# FoDA
## L23

- Clustering

- Lloyd's Algorithm
  *for k-means*

# K-Means Clustering Formulation

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

$K \in \mathbb{Z}^+ \in [1, 2, \ldots]$

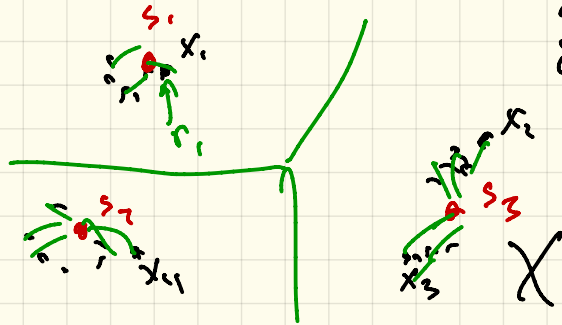$d: X \times X \to \mathbb{R}_*$   $d(x_1, x_2) = \|x_1 - x_2\|$

Goal   Find set $S = \{s_1, s_2, \ldots s_k\} \subset \mathbb{R}^d$

minimize $\sum_{i=1}^{n} d\left(x_i, \phi_S(x_i)\right)^2$

↳ maps $x_i$ to $\underset{s_j \in S}{\operatorname{argmin}} \|x_i - s_j\|$

$k = 3$

# Lloyd's Algorithm $\rightarrow$ $X = \bigcup_{j=1}^{k} X_j$

$X_j \subset X$

(arbitrarily?)

