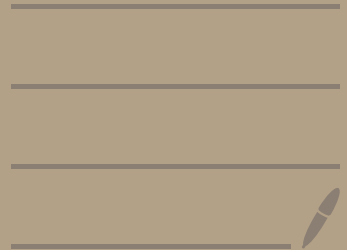


L22: K-Means Clustering w/ Lloyd's Algorithm



Assignment-based Clustering

Input $X \subset \mathbb{R}^d$

value $k > 1 = \# \text{ clusters}$

distance $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$

$$d(x_1, x_2) = \|x_1 - x_2\|$$

Goal: Set of k sites $S = \{s_1, s_2, \dots, s_k\} \subset \mathbb{R}^d$

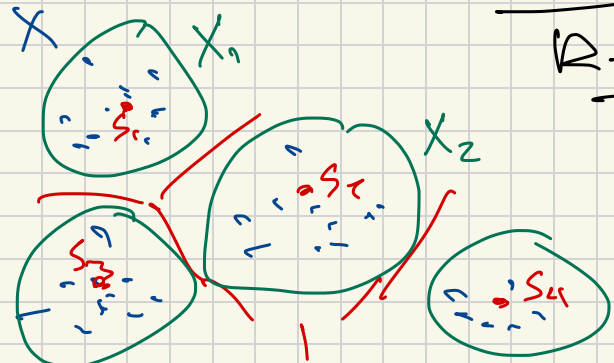
$$X = \bigcup_{j=1}^k X_j$$

$$X_j \subset X$$

$$X_j = \{x \in X \mid \phi_S(x) = s_j\}$$

hard clustering

$$X_j \cap X_{j'} = \emptyset \quad j \neq j'$$



k-means

$$S^X = \operatorname{argmin}_{|S|=k}$$

$$\sum_{i=1}^n \|x_i - \phi_S(x_i)\|^2 = \operatorname{cost}(X, S)$$

$$\phi_S(x) = \operatorname{argmin}_{s_j \in S} \|x - s_j\|$$

Lloyd's Algorithm

0: Initialize Choosing k pts $S \subset X$
 $S = \{s_1, s_2, \dots, s_k\}$

1: repeat

