

Problem 1

Part (a)

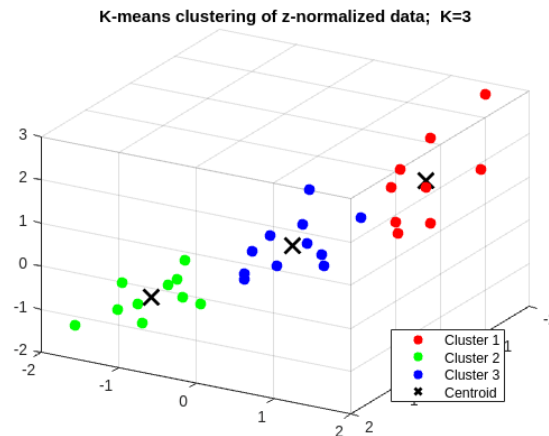


Figure 1: Clusters formed using K-means Clustering

K-means is sensitive to the initial placement of the centroids. Running the algorithm with different initializations and selecting the best result can help mitigate the risk of getting stuck in a local minimum. The algorithm is run multiple times with different initial centroid positions, and the solution with the lowest within-cluster sum of squared distances (WCSS) is chosen.

Part (b)

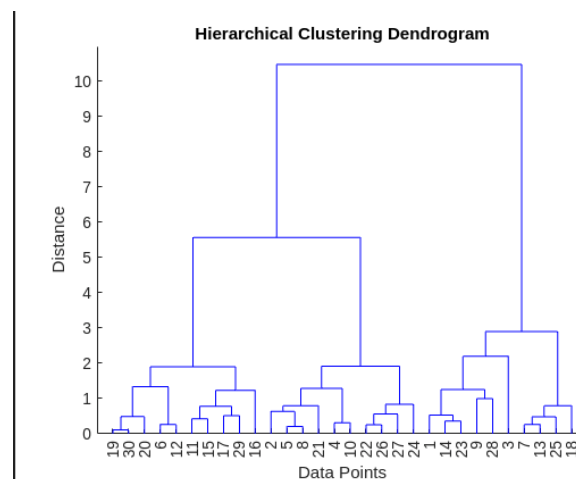


Figure 2: Dendrogram created by Hierarchical Clustering

The dendrogram can be used to visually understand the distance between the means of different datapoints. We start with the two data points which are the closest in terms of euclidean distance i.e 19 and 30. These two are clubbed and their mean is then treated as a datapoint to calculate distances with other data points. In this way we can see that at each stage, the two data points with the smallest euclidean distance are paired up together to form a cluster. It is essential to note that repeating the clustering multiple times has no effect on the resultant clustering if the datapoints and the distance function remain the same. Hence, multiple iterations of this type of clustering can not be carried out in hopes of achieving a global optima. The clustering result obtained can not be assured to be optimal.

Comparison

The clustering obtained by hierarchical clustering and K-means clustering differs under certain scenarios.

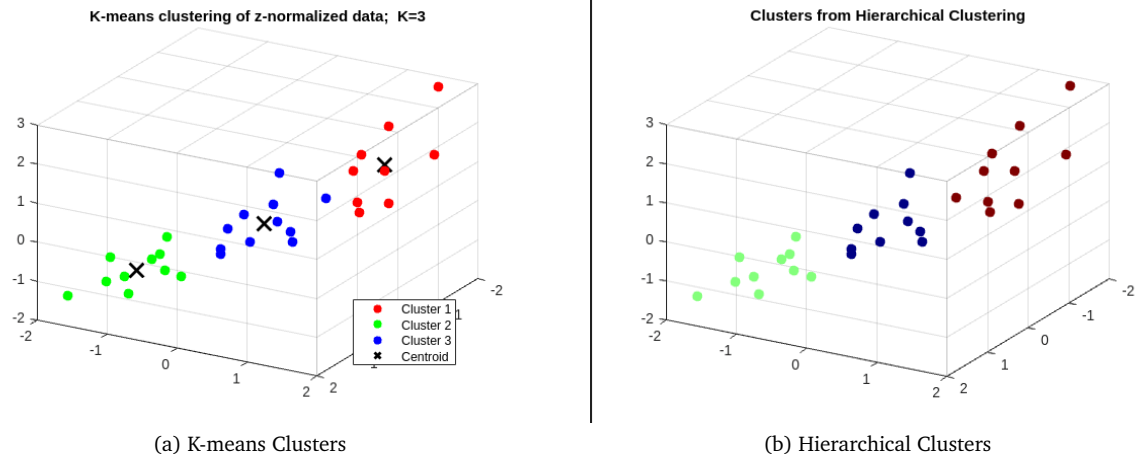


Figure 3: Comparison of results obtained from different clustering schemes

It can be seen that there is a difference in the allocation of cluster to a point that is almost evenly far from the the centroids of the K-means cluster.

