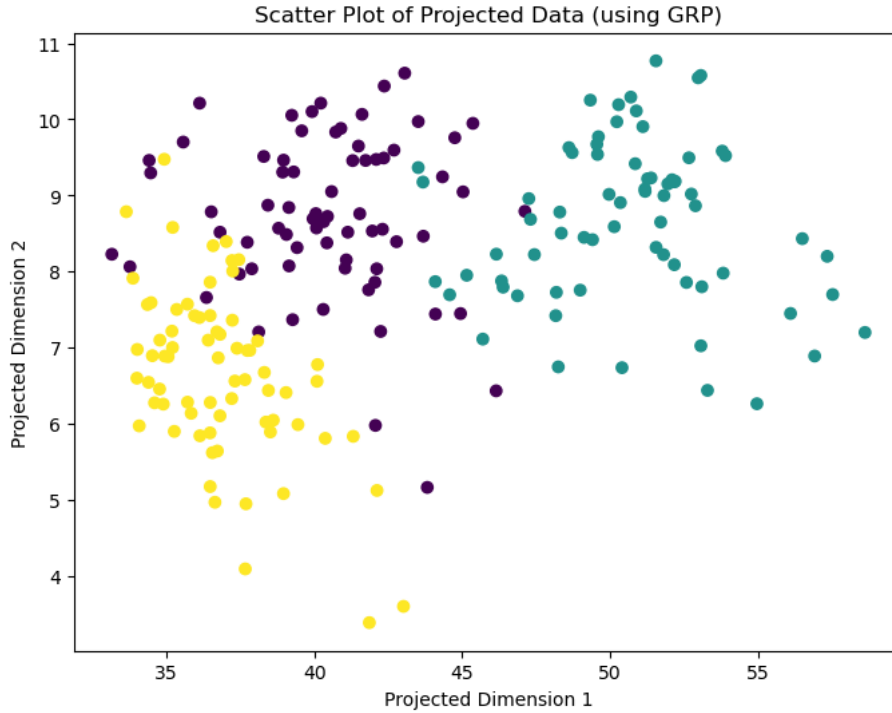


Figure 2: Gaussian Random Projection

Our source code can be found in Appendix A of this document.

Problem 2: Unsupervised Learning: Determining the appropriate number of clusters

(A) To begin, we read the dataset into a pandas DataFrame. Then, we remove the class label column from the DataFrame, as it is not needed for the clustering task. Next, we apply normalization to the remaining features. To normalize the dataset, we use the StandardScaler class from the scikit-learn library. This scaler standardizes the features by subtracting the mean and dividing by the standard deviation. This process transforms the features to have a mean of zero and a standard deviation of one, effectively bringing them to a comparable scale. Finally, we converted the normalized data back to a pandas Dataframe.

(B) After applying normalization, we proceed with the clustering task. We use the K-means algorithm for this purpose. K-means is an iterative algorithm that aims to partition the data into K clusters, where K is a predefined number. In our case, we explore different values of K to assess the optimal number of clusters for the dataset. During the clustering process, we calculate the inertia for each value of K. Inertia represents the sum of squared distances of samples to their closest cluster center. It is a measure of how well the data points within

