

SENG 474: Data Mining

Assignment 3

Clustering Experiments: *Lloyd's Algorithm and Hierarchical Agglomerative Clustering*

Sean McAuliffe, V00913346 (A02)

March 23rd, 2023 AD

Contents

1	Introduction	3
2	Implementation	3
2.1	Provided Datasets	3
3	Lloyd's Algorithm	4
3.1	Initialisation	6
3.2	The Elbow Method	6
3.3	Experimental Results	7
3.3.1	Lloyd's Algorithm on Dataset 1	7
4	Hierachical Agglomerative Clustering	8

1 Introduction

This report summarizes the results obtained from applying two clustering algorithms, Lloyd's Algorithm and Agglomerative Clustering to two provided datasets. The first dataset contains 3500 points in 2 dimensions, and the second dataset contains 14801 points in 3 dimensions. The purpose of this exercise is to gain familiarity with implementing clustering algorithms, and to gain an understanding of the differences between approaches; in order to be able to interpret their results, choose an appropriate number of clusters for a given problem and apply appropriate solutions for different situations.

Clustering is an unsupervised learning technique with the goals of categorizing data into groups based on their similarity. In this way, labels can be learned from the data itself; without the need for pre-labelled data. Clustering is a useful technique in data mining which can be used to find patterns in datasets, perhaps as a pre-processing step for further supervised learning methods. The notion of a cluster is not easy to formalise (one recognises a cluster when one sees it). Accordingly, many algorithms exist for producing clusterings from data.

2 Implementation

The experiments conducted for this report were composed in an accompanying Jupyter notebook running on a *Python 3.10.6* Kernel. An implementation of Lloyd's algorithm was written for this assignment, according to the description provided in the lecture 23 slides. An implementation of agglomerative clustering was imported from the *sci-kit learn* library. The *sklearn.cluster.AgglomerativeClustering* class was used, which provides the necessary connectivity metrics for determining the distance between clusters.

All code for implementation of the algorithms, use of external libraries, data loading, and graph generation is contained in the accompanying Jupyter notebook.

2.1 Provided Datasets

In order to build an intuition for the behaviour of the algorithms in later experiments, the provided data were plotted in 2 and 3 dimensions respectively.

This is shown in Figure 1.

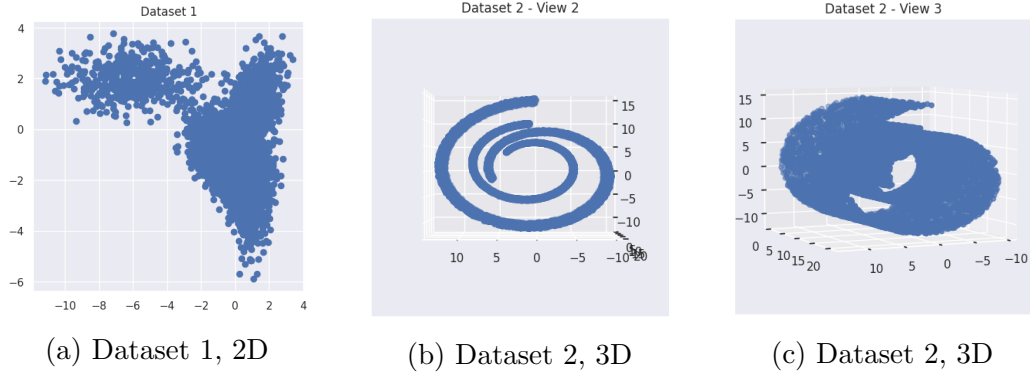


Figure 1: The shape of the provided datasets in 2 and 3 dimension

The first dataset seems like it ought to be separable into roughly 2 clusters. The clusters seem to be globular, but clearly oblong. The second dataset seems to also be separable into 2 clusters, however it is highly non-globular. Rather, it is composed of two intertwined spirals having the appearance of a dessert.

3 Lloyd's Algorithm

Lloyd's algorithm is a k-means clustering algorithm in which clusters are defined by the mean of the points in the cluster. Lloyd's algorithm is an iterative algorithm which converges on a clustering of minimum cost by repeatedly assigning points to their nearest clusters, and then recomputing the cluster means, until the assignments of points to clusters no longer changes. Through this process, the means which describe a clustering are learned.

