

Pattern Recognition Practical 6

Group 24: Maikel Withagen (s1867733) Steven Bosch (s1861948)

October 22, 2015

Assignment 1

1

When we take $k = 2$ the Minkowski metric is the same as the Euclidean distance between the points, which is used as error function in other clustering methods such as K-means clustering.

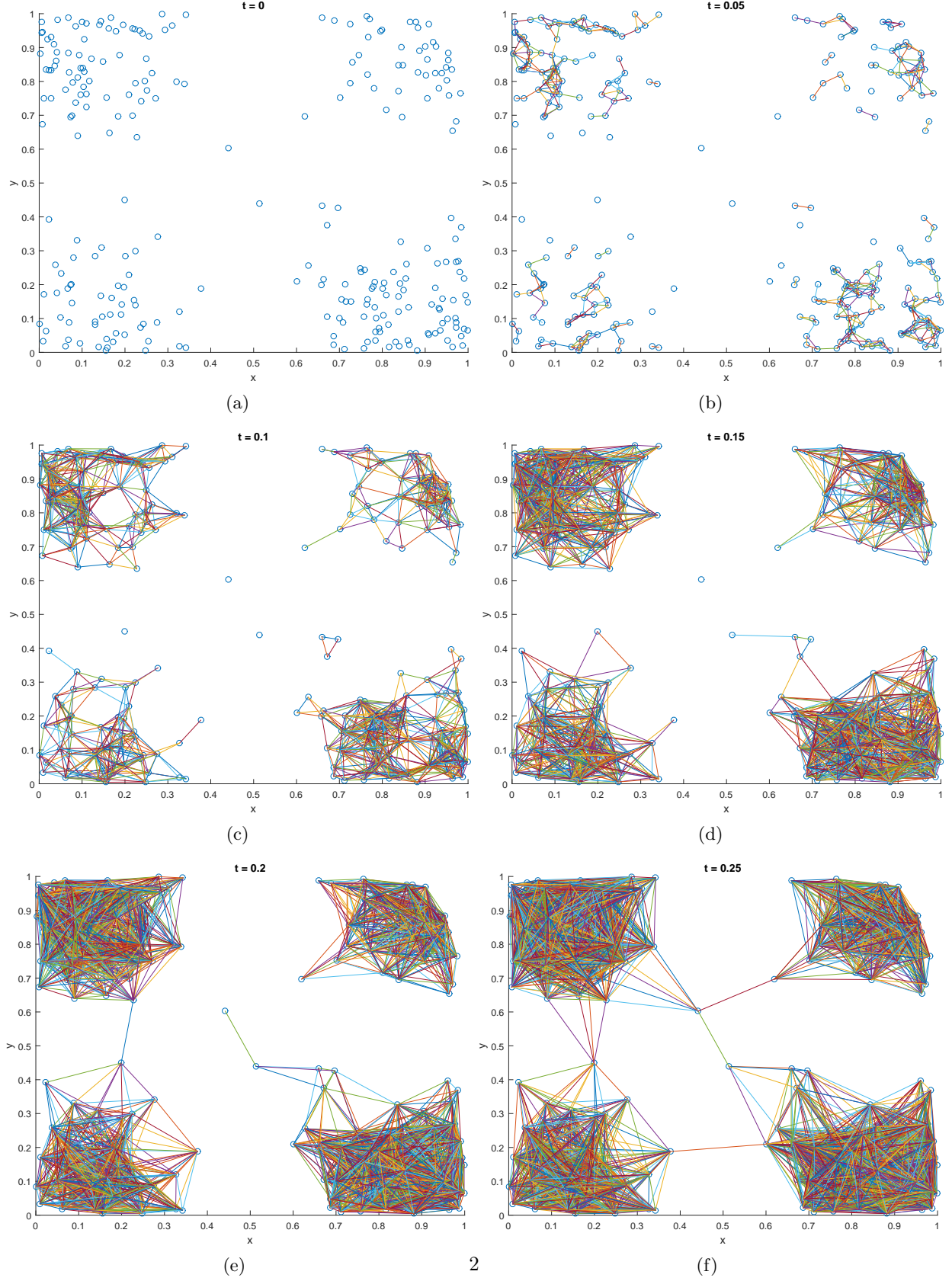
2

See section A in the appendix for our implementation.

