Last time: K-means clustering

If $X$ is a $N \times P$ mtx of numeric covariates and I define my dissim mtx $D$ so that

$$D_{ii'} = \| x_i - x_{i'} \|^2$$

Combinatorical Clustering Problem:

Choose $G$s $(G_1, ..., G_K)$ so that

$$W = W(G_1, ..., G_K) = \frac{1}{2} \sum_{k=1}^{K} \sum_{i,i' \in G_k} D_{ii'}$$

is minimized, i.e.

$$\hat{G}_1, \hat{G}_2, ..., \hat{G}_K = \underset{G_1, ..., G_K}{\arg\min} \; W(G_1, ..., G_K)$$

Can show: Equivalent to

\# obs in $K^{th}$ cluster

$$W = \sum_{k=1}^{K} N_k \sum_{i \in G_k} \| x_i - \bar{x}_k \|^2$$

$\bar{x}_k \in \mathbb{R}^P$

$$W = \frac{1}{2} \sum_{k=1}^{K} \sum_{i,i' \in G_k} D_{ii'} = \frac{1}{2} \sum_{k=1}^{K} \sum_{i,i' \in G_k} \| x_i - x_{i'} \|^2$$

$$W = \frac{1}{2}\sum_{k=1}^{K}\sum_{i,i'\in G_k} D_{ii'} = \frac{1}{2}\sum_{k=1}^{K}\sum_{i,i'\in G_k} \|X_i - X_{i'}\|^2$$

$$\|a\|^2 = a^T a$$

$$\|X_i - X_{i'}\|^2 = \|X_i - \overline{X}_k + \overline{X}_k - X_{i'}\|^2$$

$$= \left((X_i - \overline{X}_k) + (\overline{X}_k - X_{i'})\right)^T(\cdots)$$

$$(X_i - \overline{X}_k)^T(X_i - \overline{X}_k) \underbrace{}_{\|X_i - \overline{X}_k\|^2} + (\overline{X}_k - X_{i'})^T(\overline{X}_k - X_{i'})$$
$$\underbrace{}_{\|X_{i'} - \overline{X}_k\|^2}$$
$$+ 2(X_i - \overline{X}_k)^T(\overline{X}_k - X_{i'})$$

$$= \frac{1}{2}\sum_{k=1}^{K}\left[\sum_{i,i'\in G_k}\|X_i - \overline{X}_k\|^2 + \|X_{i'} - \overline{X}_k\|^2 + 2\underbrace{\sum_{i,i'\in G_k}(X_i - \overline{X})^T(\overline{X} - X_{i'})}_{0}\right]$$

$$\sum(X_i - \overline{X}) = 0$$

$$= \frac{1}{2}\sum_{k=1}^{K}\sum_{i\in G_k}\sum_{i'\in G_k}\left(\|X_i - \overline{X}_k\|^2 + \|X_{i'} - \overline{X}_k\|^2\right)$$

$$= \frac{1}{2}\sum_{k=1}^{K}\left[\left(\sum_{i}\sum_{i'}^{N_k}\|X_i - \overline{X}_k\|^2\right) + \left(\sum_{i'}\sum_{i}^{N_k}\|X_{i'} - \overline{X}_k\|^2\right)\right]$$

$$= \frac{1}{2}\sum_{k=1}^{K}\left[\sum_{i}N_k\|X_i - \overline{X}_k\|^2 + \sum_{i}\|X_i - \overline{X}_k\|^2\right]$$

$$\boxed{= \sum_{k=1}^{K} N_k \sum_{i\in G_k}\|X_i - \overline{X}_k\|^2} = TWCSS$$

# Proposed Lloyd's Algorithm:

$$t = 1, 2, 3, \ldots$$

Update ① $\mu_k^{(t)} = \dfrac{1}{N_k} \sum_{i \in G_k^{(t)}} X_i$

Assign ② $G_k^{(t)} = \left\{ X_i \,\middle|\, \|X_i - \mu_k^{(t)}\| \leq \|X_i - \mu_{k'}^{(t)}\| \quad \forall k' \right\}$

## Why does Lloyds work?

Problem: $\hat{G}_1, \ldots, \hat{G}_K = \underset{G_1, \ldots, G_K}{\arg\min} \sum_k N_k \sum_{i \in G_k} \|X_i - \overline{X}_k\|^2$

Generalize Problem

$\underset{G_1, \ldots, G_K, \underbrace{m_1, \ldots, m_k}_{\text{Centers}}}{\min} \quad \boxed{\underbrace{\sum_k N_k \sum_{i \in G_k} \|X_i - m_k\|^2}_{W}}$