The algorithm is as follows:

1. Initialise the means of the clusters (using one of the below methods)
2. Assign each point to the cluster with the nearest mean.
3. Recompute the means of the clusters.
4. Repeat steps 2 and 3 until the assignments of points to clusters no longer changes.

The goal is to learn a minimum-cost clustering of the dataset into k clusters. The cost of a k -clustering, C is defined as:

$$W(C) = \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2 \quad (1)$$

This assignment required Lloyd's algorithm to be implemented, the implementation is shown in the code snippet below.

```

1 def lloyd(X, k, init):
2     # Initialize the clusters
3     C = init(X, k)
4     # Repeat until the centres do not change
5     previous_centres = get_centres(C)
6     i = 0
7     while True:
8         i += 1
9         # Assign each data point to the cluster with the
10        closest centre
11        for cluster in C:
12            cluster.clear_data()
13            assign_to_closest(X, C)
14        # Update the cluster centres
15        for cluster in C:
16            ideal_centre = np.mean(cluster.data, axis=0)
17            # Find the point in the data closest to the
18            idealized centre
19            min_dist = np.inf
20            min_point = None
21            for point in cluster.data:
22                dist = euclidean_distance(point, ideal_centre)
23                if dist < min_dist:
24                    min_dist = dist
25                    min_point = point
26            cluster centre = deepcopy(min_point)
27        # Check if the centres have changed
28        new_centres = get_centres(C)
29        if all_centres_equal(previous_centres, new_centres):
30            break
31        previous_centres = new_centres
32    return i, C

```

3.1 Initialisation

The algorithm is probabilistic in that the initialisation of the means is probabilistic, and so the algorithm may converge to a different clustering each time it is run, and each run may require a variable number of iterations to converge. To find a low cost clustering it may be necessary to run the algorithm multiple times, and to choose the clustering with the lowest cost.

Uniform random initialisation is a method of initialising the means of the clusters by randomly sampling, without replacement, points from the dataset with equal probability. This method is simple to implement, however it can lead to poor results and poor runtimes if the initial means are poorly chosen.

K-Means++ Initialisation proceeds from the notion that cluster centres ought not to be near to each other, as such points are likely to be assigned to the same cluster. K-Means++ initialisation can reduce runtime by preferring to sample points which are far from the current cluster centres. The first cluster centre is chosen uniformly at random from the dataset, subsequent centres are chosen with probability proportional to the squared distance from the nearest cluster centre. This method is more likely to produce good results than uniform random initialisation, resulting in lower expected runtimes.

This assignment involved implementing k-means++ initialisation, a snippet showing part of the implementation is below.

```
1 for x in X:
2     min_dist = np.inf
3     for cluster in clusters:
4         dist = cluster.distance(x)
5         if dist < min_dist:
6             min_dist = dist
7             dists.append(min_dist)
8     probability_vector = dists/np.sum(dists)
9     centre = X[np.random.choice(len(X), p=probability_vector)]
```

3.2 The Elbow Method

Determining how many clusters are present in a dataset is a difficult problem to formalise. Instead, for this assignment the elbow method was used to determine the proper number of clusters to use. The elbow method is a heuristic method in which the number of clusters, k is incremented, and

a clustering algorithm is run to produce a clustering for each k . The cost of each k -clustering is measured (for each k the algorithm is run multiple times and the lowest cost clustering is chosen as the result). The cost of each k -clustering is plotted as a function of k . The resulting plot should contain a steadily decreasing curve, ideally with a sharp bend at the point of diminishing returns. The point at which the curve bends is the optimal number of clusters to use.

3.3 Experimental Results

3.3.1 Lloyd's Algorithm on Dataset 1

The number of clusters was varied from $[2, 10]$ and for each value of k the algorithm was run 10 times, and the lowest cost clustering was chosen as the result. For each of the chosen clusterings, its cost and the number of iterations required to converge were recorded. The results are shown in 2