(a) $\forall x_i \in X$: assign x_i to x_j so $\phi_S(x_i) = s_j$

j th subset of X
"cluster"
associated w/ s_j

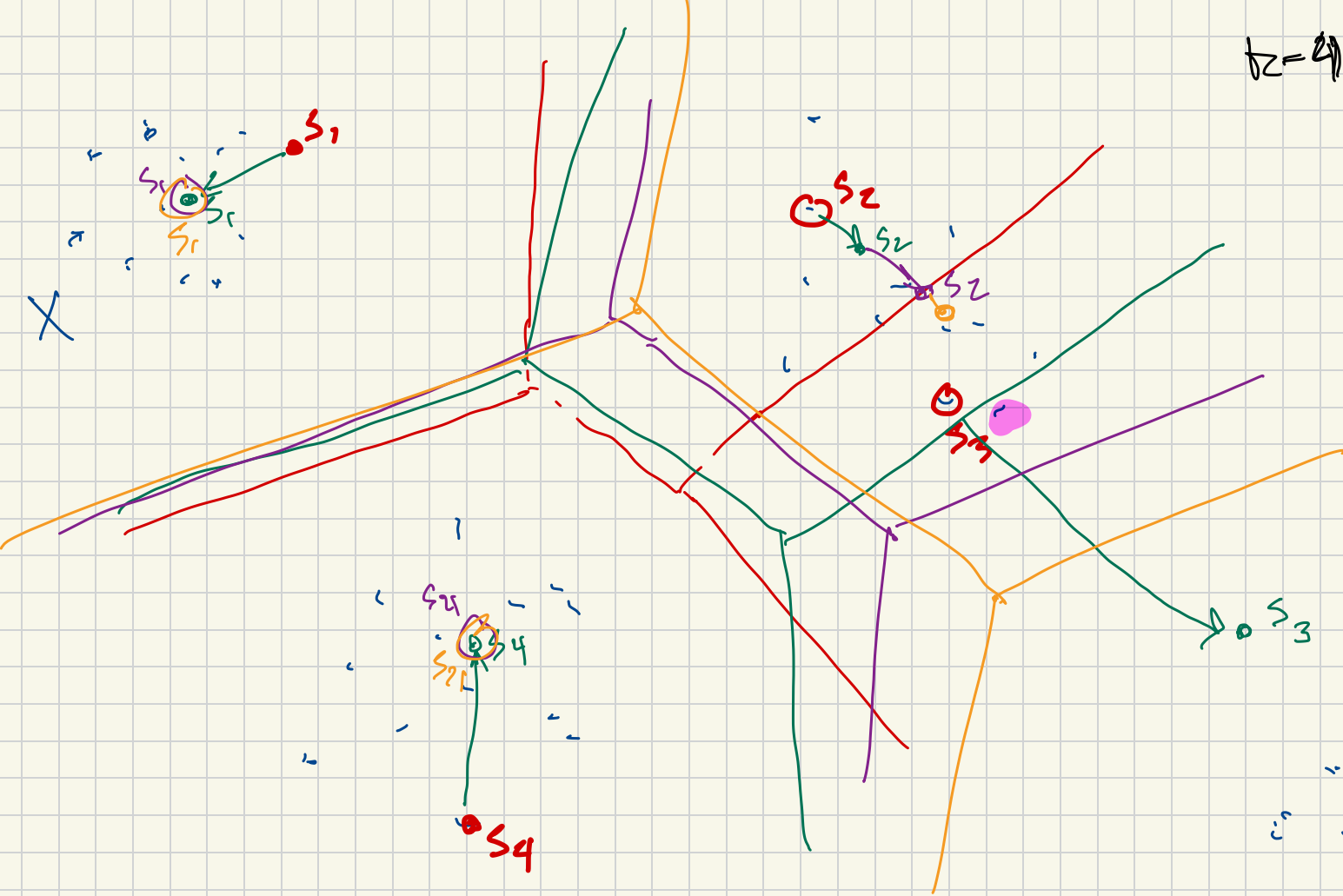
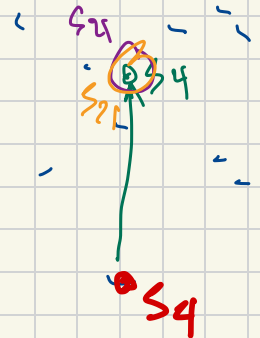
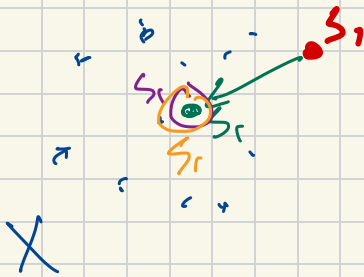
(b) $\forall s_j \in S$: update $s_j = \frac{1}{|X_j|} \sum_{x \in X_j} x = \text{average}(X_j)$

until ("converged" or \bar{K} steps)

\uparrow 10 or 20

like "Expectation - Maximization" (EM)

$t_2 = 21$



Does Lloyd's Also Always Converge? Yes

$$\text{Cost}(X, S) = \sum_{x_i \in X} \|x_i - \phi_S(x_i)\|^2 \quad \text{step (a)}$$

$$= \sum_{s_j \in S} \left(\sum_{x \in X_j} \|s_j - x\|^2 \right) \quad \text{step (b)}$$

not increase

Claim each step (a) or (b) decreases $\text{Cost}(X, S)$

step (a): $\forall x \in X$ | assign to X_j so $\phi_S(x) = s_j$
closest site.

step (b): $\forall s_j \in S$ | $s_j = \frac{1}{|X_j|} \sum_{x \in X_j} x$ ← set to average

