# DIT407 Introduction to data science and AI Assignment 5

Reynir Siik
reynir.siik@gmail.com

Franco Zambon
guszamfr@student.gu.se

2024-02-25

## Abstract

This text is from the file `abstract.tex` and will be included in the abstract.

## 1  Problem 1: Preprocessing the dataset

The data in the file have different metrics and units. To be able to do a meaningful clustering the data needs to be normalized. A method that distorts the shape of the data as little as possible is the Z-score method. The z-score method is chosen on those premises.

## 2  Problem 2: Determining the appropriate number of clusters

With k-means clustering and calculation of inertia the elbow method is used to find the number of clusters.1 While there is a big difference between 1-2 and 2-3 clusters, there are almost equal distance between the following clusters. The number of clusters the model predicts seems to be three.

When applying linear regression to the "lower arm" part of the elbow diagram and calculating Mean Squared Error (MSE) for 2, 3 and 4 clusters the following approximate numbers will be the results:

When looking at the table a clear elbow effect can be observed when going from 3 to 4 clusters. This would support the presumption that there is three clusters.

| Clusters | MSE |
|:---:|:---:|
| 2 | 30 600 |
| 3 | 9 500 |
| 4 | 6 400 |
| 5 | 3 500 |

Table 1: MSE resulting from various number of clusters
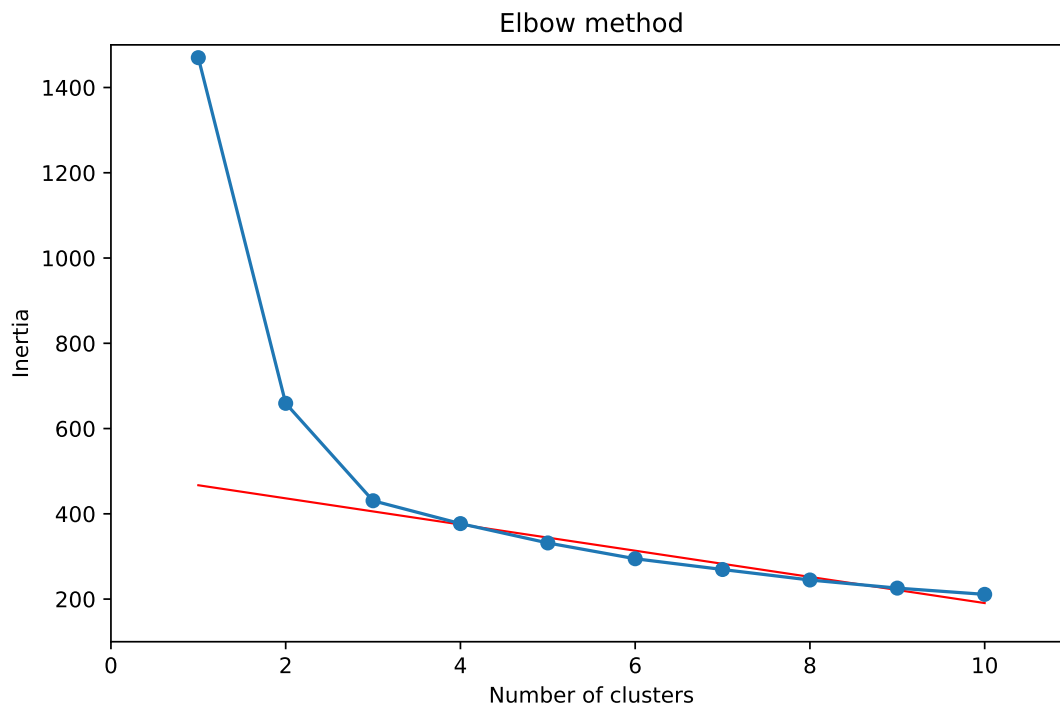
Elbow method



Figure 1: Cluster analysis

# 3    Problem 3: Visualizing the classes

## 3.1    Pair Projection

Plotting each pair of features we can see how the different classes of seed are well differentiable in most of the plots, as it's visible in figure 2. The most interesting plot seems to be the one between C4 and C7, since it is the one where the classes overlap the less.