3

Looking at figure 1a we can clearly distinguish four main clusters within the data. The figures show that for higher values of t more connections are plotted. This is logical since a higher threshold permits higher distances between two points to be plotted, which overall results in more plotted connections. As for the optimal value of t , 0.05 is clearly too low, because we can see some connections within the clusters, but a lot of points that clearly belong to the clusters are left out because their distance to the other points is too high (see figure 1b). We can still see this for a t of 0.1, but on a smaller scale (see figure 1c). On the other hand a t of 0.25 is clearly too high, because multiple different clusters get connected through outliers, causing multiple clusters to be clustered together (see figure 1f). The same thing happens with a t of 0.2, where the two left clusters are connected (see figure 1e). Since this does not happen for a t of 0.15 and almost all of the points are assigned to a cluster (see figure 1d), this seems to be the optimal value of t . There is one point that does not get assigned to a cluster, so we will have to accept this as an outlier that does not belong to any cluster.

Figure 1: Minkowski clustering for different threshold values.



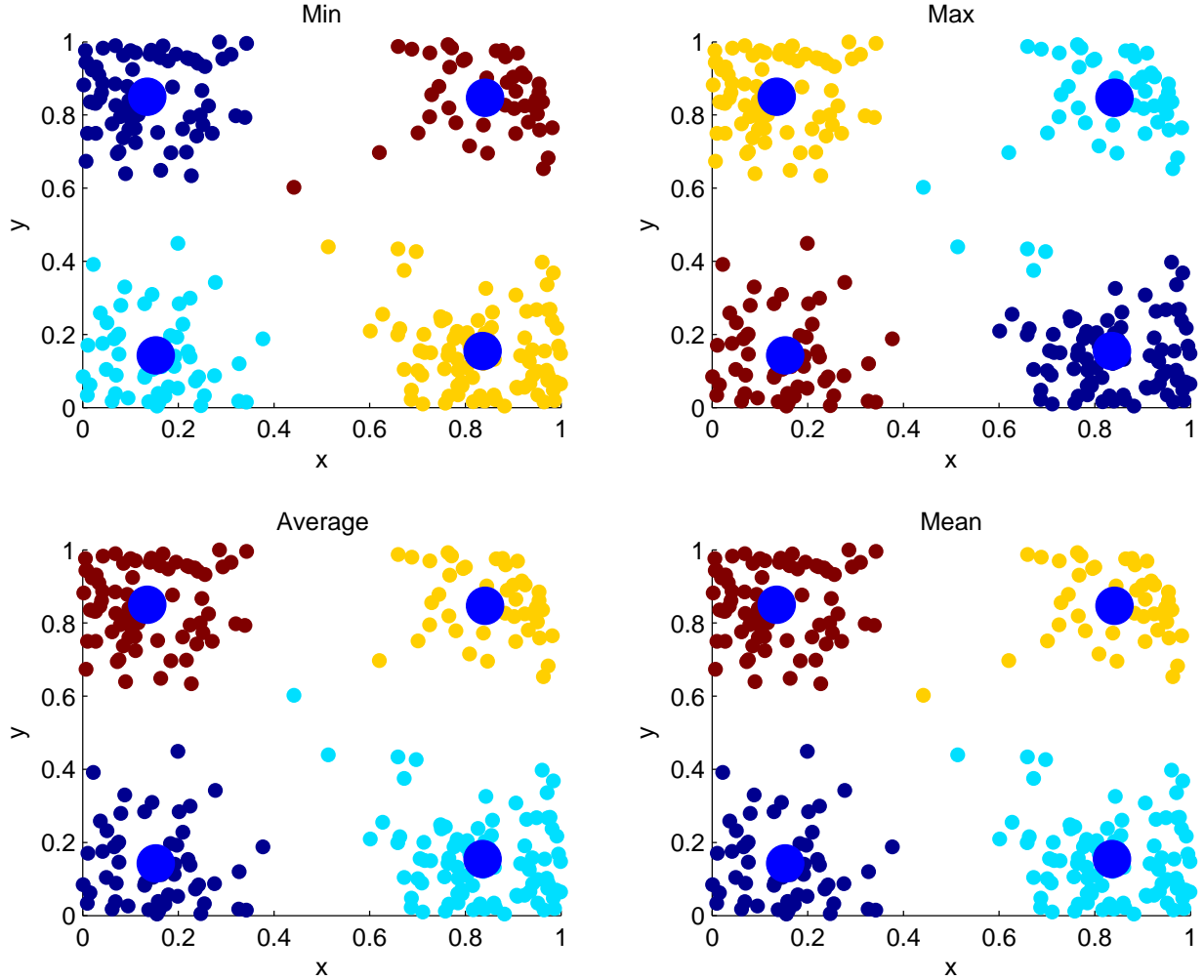


Figure 2: Agglomerative hierarchical clustering using different distance functions. The big circles are the centroids per cluster.

Assignment 2

1

See section B in the appendix for our implementation. Figure 2 shows the different clusterings the agglomerative hierarchical clustering method yields for different distances. As we can see, the different distances yield similar clusters but not entirely equal. The data points that are in the middle of the graphs get attributed to different clusters, depending on the distance type used. When we interpret the different clusterings by just looking at the plots, it seems like the minimum and the mean distance yield the best solutions out of the four. But to know for certain an error could be computed to see which distance yields the best results.

2

Figure 3 gives the dendrograms for the agglomerative hierarchical clustering using different distance functions. Let us assume four main clusters again. We can extrapolate information about the distances between

