# L9 Assignment-Based Clustering

<u>Input:</u>  1) $X \subset \mathbb{R}^d$

   • dale point $x_i \in X$

   2) Distance  $d: \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$
      metric

<u>Goal:</u>  $S = \{ S_1, S_2, S_3, \ldots S_k \}$

   • $S_i \subset X$

   • $S_i \cap S_j = \emptyset$

   • $\cup S_i = X$

## Assignment - Based Clustering

• Clusters  $S_1, S_2, S_3, \ldots, S_k$

• each with a "centroid"
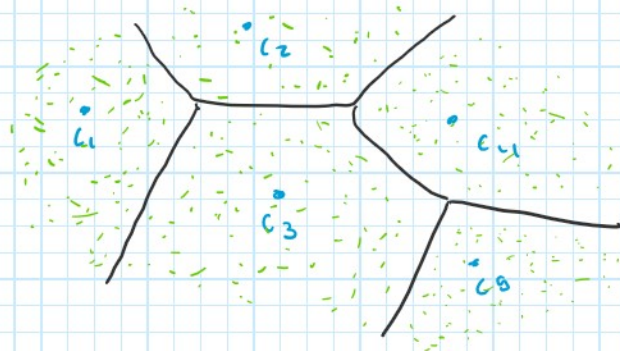
   $G = \{ c_1, c_2, \ldots, c_k \}$    ⎫  restricts the type of clusters
                                        ⎭  being considered

   ↳ representation points

Nearest Neighbor Function

$$\phi_G : \mathbb{R}^d \to G$$



• In particular...

   $\phi_G(x) = \underset{c_i \in G}{\arg\min}\ d(x, c_i)$

<u>Goal:</u> Find $G = \{ c_1, c_2, \ldots c_k \}$ given some $k$.

Formulations:

• K-means :  minimize $\sum_{x \in X} d(x, \phi_G(x))^2$  ← most popular

   Loyds                    $d =$ Euclidean

• K-center : minimizes $\underset{x \in X}{\max}\ d(x, \phi_G(x))$

   Gonzalez

• K-median : minimize $\sum d(x, \phi_G(x))$

• K-mediod : minimize $\sum d(x, \phi_{G'}(x))$
                                        $G' \subset X$

# Gonzalez Algorithm For K-Center

Main Idea: Build $G$ incrementally.

$$G_1 \rightarrow G_2 \rightarrow G_3 \rightarrow \cdots \rightarrow G_K = G$$

ie, build these in K steps where $|G_i| = i$

S1. Choose some random $c_1$ arbitrarily $\Rightarrow G_1 = \{c_1\}$

S2. for $j = 2$ to $K$:

set $c_j = \underset{x \in X}{\arg\max} \; d(x_1, \phi_{G_{j-1}}(x))$

ie, the distance between $x$ and it's assigned max!

if $d$ is a **metric**, then:

↳ output 2-approximation of the optimal.

Example:



$O(Kn)$ time
- on each round (K of them)
- check re-assignment to cluster
- keep track of the max
  - the max becomes the new center

- whose assignment has the greatest cost?
- make them the next $c_i$

- implementation tip:
  - maintain an array of size equal to the # of points
  - for each point, store the index of the closest center seen so far.
  - update the array
  - This done in linear time each round, where there are K rounds.

# Lloyd's Algorithm For K-Means: (Minimizing the Sum of squared errors)

S1. choose K points $\rightarrow G$      $d = $ euclidean

S2. repeat

S2a) For all $x \in X$, find $\phi_G(x) \rightarrow S_1, S_2, \ldots, S_K$ where each $S_i = \{x \in X \mid \phi_G(x) = c_i\}$

S2b) For all $i \in 1 \ldots K$, let $c_i = \text{average}(S_i)$

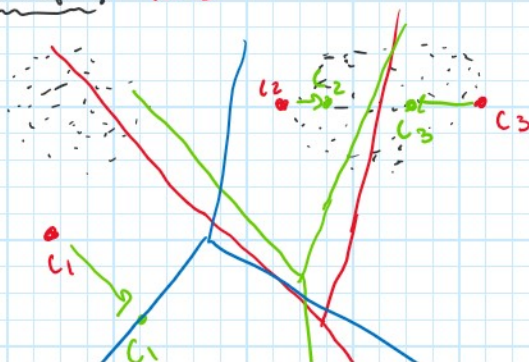S3. until ( $S$ unchanged or change is small)

\* usually $< 20$ steps

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

- add up all values of $x$ assigned to the center
- divide the number of points

$$\underset{z \in \mathbb{R}^n}{\arg\min} \sum_{x \in S_i} \|z - x\|^2$$
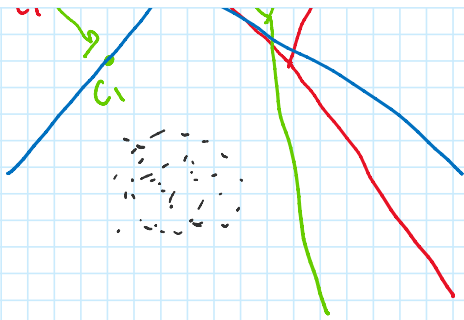
(the point which minimizes the squared errors)

Example:    K = 3



- re-assign centers (cost ↓ each round)
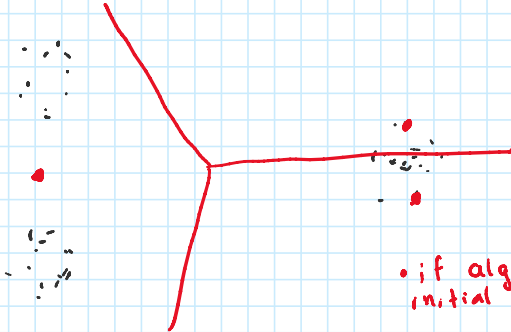- re-assign the points (cost ↓ each round)

↳ thus, we are converging

↳ thus, we are converging

But Loyd's can get stuck at a local min

ex

• if algo gets stuck, just start fresh with new initial points.

## How do we choose the initial K centers?

1. Pick random subset $G \subset X$ (random indices from the set of $x$)

   • issue: small clusters are unlikely to receive a centroid
     ↳ the algorithm will only divide the center, which is sub-optimal

2. Use the Gonzalez Algorithm to initialize points

   • issue: sensitive to outliers. Also deterministic.

3. k-means++

   s1. choose $c_1$ arbitrarily where $c_1 \in X$

   s2. for $i = 2$ to $k$:

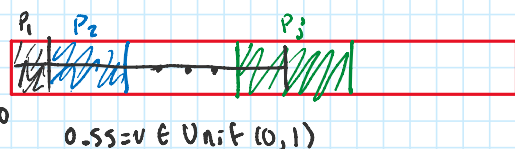        Choose $c_i$ from $X$ w/probability proportional to

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

   (e, instead of picking longest distance, we will calculate all the distances, square them, and then sample from the datapoints were we're more likely to pick the points with a larger squared difference.

   $$V = \sum_{j=1}^{n} v_j \implies prob(x_j) = \frac{v_j}{V} = p_j$$

   most important to implement of k-means++

   $p_1$ $p_2$     $p_3$

   $(\Sigma p_j = 1)$

   0     $0.55 = v \in Unif(0,1)$     1

   (each time we pick a center, we generate a probability distribution and pick an element from it)