

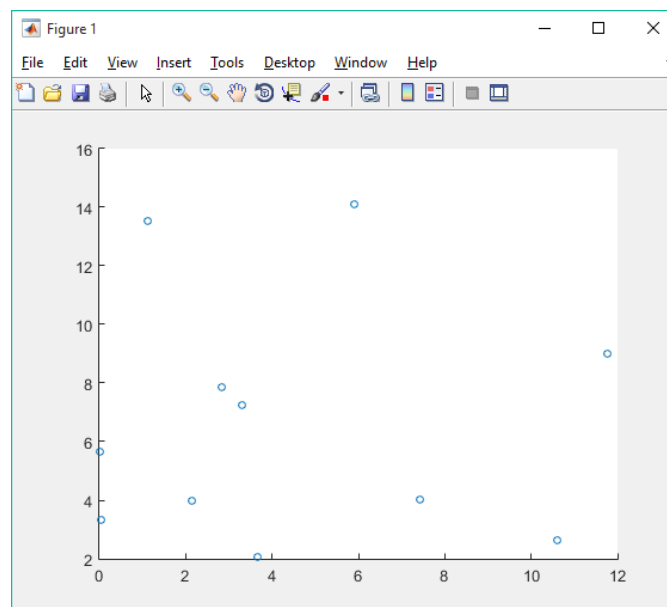
Hierarchical Clustering

1. Load Data and Visualize

- hierarchicalMain.m

```
hierarchicalMain.m* Variables - pairwise_dist
1 %% Load the Data
2 hierarchical_1 = csvread('Clustering_DataSets/Hierarchical/Hierarchical_1.csv');
3 X = hierarchical_1;
4
5 %% Visualize the Data
6 scatter(X(:,1),X(:,2),20);
7
8
```

- Output



2. 4 Function to measure dissimilarity

- Single-Link

```
-Inside\Documents\Machine Learning\Assignment 4.2 Hierarchical Clustering\dSingleLink.m
1 function distance = dSingleLink(X , Y)
2     % this function is to calculate minimum distance between two cluster
3     distance = min(X,Y);
4 end
```

- Complete –Link

```
o-Inside\Documents\Machine Learning\Assignment 4.2 Hierarchical Clustering\dCompleteLink.m
1 function distance = dCompleteLink( X , Y )
2     % this function is to choose the max from 2 cluster
3     distance = max( X , Y );
4 end
```

- Group Average

```
p-Inside\Documents\Machine Learning\Assignment 4.2 Hierarchical Clustering\dAverage.m
1  function distance = dAverage( X , Y)
2      % this function is to calculate the average between to point
3  -      distance = (X + Y) / 2;
4  -  end
```

iv. Centroid-based

```
p-Inside\Documents\Machine Learning\Assignment 4.2 Hierarchical Clustering\dCentroid.m
1  function distance = dCentroid( X,Y,N,R,i,j,v )
2      % this function is to calculate the distance using centroids approach
3
4  -      K = N(R(i))+N(R(j));
5  -      distance = (N(R(i)).*X+N(R(j)) .*Y - (N(R(i)).*N(R(j))*v) ./K) ./K;
6
7  -  end
```

3. Implement Hierarchical Clustering

```

Inside\Documents\Machine Learning\Assignment 4.2 Hierarchical Clustering\agglomerative.m
1 function result = agglomerative(Y, method)
2 % this function is hierarchical clustering implementation
3 % Y is distance from pdist function
4 % method is alogrithm approach to merge cluster
5
6 n = size(Y,2); % length of Y (pdist)
7 m = ceil(sqrt(2*n)); % total of class
8 result = zeros(m-1,3); % allocate the output matrix.
9
10 N = zeros(1,2*m-1); % N how many points are contained in each cluster.
11 N(1:m) = 1;
12 n = m; % since m is changing, we need to save m in n.
13 R = 1:n; % is a index of class
14
15 % Square the distances so updates are easier.
16 if any(strcmp(method,'centroid'))
17     Y = Y .* Y; % only on centroids method
18 end
19
20 % repeat until cluster = 1
21 for s = 1:(n-1)
22     [v, k] = min(Y); % search the smallest values and return the index and it
23
24     i = floor(m+1/2-sqrt(m^2-m+1/4-2*(k-1))); % Search the row to merge from
25     j = k - (i-1)*(m-i)/2+i; % Search the coloumn to merge from index
26
27     result(s,:) = [R(i) R(j) v]; % get the class from index of i and j and al
28
29     % Update Y.
30     I1 = 1:(i-1);
31     I2 = (i+1):(j-1);
32     I3 = (j+1):m; % these are temp variables
33     U = [I1 I2 I3]; % remaining point not yet clustered.
34     I = [I1.*(m-(I1+1)/2)-m+i i*(m-(i+1)/2)-m+I2 i*(m-(i+1)/2)-m+I3]; %index
35     J = [I1.*(m-(I1+1)/2)-m+j I2.*(m-(I2+1)/2)-m+j j*(m-(j+1)/2)-m+I3]; %inde
36
37     switch method
38     case 'single' % single linkage
39         Y(I) = dSingleLink(Y(I),Y(J));
40     case 'complete' % complete linkage
41         Y(I) = dCompleteLink(Y(I),Y(J));
42     case 'average' % weighted average linkage
43         Y(I) = dAverage(Y(I), Y(J));
44     case 'centroid' % centroid linkage
45         Y(I) = dCentroid(Y(I),Y(J),N,R,i,j,v);
46     end
47     J = [J i*(m-(i+1)/2)-m+j];
48     Y(J) = []; % no need for the cluster information about j.
49
50     % update m, N, R
51     m = m-1;
52     N(n+s) = N(R(i)) + N(R(j));
53     R(i) = n+s;
54     R(j:(n-1))=R((j+1):n);
55 end
56
57 if any(strcmp(method,'centroid'))
58     result(:,3) = sqrt(result(:,3));
59 end
60
61 result(:,[1 2])=sort(result(:,[1 2]),2); % sort from dissimilarity
62 end

