

# L9: Assignment-based Clustering (e.g. k-means)

Sep 17, 2025  
Data Mining

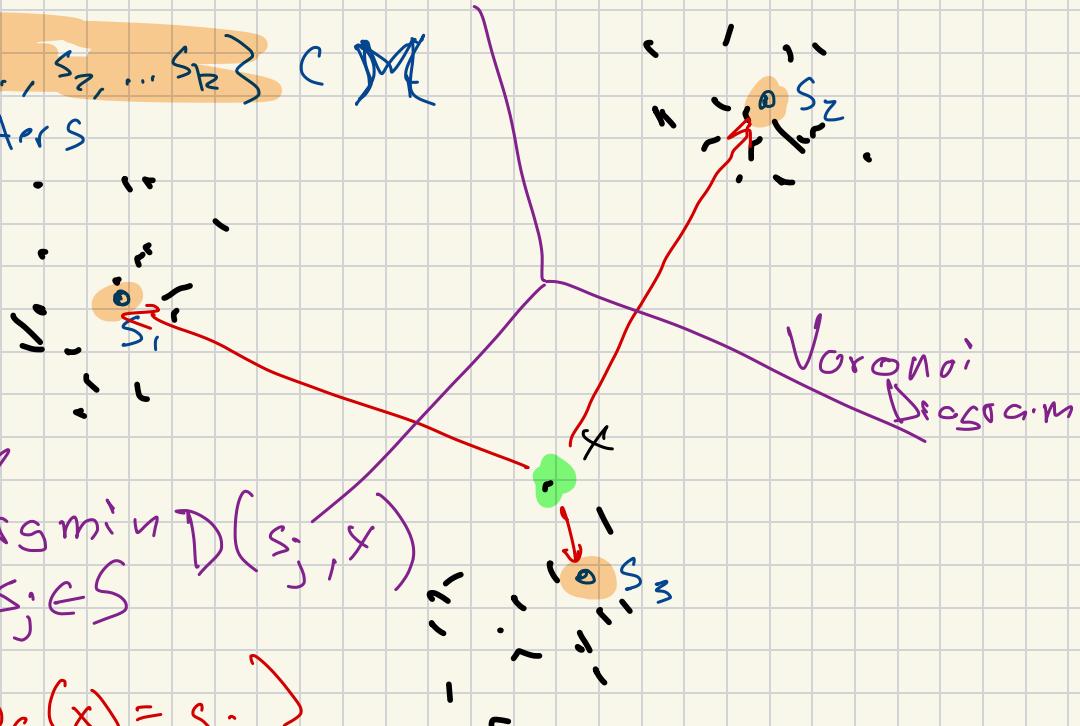
Jeff M. Phillips

Input : data  $X \subset \mathbb{M}$

$$k=3$$

distance  $D: M \times M \rightarrow \mathbb{R}_{\geq 0}$

Output : set  $S = \{s_1, s_2, \dots, s_k\} \subset M$   
sites / centers



Mappings Function

$$s_j^* = \phi_S(x) = \underset{s_j \in S}{\operatorname{arg\,min}} D(s_j, x)$$

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

defined implicitly

## Variants Assignment-based clustering

Define cost of clustering, want to minimize.

- k-means  $\text{Cost}_z(X, S) = \frac{1}{|X|} \sum_{x \in X} D(x, \phi_S(x))^z$

[Lloyd's Also,  $X \subset \mathcal{X} = \mathbb{R}^d$ ,  $D(x, s) = \|x - s\|_c$ ]

- k-median  $\text{Cost}_1(X, S) = \frac{1}{|X|} \sum_{x \in X} D(x, \phi_S(x))$

better w/ outliers, / not great also.

- k-medoid  $\text{Cost}_1(X, S)$ , but  $S \subset X$ .

- k-center  $\text{Cost}_{\infty}(X, S) = \max_{x \in X} D(x, \phi_S(x))$

outliers, Gonzalez Also.

