

Homework 3

Posted: March 18, 2019

Due: April 1, 2019

1. Let $X = \{-10, -9, -7, -6, -1, 0, 2, 6, 7, 9\}$ be a set of 10 numbers in \mathbb{R} . Draw the single-link dendrogram for this set.
2. Let $S = \{x_1, \dots, x_n\}$ be a set of n points in \mathbb{R} . Prove that the largest height of a dendrogram constructed for S is $n - 1$ and the smallest is $\lceil \log_2 n \rceil$.
3. Consider the set of points that consists of two “entangled spirals” shown in Figure 1. Extract the coordinates of the points from Figure 1 and compute the `dist` object using the Euclidean metric. Show that the sets of points of the two curves can be separated by using the single link method; in other words, it is possible to cut the dendrogram obtained by this method such that the points of the two curves belong to two different clusters. Verify that this is not possible if we use the complete link and explain why is this the case,

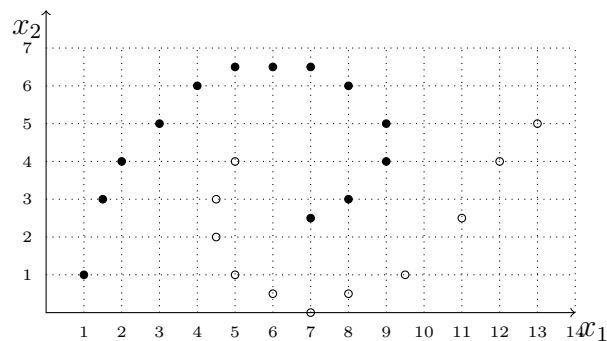


Figure 1: Two entangled curves

If m is a matrix of dissimilarities between objects, then m is a symmetric matrix that has 0s as its diagonal elements. To create a `dist` object as required by the `hclust` function we can use the coercion `as.dist`. For example, starting from a complete matrix of distances m :

```
m
      [,1] [,2] [,3] [,4] [,5] [,6]
[1,]    0    1    5    8    3    4
[2,]    1    0    2    2    8    6
[3,]    5    2    0    1    7    9
[4,]    8    2    1    0    4    2
[5,]    3    8    7    4    0    7
[6,]    4    6    9    2    7    0
```

we can write `d <- as.dist(m)`.

This would yield the lower half of m :

```
> d

  1 2 3 4 5
2 1
3 5 2
4 8 2 1
5 3 8 7 4
6 4 6 9 2 7
```

4. Starting from the matrix m given above construct a single-link clustering and a complete-link clustering. Note that the dendrograms of these clusterings are distinct. What is the least number of entries of m that you need to change to obtain the same dendrogram using these two distinct methods?
5. Starting from the following dissimilarity matrix defined on a set of objects $\{x_1, x_2, x_3, x_4, x_5\}$

	x_1	x_2	x_3	x_4	x_5
x_1	0	2	4	2	5
x_2	2	0	3	4	5
x_3	4	3	0	2	4
x_4	2	4	2	0	6
x_5	5	5	4	6	0

apply the DIANA divisive algorithm. Compare the results with the dendrograms obtained using the single-link and complete link methods from the **R** `cluster` package. Note that the dissimilarity matrix will have to be converted to a `dist` object using `as.dist` before the function `hclust` can be applied.