


FoDA L23

K-means clustering

↳ Lloyd's Algorithm

Nov 15, 2022



Assignment-based Clustering

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

value k
clusters

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

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

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

k-means

$$S^* = \underset{|S|=k}{\text{arg min}}$$

$$\sum_{i=1}^n \|x_i - \phi_S(x_i)\|^2$$

$$= \text{cost}(X, S)$$

$$\begin{aligned} \phi_S(x) &= \text{closest site } s_j \text{ to } x \\ &= \underset{s_j \in S}{\text{arg min}} \|s_j - x\| \end{aligned}$$



Lloyd's Algorithm

0. Initialize: Choose k points $S \subset X$

1. repeat

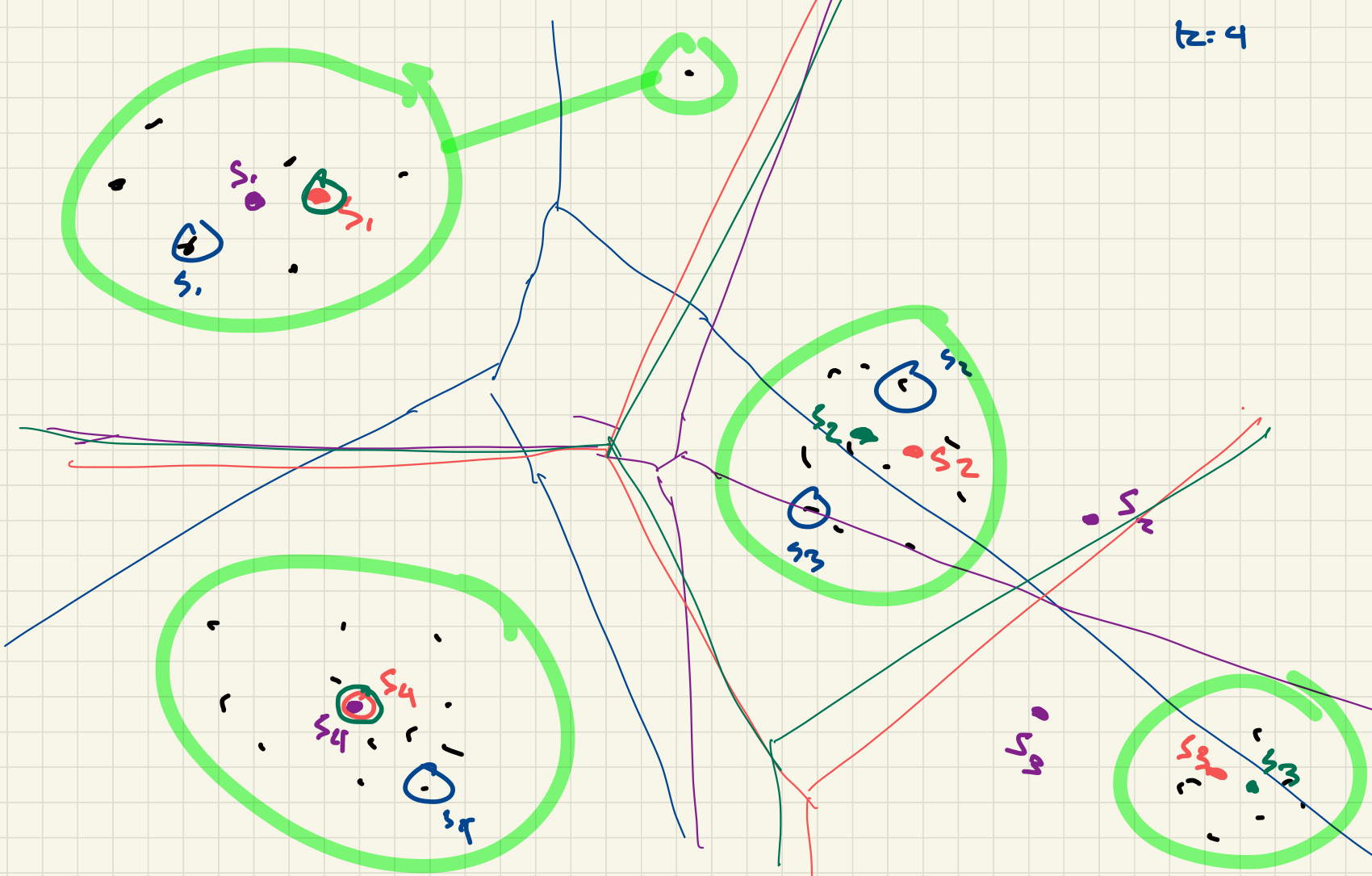
(a) $\forall x_i \in X$: assign x_i to x_j so $\phi_S(x_i) = s_j$

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

until ('converged' or K steps)
so $R=10$ or ∞

with subset clusters $\subset X$

t: 4



Does Lloyd's Algo Converge?