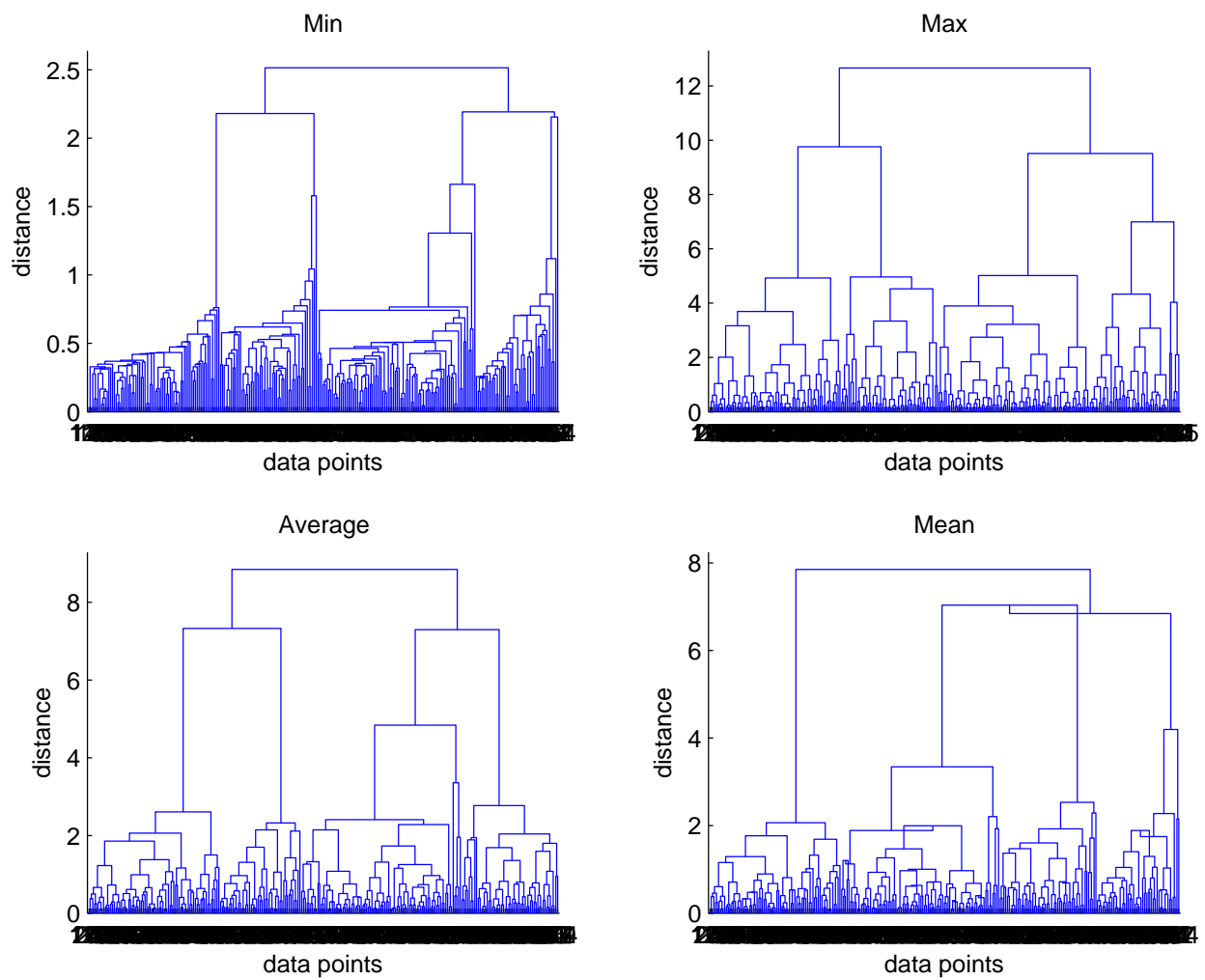


Figure 3: Dendroids for the agglomerative hierarchical clustering using different distance functions.

members of these four clusters from these graphs. For example, the minimum distance dendrogram shows us the minimum distances between members of the four main clusters: the nearest members of the first two main clusters shown in the graph (from left to right) lie about 0.75 units apart, whereas the nearest members of the last two main clusters lie very close (in the order of 0.05 units) to each other. The nearest members between the first two clusters and the last two clusters are about 0.25 units apart.

Similarly we can deduce from the maximum distance graph that the farthest members from the first two main clusters lie about 5 units apart and from the last two main clusters about 3 units. Between the first and last two clusters the farthest members lie about 3 units apart. The average member of the first two clusters lie about 4 units apart and those of the last two clusters about 2 units. The average members between the first and last two clusters lie about 2 units apart.

Finally in the mean graph we can see something different with respect to the previous three. Here the centroid of one cluster is farther away from the three other clusters than their centroids are to each other. So there is one cluster that lies further apart from the other three, whereas the previous three graphs did not show that. This is probably due to the fact the last method uses distances between centroids, which is more robust towards outliers than the other three methods. Therefore the first three graphs do not show this, while the last one does.

Assignment 3

1

Using the code given in section C in the appendix, we computed the J-values for the different clusterings, shown in figure 1

Table 1: J_e -values for different clusters

Clustering	J_e
$\{\{x1, x2, x3\}, \{x4, x5\}\}$	13.1667
$\{\{x2, x3, x5\}, \{x1, x4\}\}$	20.6667
$\{\{x4\}, \{x1, x2, x3, x5\}\}$	17.7500
$\{\{x4, x5\}, \{x1, x2, x3\}\}$	13.1667
$\{\{x3, x5\}, \{x1, x2, x4\}\}$	22.6667

2

Table 1 shows that the clusterings $\{\{x1, x2, x3\}, \{x4, x5\}\}$ and $\{\{x4, x5\}, \{x1, x2, x3\}\}$ minimize the sum-of-squared error. The order of the clusters does not matter (different orderings of a sum yield the same result), as long as they have the same members. That is why these two clusterings yield the same result.