Fact: $\overline{X} = \underset{m}{\arg\min} \|X_i - m\|^2$
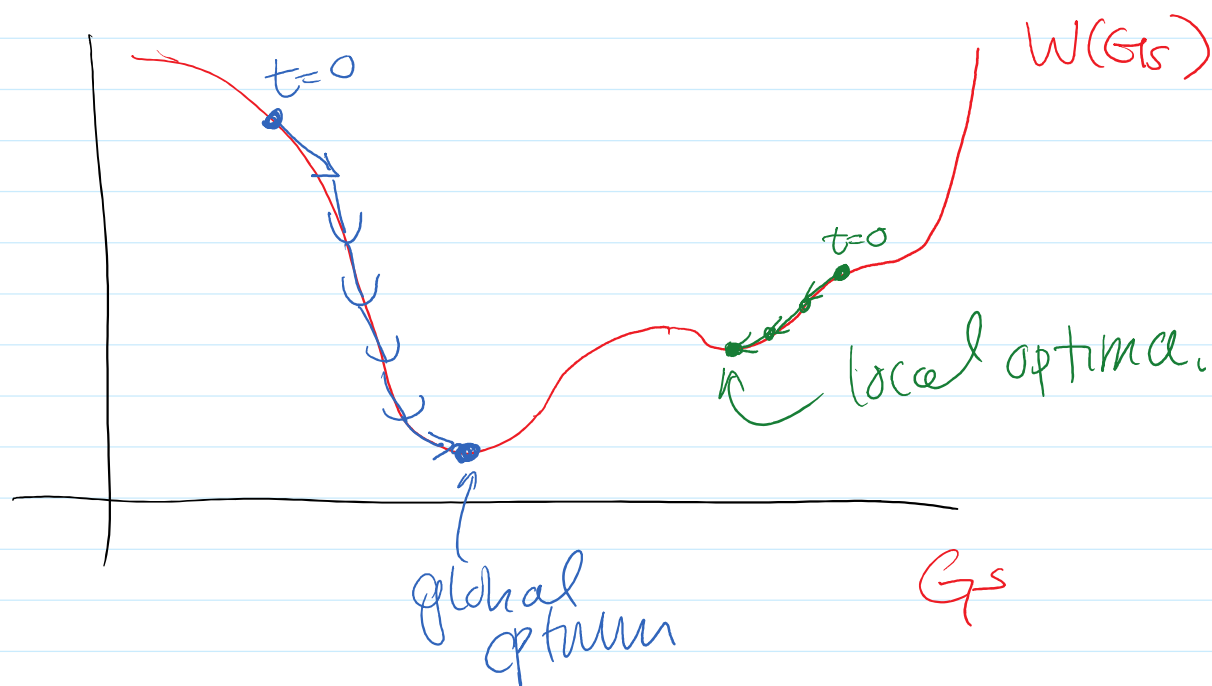
Step ① Given $G$s, set $m_k = \overline{X}_k$

this reduces W

Step ② Given Ms, choose Gs as pts closest
to Ms, this reduces W

So at each step this makes W smaller.

HOWEVER
We don't nec. get to globally optimum soln.



Soln: Try several random intializations
take soln w/ lowest W.

## What about non-numeric data?

Just have $D$?     $\curvearrowleft$ non-euclidean dist/dissim.

## K-Mediods

<u>Soln:</u> Replace step ① (update means)
   by setting $m_k$ = pt in cluster closest to everything else in cluster

Step ①* Find obs in cluster closest to others

$\color{red}\longrightarrow$ $i_k^* = \underset{i \in G_k}{\text{argmin}} \sum_{i' \in G_k} D_{ii'}$

$\color{red}\curvearrowleft$ explicit optim problem

② <u>Assignment step.</u>

Object $i$ goes in Group $G_k$ if

$$D_{ii_k^*} \le D_{ii_{k'}^*} \quad \forall k'$$

Nice fact: No $X$, just $D$.

Nice fact: No X, just D.