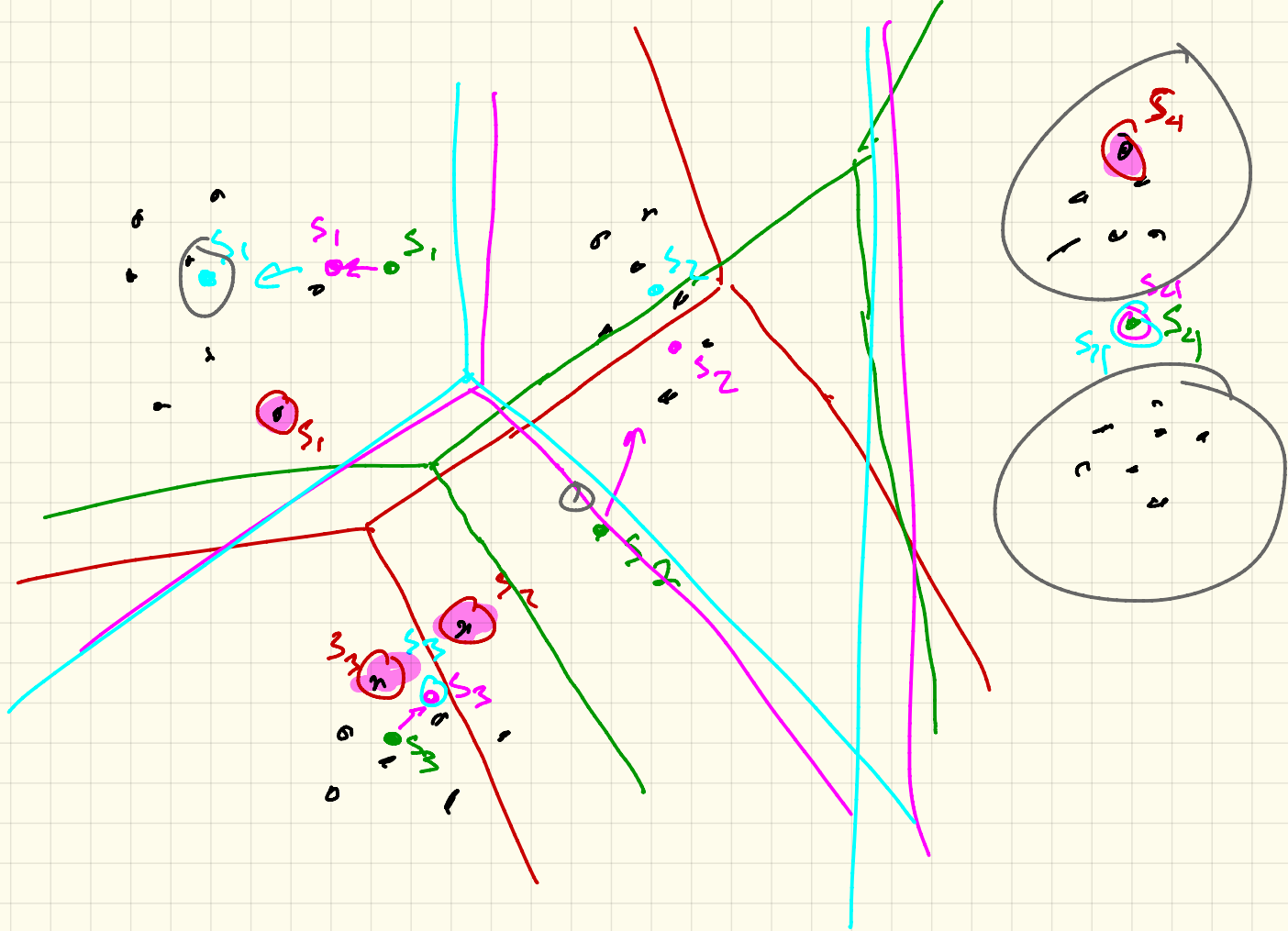0. Initialize $k$ points $S \subset X$

1. repeat
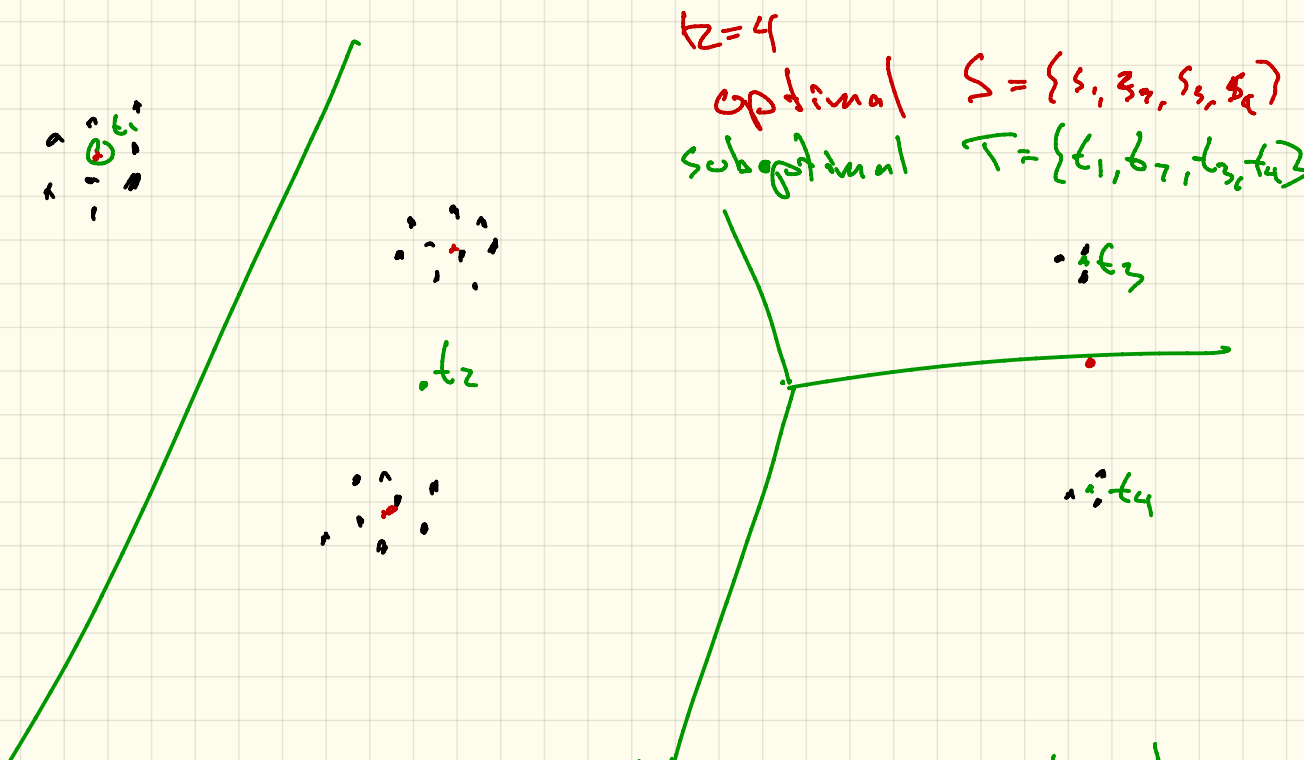
   **assignment / Voronoidy** (a.) for all $x_i \in X$ : assign $x_i$ to $X_j$ so $\phi_S(x_i) = s_j$

   **recenter / average** (b.) for all $s_j \in S$ : update $s_j = \frac{1}{|X_j|} \sum_{x \in X_j} x = \text{average}(X_j)$

   until ( the set $S$ unchanged
            or change is small )

$t_2 = 4$

optimal $S = \{s_1, s_2, s_3, s_4\}$

suboptimal $T = \{t_1, t_2, t_3, t_4\}$

$c_i$

$\cdot t_3$

$\cdot t_2$

$\cdot t_4$

Lloyds Algo is stuck

Most of the time Lloyd's works well
(with a little help)

