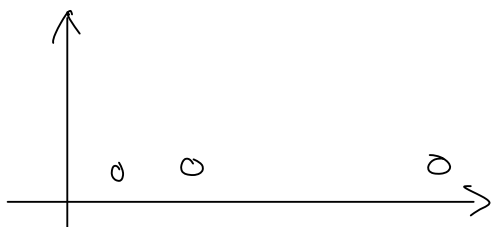# K-MEANS

## Main Idea:
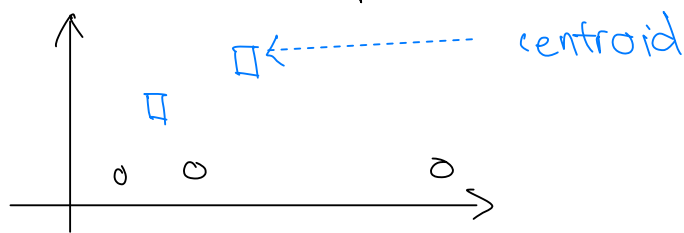
Separating data points into K clusters by alternatively set the assignments and centroids.

For example:

- Given some data points:
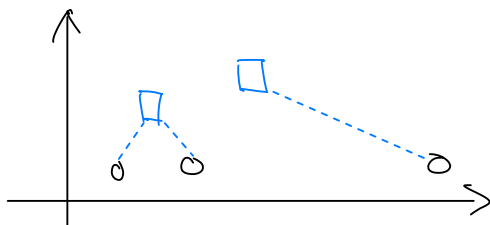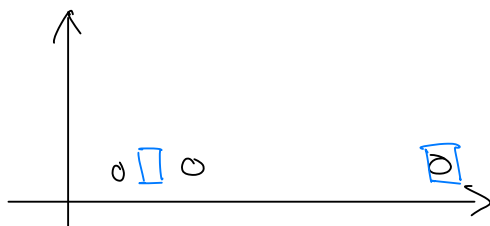
- Initialize random centroids:

  centroid

- Iterate through these steps:

  - Assignment: Assign each point to the nearest centroid

  - Centroid: Update position of centroids (mean of assignments)

## K-means algorithm

- **Input:** dataset $\{x_i\}_{i=1}^{n}$ and a number $K$ of clusters

- **Algorithm:**

  - Initialization: Randomly placed $K$ centroids

  - Iterate till converge:

    i) Assignment: For each $x_i$, assign it to closet centroid

    $$z_i = \underset{k=1,\dots,K}{\arg\min} \| x_i - \mu_k \|_2$$

    centroid index assigned to $x_i$ ← $z_i$

    $\mu_k$ → $k^{th}$ centroid

    $x_i$ → $i^{th}$ data point

    ii) Centroids: For each centroid, update based on new assignments

    $$\mu_k = \frac{1}{|S_k|} \sum_{i \in S_k} x_i \quad , \text{ where } S_k = \{ i : z_i = k \}$$

    $\mu_k$ → centroid of $k^{th}$ cluster

    $S_k$ → $k^{th}$ cluster

# K-means as Optimization

K-means can be viewed as a "==coordinate descent==" algorithm for ==optimizing== a ==objective function== of centroids and assignments.

## Problem definition:

- Given: $\{x_i\}_{i=1}^n$

- Define objective function:

$$L(\mu, z) = \sum_{i=1}^n \| x_i - \mu_{z_i} \|^2$$

$\longrightarrow$ centroid assigned to $x_i$

- We try to minimize that function:

$$\min_{\mu, z} L(\mu, z)$$

$\searrow$ Tricky: - Mixed optimization (not convex optimization)

Because: $\begin{cases} \circ \ \mu \ \text{is continuous} \\ \circ \ z \ \text{is discrete} \end{cases}$

$\Rightarrow$ Many local optima

## Solution to mixed optimization problem:

Use coordinate descent to find the local optima

- <u>Coordinate descent:</u>

  - Initialize $\mu_0$

  - Repeat: (at iteration $t^{th}$)        $\dashrightarrow$ actually K-means

    i) Update $z$, with fixed $\mu$                                algorithm

    vector $\leftarrow z^t = \underset{z}{\arg\min} \ L(\mu^t, z)$        (Proof next page)

    ii) Update $\mu$, with fixed $z$

    vector $\leftarrow \mu^{t+1} = \underset{\mu}{\arg\min} \ L(\mu, z^t)$

- <u>To find the global optima</u>, run many coordinate descent with different initialization values and return the best optima value

# Proof coordinate descent on $L(\mu, z)$ is actually K-means algorithm

○ We know that: $L(\mu, z) = \sum\limits_{i}^{n} \| x_i - \mu_{z_i} \|^2$

○ Coordinate descent says that:

  ○ Update $z$, fixed $\mu$

$$\underbrace{z^t}_{\text{Vector}} = \arg\min_{z} L(\mu^t, z)$$

$$= \arg\min_{z} \sum\limits_{i}^{n} \| x_i - \mu_{z_i} \|^2$$

$$\Rightarrow \underbrace{z_i}_{\text{scalar}} = \arg\min_{z_i} \| x_i - \mu_{z_i} \|^2$$

   This is __assignment step__ in K-means

  ○ Update $\mu$, fixed $z$

$$\underbrace{\mu^{t+1}}_{\text{Vector}} = \arg\min_{\mu} L(\mu, z^t)$$

$$= \arg\min_{\mu} \sum\limits_{i}^{n} \| x_i - \mu_{z_i} \|^2$$

$$\underbrace{\mu_k}_{\text{Scalar}} = \arg\min_{\mu_k} \sum\limits_{i \in S_k} \| x_i - \mu_k \|^2 \quad , \text{ where } S_k = \{ i : z_i = k \}$$

$$= \arg\min \left[ \underbrace{\sum \| x_i \|^2}_{\text{const}} - 2 \left( \sum x_i \right) \mu_k + |S_k| \cdot \| \mu_k \|^2 \right]$$

$$\nabla_\mu f(\mu_k) = 0$$

$$\boxed{\Rightarrow \mu_k = \frac{1}{|S_k|} \sum\limits_{i \in S_k} x_i}$$

   This is __centroids step__ in K-means

○ __Conclusion__ :

   Coordinate descent on $L(\mu, z) \iff$ K-means algorithm