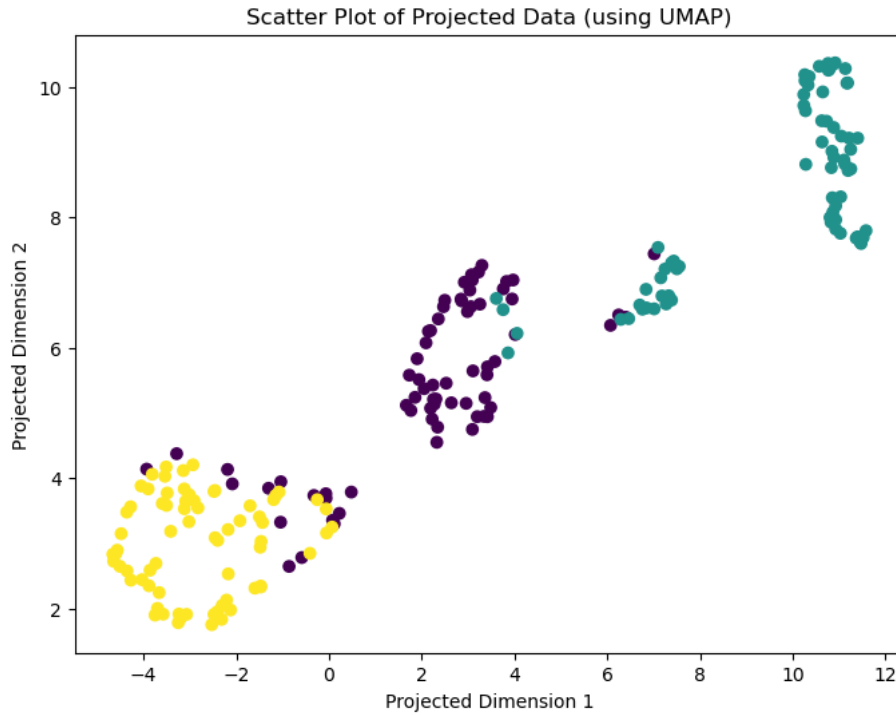


Figure 3: UMAP Projection

each cluster are grouped together. Finally, we plot the inertia values as a function of K to analyze the results (see **Figure 4**). By visually examining the graph, we can observe a noticeable change in the inertia values starting from $K = 3$, which also corresponds to the elbow point. Based on this observation, we can determine that the appropriate number of clusters is 3.

Problem 3: Evaluating clustering

In this task, we evaluate the quality of the clustering results obtained using k-means by computing the Rand Index and accuracy. But this time, we applied k-means to our normalised data so that k matches the number of underlying species of wheat.

Compute Rand Index: The `adjusted_rand_score` function from `scikit-learn` is used to calculate the Rand Index between the cluster labels (`cluster_labels`) obtained from k-means clustering and the original class labels (`class_labels`). The resulting value of Rand Index is 0.77 (2sf).

Compute accuracy: To compute the accuracy of the clustering, our code finds the permutation based on the number of all unique labels. It iterates through each permutation and calculates the accuracy score by assigning the permuted labels to the cluster labels. If the accuracy is higher than the current

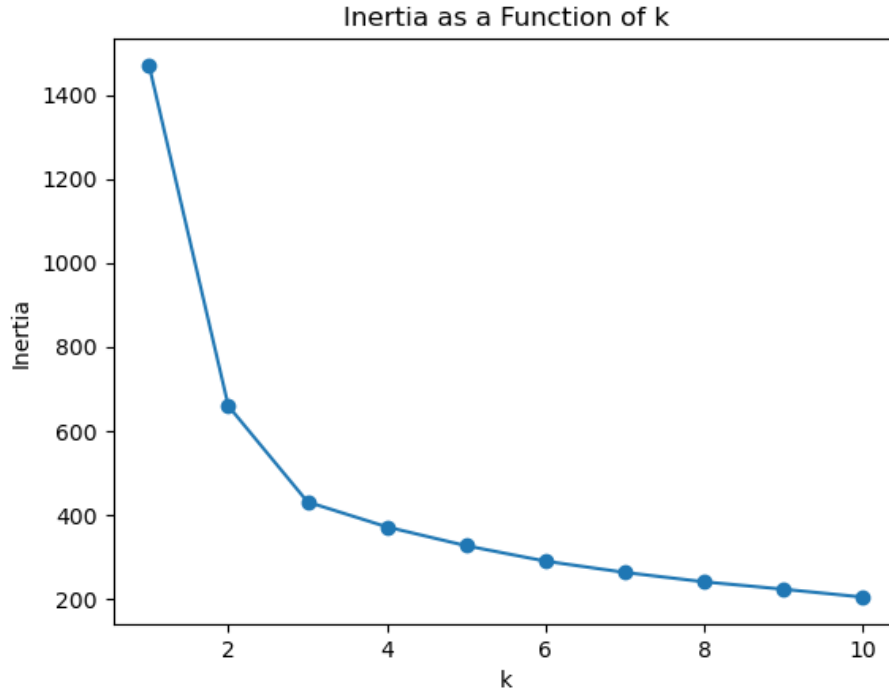


Figure 4: Inertia as a function of k.

best accuracy, the variable `best_accuracy` is updated. The resulting accuracy value is 0.92 (2sf).

Problem 4: Supervised Learning: K-nearest neighbour

(A) When choosing a value, our approach would be to start with a small value and gradually increase it. While doing so, we also plan to calculate the accuracy score against the test sets for several expected values of the parameter `n_neighbour`. Additionally, plotting a graph to visualize this would help us to compare the accuracy score at different values. Our goal is to find a balance between underfitting (too small `n_neighbors`) and overfitting (too large `n_neighbors`).