Assignment 4

	S Holland	N Holland	Utrecht	Limburg	N Brabant	Gelderland	Overijssel	Flevoland	Groningen	Zeeland	Friesland	Drenthe
S Holland	0	1.57	3.69	3.99	3.12	3.95	4.41	5.72	4.99	5.34	5.22	5.41
N Holland	1.57	0	2.5	2.45	2.38	2.93	2.95	4.25	3.57	3.81	3.76	3.91
Utrecht	3.69	2.5	0	2.38	3.83	4.21	3.26	3.96	2.17	3.06	3.85	3.78
Limburg	3.99	2.45	2.38	0	3.02	2.64	1.17	1.93	2.11	1.4	1.7	1.6
N Brabant	3.12	2.38	3.83	3.02	0	1.34	2.54	4.58	3.57	4	3.24	3.72
Gelderland	3.95	2.93	4.21	2.64	1.34	0	1.78	3.77	3.58	3.43	2.29	2.84
Overijssel	4.41	2.95	3.26	1.17	2.54	1.78	0	2.13	2.34	1.66	0.84	1.2
Flevoland	5.72	4.25	3.96	1.93	4.58	3.77	2.13	0	3.18	1.3	1.77	1.16
Groningen	4.99	3.57	2.17	2.11	3.57	3.58	2.34	3.18	0	1.9	2.63	2.61
Zeeland	5.34	3.81	3.06	1.4	4	3.43	1.66	1.3	1.9	0	1.52	1.07
Friesland	5.22	3.76	3.85	1.7	3.24	2.29	0.84	1.77	2.63	1.52	0	0.64
Drenthe	5.41	3.91	3.78	1.6	3.72	2.84	1.2	1.16	2.61	1.07	0.64	0

