

Data Mining

L9 - Assignment-based Clustering

Input • $X \subset \mathbb{R}^d$

data point $x_i \in X$

• distance metric $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Goal: $\mathcal{S} = \{S_1, S_2, \dots, S_k\} \leftarrow \text{clusters}$

$$S_i \subset X$$

$$S_i \cap S_j = \emptyset$$

$$\bigcup S_i = X$$

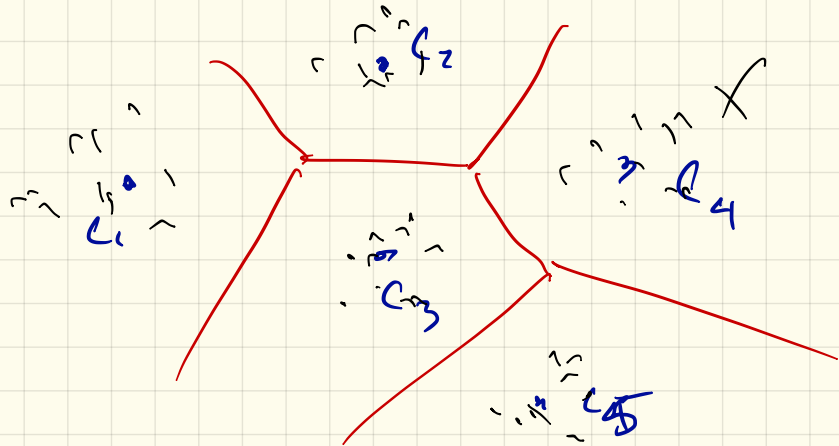
Clusters S_1, S_2, \dots, S_k

center $C = \{c_1, c_2, \dots, c_k\}$ ← representative pts.

Nearest-Neighbor
function

$$\phi_C: \mathbb{R}^d \rightarrow C$$

$$\phi_C(x) = \arg \min_{c_i \in C} d(x, c_i)$$



Goal: Find $\mathcal{C} = \{c_1, c_2, \dots, c_k\}$ k given

Formulations

k -means: minimize $\sum_{x \in X} d(x, \phi_{\mathcal{C}}(x))^2$
Llogds
 $d = \text{Euclidean}$

k -center: minimize $\max_{x \in X} d(x, \phi_{\mathcal{C}}(x))$
Gonzalez

k -median: minimize $\sum_{x \in X} d(x, \phi_{\mathcal{C}}(x))$

k -medoid: minimize $\sum_{x \in X} d(x, \phi_{\mathcal{C}}(x))$
 $\mathcal{C} \subset X$

Gonzalez Alg. for K-center

Build C incrementally $C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_k$

$$|C_i| = i$$

$$C = \bigcup_{i=1}^k C_i$$

0. choose c_1 ← arbitrarily $\Rightarrow C_1 = \{c_1\}$

1. for $j=2$ to k

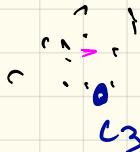
$$\text{Set } c_j = \arg \max_{x \in X} d(x, \Phi_{C_{j-1}}(x))$$

if d
metric

↳ output

2-approx

of
optimal!



$O(kn)$ time

Lloyd's Algo for k -means

$d = \text{Euclidean}$

0. Choose k pts $\rightarrow C_i$

repeat

1a. For all $x \in X$, find $\phi_c(x) \rightarrow S_1, S_2, \dots, S_k$

1b. For all $i \in 1 \dots k$, let $C_i = \text{average}(S_i)$

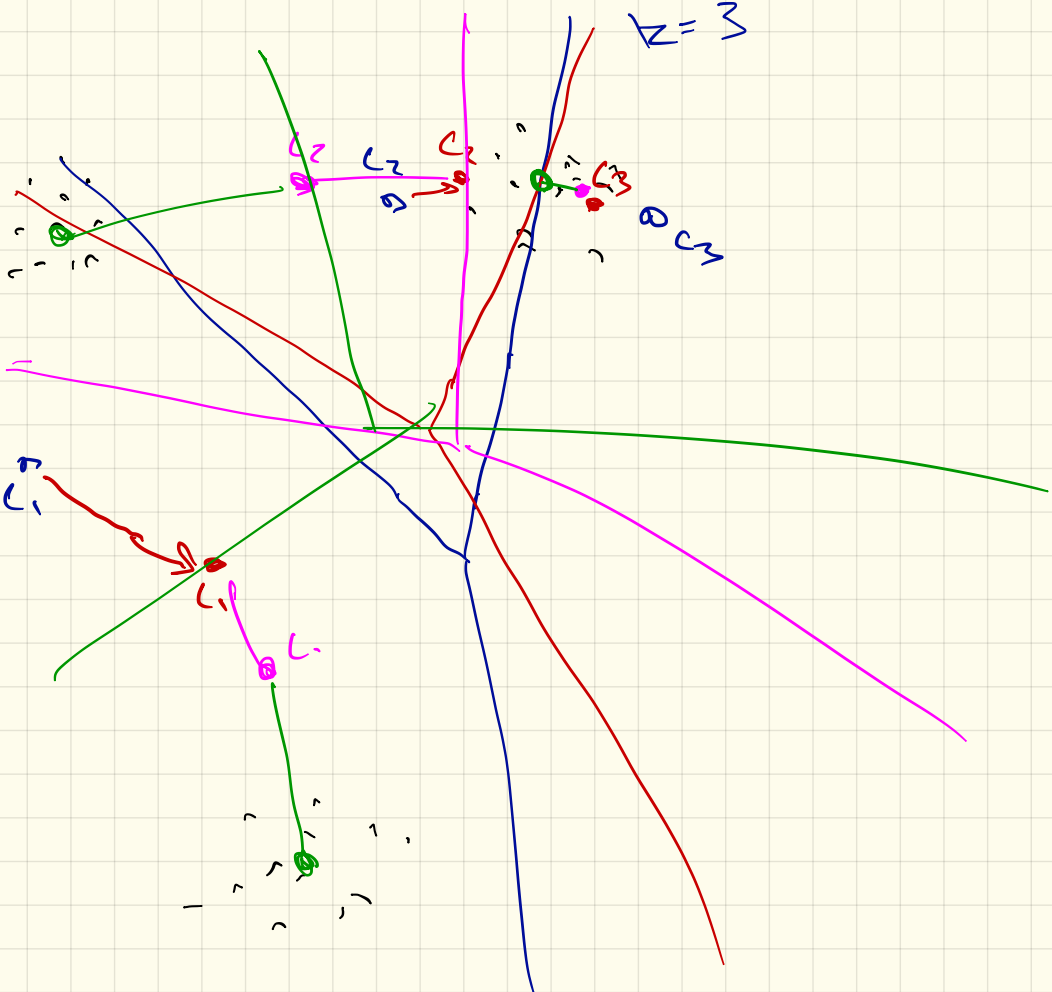
2 until (S unchanged
or
change is small)