(B) Our chosen value for `n_neighbour` is 5. From our plot, we can determine that at `n_neighbour = 5` the accuracy score complies with both the training set and testing set (see Figure 5). We also experimented for values up to 150, although we got similar accuracy scores, we decided to not use a higher value to avoid overfitting our model. The code is included in the Appendix A of this document.

(C) For our chosen value, the accuracy on the test set is 0.93 (2sf).

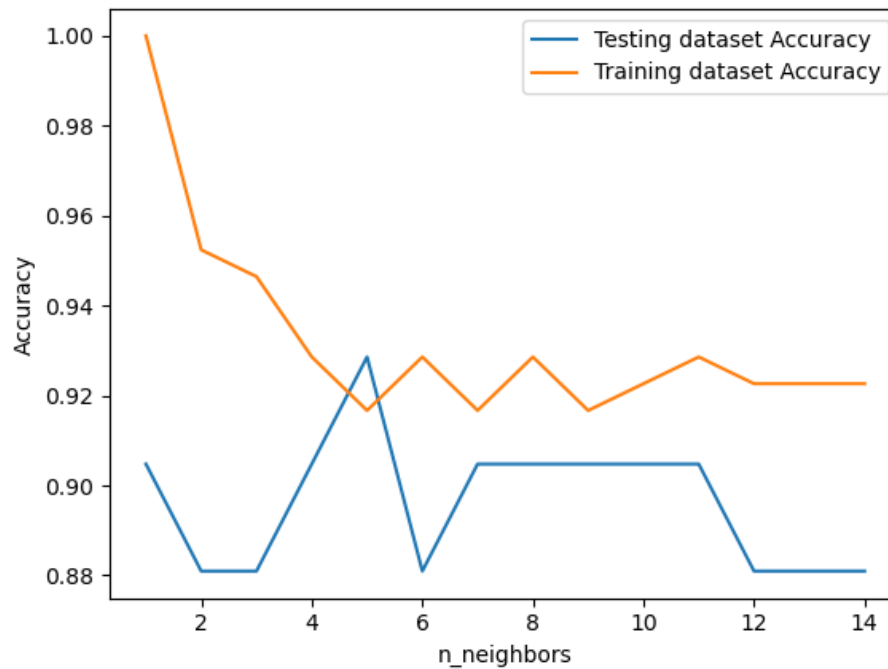


Figure 5: Accuracy Vs. n_neighbors

(D) Using a different split for train-test would result in different accuracy values, which can vary from the results obtained with the original split. The randomness in the train-test split introduces variability in the composition of the training and test sets. This variability can lead to different outcomes, particularly if the new training split is imbalanced, meaning it contains a higher proportion of data from a specific seed species. In such cases, the model may become biased towards the overrepresented seed species, resulting in lower accuracy when evaluated on the test set.

A Python code

```

1  """
2  Assignment 5
3  """
4  import numpy as np
5  import pandas as pd
6  import matplotlib.pyplot as plt
7  from sklearn.random_projection import
    GaussianRandomProjection
8  from umap import UMAP

```

```

9  from sklearn.cluster import KMeans
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.metrics import adjusted_rand_score
12 import itertools
13 from sklearn.metrics import accuracy_score
14 from sklearn.model_selection import train_test_split
15 from sklearn.neighbors import KNeighborsClassifier
16 """
17 1. Load and Visualizing the classes in the data
18 """
19 data = pd.read_csv("seeds.tsv", delimiter="\t", header
20                    =None)
21 column_names = [
22     "Area",
23     "Perimeter",
24     "Compactness",
25     "Length_of_Kernel",
26     "Width_of_Kernel",
27     "Asymmetry_Coefficient",
28     "Length_of_Kernel_Groove",
29     "Class_Label"
30 ]
31
32 data.columns = column_names
33 print(data.columns)
34
35 # a.
36 feature_columns = data.columns[:-1] # Exclude the
37                                     # last column (Class Label)
38 class_labels = data["Class_Label"].unique()
39
40 for i in range(len(feature_columns)):
41     for j in range(i + 1, len(feature_columns)):
42         plt.figure(figsize=(6, 5))
43         for label in class_labels:
44             subset = data[data["Class_Label"] == label
45                             ]
46             plt.scatter(subset[feature_columns[i]],
47                         subset[feature_columns[j]], label=label
48                         )
49             plt.xlabel(feature_columns[i])
50             plt.ylabel(feature_columns[j])
51             plt.legend()
52             plt.title(f"Scatter_Plot_of_{feature_columns[i]}
53                       _vs_{feature_columns[j]}")
54             plt.show()
55
56 # b.
57 grp = GaussianRandomProjection(n_components=2)

