

# Hierarchical Agglomerative Clustering

February 5, 2018

When data is easily "clusterable",  
most clustering algorithms work quickly  
and well.

When data is not easily "clusterable",  
then no algorithm will find good  
clusters

# What is clustering?

Input Data set  $X = \{x_1, x_2, \dots, x_n\}$

distance  $d: X \times X \rightarrow \mathbb{R}$

↑  
input

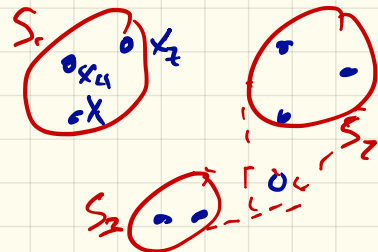
$X \subset M = \mathbb{R}^d$   $d: M \times M \rightarrow \mathbb{R}$

Output  $\mathcal{C}(X) = \{S_1, S_2, \dots, S_k\}$

(1)  $S_i \subset X$  ← ith "cluster"

(2)  $S_i \cap S_j = \emptyset$   $i \neq j$  "hard clustering" or "soft"

(3)  $\bigcup_i S_i = S_1 \cup S_2 \dots \cup S_k = X$



Usually some objective function

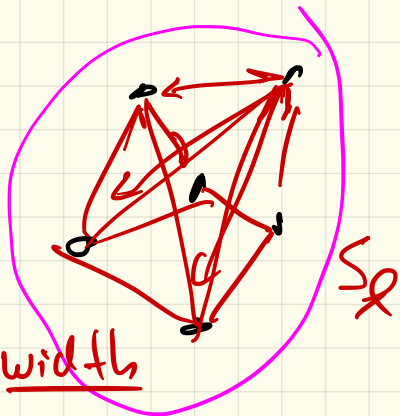
# Clustering Objective

large:  $\frac{\text{split}}{\text{width}}$  or split-width

$$X \subset \mathbb{R}^2$$

$$d = \|\cdot - \cdot\|$$

Euclidean

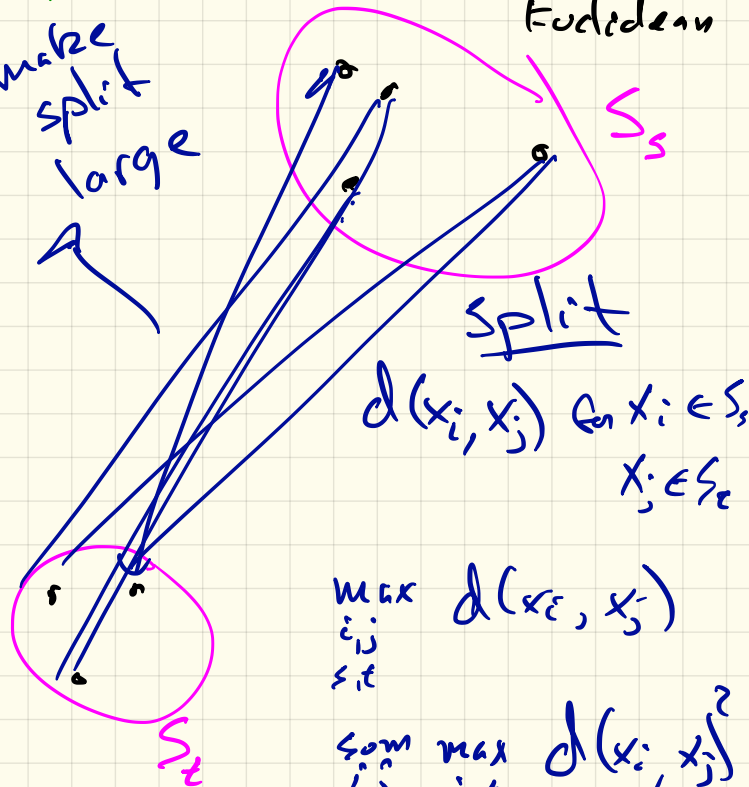


$$d(x_i, x_j) \text{ for } x_i, x_j \in S_e$$

$$\sum_{i,j} d(x_i, x_j)$$

$$\left. \begin{array}{l} \max_{i,j} d(x_i, x_j) \\ \sum_{i,j} d(x_i, x_j)^2 \end{array} \right\} \rightarrow \text{make width small}$$