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

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

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

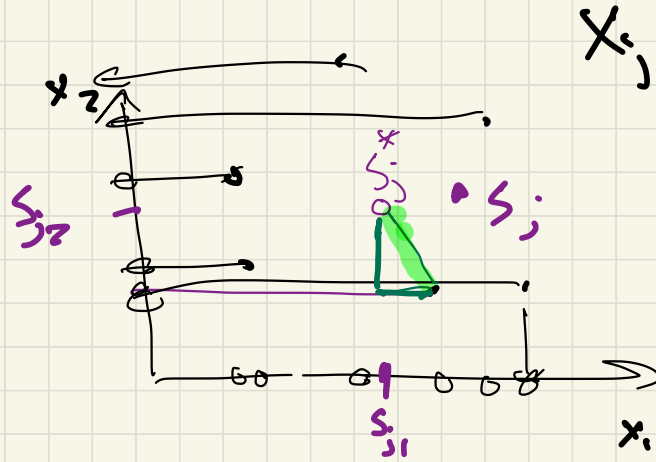
every step (a) or (b) decreases $\text{Cost}(X, S)$

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

closest site

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

$s_j^* = \arg \min_s \sum_{x \in X_j} \|x - s\|^2$

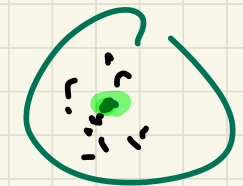
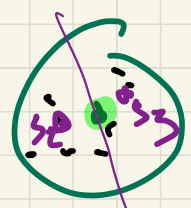


$$\|x_i - s_j\|^2 = (x_{i1} - s_{j1})^2 + (x_{i2} - s_{j2})^2$$

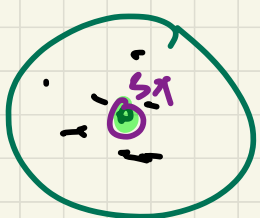
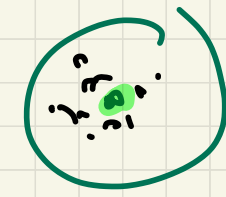
Lloyd's Algo

$k=4$

not guaranteed
to find optimal
solution.



s_4



optimal

Address non-optimal convergence.

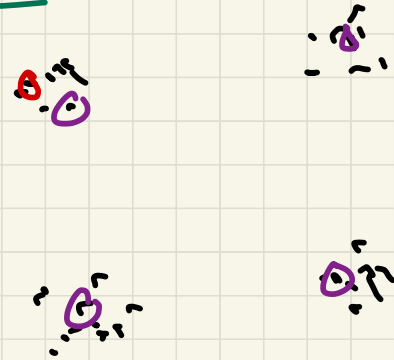
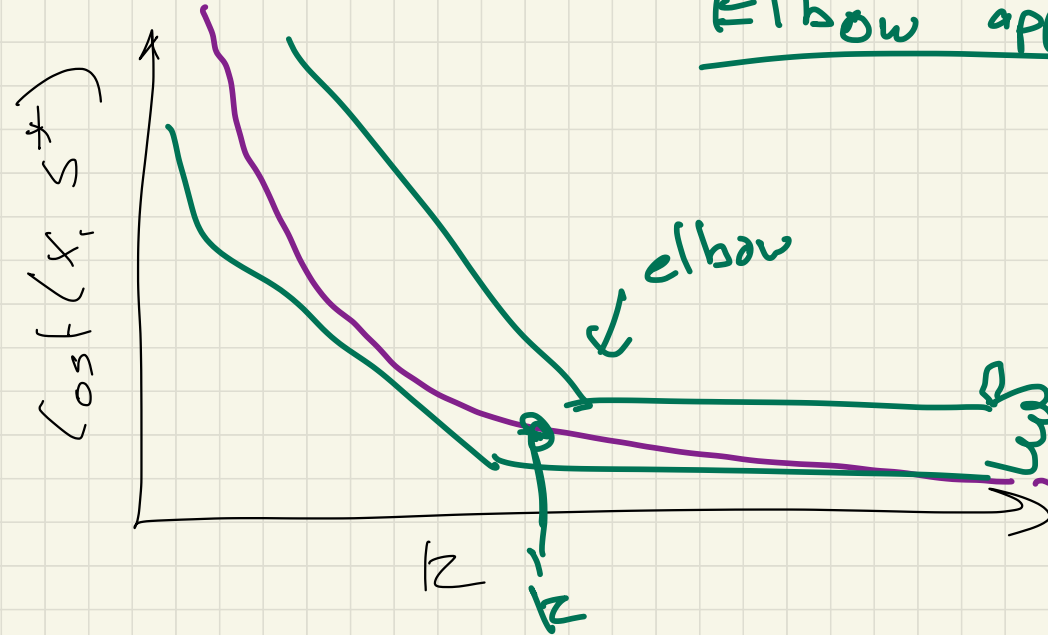
- ① Try more \downarrow random initializations "random restart"
- ② Better initialization
 - (a) Gonzalez k -centers
 - (b) k -means++ (randomized)
- ③ Randomize averaging step
 \rightarrow simulated annealing
- ④ Use more clusters $k' \gg k \rightarrow$ then merge.
 $k' = k \log n$

How to choose k ?



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

Elbow approach



$k=2 \Rightarrow Cost(x, S^*) = 0$

