

# K-MEANS CLUSTERING

Vishal S  
MM14B048

## Data Set

The data is visualised using MATLAB. The scatter plot is depicted in Fig. below.

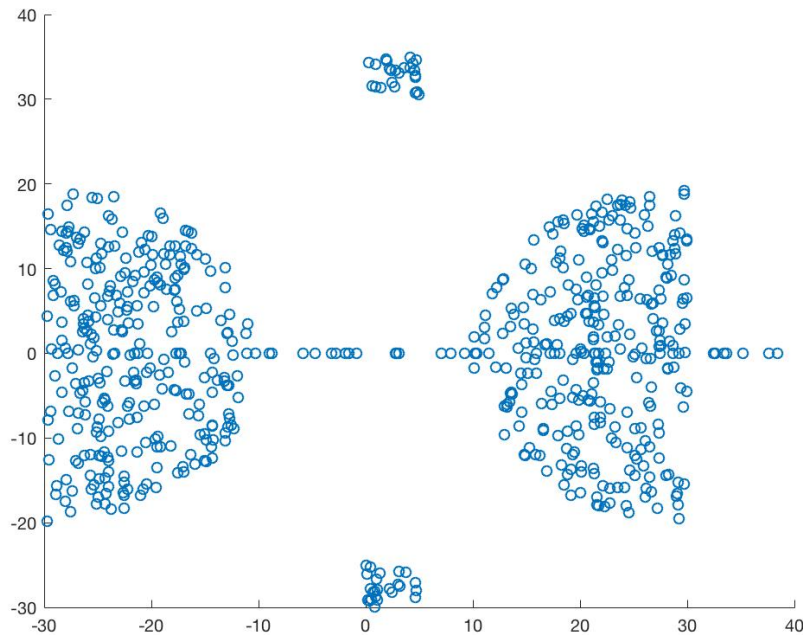


Figure 1: Outlier.txt

## Outlier.txt

The data has 4 clusters and some noise points.

| Cluster Id | No of Points |
|------------|--------------|
| Cluster 1  | 276          |
| Cluster 2  | 276          |
| Cluster 3  | 24           |
| Cluster 4  | 24           |
| Noise      | 50           |
| Total      | 650          |

## K-means algorithm

A MATLAB code was written to cluster the data based on K-means algorithm. First the data is read into the code using the code snippet given below.

```
1 %% Load data
2 clear all
3 clf
4 clc
```

```

5 hold on
6 filename = 'outlier.txt';
7 % read from the file
8 fileID = fopen(filename, 'r');
9 formatSpec = '%f,%f,%f';
10 sizeA = [3 Inf];
11 A = fscanf(fileID, formatSpec, sizeA);
12 fclose(fileID);

```

Once the data is loaded, the K-means algorithm is implemented using the code snippet given below.

```

1 %% Kmeans algorithm
2 % initialising the clusters' means
3 for i = 1:no_of_clusters
4     for j = 1:2
5         mean_old(j, i) = A(j, initial_mean(i));
6     end
7 end
8
9 while ((convergence > convergence_limit)&&(iteration < max_iteration))
10     X = ['Iteration = ', num2str(iteration)];
11     disp(X);
12     % Assigning each point to a cluster
13     dist = 0.0;
14     for i = 1:max_len
15         dist_x = (A(1, i) - mean_old(1, 1))^2;
16         dist_y = (A(2, i) - mean_old(2, 1))^2;
17         dist = dist_x + dist_y;
18         index = 1;
19         for j = 2:no_of_clusters
20             dist_x = (A(1, i) - mean_old(1, j))^2;
21             dist_y = (A(2, i) - mean_old(2, j))^2;
22             if (dist > dist_x + dist_y)
23                 dist = dist_x + dist_y;
24                 index = j;
25             end
26         end
27         cluster_id(i, 1) = index;
28     end
29
30     for i = 1:no_of_clusters
31         sum(1, i) = 0.0;
32         sum(2, i) = 0.0;
33         no_of_points(1, i) = 0.0;
34     end
35     % Calculating new mean
36     for i = 1:max_len
37         sum(1, i) = 0.0;
38         sum(2, i) = 0.0;
39         sum(1, cluster_id(i, 1)) = sum(1, cluster_id(i, 1)) + A(1, i);
40         sum(2, cluster_id(i, 1)) = sum(2, cluster_id(i, 1)) + A(2, i);
41         no_of_points(1, cluster_id(i, 1)) = no_of_points(1, cluster_id(i, 1)) + 1;
42     end
43
44     % Calculating the new means
45     for i = 1:no_of_clusters
46         for j = 1:2
47             mean_new(j, i) = sum(j, i) / no_of_points(1, i);
48         end
49     end
50
51     % Calculating SSE
52     sse_new = 0.0;
53     for i = 1:max_len
54         for j = 1:2
55             sse_new = sse_new + (A(j, i) - mean_new(j, cluster_id(i)))^2;