make split large



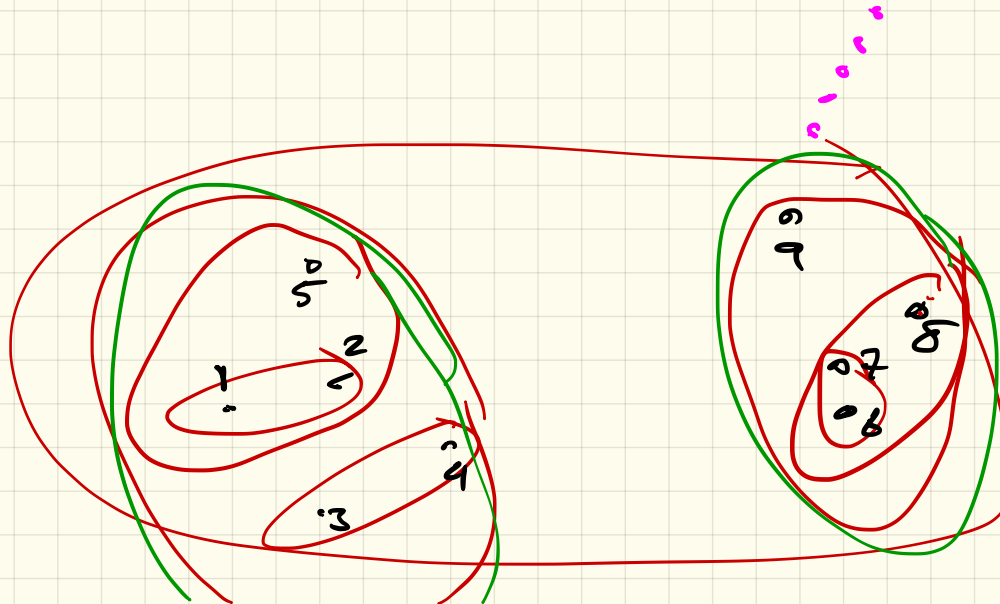
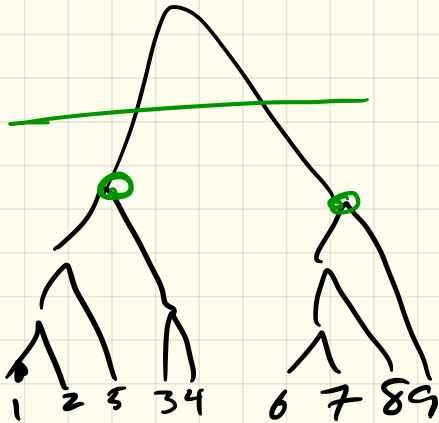
$$d(x_i, x_j) \text{ for } x_i \in S_s, x_j \in S_t$$

$$\max_{\substack{i,j \\ S_t}} d(x_i, x_j)$$

$$\sum_{i,j} \max_{S_t} d(x_i, x_j)^2$$

# Hierarchical Agglomerative Clustering

Algo: If 2 points (clusters) are close enough, put them in same cluster.