Bad fact: more comp. intensive

---

## How do I choose K?

W | as $K \uparrow \Rightarrow W \downarrow$ (monotonically)

like KNN can't choose K in this way.
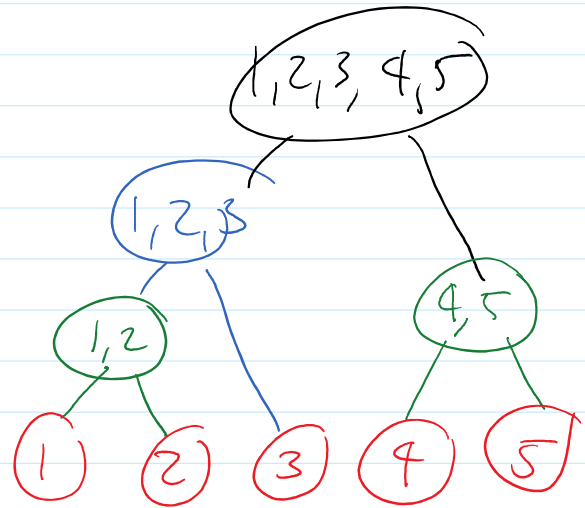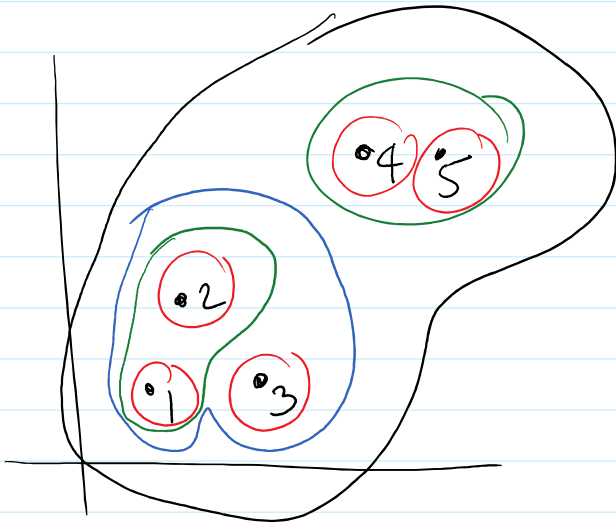
K

One way!

W

Kink, Knee

$K^*$      K

---

## Hierarchical Clustering

Build a collection (hierarchy) of nested clusters.

① Agglomerative Clustering

(i) start w/ clusters as individual pts

(ii) merge clusters that a "similar" or "close"

(iii) recursively repeat until everytls is in one cluster.



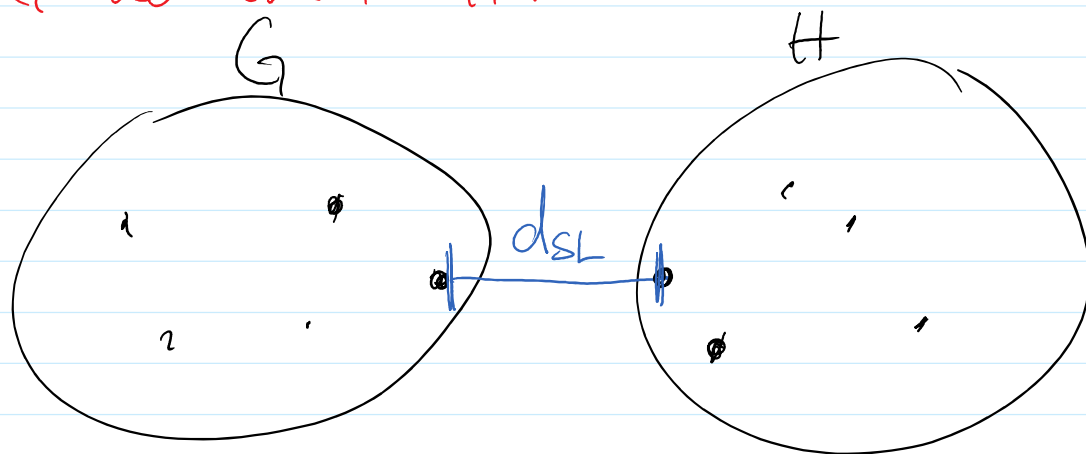(2) <u>Divisive Clustering</u> (opposite)

① start w/ 1 large cluster

② recursively break into smaller clusters

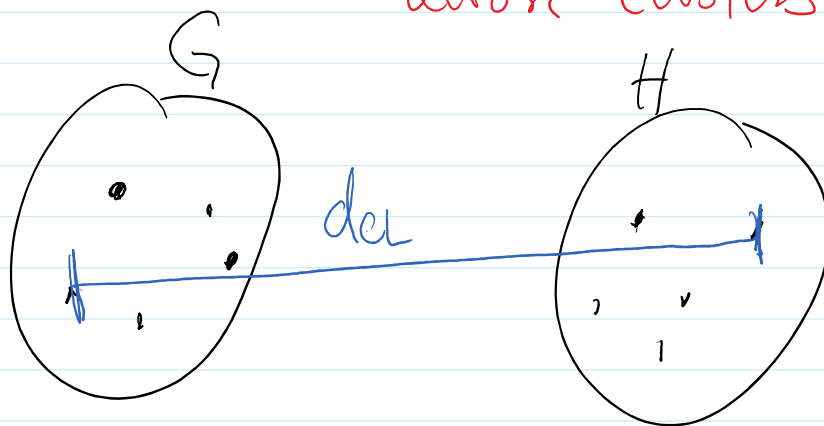To do <u>agglomerative clustering</u> we need a measure of "closeness" between clusters.
"Similarity"

