

FODA L23

K-Means

Clustering

Assignment-Based Clustering

Input $X \subset \mathbb{R}^d$ $X = \{x_1, x_2, \dots, x_n\}$

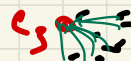
$k=3$

distances $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$



value k

$$d(x, p) = \|x - p\|$$



Goal Set of sites $S = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

minimize

$$f_{X \in X} (d(x, \phi_S(x)))$$

$$\sum_{x \in X} \|x - \phi_S(x)\|^2 \quad k\text{-means}$$

residual

maps x to closest $c_j \in S$.

$$\sum_{x \in X} \|x - \phi_S(x)\| \quad k\text{-median}$$

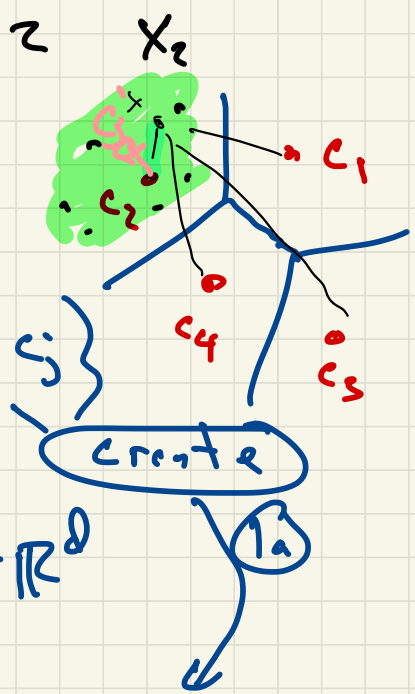
k -median
 k -medial SCX

$$\max_{x \in X} \|x - \phi_S(x)\| \quad k\text{-center}$$

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

Logd's Algorithm

$$X_j = \underbrace{\{x \in X \mid \phi_S(x) = c_j\}}_{\text{sets}}$$



①. Init: Choose k points $S \subset X$

$$S = \{c_1, \dots, c_k\} \subset \mathbb{R}^D$$

sites

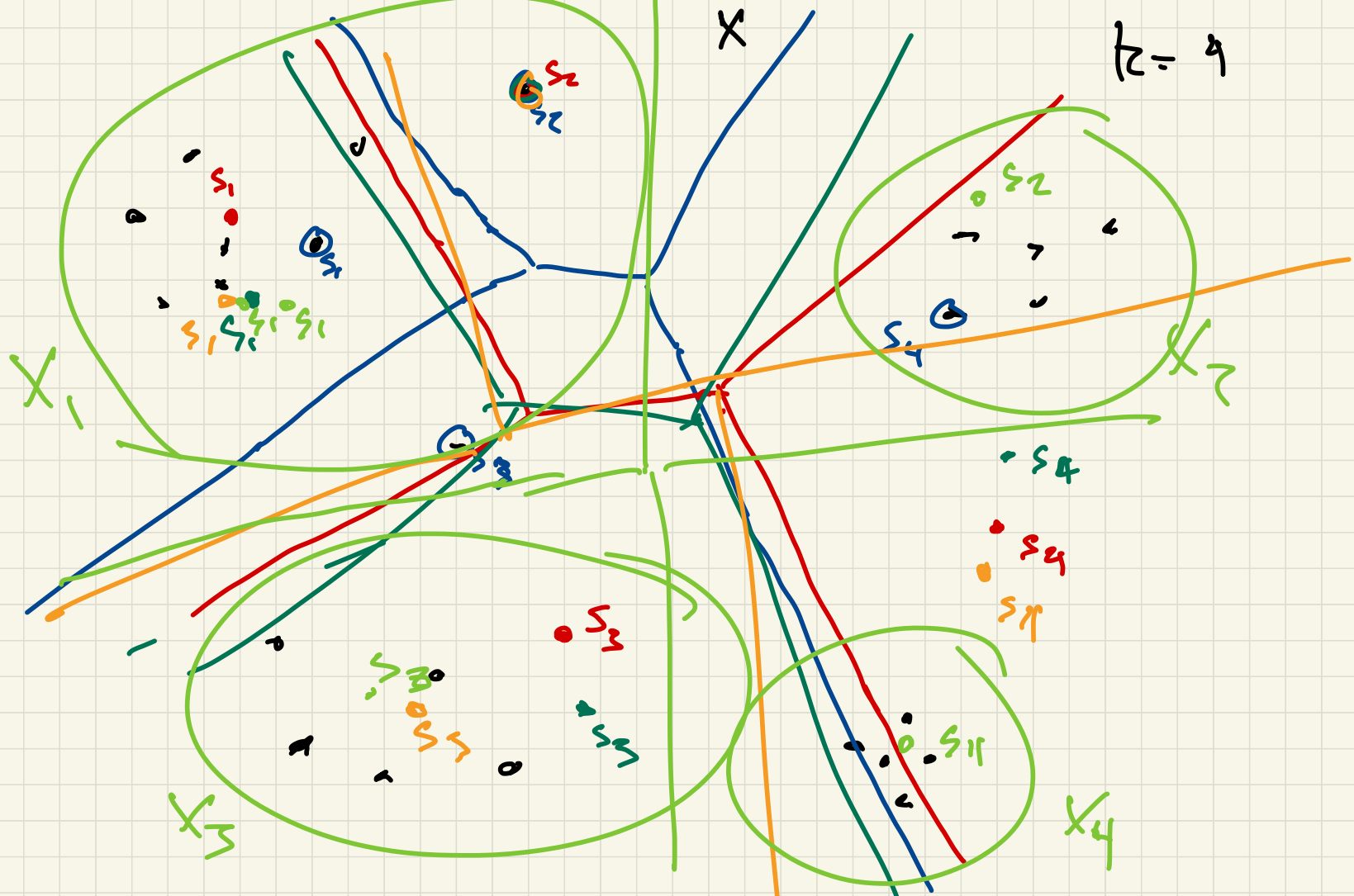
1. repeat

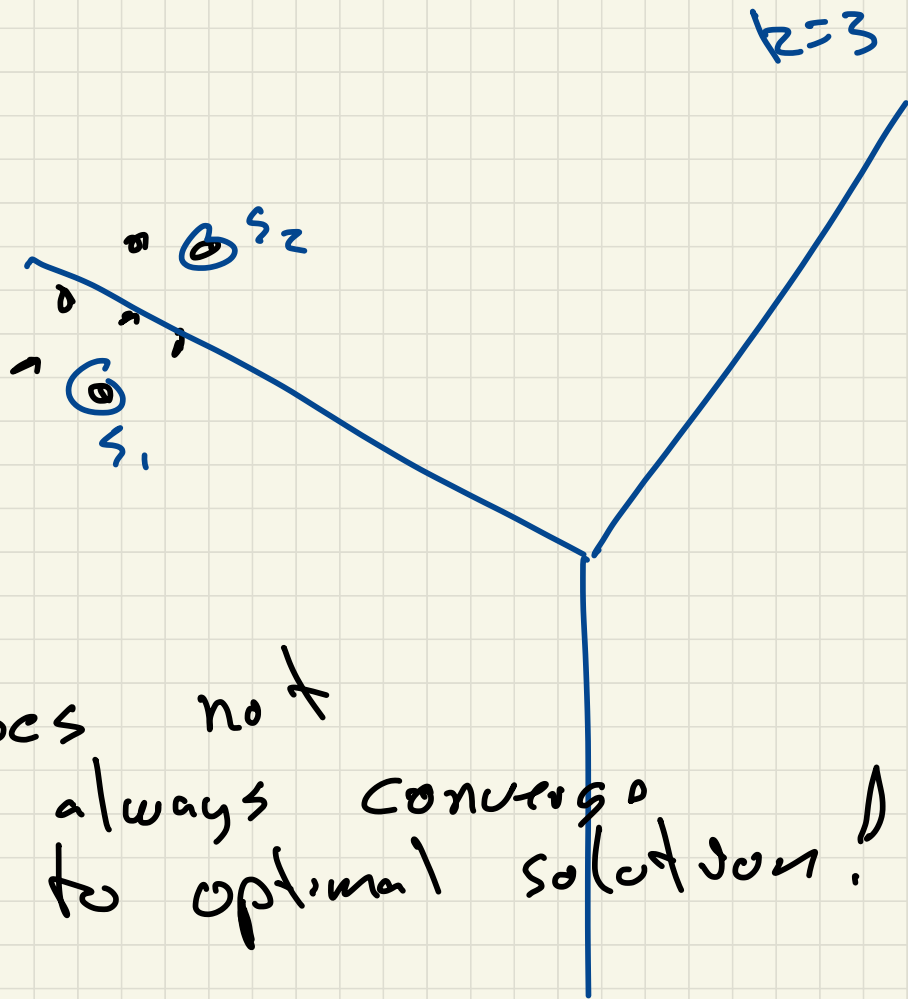
1a for all $x \in X$: assign x to X_j s.t. $\phi_S(x) = c_j$ update sets