In this case, it is observed that the datapoint is clustered along with the cluster of the nearest data point. This is because, the hierarchical clustering is agglomerative and hence if two data points, which should ideally belong to different clusters in the global optima, are near each other, they are clustered together. This can even be interpret as hierarchical clustering being more robust to outliers, while K-means clustering is sensitive to outliers as they can significantly impact the position of centroids.

It even suggests that K-means clustering is more suited to work well with clusters that are spherical, equally sised and have similar density.

APPENDIX A

MATLAB Codes

Problem 1

```
1 clear all
2 close all
3 clc
4
5 %% Load the data
6 X = load('data.mat').X;
7 % 30x3 data matrix for 30 different
8 % Column1 : strength [MPa]
9 % Column2 : strain to failure [%]
10 % Column3 : elastic modulus [GPa]
11
12 %% Normalize dataset
13 mean = mean(X);
14 stdDev = std(X);
15 zNormalData = normalize(X);
16 %X_reconstructed = zNormalData.*stdDev + mean;
17 %X_error = X_reconstructed - X;
18
19 %% Part (a) : K-means clustering
20 K = 3;
21 numInitializations = 1000;          % define the number of random
    initializations
22 choiceMultiple = 10;
23 % random centroids generated are numInitializations*choiceMultiple
24 % this is done to further increase randomness by using datasample
25 paddedNumInits = choiceMultiple*numInitializations;
26
27 % Variables to store the results
28 bestIdx = [];
29 bestCentroids = [];
30 bestWCSS = Inf;
31 randomCentroids = [randn(paddedNumInits, 1), randn(paddedNumInits, 1),
    randn(paddedNumInits, 1)];
32
33 for i = 1:numInitializations
34     % Randomly initialize cluster centroids
35     initialCentroids = datasample(randomCentroids, K, 'Replace', false
        );
36
37     % Perform K-means clustering with these initial centroids
38     [idx, centroids, sumd] = kmeans(zNormalData, K, 'Start',
        initialCentroids);
39
40     % Calculate the total within-cluster sum of squares (WCSS)
41     wcss = sum(sumd);
42
43     % Check if this initialization results in a lower WCSS
44     if wcss < bestWCSS
45         bestIdx = idx;
46         bestCentroids = centroids;
47         bestWCSS = wcss;
48         %disp(bestWCSS);
49     end
50 end
51
```

```

52
53 % Plot the best clustered data
54 colors = ['r', 'g', 'b', 'c', 'm', 'y', 'k'];
55 figure;
56 hold on;
57 for i = 1:K
58     cluster_points = zNormalData(bestIdx == i, :);
59
60     scatter3(cluster_points(:, 1), cluster_points(:, 2),
61             cluster_points(:, 3), 50, colors(i), 'filled', 'DisplayName', ['
Cluster ', num2str(i, '%d') ]]);
62 end
63 colormap('jet'); % Set the colormap as needed
64 scatter3(bestCentroids(:, 1), bestCentroids(:, 2), bestCentroids(:, 3)
65         , 200, 'k', 'Marker', 'x', 'LineWidth', 2, 'DisplayName', 'Centroid'
66         );
67 title(['K-means clustering of z-normalized data; K=', num2str(K, '%d')
68         ]);
69 legend('Location', 'Best');
70 grid on;
71 hold off;
72
73 % Set the perspective view
74 view(-30, -30); % Azimuth: 30 degrees, Elevation: 45 degrees
75
76 %% Part (b) : Hierarchical clustering
77
78 % Perform hierarchical clustering
79 Z = linkage(zNormalData, 'ward'); % You can choose different linkage
80     methods
81
82 % Visualize the hierarchical clustering using a dendrogram
83 figure;
84 dendrogram(Z);
85
86 title('Hierarchical Clustering Dendrogram');
87 xlabel('Data Points');
88 ylabel('Distance');
89
90 % Use Hierarchical Clustering to form Clusters
91 T = cluster(Z, "maxclust", 3);
92
93 % Display the clusters using different colors in a scatter plot
94 figure;
95 scatter3(zNormalData(:, 1), zNormalData(:, 2), zNormalData(:, 3), 50,
96         T, 'filled');
97 title('Clusters from Hierarchical Clustering');
98 colormap('jet');
99
100 % Set the perspective view
101 view(-30, -30); % Azimuth: 30 degrees, Elevation: 45 degrees

```