① <u>Single Linkage</u>: dist. btwn clusters G and H is the min dissim btwn any 2 pts:

are in $G$ and are in $H$.



$$d_{SL}(G, H) = \min_{\substack{i \in G \\ i' \in H}} D_{ii'}$$
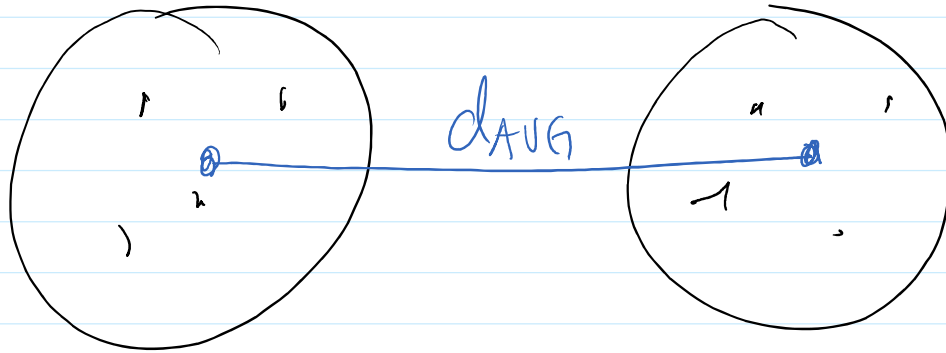
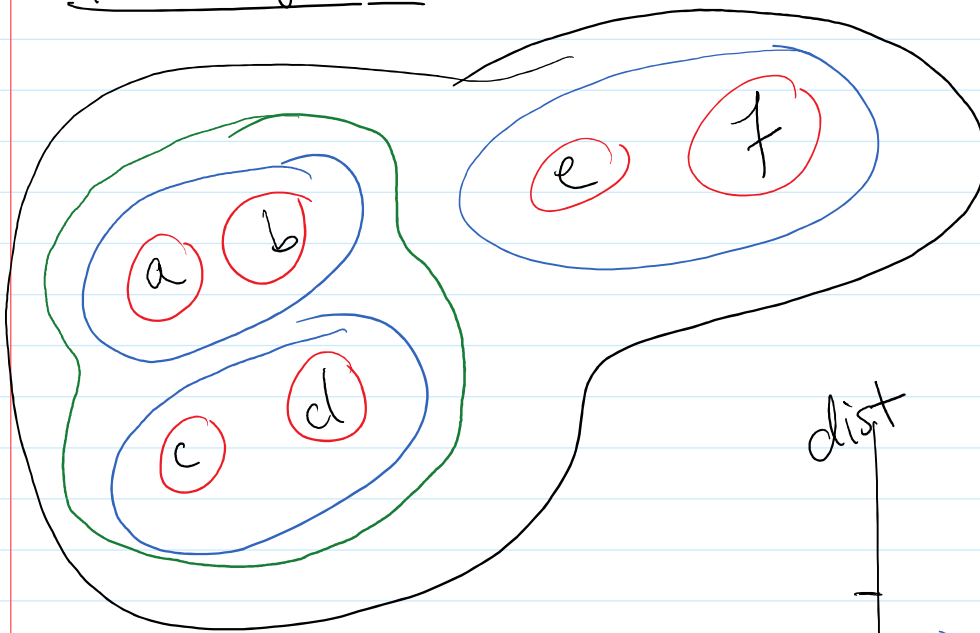② <u>Complete Linkage</u>:  max dissim btwn 2 pt
across clusters



$$d_{CL}(G, H) = \max_{\substack{i \in G \\ i' \in H}} D_{ii'}$$

③ <u>Average Linkage</u>:   avg dissim btwn clust.

$$d_{AVG}(G, H) = \frac{1}{N_G N_H} \sum_{i \in G} \sum_{i' \in H} D_{ii'}$$



$d_{AVG}$

## Dendogram



dist      dendogram

height =
dist btwn
clusters