```

- hierarchicalMain.m

```
9 %% do hierarchical clustering
10 % single-link method
11 - Y = pdist(X); %calculate euclidean distance
12 - Z = agglomerative(Y,'single');
13 - disp(Z);
14
15 % complete-link method
16 - Y = pdist(X); %calculate euclidean distance
17 - Z = agglomerative(Y,'complete');
18 - disp(Z);
19
20 % average group method
21 - Y = pdist(X); %calculate euclidean distance
22 - Z = agglomerative(Y,'average');
23 - disp(Z);
24
25 % centroid method
26 - Y = pdist(X); %calculate euclidean distance
27 - Z = agglomerative(Y,'centroid');
28 - disp(Z);
```

Output :

- Single-Link

```
Command Window
>> %% do hierarchical clustering
% single-link method
Y = pdist(X); %calculate euclidean distance
Z = agglomerative(Y,'single');
disp(Z);
    8.0000    10.0000    0.7768
    7.0000    11.0000    2.1869
    9.0000    13.0000    2.3226
    3.0000    14.0000    2.4257
    1.0000     2.0000    3.4645
   12.0000    15.0000    3.4712
   16.0000    17.0000    4.2426
    5.0000     6.0000    4.8313
   18.0000    19.0000    5.9168
    4.0000    20.0000    6.4684
```

- Complete-Link

```
Command Window
>> % complete-link method
Y = pdist(X); %calculate euclidean distance
Z = agglomerative(Y,'complete');
disp(Z);
    8.0000    10.0000    0.7768
    7.0000    11.0000    2.1869
    9.0000    13.0000    2.6952
    1.0000     2.0000    3.4645
    5.0000     6.0000    4.8313
    3.0000    14.0000    5.0876
   12.0000    17.0000    5.8343
    4.0000    15.0000    6.5989
   16.0000    18.0000   12.2626
   19.0000    20.0000   14.4433
```

- Average

```
Command Window

>> % average group method
Y = pdist(X); %calculate euclidean distance
Z = agglomerative(Y,'average');
disp(Z);
    8.0000    10.0000    0.7768
    7.0000    11.0000    2.1869
    9.0000    13.0000    2.5089
    1.0000     2.0000    3.4645
   12.0000    14.0000    4.0415
    3.0000    16.0000    4.8028
    5.0000     6.0000    4.8313
    4.0000    15.0000    6.5337
   17.0000    19.0000    8.6695
   18.0000    20.0000   10.5188
```

- Centroid

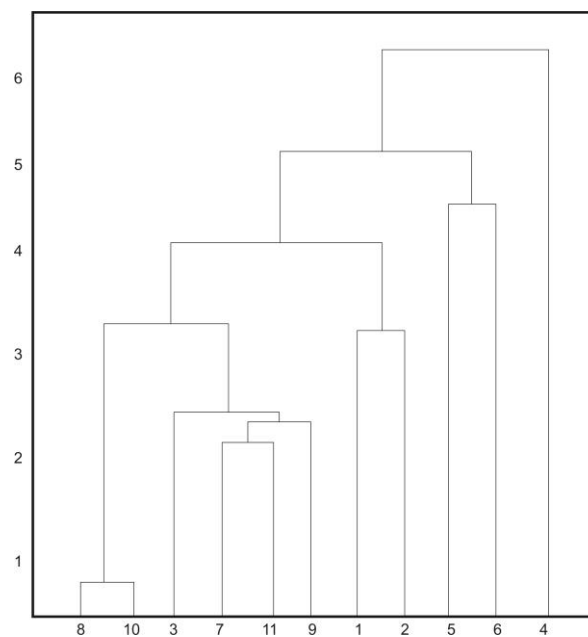
```
Command Window

>> % centroid method
Y = pdist(X); %calculate euclidean distance
Z = agglomerative(Y,'centroid');
disp(Z);
    8.0000    10.0000    0.7768
    7.0000    11.0000    2.1869
    9.0000    13.0000    2.2658
    1.0000     2.0000    3.4645
    3.0000    14.0000    3.6727
   12.0000    16.0000    4.1225
    5.0000     6.0000    4.8313
    4.0000    15.0000    6.3002
   17.0000    19.0000    7.9334
   18.0000    20.0000    8.7954
```