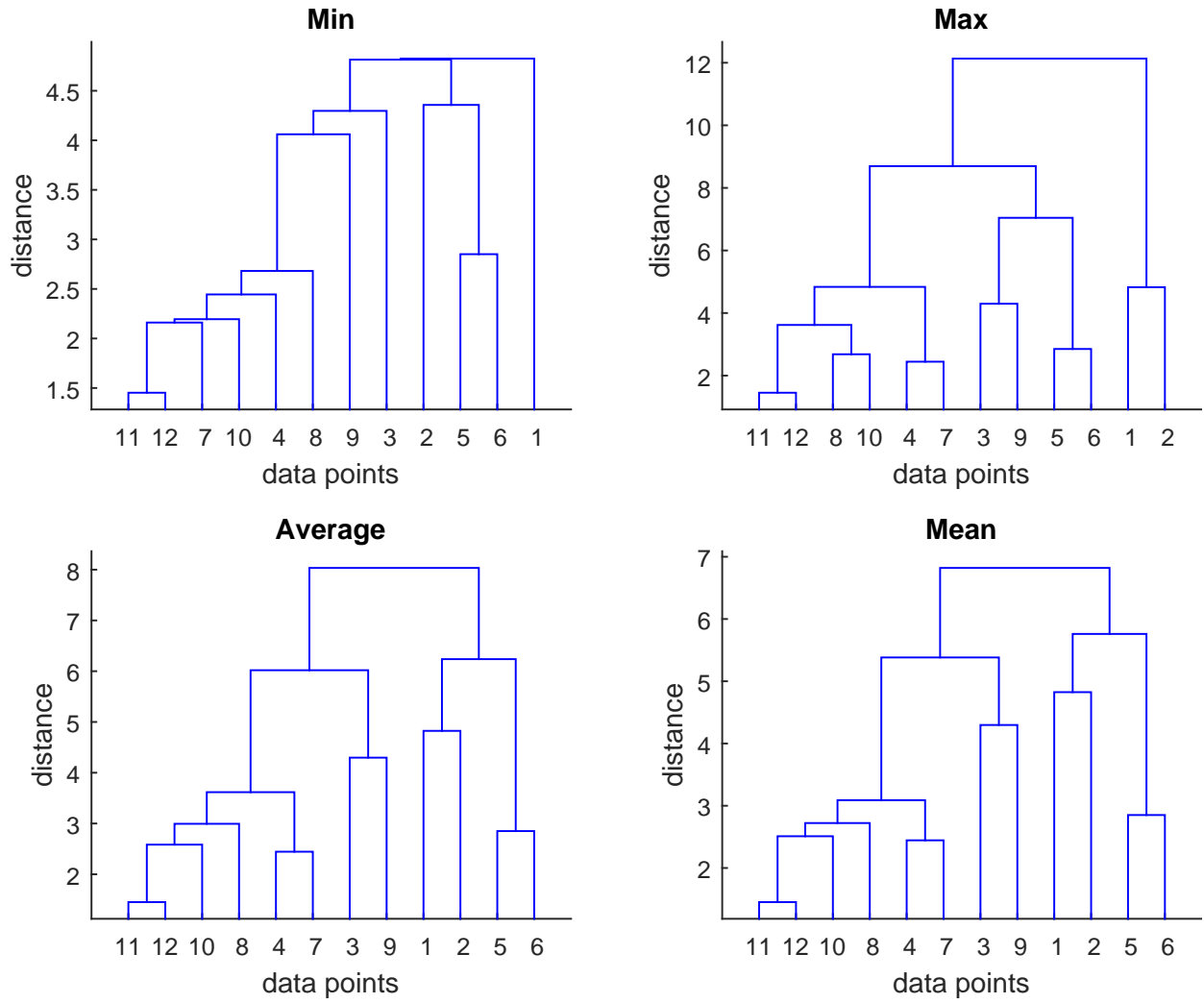


Figure 4: Dendroids for the agglomerative hierarchical clustering of the provinces using different distance functions.

Appendix

A Assignment 1

../Code/Ass1.m

```
1 close all;
2 load('cluster_data.mat', 'cluster_data');
3 dat = cluster_data;
4 k = size(dat, 2);
5 % Calculate the minkowski distances between the points for k dimensions
6 dist = pdist2(dat, dat, 'minkowski', k);
7
8 % Make a new figure for every t-value and plot the relevant connections
9 for t = 0.00 : 0.05 : 0.25
10     figure;
11     hold on;
12     plot(dat(:,1), dat(:,2), 'o');
13     % Loop over all points and plot the connections when the distance
14     % between two points is smaller than t
15     for point = 1 : length(dat)
16         for point2 = point+1 : length(dat)
17             if dist(point, point2) < t
18                 plot([dat(point,1) dat(point2,1)], [dat(point,2) dat(point2,2)])
19             end
20         end
21     end
22     xlabel('x'); ylabel('y'); title(['t = ' num2str(t)]);
23     print(sprintf(['../Report/Ass1-' num2str(t*100)]), '-depsc');
24     hold off;
25 end
```

B Assignment 2

../Code/Ass2.m

```
1 X = cluster_data;
2
3 % Calculate the clusters using linkage with different distance measures
4 c_min = cluster(linkage(squareform(pdist(cluster_data))), 'single'), 'maxclust', 4);
5 c_max = cluster(linkage(squareform(pdist(cluster_data))), 'complete'), 'maxclust', 4);
6 c_avg = cluster(linkage(squareform(pdist(cluster_data))), 'average'), 'maxclust', 4);
7 c_mean = cluster(linkage(squareform(pdist(cluster_data))), 'centroid'), 'maxclust', 4);
8
9 % Plot the clusters and their centroids
10 hold on;
11 figure();
12 subplot(2,2,1)
13 scatter(X(:,1), X(:,2), [], c_min, 'filled');
