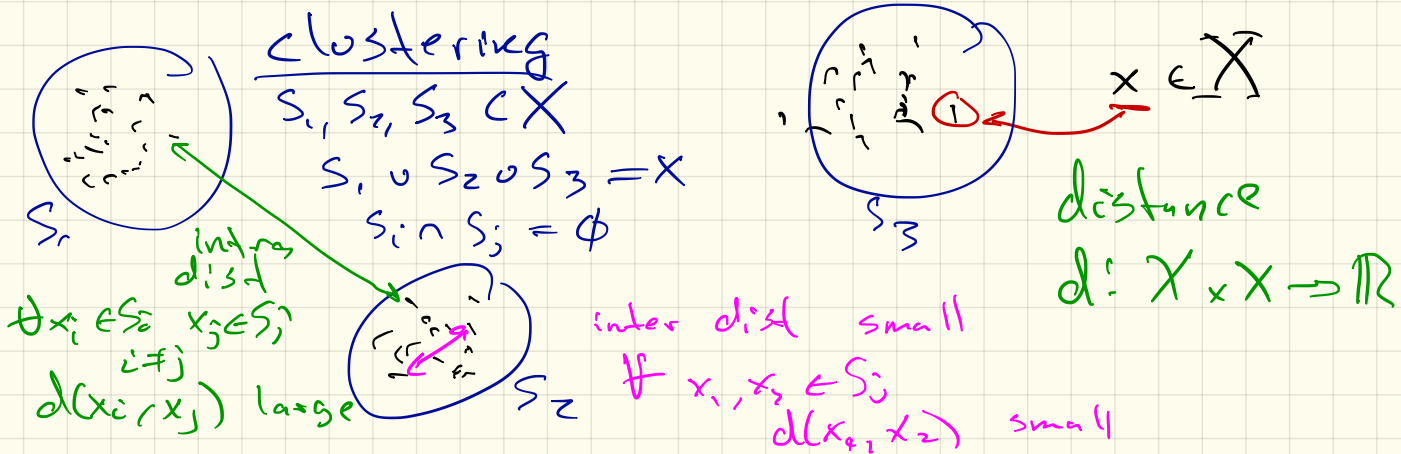# CLUSTERING

## Hierarchical Agglomerative Clustering

1000s of types of clustering!

- HAC 
- Assignment-based clustering
- Spectral

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

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

## clustering

$S_1, S_2, S_3 \subset X$

$S_1 \cup S_2 \cup S_3 = X$

$S_i \cap S_j = \phi$

$S_1$

intra dist

$\forall x_i \in S_i, \ x_j \in S_j$
$i \neq j$
$d(x_i, x_j)$ large

$S_2$

inter dist small
$\forall x_1, x_2 \in S_i$
$d(x_1, x_2)$ small

$S_3$

$x \in X$

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

common data

# Hier Agg Clust

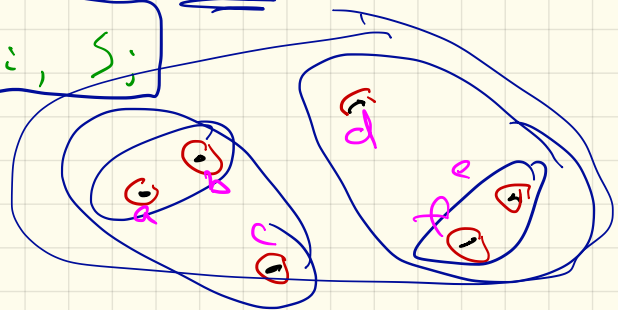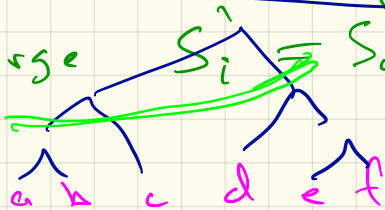Input $X$, dist $\boxed{d.}$ $\quad d: X \ast X \rightarrow \mathbb{R}$

Algo: If two points (or clusters) are close enough, put them in the same cluster. Repeat.

0. Each $x_i \in X \rightarrow$ put in separate cluster $S_i$

1. While (two clusters are close enough)

   1a. Find closest pair $S_i, S_j$

   1b. Merge $S_i = S_i \cup S_j$

# Distance between pair clusters $S_1, S_2$

- find "center" $c_1$ of $S_1$, $c_2$ of $S_2$

$$D(S_1, S_2) = d(c_1, c_2)$$

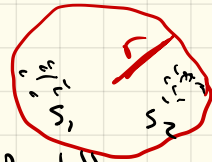  - $c_i = \text{average}(S_i)$     — $c_i = \text{random}(S_i)$
  - $c_i = \arg\min_{x \in S_i} \|x\|$
  - $c_i = \text{median}(S_i) = \arg\min_{c \in S_i} \sum_{x \in S_i} d(x, c)$

- $D(S_1, S_2) = \text{radius}\left(\text{min enclosing ball}(S_1 \cup S_2)\right)$

- $D(S_1, S_2) = \underbrace{\overbrace{\{\min\}}^{\text{average}}}_{x_1 \in S_1, \ x_2 \in S_2} d(x_1, x_2)$   max

  "single link"



- Build generative models
  compare likelihoods $L(S_1), L(S_2), L(S_1 \cup S_2)$



two moons

# Which Variant?

- gives
  proper clostering; "smallest error"

$$\hookrightarrow Dist \quad D \quad linked$$
$$to \quad error \quad eval.$$

- Computational Complexity.
  1. $O(n^2)$ pairwise dist.
  2. $O(n)$ nodes in hierarchy.
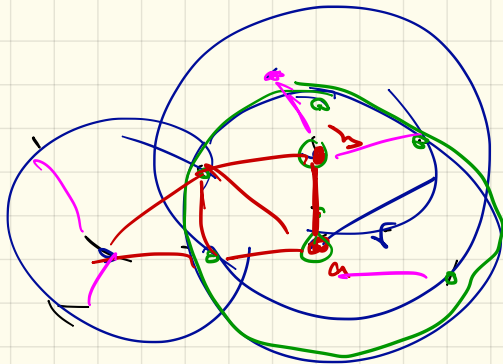
$O(n^3) \longrightarrow O(n^2 \log n)$ (w/ PQ)
   2. $O(n \log n)$ or $O(nk)$

# DB Scan   $O(n \log n)$?

$X, d,$   param: $r$   $\tau$
radius   threshold density

each $B_\sigma(x)$

$|B_r(x) \cap X| = \delta$

if $\geq \tau$

$x \to$ "core point"

link core points
$d(x, x') < r$
↳ cluster.

# K-Center Cluster       $O(nk)$
## Gonzalez Algo.

In $X, d, k$    any $d$   metric

Out: $c_1, c_2, \dots c_k$ ← centers of clusters.

$c_1 \leftarrow$ arbitrarily $(X)$

for $j = 2$ to $k$

$$c_j = \underset{x \in X}{\arg\max} \left( \underset{c \in \{c_1, \dots c_{j-1}\}}{\min} \| x - c \| \right)$$