

L8:

Hierarchical Agglomerative Clustering

Sep 15, 2025
Data Mining



Jeff M. Phillips

When data is well clusterable,
most clustering algorithms/formulations
will work well.

When data is not well clusterable,
no clustering algorithm/formulation/
will find a good clustering.

Clustering

Output

$$C(X) = \{S_1, S_2, \dots, S_{|C|}\}$$

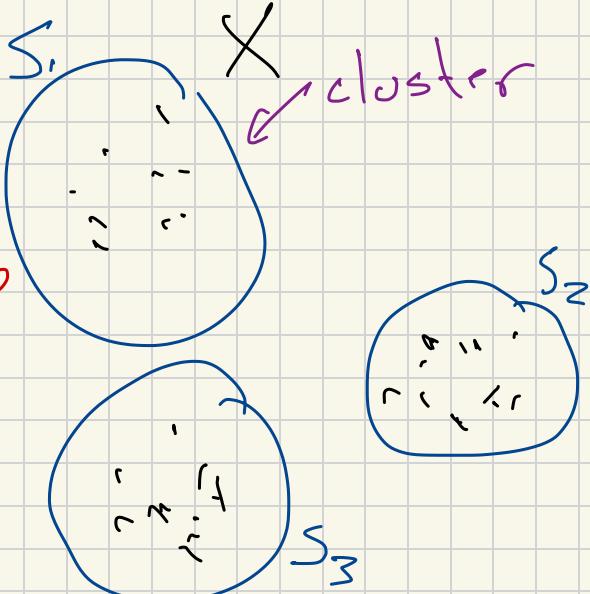
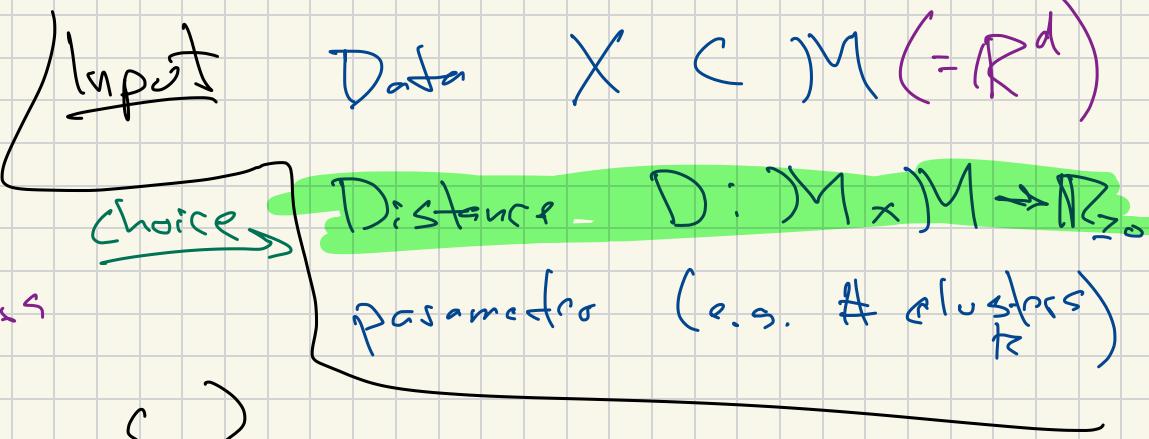
↑ clusters

= each $S_i \subset X$

- each pair $S_j \cap S_{j'} = \emptyset$ hard clustering
 $j \neq j'$

$$= \bigcup_{j=1}^{|C|} S_j = X$$

(outliers)



What makes a good clustering?