1b for all $c_j \in S$: update $c_j = \frac{1}{|X_j|} \sum_{x \in X_j} x$ update sites

until (k rounds or X_j unchanged)

2. Return S

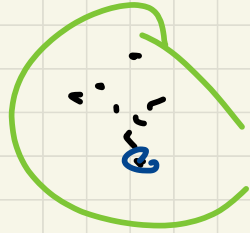




Localization is important

Goal: 1 site in each free cluster.

Option 1: plot, choose, optimize

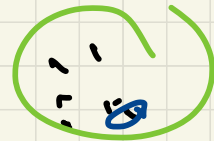


option 2: ensure all dist $\|c_j - c_{j+1}\|$ large

Gonzalez alg.

k-means++

option 3: Random: choose $S \subset X$
at random



option 4: Use large $k' = 10 \cdot k \rightarrow$ then merge.

Random Restarts

best-score = ∞

for $j=1 \dots T$ steps

0. Random init S_j

1. Run Lloyd's on (X, S_j)

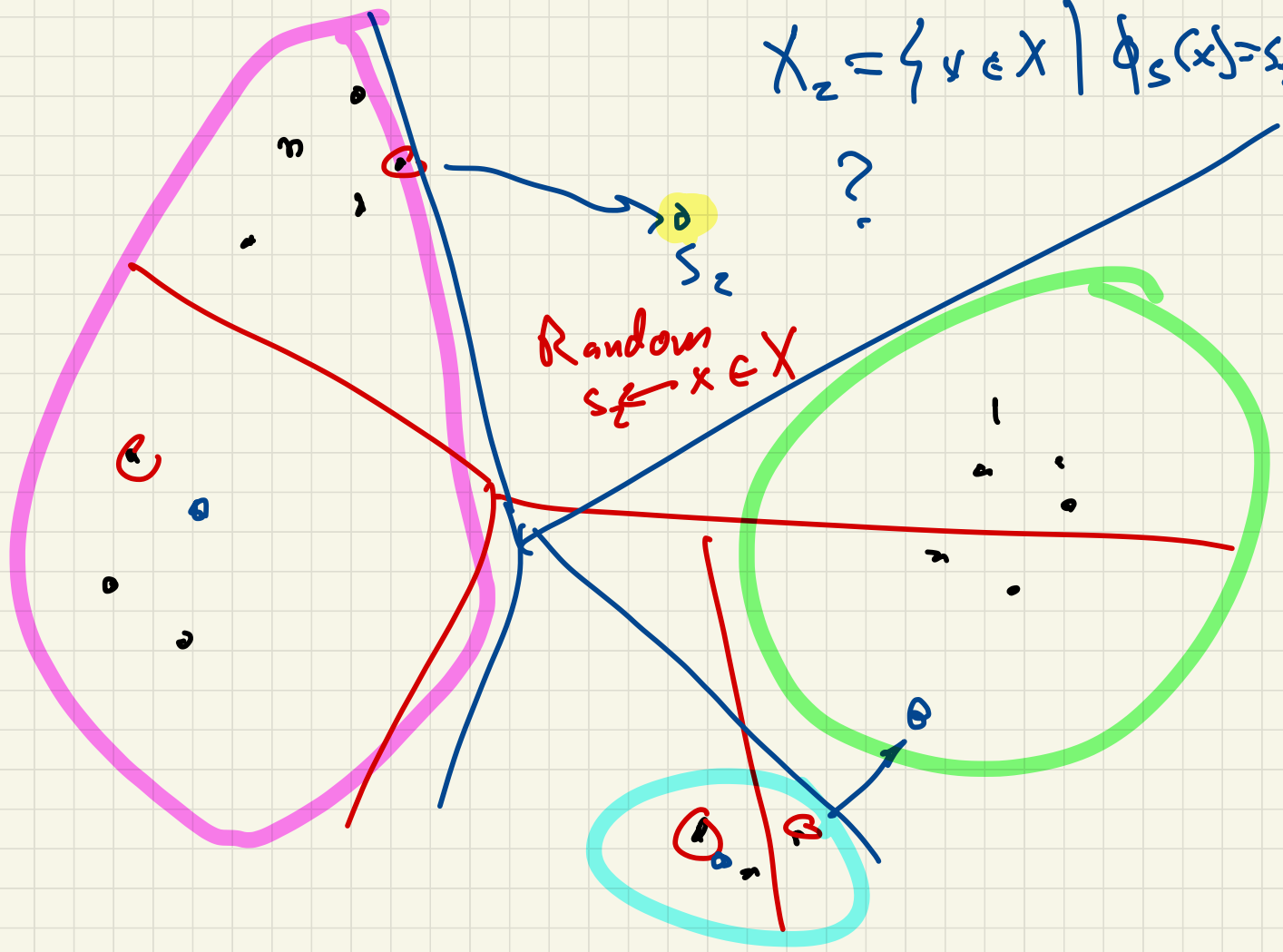
2. if $\text{cost}(X, S_j^*) < \text{best-score}$