## 3.2    Gaussian Random Projection

Another interesting result is obtained with the Gaussian Random Projection (figure 3). In this case we distinguished 3 classes (red, green, blue) and two other samples (grey) that seem to be outliers, even though they belong to class 2.

## 3.3    UMAP Projection

In the UMAP projection (figure 4) the 2 outliers are even more evident (they would belong to the blue class). Probably the best thing to do would be to delete these samples from the dataset.

# 4    Problem 4: Evaluating clustering

Since there is no obvious linear border between the clusters a "best effort" approach will be executed. There is always the possibility that some data have flaws, and that may add to the complexity of the problem. There seems to be three clusters, but then there are outliers that blur the picture. The K-Means clustering function does a fairly good job in allocating
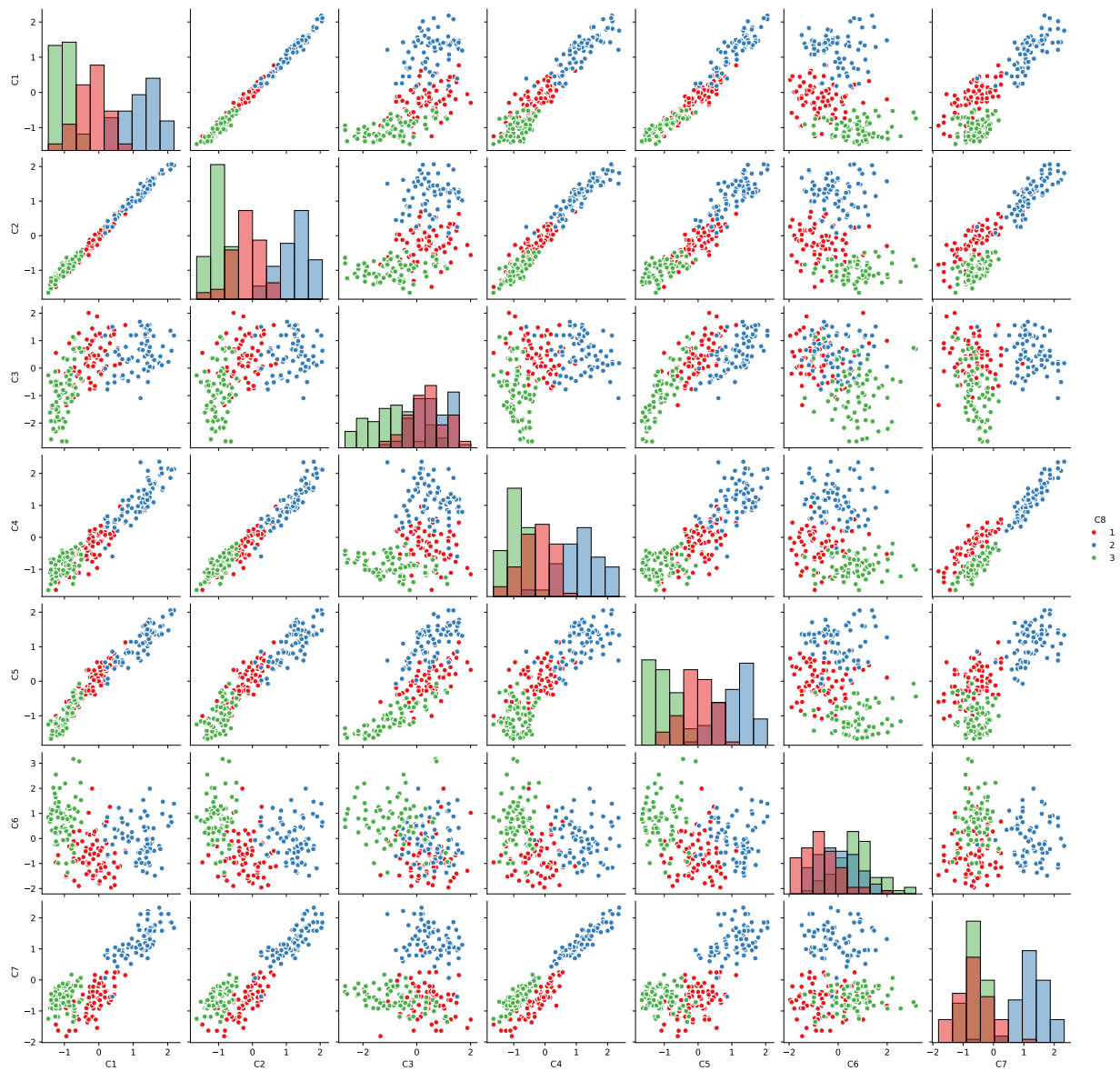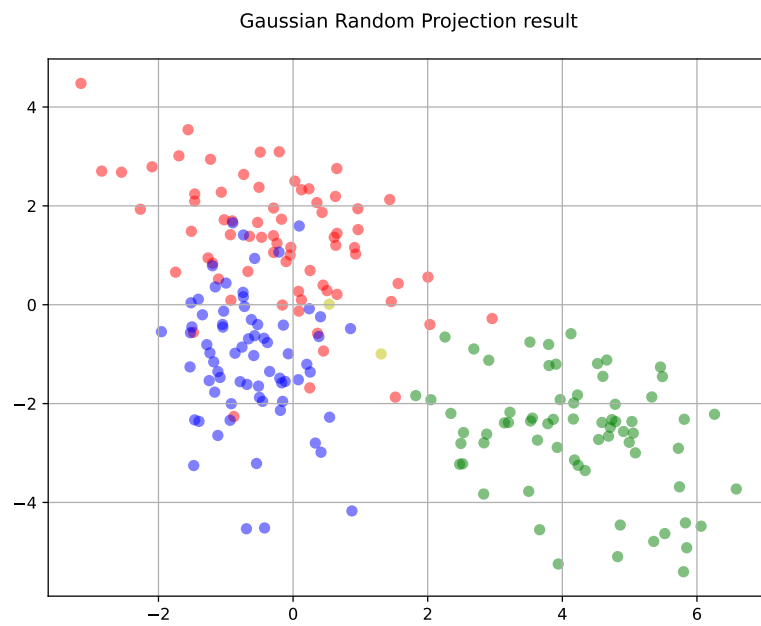
Figure 2: Pair Projection
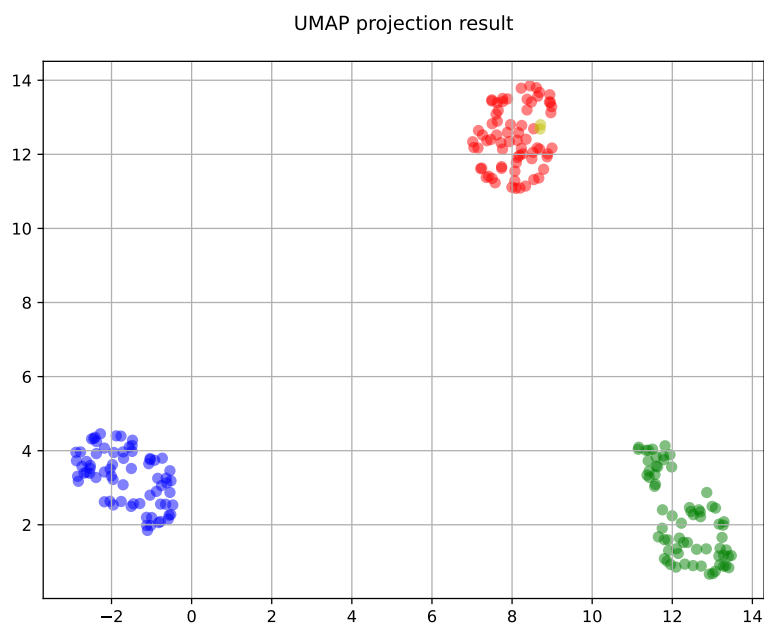
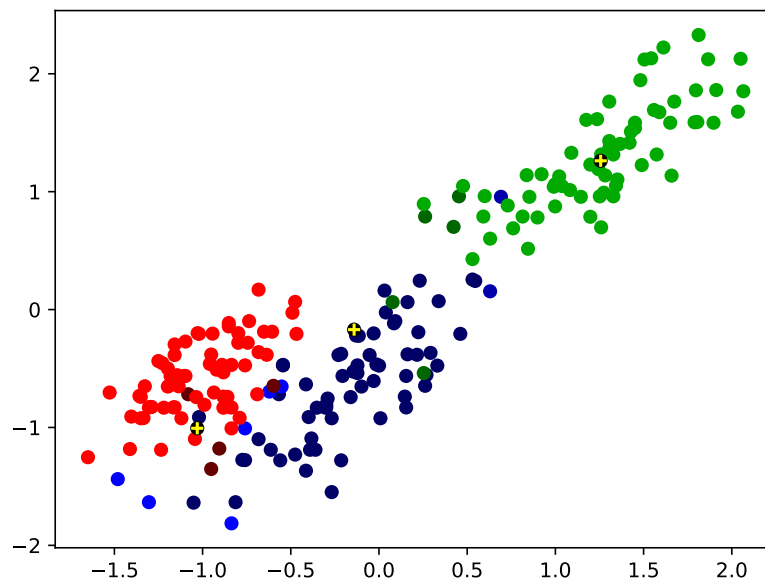Figure 3: Gaussian Random Projection