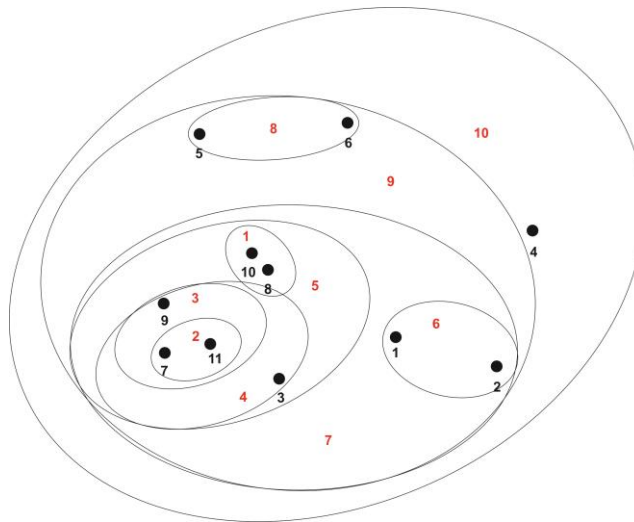
4. Visualize

- Single-Link

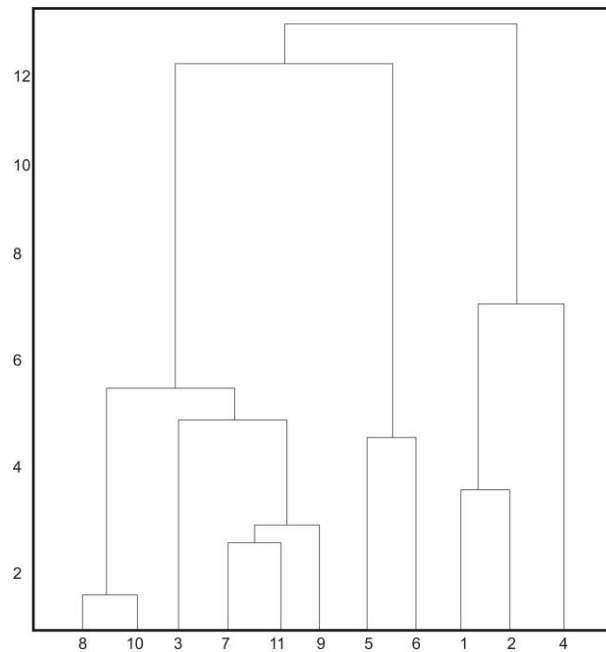
- Dendrogram



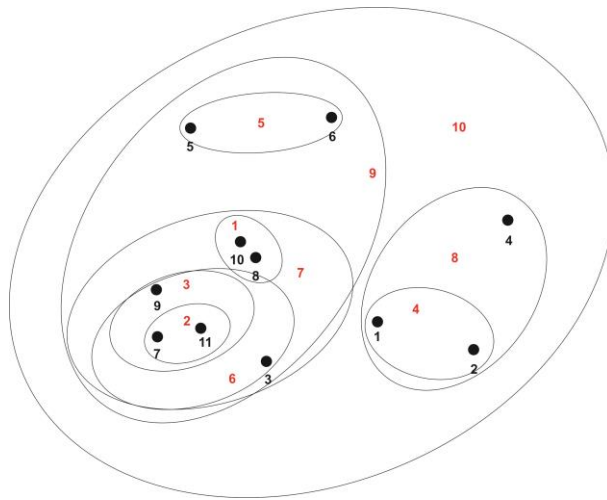
- Nested Cluster



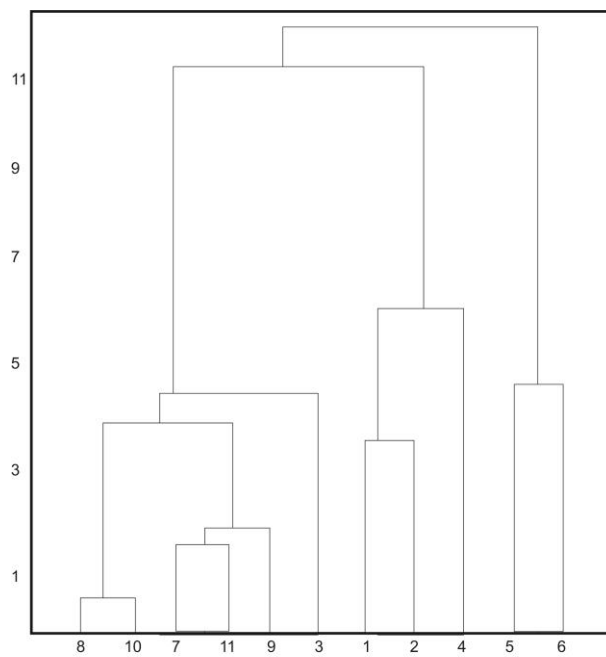
- Complete-Link
 - Dendrogram



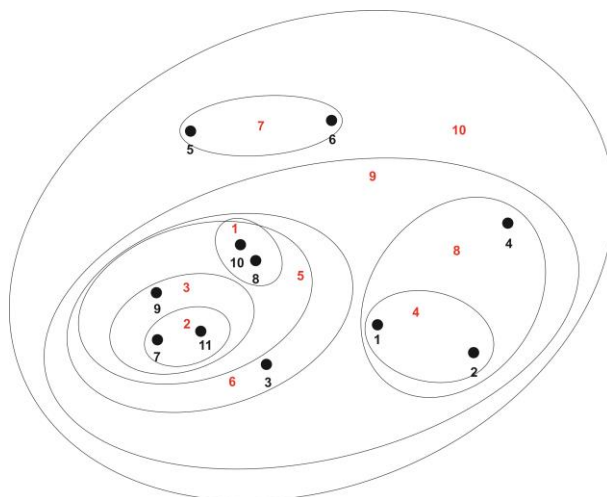
- Nested Cluster



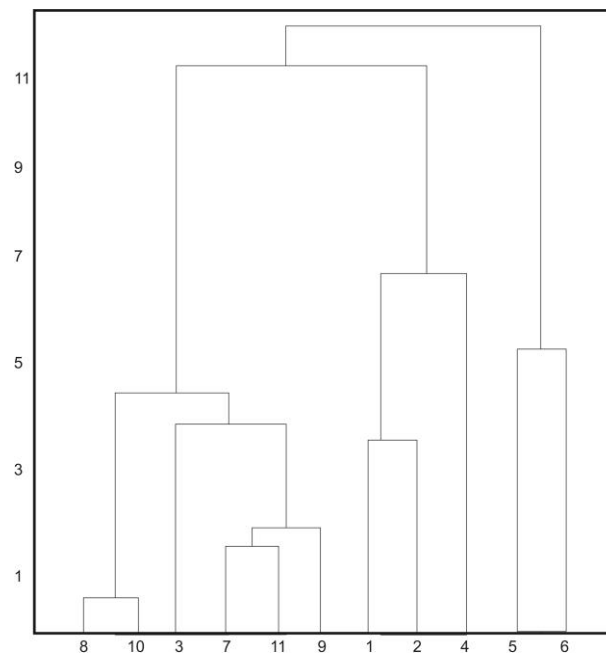
- Average
 - Dendrogram



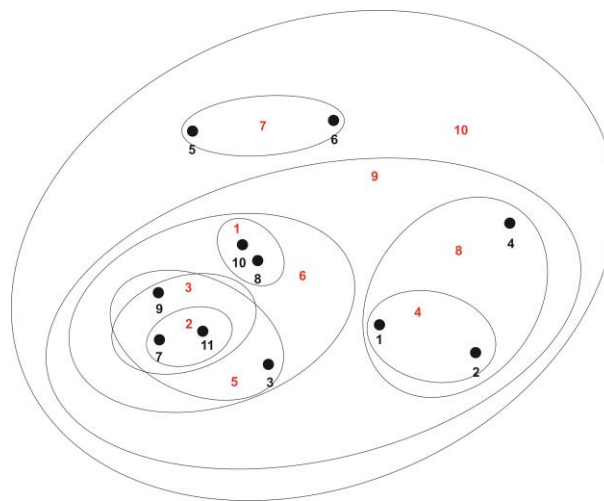
- Nested Cluster



- Centroid
 - Dendrogram



- Nested Cluster



5. Compare Cluster

Analisis :

Untuk single-link list sangat cocok untuk cluster yang benar-benar terpisah. Tetapi jika saling berdekatan cenderung mengidentifikasi secara berurutan. Adapun untuk complete-link pengclusteran cenderung memiliki diameter yang sama karena menghitung jarak terjauh dari cluster sehingga akan sensitive terhadap data percilan. Kemudian untuk average mempertimbangkan nilai kesluruhan anggota cluster sehingga lebih sentral untuk mengelompokkan cluster. Dan method centroid sama seperti average namun perhitungan melibatkan centroid pada cluster.