width

For each S_j , $x, x' \in S_j$, $D(x, x')$ small

split

For S_i, S_j , (most) $x \in S_i, x' \in S_j$

$D(x, x')$ large

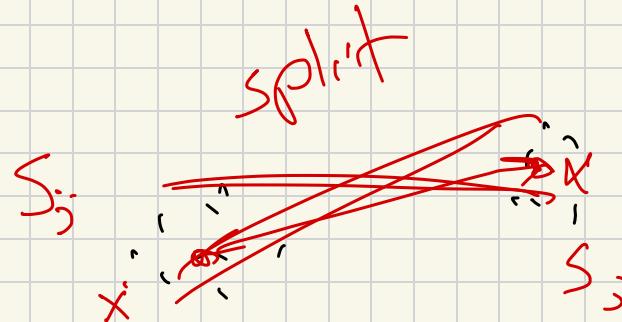
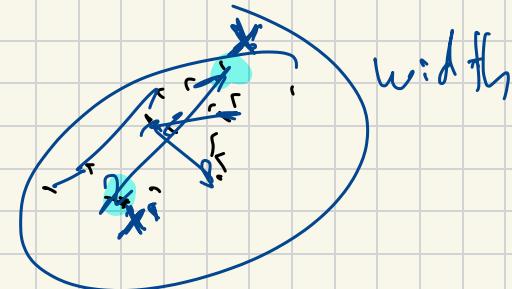
overall

width - split

as

$\frac{\text{width}}{\text{split}}$

small



Hierarchical Agglomerative Clustering (HAC)

If two clusters (^{points}) are close, → merge them

HAC

0. Each $x_i \in \mathcal{X}$ in own cluster S_i :

1. while (two clusters are close enough)

Find closest pair S_i, S_j

Merge $S_i, S_j \rightarrow S = S_i \cup S_j$

$$D(S, S')$$

(1) $\mu = \text{represent}(S)$
 $\mu' = \text{mean}(S')$

$$D(S, S') = D(\mu, \mu')$$

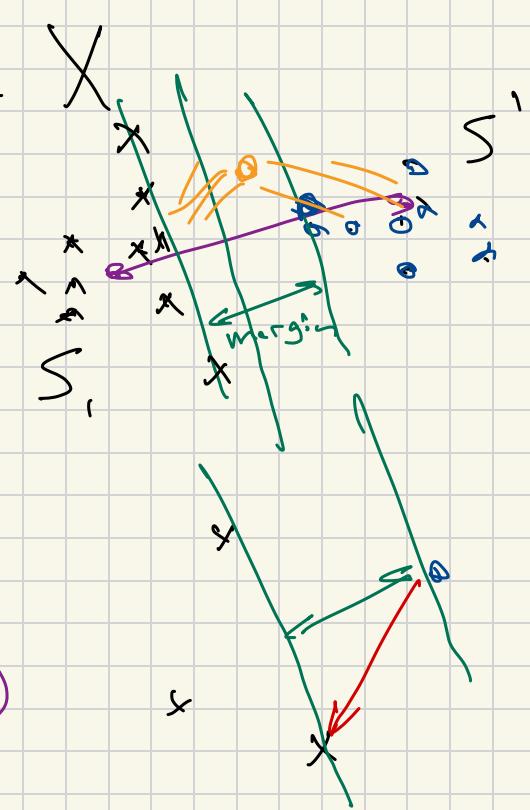
(2)
$$\frac{1}{|S|} \sum_{x \in S} \left(\frac{1}{|S'|} \sum_{x' \in S'} D(x, x') \right)$$