$$S^* = S_j^*$$

$$\text{best-score} = \text{cost}(X, S_j^*)$$

Return S^*

$$X_2 = \{x \in X \mid \phi_S(x) = s_2\}$$



Random $s \leftarrow X$

?

c

a

b

3

2

1

2

3

b

1

2

Convergence Lloyd's Algo.

$$\text{Cost}(X, S) = \sum_{x \in X} \|x - \phi_S(x)\|^2 \quad (1)$$

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

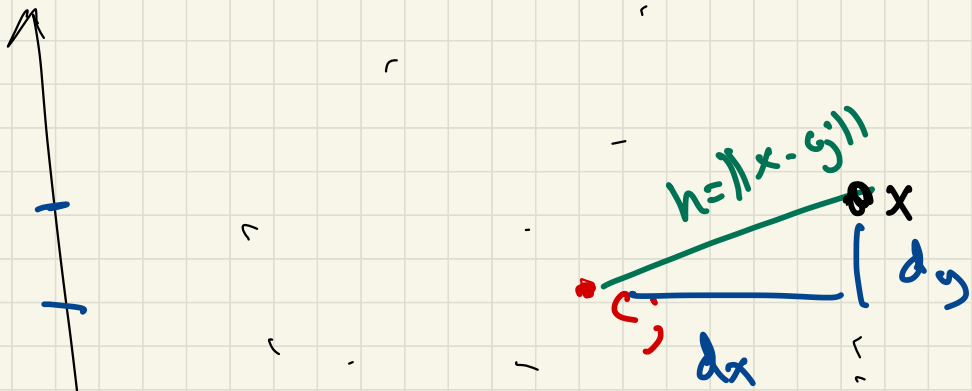
Lloyd's

1a optimize $X_j = \{x \in X \mid \phi_S(x) = s_j\}$

1b optimize $s_j = \frac{1}{|X_j|} \sum_{x \in X_j} x$

fix X_j (not implicit w/ ϕ_S)

proved in $d=1$
average minimizer $\bar{x}(-)$



x_j

$$h^2 = dx^2 + dy^2$$

$$\|x - c_j\|^2 = dx^2 + dy^2$$

• optimize
independently

(ii)

Choose k ?

