# L12: Assignment-based Clustering

k-means, k-center, k-median, k-medioid

Feb 19, 2025

Jeff M. Phillips

Input : 1. $X \subset \mathcal{X} = \mathbb{R}^d$    $X = \{x_1, x_2, \dots x_n\}$

data point $x_i \in \mathbb{R}^d$

2. Distance $D : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_{\geq 0}$

metric

Lloyds requires $\mathcal{X} = \mathbb{R}^d$, $D = \|\cdot - \cdot\|_2$

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

$S_j < X$     subsets

$S_i \wedge S_j = \phi$ (hard clustering)

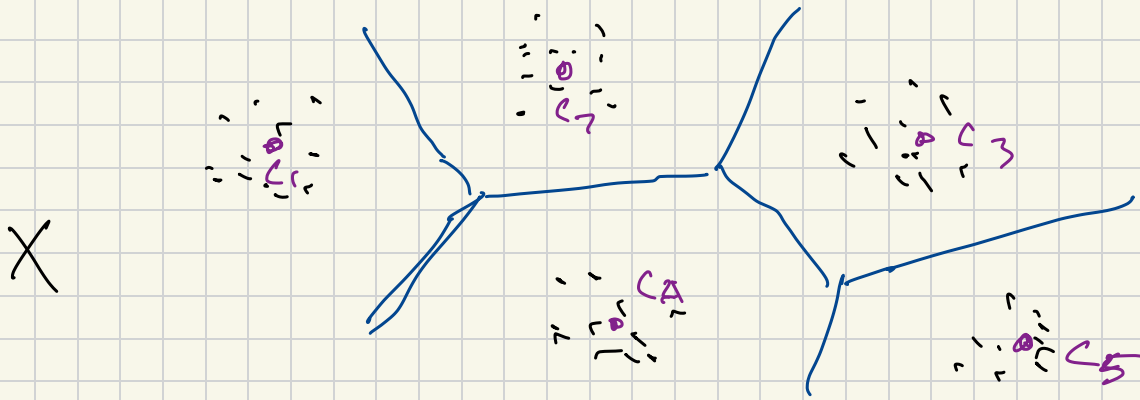$\bigcup_j S_j = X$

# Assignment-based Clustering

Clusters $S_1, S_2, \ldots S_k$

Centers $C = \{c_1, c_2, \ldots c_k\} \subset \mathbb{R}^d$

$\uparrow c_i$ representative of $S_i$

Nearest Neighbor: $\phi_C : \mathbb{R}^d \to C$

$$\phi_C(x) = \underset{c_j \in C}{\arg\min} \; D(x, c_j)$$



$X$

$$\underline{\text{Goal}} \qquad \text{Find} \qquad C = \{ c_1, c_2, \dots c_k \}$$

## Formulations

$\underline{\text{k-means}}$ : minimize $\sum\limits_{x \in X} D(x, \phi_C(x))^2$

$\quad$ $\color{red}{\text{Lloyds}}$ $\qquad \color{red}{D = \text{Euclidean}}$

$\underline{\text{k-center}}$ : minimize $\max\limits_{x \in X} D(x, \phi_C(x))$

$\qquad \color{red}{\text{Gonzalez}}$

$\underline{\text{k-median}}$ : minimize $\sum\limits_{x \in X} D(x, \phi_C(x))$

$\underline{\text{k-medioid}}$ : minimize $\sum\limits_{x \in X} D(x, \phi_C(x))$

$C \subset X$

# Gonzalez Algo for k-center

Build centers $C_i$ incrementally $C_1 \to C_2 \to \dots \to C_k = \{c_i\}$

$|C_j| = j$

Gives 2-approx. of optimal.

0. Choose $c_1$ arbitrarily $\Rightarrow C_1 = \{c_1\}$

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

   Set $c_j = \displaystyle\arg\max_{x \in X} D(x, \phi_{C_{j-1}}(x))$

NP-hard to do better then 2-approx

"Centers" on boundary

$k = 5$

$c_1$

$c_2$

$c_3$

$c_5$

$c_4$

# Lloyd's Algorithm for k-means

$$D = \text{Euclidean}$$

0. $\boxed{\text{Choose } k \text{ pts} \Rightarrow C_j}$ Do something / Random

1. repeat

   1a. For all $x \in X$, find

   1b. For all $j \in 1 \ldots k$, let

2. until ( $S$ unchanged

   or change is small )

in practice $k = 20$ iterations

$\boxed{\phi_q(x) \rightarrow c_j \Rightarrow x \rightarrow S_j}$ implicit

$$\boxed{c_j = \text{average}(S_j) = \frac{1}{|S_j|} \sum_{x \in S_j} x}$$

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

Not Global opt
for $k$-means

$c_1$

$c_2$

$c_3$

Lloyd's Also
stops!

# Initialize Lloyd's Algo for k-Means
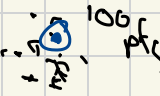
1. **Init** Choose $C$ at random.

2. **Init** $C \leftarrow$ Gonzalez $(X)$

always choose

3. k-means++

# K-means++

0. Choose $c_1$ arbitrarily $c_i \in X$

1. for $j = 1$ to $k$

   Choose $c_j$ from $X$    w/ probability proportional to

   $x_i \in X$

   $$V_i = D(x_i, \phi_{c_j}(x_i))^2$$

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

Select $x_i \sim X$

w.p. $\text{prob}(x_i) = \dfrac{V_i}{V}$

Implement $\boxed{U \sim Unif(0,1]}$

$V_i = \sum_{i=1}^{i} v_i$

$U = 0.43$

0

$V_1$ $V_2$ $V_i$ $V_c$

$V_{i<}$ $V_i$

1