$$s_j^x = \underset{\text{set } S}{\text{argmin}} \sum_{x \in X_j} \|x - s\|^2$$

Must be $S_t \rightarrow S_{t+1} \rightarrow S_{t+2}$

claims $S_t \neq S_{t+1}$ $\text{cost}(X, S_t) > \text{cost}(X, S_{t+1})$

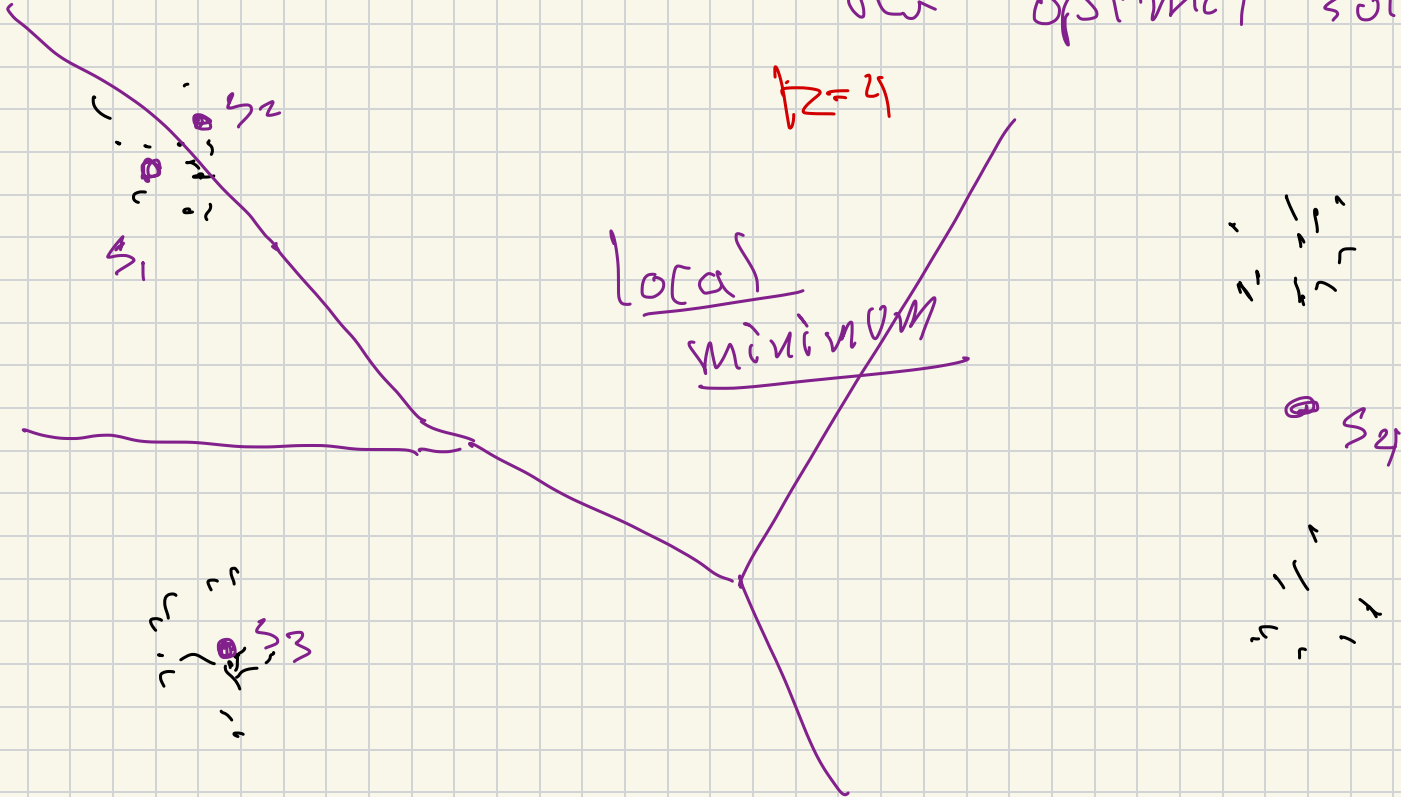
$\text{cost}(X, S_{t+1}) > \text{cost}(X, S_{t+2})$

$\text{cost}(X, S_t) > \text{cost}(X, S_{t+2})$

$\Rightarrow S_t \neq S_{t+2}$


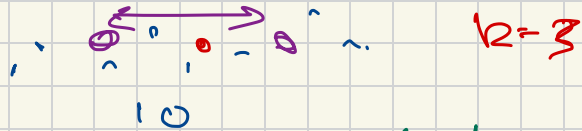
Lloyd's Also

May not converge to
the optimal solution!

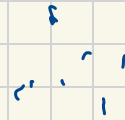
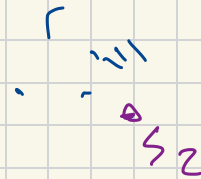
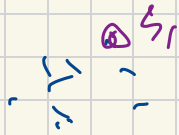


How to Address non-optimal convergence?

① Try more initializations of random
step 0: select $S \subset X$ of random.
→ Initialize each s_j is course. $X_j \subset X$ of random.

