

The shape of a data set  $S \subseteq \mathbb{R}^n$  can be analysed using cluster analysis.

Compute

$$c = \frac{1}{|S|} \sum_{v \in S} v \quad (\text{Centre})$$

and for  $r > 0$

$$S_r = \{v \in S : \|v - c\| \geq r\}$$

and for  $\Sigma > 0$  the graph

$$G_\Sigma(S_r) : \begin{array}{l} \text{vertices} = S_r \\ \text{edge } \overset{u}{\bullet} \text{---} \overset{v}{\bullet} \text{ if } \|u - v\| \leq \Sigma. \end{array}$$

For fixed  $\Sigma > 0$  and for

$$r_1 > r_2 > r_3 > \dots$$

we get a chain of inclusions of graphs

$$G_\Sigma(S_{r_1}) \subseteq G_\Sigma(S_{r_2}) \subseteq G_\Sigma(S_{r_3}) \subseteq \dots$$

For  $\Sigma$  in a "stable range" the resulting barcode (for connected components) may be informative.

# Hierarchical Clustering Algorithms

Given  $n$  items to be clustered, and an  $n \times n$  distance matrix, do:

- ① Start with  $n$  clusters, each containing 1 item. Let the distances between clusters equal to the distances between items.
- ② Find the closest pair of clusters and merge them, so that you have one fewer clusters.
- ③ Compute "distances" between the new cluster and each of the remaining old clusters.
- ④ Repeat ② and ③ until there is just a single cluster with  $n$  items.
- ⑤ Return the output as a dendrogram or barcode.

There are various ways to define distance  $d(X, Y)$  between clusters  $X, Y$  in step ③.

Single-linkage:  $d(x, y) = \min_{\substack{x \in X \\ y \in Y}} d(x, y)$

Complete-linkage:  $d(x, y) = \max_{\substack{x \in X \\ y \in Y}} d(x, y)$

Average-linkage:

$$d(x, y) = \frac{1}{|X||Y|} \sum_{\substack{x \in X \\ y \in Y}} d(x, y)$$

# Algorithm for single linkage

① set  $m := 0$   
 $L[0] = 0$  (we'll say the level is  $L[m]$  at stage  $m$ )

$D_{(0)} = (d_0(i, j))$  The initial  $n \times n$  matrix of distances

② while  $m < n$  do:

③ find  $1 \leq t, s \leq n-m$  such that

$$d_m(t, s) = \min_{1 \leq i, j \leq n-m} d_m(i, j).$$

set  $L[m+1] = d_m(t, s)$

④ Let  $D_{(m+1)} = (d_{m+1}(i, j))$  be the  $(n-m-1) \times (n-m-1)$  matrix obtained

$D_{(m)}$  by

— delete row  $s$ , row  $t$ , column  $s$ , column  $t$

— re-index so that old row  $i'$  / column  $i'$  becomes the new row  $i$  / column  $i$

— add a new final row and column  $n$  with

$$d_{m+1}(n-m-1, j) = d_{m+1}(j, n-m-1) = \min(d_m(s, j'), d_m(t, j')).$$

end do

⑤ Return the list  $L[0], L[1], \dots, L[n-1]$  as a barcode.

### Example

$D_{(0)} =$

	1	2	3	4	5
1	0	11	10	14	22
2	11	0	3	13	21
3	10	3	0	12	20
4	14	13	12	0	16
5	22	21	20	16	0

$n=0$

$L=[0]$

merge 2 & 3

	1	4	5	6
1	0	14	22	10
4	14	0	16	12
5	22	16	0	20
6	10	12	20	0

$n=1$

$L=[0, 3]$

merge 1 & 6

	4	5	7
4	0	16	12
5	16	0	20
7	12	20	0

$n=2$

$L=[0, 3, 10]$

merge 4 2 7

$n = 3$

	5	8
5	0	16
8	16	0

$L = [0, 3, 10, 12]$

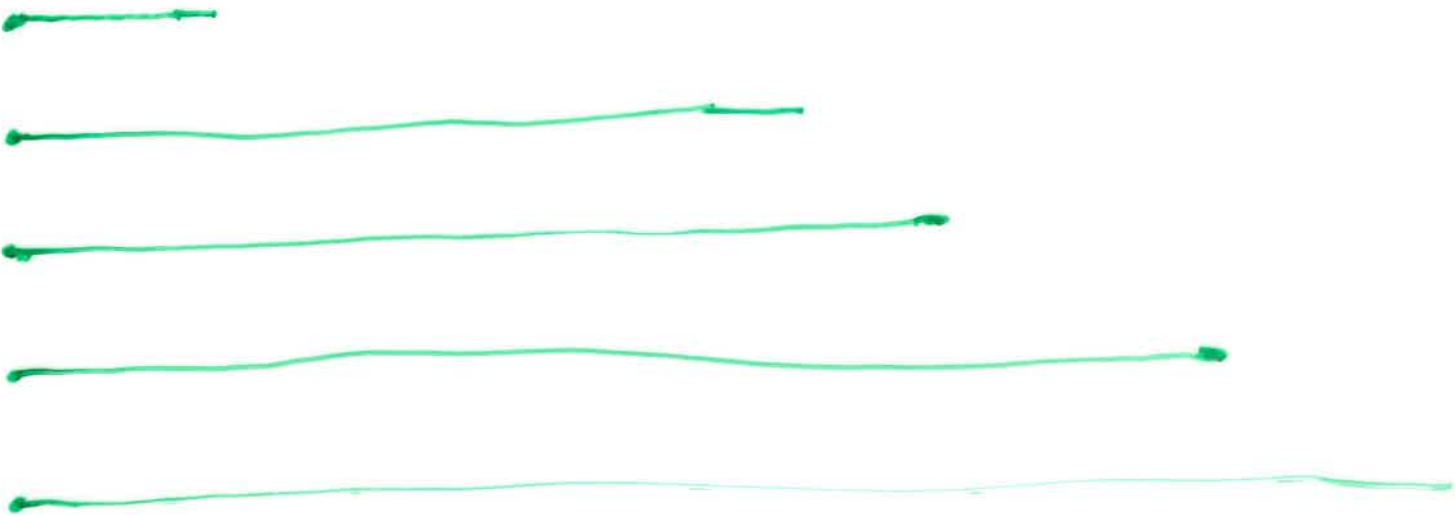
merge 5 2 8

$n = 4$

9
0

$L = [0, 3, 10, 12, 16]$

Return the bar code



0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18