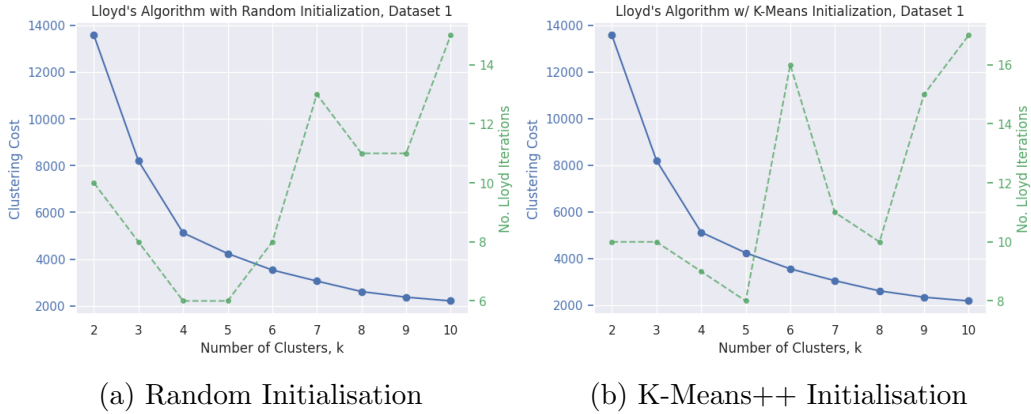


Figure 2: The elbow method graph for dataset 1

One thing to notice is that the results are very similar. Indeed, the cost of the clusterings chosen for each value of k are almost identical. In both plots, there appears to be an elbow that appears sharpest at $k = 4$. This choice is somewhat arbitrary, and sensitive to such things as the presented scale of the axes.

Another thing to notice is that the number of iterations required to converge increases sharply after the elbow-point. And, that k-means++ initialisation

did not appear to significantly reduce the number of iterations required to converge. This is likely due to the small size of the dataset, with so few points and clusters, the initialisation is very sensitive to the choice of initial means.

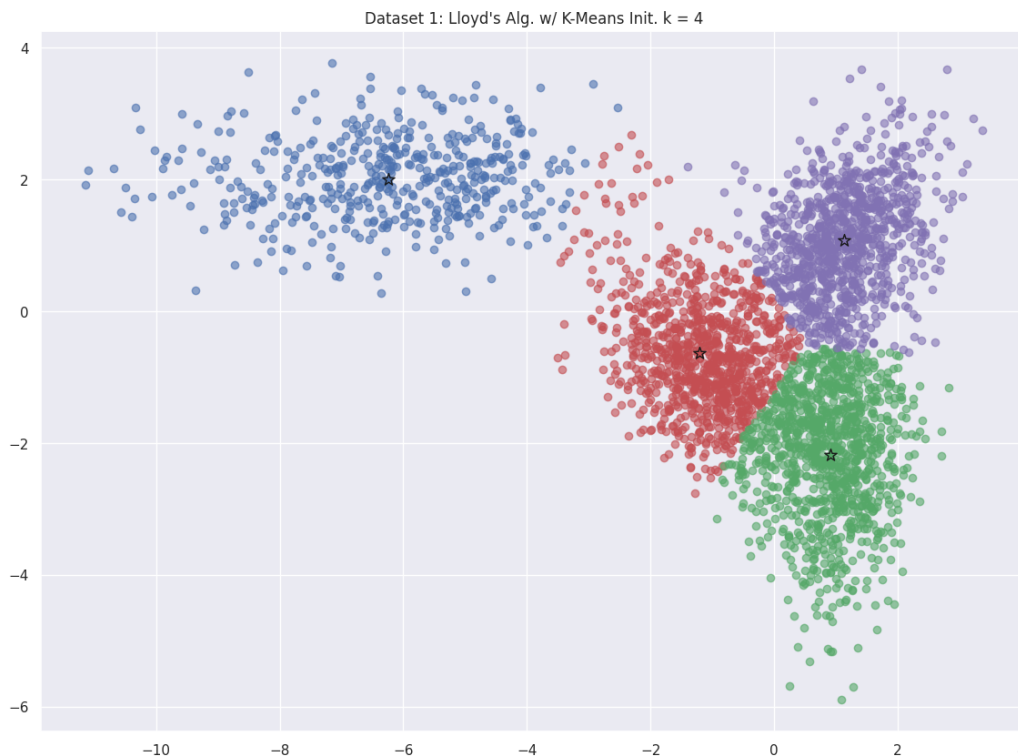


Figure 3: Clustering produced by Lloyd's Algorithm w/ $k = 4$

4 Hierarchical Agglomerative Clustering

Hierarchical agglomerative clustering is a connectivity-based algorithm in which clusters are defined by their member points (depending on the connectivity, or linkage, metric used). The algorithm begins by assigning each point to its own cluster. Then, at each iteration, the two clusters with the smallest distance between them are merged together. This process is repeated until the desired number of clusters is reached.