Figure 4: UMAP projection

Figure 5: Clustering, with markings for cluster center

| | |
|---|---|
| Accuracy score: | **0.9190** |
| Rand index: | **0.8997** |

Table 2: Metrics for K-Mean clustering result

points to clusters. Given the class label data, and comparing it to the clusters there is an accuracy of **0.9190**, which counts as a quite good value. Trying to elevate that value might mean adapting to noise.

# 5 Problem 5: Agglomerative clustering

We used agglomerative clustering to compute a hierarchical clustering for the data and we then tried different linkage options.
The result is visible in figure 7. Looking at the plots it seems that the best work is done by the ward, the average and the complete linkage options.
To determine which one is the best one we looked at the accuracy of each linkage option. The worst is the single. The result was the following:

- Ward: 0.3926339709101015

- Complete: 0.35019845816108097

| | Data labels | | |
|---|---|---|---|
| | 65 | 0 | 5 |
| Predicted labels | 0 | 66 | 4 |
| | 2 | 6 | 62 |

Table 3: Confusion matrix, se diagram 6

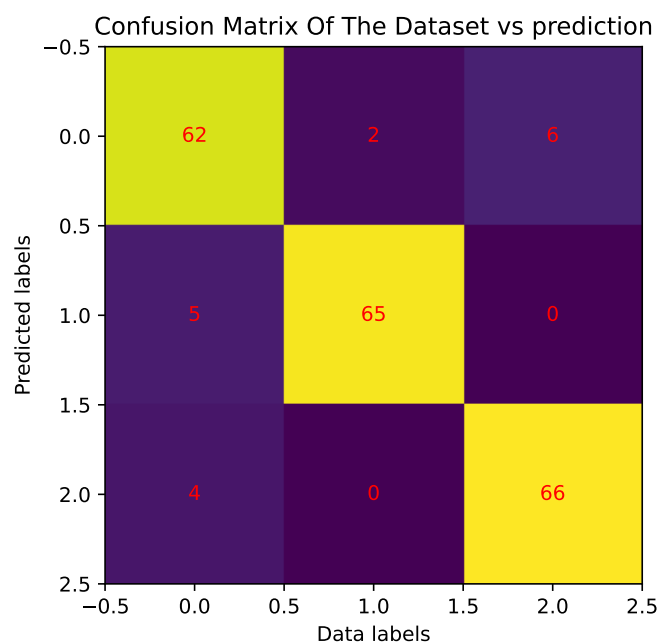Figure 6: Confusion matrix showing number of accurate classification vs faulty classifications
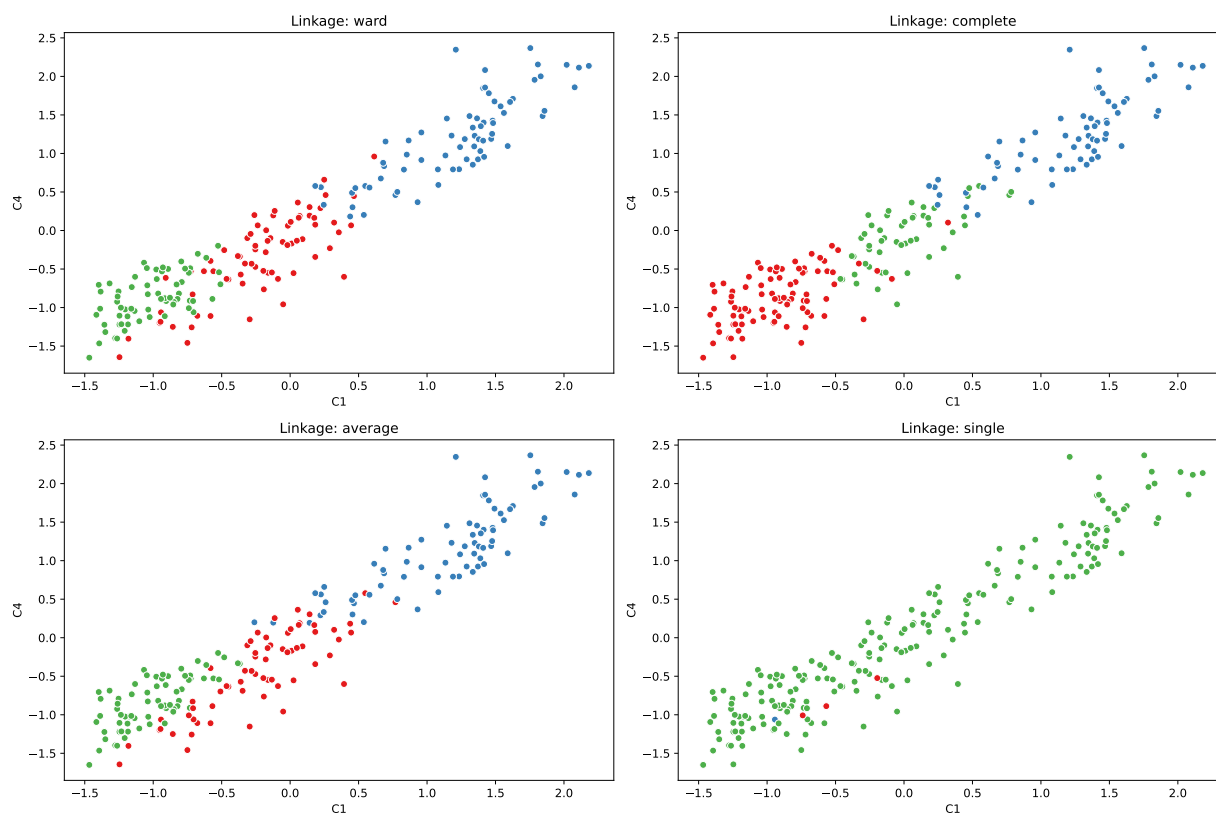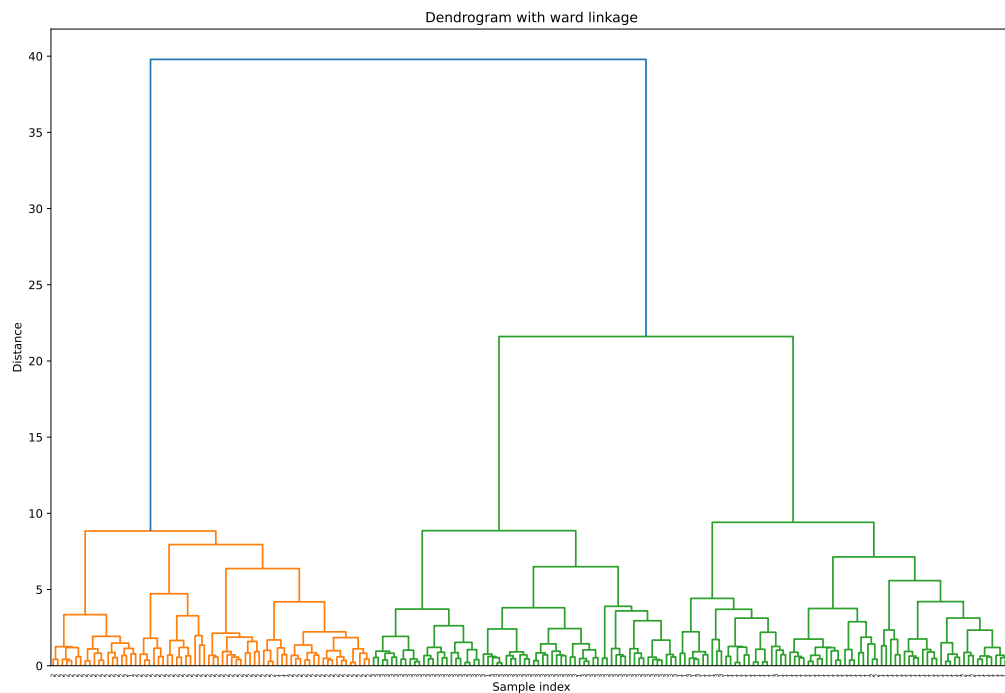