```

```

56         end
57     end
58     convergence = abs((sse_new - sse_old)) / sse_new ;
59     % update the variables
60     iteration = iteration + 1;
61     sse_old = sse_new;
62     mean_old = mean_new;
63 end

```

Finally, the results are plotted.

```

1 %% Plot results
2 k=max(cluster_id);
3 Colors=hsb(k);
4 Legends = {};
5 for i=0:k
6     A_i=A(:,cluster_id==i);
7     if i~=0
8         Style = 'x';
9         MarkerSize = 8;
10        Color = Colors(i,:);
11        Legends{end+1} = ['Cluster #' num2str(i)];
12    end
13    if ~isempty(A_i)
14        %scatter3(A_i(1,:),A_i(2,:),A_i(3,:),Style,'MarkerSize',MarkerSize,'
15        MarkerFaceColor',Color);
16        scatter(A_i(1,:),A_i(2,:),Style,'MarkerFaceColor',Color);
17    end
18    hold on;
19 end
20 Style = 'o';
21 MarkerSize = 10;
22 Color = [0 0 0];
23 Legends{end+1} = 'Centroids';
24 scatter(mean_new(1,:),mean_new(2,:),Style,'MarkerFaceColor',Color);
25 hold off;
26 grid on;
27 legend(Legends);
28 legend('Location','NorthEastOutside');

```

## K-means Implementation

### Outlier.txt

#### No of clusters = 2

The above code was executed with the number of clusters equal to two and the following output was obtained. The centroids and the SSE obtained for each case is given below.

```

1 SSE =138047.9529
2 Centroids
3     Centroid #1=[18.9961 0.78908]
4     Centroid #2=[-21.0431
5     -0.27064]

```

Listing 1: Case 1

```

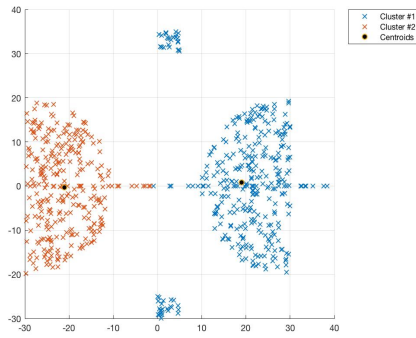
5 SSE =137470.4669
6 Centroids
7     Centroid #1=[20.122 2.4939]
8     Centroid #2=[-19.5722 -2.0021]

```

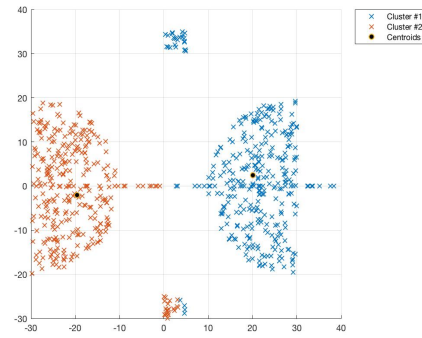
Listing 2: Case 2

#### No of clusters = 4

The above code was executed with the number of clusters equal to four and the following output was obtained. The centroids and the SSE obtained for each case is given below.

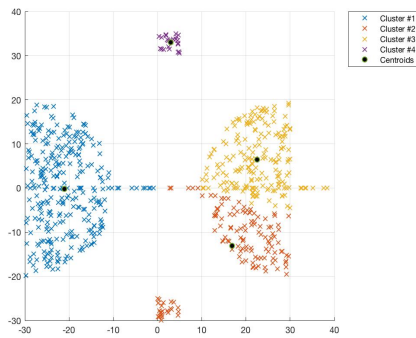


(a) Case 1

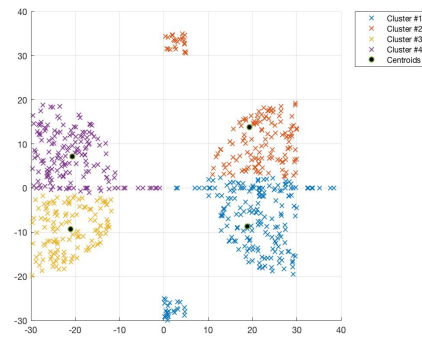


(b) Case 2

Figure 2: K-means with K=2



(a) Case 1



(b) Case 2

Figure 3: K-means with K=4

```

9 SSE = 71740.5202
10 Centroids
11   Centroid #1=[-21.0418
12     -0.25116]
13   Centroid #2=[16.8663  -13.1015]
14   Centroid #3=[22.5359  6.4484]
   Centroid #4=[3.0439  32.9744]

```

Listing 3: Case 1

```

16 SSE = 74377.92
17 Centroids
18   Centroid #1=[18.8815  -8.6685]
19   Centroid #2=[19.2873  13.8065]
20   Centroid #3=[-21.094  -9.3189]
21   Centroid #4=[-20.6517  7.1059]

```

Listing 4: Case 2

## Observation

We observe that, the bigger clusters are getting classified as two different clusters. The noise points are also getting clustered.