```

```

52 projected_data = grp.fit_transform(data.iloc[:, :-1])
    # Exclude the last column (Class Label)
53 plt.figure(figsize=(8, 6))
54 plt.scatter(projected_data[:, 0], projected_data[:,
    1], c=data["Class_Label"])
55 plt.xlabel("Projected_Dimension_1")
56 plt.ylabel("Projected_Dimension_2")
57 plt.title("Scatter_Plot_of_Projected_Data_(using_GRP)"
    )
58 plt.show()
59 # c.
60 umap = UMAP(n_components=2)
61 projected_data = umap.fit_transform(data.iloc[:, :-1])
    # Exclude the last column (Class Label)
62 # Plot a scatter plot of the projected data
63 plt.figure(figsize=(8, 6))
64 plt.scatter(projected_data[:, 0], projected_data[:,
    1], c=data["Class_Label"])
65 plt.xlabel("Projected_Dimension_1")
66 plt.ylabel("Projected_Dimension_2")
67 plt.title("Scatter_Plot_of_Projected_Data_(using_UMAP)"
    )
68 plt.show()
69 """
70 2. Unsupervised Learning: Determining the appropriate
    number of clusters
71 """
72 # a.
73 feature_columns = "Class_Label"
74 class_labels = data[feature_columns]
75 data = data.drop(feature_columns, axis=1)
76 scaler = StandardScaler()
77 normalized_data = scaler.fit_transform(data)
78 normalized_data = pd.DataFrame(normalized_data,
    columns=data.columns)
79 # b.
80 inertia_values = []
81 k_values = range(1, 11) # Setting the maximum value
    of k as 10, but we can adjust it as needed
82
83 for k in k_values:
84     kmeans = KMeans(n_clusters=k, n_init=10,
        random_state=42) # Explicitly setting n_init
        to 10
85     kmeans.fit(normalized_data)
86     inertia_values.append(kmeans.inertia_)
87     cluster_labels = kmeans.labels_
88
89 plt.plot(k_values, inertia_values, marker='o')
90 plt.xlabel('k')

```



```

91 plt.ylabel('Inertia')
92 plt.title('Inertia as a Function of k')
93 plt.show()
94 '''
95 3. Evaluating clustering
96 '''
97 class_labels = data["Class_Label"]
98 kmeans = KMeans(n_clusters=len(class_labels.unique()),
99                 n_init=10, random_state=42)
100 kmeans.fit(normalized_data)
101 cluster_labels = kmeans.labels_
102 rand_index = adjusted_rand_score(class_labels,
103                                  cluster_labels)
104 print("Rand_index:", rand_index)
105 all_labels = pd.Series(cluster_labels).append(
106                      class_labels)
107 perms = itertools.permutations(range(len(all_labels.
108                                     unique()))))
109 best_accuracy = 0
110 for perm in perms:
111     perm_labels = [perm[label] for label in
112                   cluster_labels]
113     accuracy = accuracy_score(perm_labels,
114                               class_labels)
115     if accuracy > best_accuracy:
116         best_accuracy = accuracy
117         best_permutation = perm
118 print("Best_Accuracy:", best_accuracy)
119 print("Best_Permutation:", best_permutation)
120 """
121 4. Supervised Learning: K-nearest neighbour
122 """
123 X = data.iloc[:, :-1]
124 X = X.values
125 y = data.iloc[:, -1]
126 y = y.values
127
128 X_train, X_test, y_train, y_test = train_test_split(X,
129                                                       y, test_size=0.2, random_state=42)
130
131 neighbors = np.arange(1, 15)
132 test_accuracy = np.empty(len(neighbors))
133 train_accuracy = np.empty(len(neighbors))
134
135 for i, k in enumerate(neighbors):
136     knn = KNeighborsClassifier(n_neighbors=k)

```

```
134     knn.fit(X_train, y_train)
135     test_accuracy[i] = knn.score(X_test, y_test)
136     train_accuracy[i] = knn.score(X_train, y_train)
137
138 plt.plot(neighbors, test_accuracy, label = 'Testing_
dataset_Accuracy')
139 plt.plot(neighbors, train_accuracy, label = 'Training_
dataset_Accuracy')
140 plt.legend()
141 plt.xlabel('n_neighbors')
142 plt.ylabel('Accuracy')
143 plt.show()
```