

## Clustering and K-means

---

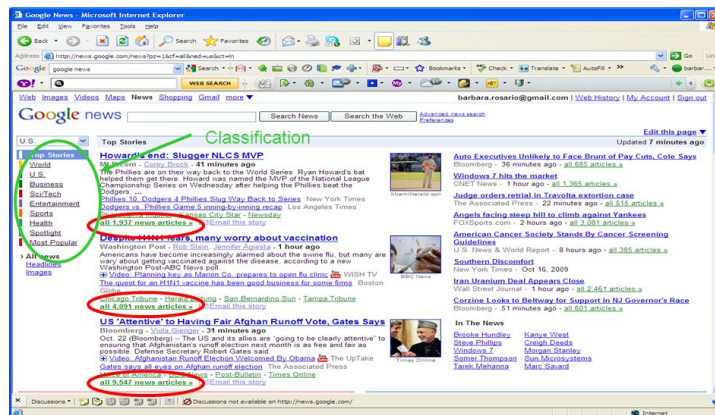
**Qiang Liu**  
UT Austin

## Clustering

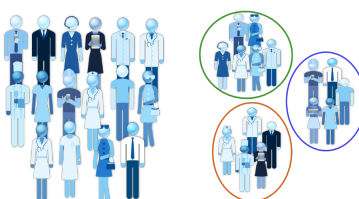
- Clustering: Partition the dataset into groups based on their similarity.



- Google News: automatic clustering gives an effective news presentation metaphor



- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs



# K-means Algorithm

- **Inputs:**  $n$  objects (data points)  $\{x_i\}_{i=1}^n$  and a number  $K$  of clusters.

- **K-means:**

- **Idea:**

Alternatively optimize

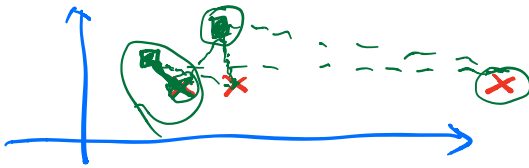
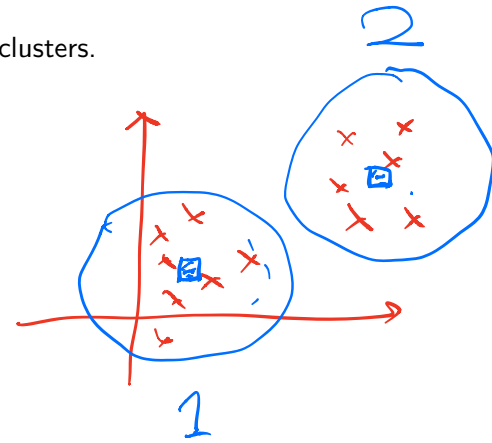
- i) the **assignment** of the data points different clusters.
    - ii) the **centroids** of the clusters.

- **Algorithm:**

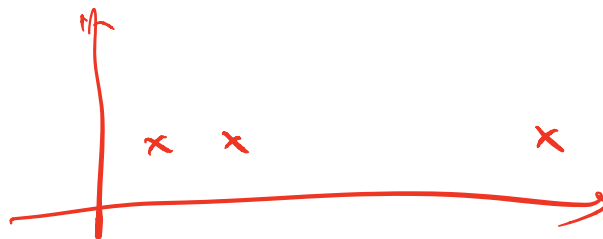
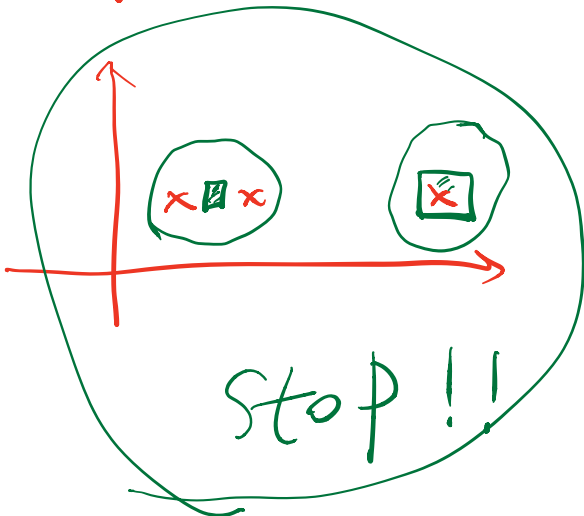
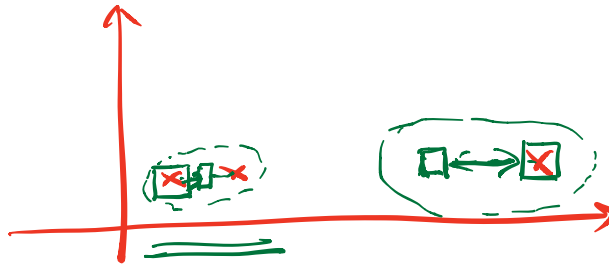
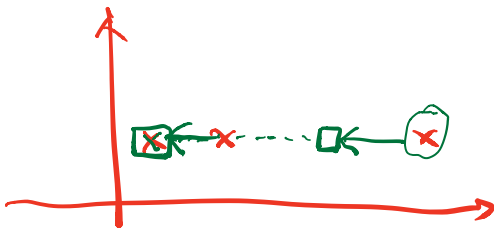
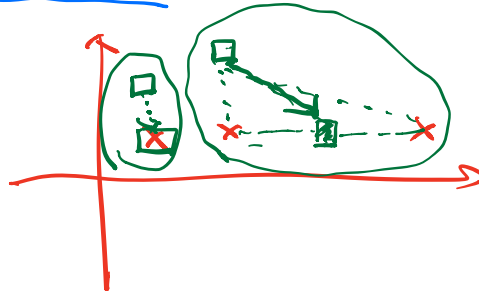
Initialize the centroids or assignments randomly;