14 hold on;
15 for group = 1:4
16     plot(mean(X(c_min == group,1)), mean(X(c_min == group,2)), 'o', 'MarkerSize', 15, '
17         MarkerFacecolor', 'b');
18 end
19 xlabel('x'); ylabel('y'); title('Min');
20 subplot(2,2,2)
21 scatter(X(:,1), X(:,2), [], c_max, 'filled')
22 hold on;
23 for group = 1:4
24     plot(mean(X(c_max == group,1)), mean(X(c_max == group,2)), 'o', 'MarkerSize', 15, '
25         MarkerFacecolor', 'b');
```

```

24 end
25 xlabel('x');ylabel('y');title('Max');
26 subplot(2,2,3)
27 scatter(X(:,1),X(:,2),[],c_avg,'filled')
28 hold on;
29 for group = 1:4
30     plot(mean(X(c_min == group,1)),mean(X(c_min == group,2)), 'o', 'MarkerSize', 15, '
        MarkerFacecolor', 'b');
31 end
32 xlabel('x');ylabel('y');title('Average');
33 subplot(2,2,4)
34 scatter(X(:,1),X(:,2),[],c_mean,'filled')
35 hold on;
36 for group = 1:4
37     plot(mean(X(c_min == group,1)),mean(X(c_min == group,2)), 'o', 'MarkerSize', 15, '
        MarkerFacecolor', 'b');
38 end
39 xlabel('x');ylabel('y');title('Mean');
40 print(sprintf(' ../ Report/ Ass2_1 '), '-depsc');
41
42 % Plot the dendograms
43 figure();
44 subplot(2,2,1); dendrogram(linkage(squareform(pdist(cluster_data)), 'single'), 270);
45 xlabel('data points');ylabel('distance');title('Min');
46 subplot(2,2,2); dendrogram(linkage(squareform(pdist(cluster_data)), 'complete'), 270);
47 xlabel('data points');ylabel('distance');title('Max');
48 subplot(2,2,3); dendrogram(linkage(squareform(pdist(cluster_data)), 'average'), 270);
49 xlabel('data points');ylabel('distance');title('Average');
50 subplot(2,2,4); dendrogram(linkage(squareform(pdist(cluster_data)), 'centroid'), 270);
51 xlabel('data points');ylabel('distance');title('Mean');
52 print(sprintf(' ../ Report/ Ass2_2 '), '-depsc');