0. Each  $x_i \in X$  a separate cluster  $S_i$   $\leftarrow$   $n$  clusters

1. while ( $2$  clusters are close enough)

la. Find closest two clusters  $S_i, S_j$

lb. Merge  $S_i, S_j \rightarrow$  new cluster  $S_k$

How to define  $d(S_i, S_j)$

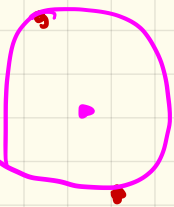
• Define center  $c_i \leftarrow S_i, c_j \leftarrow S_j$   $d(c_i, c_j)$

+ geometric median:  $c_i = \underset{c \in M}{\operatorname{argmin}} \sum_{x \in S_i} d(x, c)$

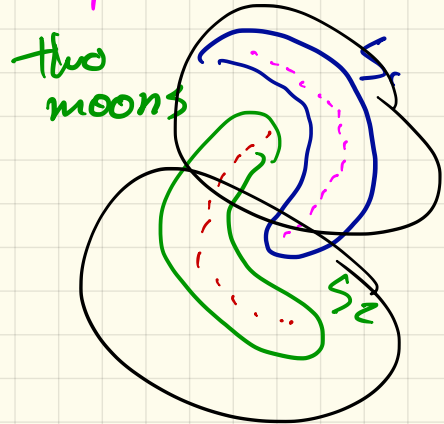
+ mean  $c_i = \underset{c \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{x \in S_i} \|x - c\|^2 = \frac{1}{|S_i|} \sum_{x \in S_i} x$

+ restrict some  $c \in S_i$

+ center of minimum enclosing ball



- Single Link :  $d(S_i, S_j) = \min_{x \in S_i, x' \in S_j} d(x, x')$   
 closest dist between pts in clusters



- Average Link

$$d(S_i, S_j) = \frac{1}{2} \sum_{x \in S_i} \sum_{x' \in S_j} d(x, x')$$

- Furthest Link

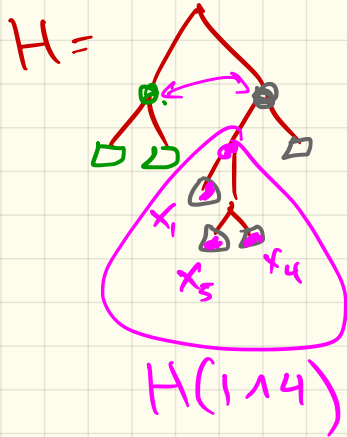
$$d(S_i, S_j) = \max_{x \in S_i, x' \in S_j} d(x, x')$$

- Define goodness of fit of model  $M$  cluster's  
 $g(M, S) > \text{threshold} = 10$   
 $d(S_i, S_j) = -g(M, S_i \cup S_j)$ 
 $M = f(S)$

Dasgupta 2016

Similarity

$$S_{ij} = S(x_i, x_j)$$



$$\text{Cost}(H) = \sum_{x_i, x_j \in X} S_{ij} \left[ \# \text{leaves } H(i, j) \right]$$

$H(i, j)$  = smallest subtree containing  $x_i, x_j$

$$\# H(i, j) = 3 = |\{x_1, x_4, x_5\}|$$

goal: Find  $H$  minimize  $\text{Cost}(H)$

# Efficiency

loop  
 $O(n)$

1. Find closest pair

$O(n^2)$   
 $(s_i, s_j)$

$O(n^3)$  time

check  
 $O(n^2)$  pairs  
 $i, j$

2. Merge  $(s_i, s_j) \rightarrow s_m$

$d(s_i, s_j)$  often update  $O(1)$

$d(s_i, s_j), d(s_i, s_e), d(s_i, s_e)$  time, is  $O(n)$  time

Merge  $s_i, s_j \rightarrow s_m$

$d(s_m, s_e) = \min(d(s_e, s_i), d(s_e, s_j))$

still  
SLOW

often w/ priority queue  $\rightarrow O(n^2 \log n)$



Why  $k$ -d trees not work  
in high dim?

battle: boxes vs. balls.

$$\text{Vol}(\text{Ball}(d)) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} \approx \frac{\pi^{d/2}}{(d/2)!}$$

$d \rightarrow \infty \rightarrow 0$  (less than 1)  
 $d > 5$

$$\text{Vol}(\text{Box}(d)) = 2^d$$

$d \rightarrow \infty \rightarrow \text{vol exponentially to } \infty$