Iterate until convergence:

- i) **Assign each data point**  $(x_i)$  to the cluster that has the closest centroid.
    - ii) **Update the centroids** of all the clusters.



$k=2$



# K-means Algorithm

● Inputs:  $n$  objects (data points)  $\{x_i\}_{i=1}^n$  and a number  $K$  of clusters.

● **K-means Algorithm:**

● Initialization: randomly place  $K$  points (as the centroids)

● Iterate until convergence:

● i) Assign each object to the group that has the closest centroid

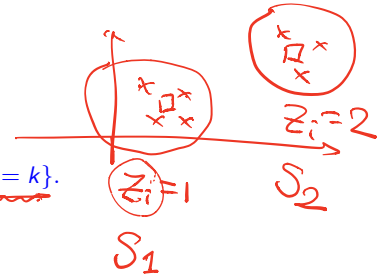
$$z_i = \arg \min_{k=1, \dots, K} \|x_i - \mu_k\|$$

$$z_i \in \{1, \dots, K\}$$

● ii) Recalculate the position of the  $K$  centroids

$$\mu_k = \frac{1}{|S_k|} \sum_{i \in S_k} x_i$$

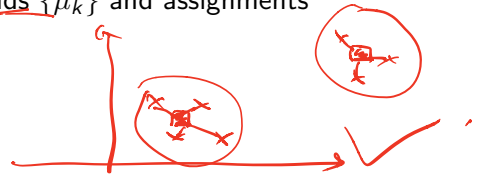
$$S_k = \{i : z_i = k\}$$



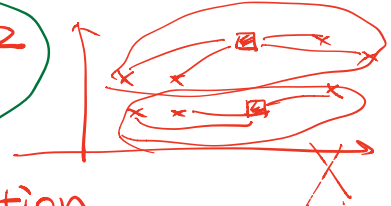
## K-means as Optimization

- K-means can be viewed as a "coordinate descent" algorithm for optimizing an objective function of the centroids  $\{\mu_k\}$  and assignments  $\{z_i\}$ .

Given  $\{x_i\}_{i=1}^n$   
Want  $\{\mu_k\}_{k=1}^K$   $\{z_i\}_{i=1}^n$   
 $\underline{\mu}$   $\underline{z}$



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



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

mixed optimization

$$\underline{\mu} \in \mathbb{R}^{d \times K} \quad \underline{z} \in \{1, \dots, K\}^n$$

Coordinate descent.

Initialize  $\mu_0$

Repeat: (t-iteration)

① update  $\underline{z}$  with fixed  $\underline{\mu}$

$$\underline{z}_t = \underset{\underline{z}}{\operatorname{argmin}} L(\underline{\mu}_t, \underline{z})$$

② update  $\underline{\mu}$  with  $\underline{z}$  fixed

$$\underline{\mu}_{t+1} = \underset{\underline{\mu}}{\operatorname{argmin}} L(\underline{\mu}, \underline{z}_t)$$

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

Find optimal  $z_i \leftarrow \underset{z_i}{\operatorname{argmin}} \|x_i - \mu_{z_i}\|^2$  (Assignment Step)

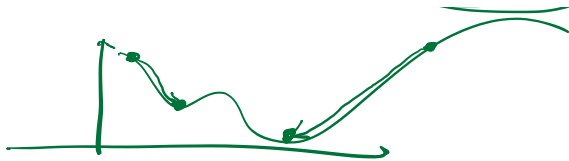
Find optimal  $\underline{\mu}_K \leftarrow \underset{\underline{\mu}_K}{\operatorname{argmin}} \sum_{i \in S_K} \|x_i - \mu_K\|^2$

$$S_K = \{i : z_i = K\}$$

$$= \sum_{i \in S_K} \|x_i\|^2 - 2 \sum_{i \in S_K} x_i^T \mu_K + |S_K| \|\mu_K\|^2$$

$$= \text{const} - 2 \left( \sum_{i \in S_K} x_i \right)^T \mu_K + |S_K| \|\mu_K\|^2$$

$$\Rightarrow \underline{\mu}_K \leftarrow \frac{\sum_{i \in S_K} x_i}{|S_K|} \quad (\text{Centroid update Step})$$



$$(u_t, z_t) \quad \text{3}$$
$$\Rightarrow L(u_{t+1}, z_{t+1}) \leq \underline{L(u_t, z_t)}$$

$$\underline{L(u_t, z_t)} \geq \underline{L(u_{t+1}, z_t)}$$
$$\geq \underline{L(u_{t+1}, z_{t+1})}$$

*[Handwritten scribbles and wavy lines]*