```

C Assignment 3

../Code/Ass3.m

```

1 x1 = [0 0];
2 x2 = [2 3];
3 x3 = [1 4];
4 x4 = [4 2];
5 x5 = [3 0];
6
7 % {{x1, x2, x3}, {x4, x5}}
8 m1 = 1/3 * (x1+x2+x3);
9 m2 = 1/2 * (x4+x5);
10 J1 = norm(x1-m1).^2+norm(x2-m1).^2+norm(x3-m1).^2+norm(x4-m2).^2+norm(x5-m2).^2
11
12 % {{x2, x3, x5}, {x1, x4}}
13 m1 = 1/3 * (x2+x3+x5);
14 m2 = 1/2 * (x1+x4);
15 J2 = norm(x2-m1).^2+norm(x3-m1).^2+norm(x5-m1).^2+norm(x1-m2).^2+norm(x4-m2).^2
16
17 % {{x4}, {x1, x2, x3, x5}}
18 m1 = x4;
19 m2 = 1/4 * (x1+x2+x3+x5);
20 J3 = norm(x4-m1).^2+norm(x1-m2).^2+norm(x2-m2).^2+norm(x3-m2).^2+norm(x5-m2).^2
21
22 % {{x4, x5}, {x1, x2, x3}}
23 m1 = 1/2 * (x4+x5);
24 m2 = 1/3 * (x1+x2+x3);
25 J4 = norm(x4-m1).^2+norm(x5-m1).^2+norm(x1-m2).^2+norm(x2-m2).^2+norm(x3-m2).^2
26

```



```

27 % {{x3, x5}, {x1, x2, x4}}
28 m1 = 1/2 * (x3+x5);
29 m2 = 1/3 * (x1+x2+x4);
30 J5 = norm(x3-m1).^2+norm(x5-m1).^2+norm(x1-m2).^2+norm(x2-m2).^2+norm(x4-m2).^2

```

D Assignment 4

../Code/Ass4.m

```

1  close all;
2  load('exercise4_data.mat', 'data');
3
4  % Calculate the z-score for every data point
5  zdata = zscore(data);
6
7  % Compute the dissimilarity matrix
8  D = squareform(pdist(zdata));
9  writetable(table(round(D,2)), '../Report/dissMatrixAss4.csv')
10
11 figure();
12 subplot(2,2,1); dendrogram(linkage(D, 'single'));
13 xlabel('data points'); ylabel('distance'); title('Min');
14 subplot(2,2,2); dendrogram(linkage(D, 'complete'));
15 xlabel('data points'); ylabel('distance'); title('Max');
16 subplot(2,2,3); dendrogram(linkage(D, 'average'));
17 xlabel('data points'); ylabel('distance'); title('Average');
18 subplot(2,2,4); dendrogram(linkage(D, 'centroid'));
19 xlabel('data points'); ylabel('distance'); title('Mean');
20
21 print(sprintf('../Report/Ass4'), '-depsc');

```