Figure 7: Agglomerative Clustering

Figure 8: Dendrogram

- Average: 0.3759568059006467

- Single: -0.005642378923309307

Looking at the accuracy we determined that the absolute best one was 'ward' and we decided to choose that.
For that we plotted the dendrogram that is possible to see in figure 8.

# A    Source code

This is the complete listing of the source code.

## A.1    Problem 1 and 2

Code for normalizing data and examine optimal number of clusters.

```python
1  # Reynir Siik, Franco Zambon
2  # DIT407 lp3 2024-02-21
3  # Assignment 5
4  # Problem 1, 2
5  # Clustering
6  ##########################################################################
7  import pandas as pd
8  import numpy as np
9  import matplotlib.pyplot as plt
10 from sklearn.cluster import KMeans
11 from scipy.stats import zscore
12 from sklearn.linear_model import LinearRegression
13 ##########################################################################
14 filelocation_in = r"Uppgift05\\00_Uppg_Data\\" # folder to get data from
15 filelocation_out = r"Uppgift05\\test_data\\" # folder to write data to
16 ## ######################################################################
17 df = pd.read_csv(filelocation_in + "seeds.tsv", sep = '\t', header=None)
18
19 ## Problem 1, normalize data, Z-score chosen
20 df_zscore = df.apply(zscore) # Calculate Z-scores for each column
21 df_zscore.drop(7, axis='columns', inplace=True) # Drop class label
22
23 ## Problem 2, find number of clusters
24 ## Use elbow method to find nbr of clusters
25 inertias = []
26
27 for i in range(1,11):
28     kmeans = KMeans(n_clusters=i)
29     kmeans.fit(df_zscore)
30     inertias.append(kmeans.inertia_)
31
32 ellbow = 3
33
34 x=[[1],[2],[3],[4],[5],[6],[7],[8],[9],[10]]
35 regr = LinearRegression()
36 regr.fit(x[ellbow-1:11], inertias[ellbow-1:11])
37 k=regr.coef_
38 m=regr.intercept_
39 print("2 Residual sum of squares: %.2f"
40       % np.mean((k*x[ellbow-1:11]+m - inertias[ellbow-1:11]) ** 2))
41
42 plt.figure(figsize=(8, 5))
43 plt.xlim(0,11)
44 plt.ylim(100,1500)
45 plt.plot(x[0:11:], k*x[0:11:]+m, color='red', linewidth=1)
46 plt.plot(range(1,11), inertias, marker='o')
47 plt.title('Elbow method')
48 plt.xlabel('Number of clusters')
49 plt.ylabel('Inertia')
50 plt.savefig(filelocation_out + 'plot_of_Elbow_method_problem2.pdf')
```