# Tricks to help Llogd's
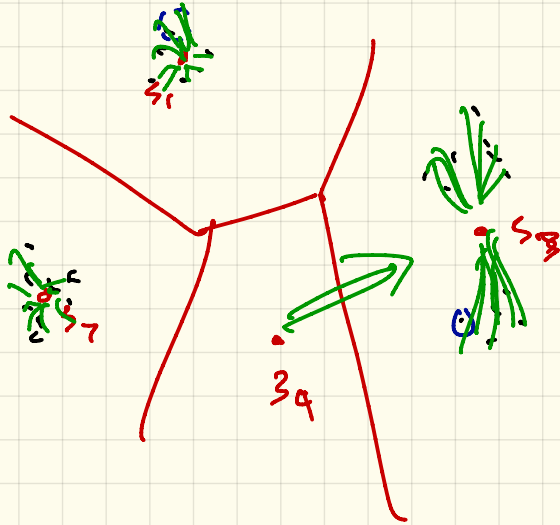
- Random Restarts
  - a. Randomly Initialize $S$.
  - $\times^{(i)}$ b. Run Llogds → compute Cost (SSE)
  - c. Repeat (a,b) say 10 times
  - d. Return final $S$ w/ lowest cost.
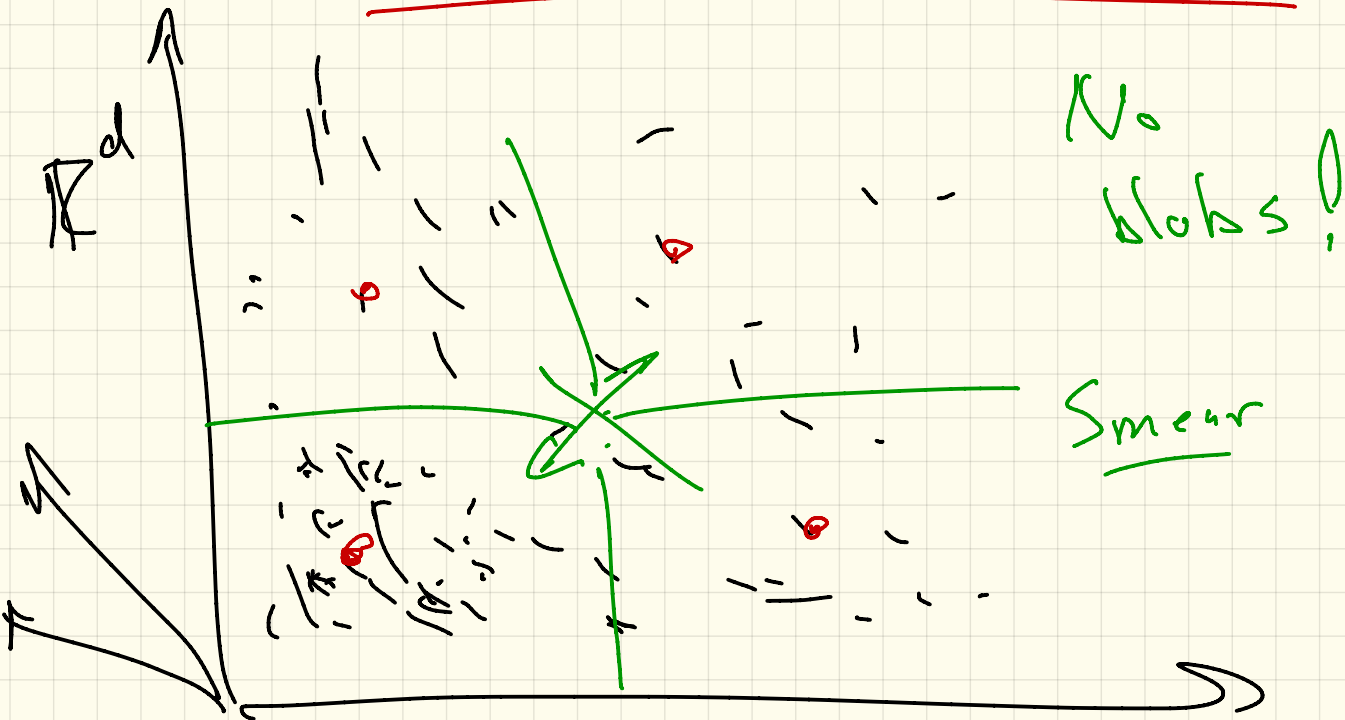
- Better Initialize
  - Gonzalez Algo
  - k-means++

# Corner Cases

- Might be site w/ no points closest to it

$S_1$

$S_3$

$0$

$S_7$

from $X_j$ w/ most variance.

$3_4$

$\Rightarrow$ Randomly assign $S_4$ (w/ no points in $X_4$) to some $x_i \in X$.

# Most Data Looks Like

$\mathbb{R}^d$

No blobs!

Smear

# Number of Clusters.  $k$?

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



$\text{Cost}(X, S)$

$k=7$

$k=3$

$k=1$

$k=4 \text{ or } 5$

where flattens out

$k=1 \quad k=7$

elbow

$k=n$

$\text{Cost}(X, S) = 0$

smear

w/ smear

w/ clusters

# Why Lloyd's Algo converges?

$$\text{cost}(X, S) = \sum_{x \in X} \| x - \phi_s(x) \|^2 \quad \leftarrow (a)$$

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

(a) assignment

(b) centering