$$= \frac{1}{|S|} \frac{1}{|S'|} \sum_{x \in S} \sum_{x' \in S'} D(x, x')$$

(3) $D(S, S') = \min_{x \in S, x' \in S'} D(x, x')$

(4) $D(S, S') = \max_{x \in S, x' \in S'} D(x, x')$

$$S, S' \subset X$$



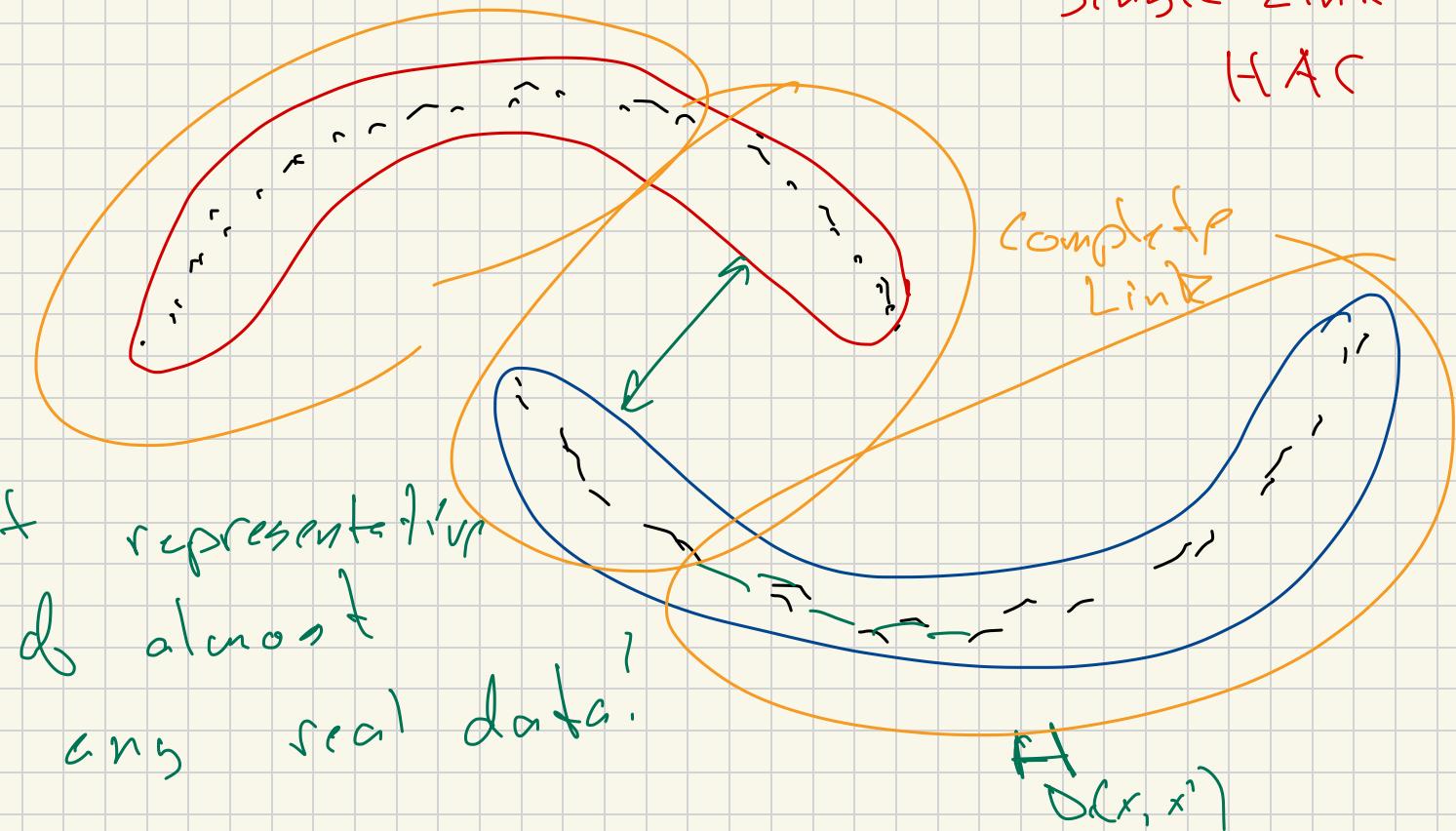
[Single Link]

[Complete Link]