## A.2    Problem 3

Code for Visualizing the classes.

```python
# Reynir Siik, Franco Zambon
# DIT407 lp3 2024-02-21
# Assignment 5
# Problem 3
# Visualizing the classes
###############################################################################
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import zscore
from sklearn.random_projection import GaussianRandomProjection
import seaborn as sns
import umap
###############################################################################
filelocation_in = r"Uppgift05\\00_Uppg_Data\\" # folder to get data from
filelocation_out = r"Uppgift05\\test_data\\" # folder to write data to
## ###########################################################################
df = pd.read_csv(filelocation_in + "seeds.tsv", sep = '\t', header=None)
df.columns = ['C1', 'C2', 'C3', 'C4', 'C5', 'C6', 'C7', 'C8']
colormap = np.array(['r', 'g', 'b', 'y'])

## Problem 1, normalize data, Z-score chosen
df_zscore = df.apply(zscore) # Calculate Z-scores for each column
df_zscore['C8'] = df['C8']
df_zscore['C8'] = np.where(df_zscore['C8']==4, 2, df_zscore['C8'])

# Select only the numerical columns for plotting
numerical_columns = ['C1', 'C2', 'C3', 'C4', 'C5', 'C6', 'C7']
data_for_plot = df_zscore[numerical_columns + ['C8']]

# Plot scatter plots between each pair of features,
# coloring the points by the class label
sns.pairplot(data_for_plot, hue='C8', palette='Set1', diag_kind='hist')
plt.savefig(filelocation_out + 'Pair Projection result.pdf')

# df_zscore.drop(['C8'], axis='columns', inplace=True) # Drop class label
# Extract numerical features from the DataFrame
X = df_zscore.values

# Apply Gaussian Random Projection to project the data to two dimensions
random_projection = GaussianRandomProjection(n_components=2,
                                             random_state=42)
X_projected = random_projection.fit_transform(X)

# Plot the projected data as a scatter plot
plt.figure(figsize=(8, 6))
scatter=plt.scatter(
    X_projected[:, 0], X_projected[:, 1],
    c=colormap[df['C8']-1], alpha=0.5)
plt.title('Gaussian Random Projection result\n')
plt.grid(True)
plt.savefig(filelocation_out + 'Gaussian Random Projection result.pdf')

## Apply UMAP
reducer = umap.UMAP(n_components=2)
```

```
56  embedding = reducer.fit_transform(X)
57  plt.figure(figsize=(8, 6))
58  plt.scatter(
59      embedding[:, 0], embedding[:, 1],          # type: ignore
60      c=colormap[df['C8']-1], alpha=0.5)
61  plt.title('UMAP␣projection␣result\n')
62  plt.grid(True)
63  plt.savefig(filelocation_out + 'UMAP␣Projection␣result.pdf')
```

## A.3   Problem 4

Code for Evaluating clustering.

```
1   # Reynir Siik, Franco Zambon
2   # DIT407 lp3 2024-02-21
3   # Assignment 5
4   # Problem 4
5   # Evaluating Clustering
6   ######################################################################
7   import pandas as pd
8   import numpy as np
9   import matplotlib.pyplot as plt
10  from sklearn.cluster import KMeans
11  from scipy.stats import zscore
12  from sklearn.metrics import rand_score
13  from sklearn.metrics import accuracy_score
14  from sklearn.metrics import confusion_matrix
15  ######################################################################
16  filelocation_in = r"Uppgift05\\00_Uppg_Data\\" # folder to get data from
17  filelocation_out = r"Uppgift05\\test_data\\" # folder to write data to
18  ## ##################################################################
19  df = pd.read_csv(filelocation_in + "seeds.tsv", sep = '\t', header=None)
20  colormap = np.array(['#000066', '#0000aa', '#0000ff', # Blue
21                       '#006600', '#00aa00', '#00ff00', # Green
22                       '#660000', '#aa0000', '#ff0000', # Red
23                       '#FFFF00'])                       # Yellow
24
25  ## Normalize data with Z-scores and extract labels
26  df_zscore = df.apply(zscore) # Calculate Z-scores for each column
27  df_zscore.drop(7, axis='columns', inplace=True) # Drop class label
28  labels = df[7] # classification labels in data
29
30  ## Problem 4, find 3 clusters
31  kmeans = KMeans(n_clusters=3, random_state=0).fit(df_zscore)
32
33  # Set appropriate permutation of labels
34  km_lbl = kmeans.labels_
35  for i in range(len(km_lbl)):
36      km_lbl[i] = (km_lbl[i] + 1) % 3 + 1
37
38  ## Calculate metrics
39  cm = confusion_matrix(labels, km_lbl)
40  print(cm)
41
42  ## Compute accuracy
43  print('Accuracy␣score:␣', accuracy_score(labels, km_lbl))
44  ## Compute Rand index
45  print('Rand␣index:␣', rand_score(labels, km_lbl))
```

```
46
47  plt.clf
48  plt.scatter(df_zscore[1], df_zscore[6],
49              c=colormap[3 * (df[7] - 1) + (km_lbl - 1)])
50  plt.scatter(kmeans.cluster_centers_[:, 0],
51              kmeans.cluster_centers_[:, 1], c='#000000', marker='o')
52  plt.scatter(kmeans.cluster_centers_[:, 0],
53              kmeans.cluster_centers_[:, 1], c='#FFFF00', marker='+')
54  plt.savefig(filelocation_out + 'P4_clusters.pdf')
55  plt.show()
56  plt.clf
57  for i in range(3):
58      for j in range(3):
59          plt.annotate(str(round(cm[i][j], 2)),
60                       xy=(j, i),
61                       ha='center', va='center', color='red')
62  plt.title("Confusion␣Matrix␣Of␣The␣Dataset␣vs␣prediction")
63  plt.xlabel("Data␣labels")
64  plt.ylabel("Predicted␣labels")
65  plt.imshow(cm,  interpolation='nearest')
66  plt.savefig(filelocation_out + 'P4_confusion_matrix.pdf')
67  plt.show()
```