# Gonzalez Algo for $k$ -center

Input:  $X \in \mathbb{R}^n$ , distance  $D$  (metric),  $k = \# \text{clusters}$

## Gonzalez

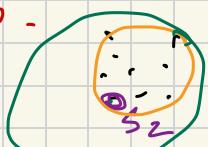
0.  $s_1 \in X$  (arbitrarily?)  $S_1 = \{s_1\}$

1. for  $j = 2 \dots k$

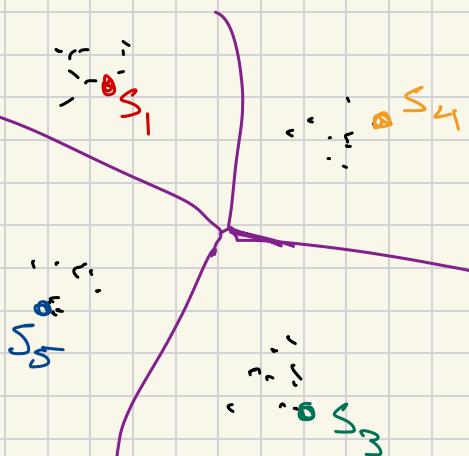
Set  $s_j = \underset{x \in X}{\operatorname{argmax}} D(x, \phi_{S_{j-1}}(x))$

2. Return  $S = S_k$

Distance to  $s_1$   
sites already chosen.



$$S_j = \{s_1, s_2, \dots, s_j\} \\ j \in [1, \dots, k]$$



# Lloyd's Algorithm for k-means

Input  $X \subset \mathbb{R}^d$ ,  $D = \| \cdot - \cdot \|_2$ ,  $k$

Lloyd's Also

o. Choose  $S$  (arbitrarily?)  $\subseteq X$

1. repeat

la. For all  $x \in X$ , set  $\phi_S(x) = \underset{s_j \in S}{\operatorname{argmin}} \|x - s_j\|$

lb. For all  $j \in [k]$ ,  $s_j = \text{average } \{x \in X \mid \phi_S(x) = s_j\}$

2. until  $(S$  unchanged)



$$s_j = \text{average} \{x \in X \mid \phi_j(x) = s_j\}$$

$X \subset \mathbb{R}^d$ ,  $D(x, s) = \|x - s\|$

$$\text{Let } X_j = \{x \in X \mid \phi_j(x) = s_j\}$$

Find  $s_j^* = \underset{s \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{x \in X_j} \|x - s\|^2$

$$= \text{average} \{x \in X_j\}$$

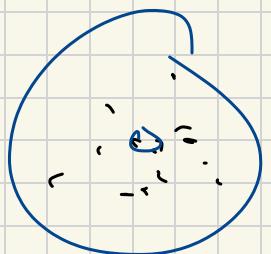
Both (1a) assignment  
(1b) coverage

$\rightarrow$  decrease  $\text{Cost}_j(x, s)$

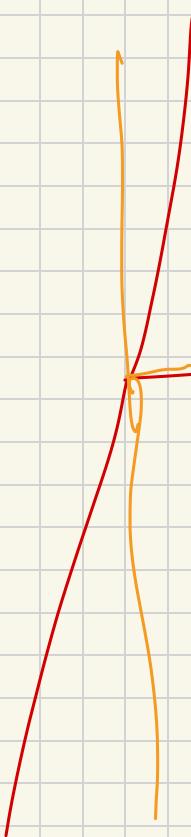
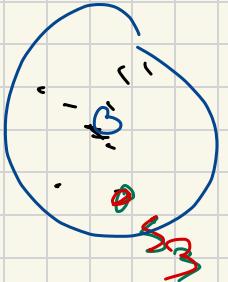
Lloyd's Algo will terminate.

$$\text{Cost}_r(x, S) = \sum_{x \in X} \sum_{S_j \in S} \left( \|x - S_j\|^2 \text{ if } q(x) = j \right)$$

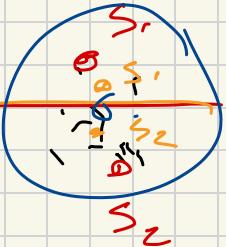
$$= \sum_{S_j \in S} \sum_{x \in X_j} \|x - S_j\|^2$$



$s_3$



$tz = 3$   
Logs can get stuck in local opt.



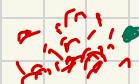
blue is better  
orange terminates.

How to initialize Lloyd's ?

1. Random restart, w  $\sim \mathcal{N}$  of random

2. Gonzalez Algo. (deterministic)

3. k-means ++



## k-means + T

0. Initialize  $s_i \in X$

1. for  $k=2$  to  $n$

choose  $s_j \sim X$  proportional to  
 $\|Q_{s_{j-1}}(x) - x\|^2$