② Force sites not too close

if stuck and points close → randomly reassign
to far point.
sometimes $X_j = \emptyset$ → 

③ Better way to initialize
→ incrementally increase k .
always start from existing sites



→ Gonzalez Also greedy

→ k -means++ (randomized)

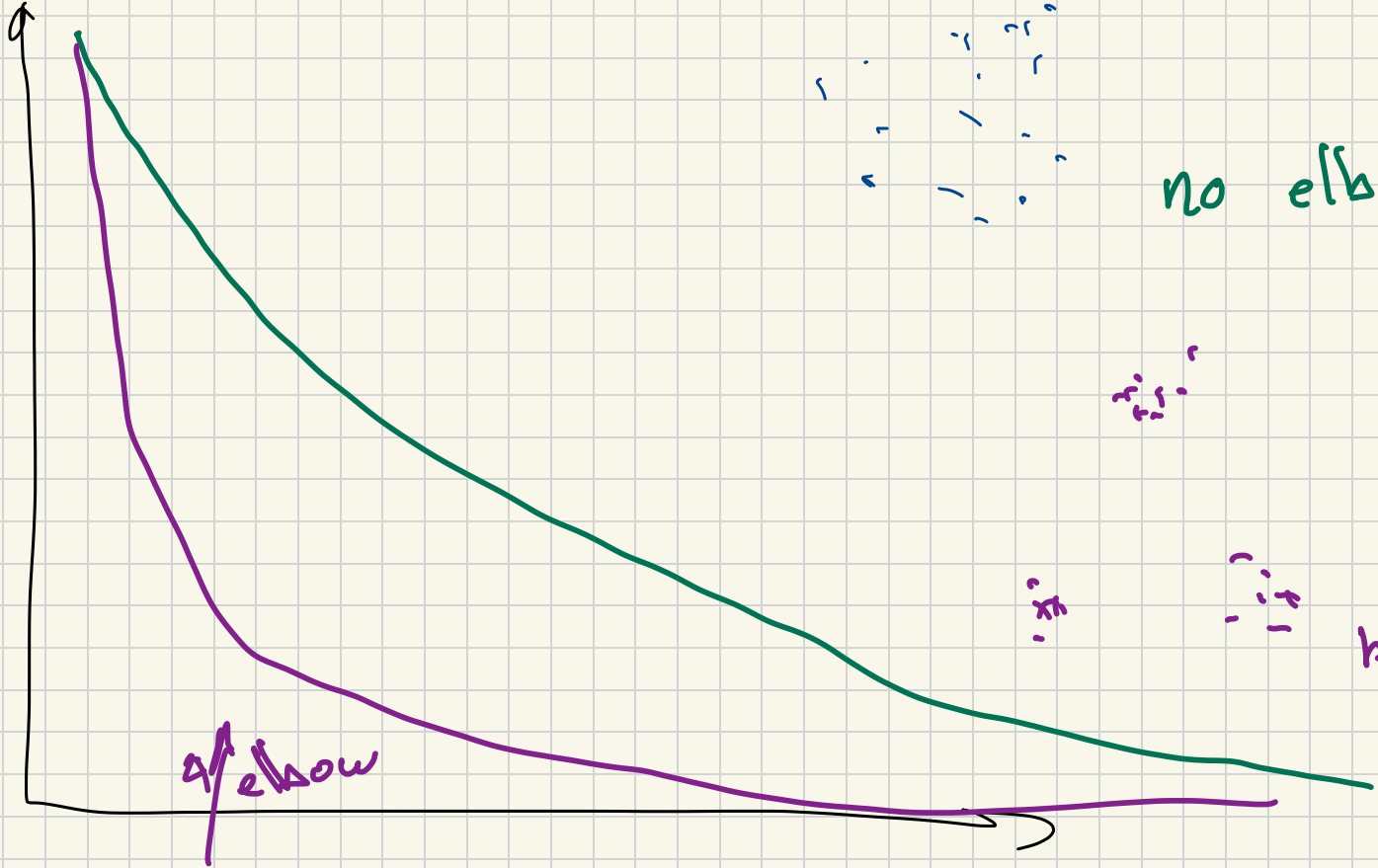
④ Run w/ much larger $k' \gg k$
↳ then combine

How to decide k ?

Goal $\text{Cost}(X, S) = \sum_{x_i \in X} \|x_i - \phi_S(x_i)\|^2$

$\text{Cost}(X, S, k)$ always decreases as k increases





blob

no elbow

elbow

no elbow

How many clusters?