## A.4  Problem 5

Code for Agglomerative clustering.

```
1   ## PROBLEM 5
2
3   import matplotlib.pyplot as plt
4   from sklearn.cluster import AgglomerativeClustering
5   import seaborn as sns
6   from sklearn.metrics import silhouette_score
7
8
9   ## Agglomerative clustering
10
11  # Extract predictors (features) and dependent variable from df_zscore
12  X = df_zscore[['C1', 'C2', 'C3', 'C4', 'C5', 'C6', 'C7']]
13  y = df_zscore['C8']
14
15  # Initialize a dictionary to store linkage options and corresponding clusterings
16  linkage_options = {'ward': AgglomerativeClustering(n_clusters=3, linkage='ward'),
17                     'complete': AgglomerativeClustering(n_clusters=3, linkage='complet
18                     'average': AgglomerativeClustering(n_clusters=3, linkage='average
19                     'single': AgglomerativeClustering(n_clusters=3, linkage='single')]
20
21  # Perform hierarchical clustering with different linkage options
22  plt.figure(figsize=(15, 10))
23  for i, (linkage, clustering) in enumerate(linkage_options.items(), 1):
24      plt.subplot(2, 2, i)
25      clustering.fit(X)
26      sns.scatterplot(x='C1', y='C4', hue=clustering.labels_, data=df_zscore, palette=
27      plt.title(f'Linkage:␣{linkage}')
28      plt.xlabel('C1')
29      plt.ylabel('C4')
30  plt.tight_layout()
31  plt.savefig('AgglomerativeClustering.pdf')
```

```python
32  plt.show()
33
34
35  ## Calculate the accuracy for each linkage option
36
37  # Initialize a dictionary to store silhouette scores for each linkage option
38  silhouette_scores = {}
39
40  # Compute silhouette score for each linkage option
41  for linkage, clustering in linkage_options.items():
42      clustering.fit(X)
43      silhouette_scores[linkage] = silhouette_score(X, clustering.labels_)
44
45  # Print silhouette scores for each linkage option
46  for linkage, score in silhouette_scores.items():
47      print(f"Silhouette Score for {linkage}: {score}")
48
49
50  ## Plot the dendrogram
51
52  clustering = AgglomerativeClustering(n_clusters=3, linkage='ward')
53  clustering.fit(X)
54
55  Z = linkage(X, method=best_linkage)
56
57  plt.figure(figsize=(15, 10))
58  dendrogram(Z, labels=y.values)
59  plt.title(f'Dendrogram with {best_linkage} linkage')
60  plt.ylabel('Distance')
61  plt.savefig('Dendrogram.pdf')
62  plt.show()
```