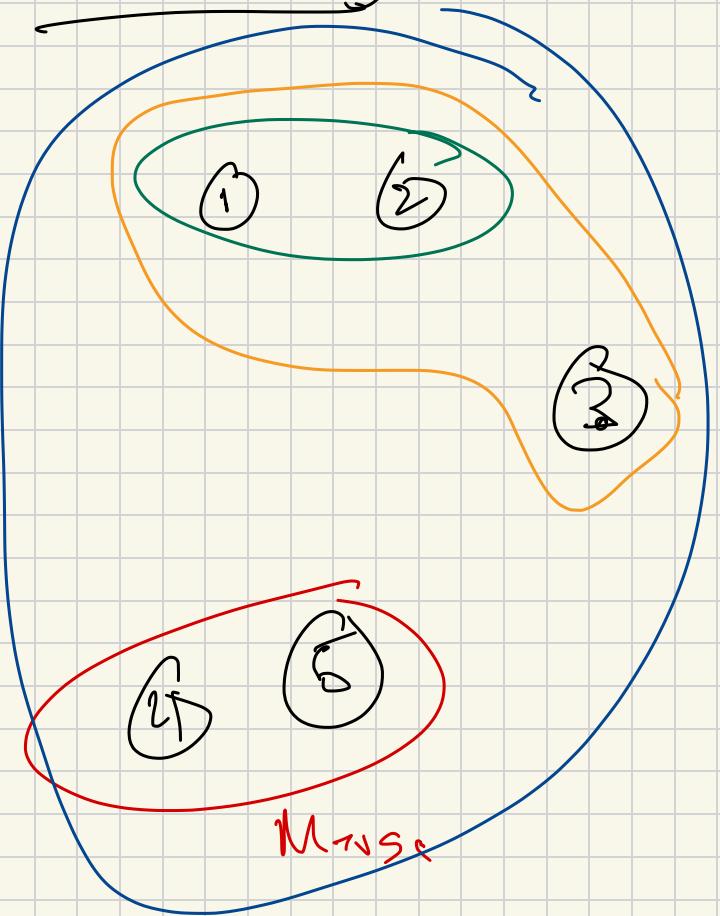
Two Moons Data Set

X

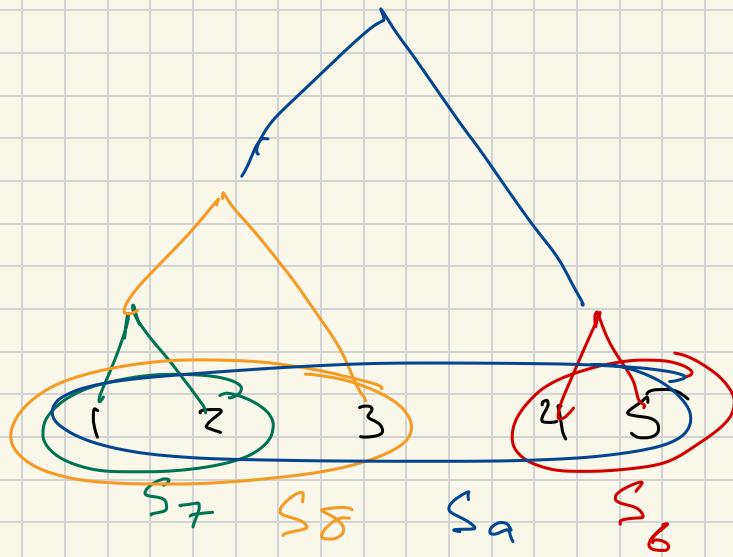
Not
representative
of almost
any real data.



Hierarchy



Single Link Hierarchy



Phylogenetic Trees

How efficient is HAC?

n rounds of while loop.

- $D(S_j, S_{j'})$ $O(n^2)$ if $(S_j, |S_{j'}|) = \alpha_n$

maybe $O(n^2)$ pairs

$$O(n^2)$$

overall
 $O(n^3)$

representative-based approach

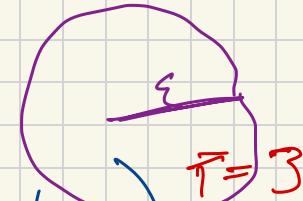
$$\hookrightarrow O(n^2 \log n)$$

Density-based clustering (DBScan)

Two Parameters

ε = radius (what is close)

τ = Min Pts (density parameter)

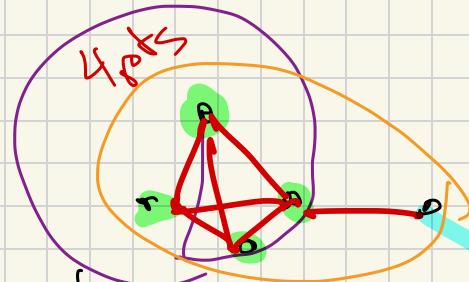


① Find core pts

(w/ T pts in radius ε)

outliers

② Edges
 $\text{dist}(x, x') < \varepsilon$
 $\rightarrow x \text{ or } x' \text{ is core}$



③ clusters are connected components