$$S_i = \{x \in X \mid \phi_c(x) = c_i\}$$

$$C_i = \frac{1}{|S_i|} \sum_{x \in S_i} x$$

\uparrow argmin $\sum_{x \in S_i} \|z - x\|^2$
 $z \in \mathbb{R}^d$

usually
2-20 steps



$k=3$

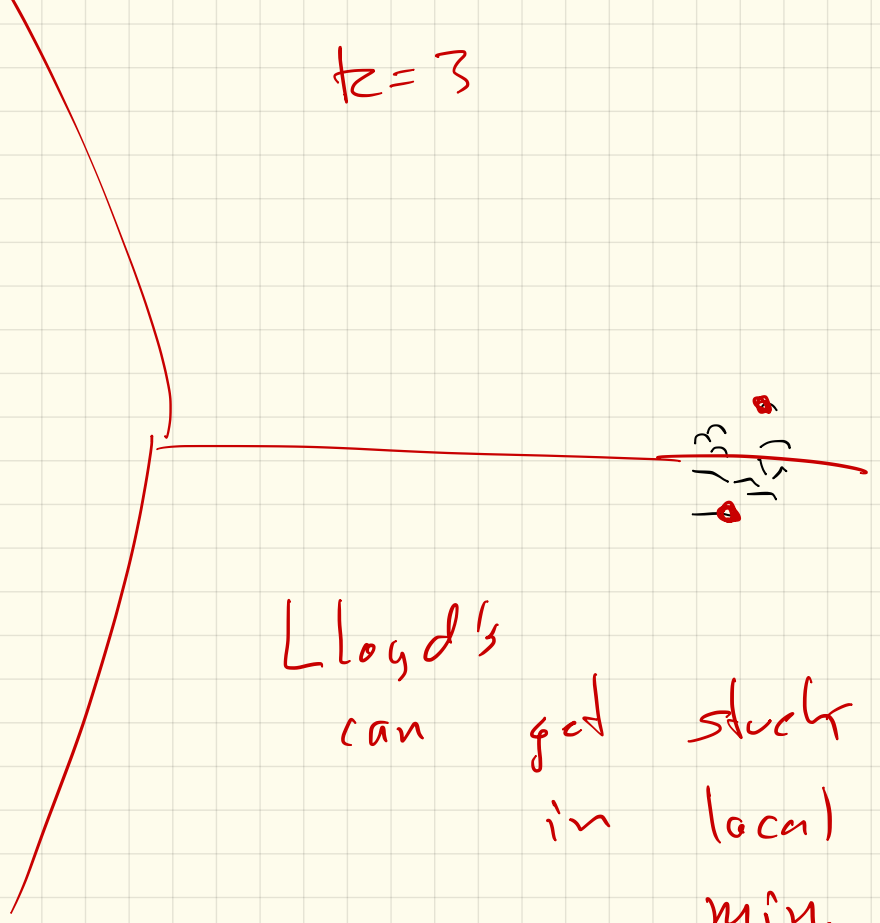
100
100
100

⊖

100
100
100

100
100
100

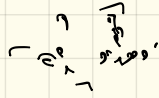
Lloyd's
can get stuck
in local
min.



Choose initial k Centers

1. Pick random subset $C \subset X$

x_1



x_2

2. Gonzalez Algo.



x_3

x_4

3. k -means ++

k-means ++

0. Choose c_i arbitrarily, $c_i \in X$

1. for $i = 2$ to k

Choose c_i from X w/ prob proportional

$$v_j = d(x_j, \phi_{c_{i-1}}(x_j))^2$$

$$V = \sum_{j=1}^n v_j$$

$$\text{prob}(x_j) = \frac{v_j}{V} = p_j$$

$\sum_j p_j = 1$

