

Clustering (K-means)

Given dataset $D = \{\vec{x}_1, \dots, \vec{x}_n\}$ $\vec{x}_i \in \mathbb{R}^d$
 $k \in \mathbb{Z}^+$

Partition D in k partitions (clusters) $C = \{C_1, \dots, C_k\}$

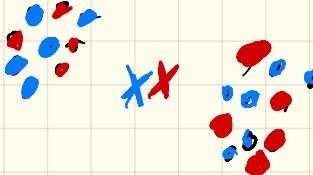
Representative clustering

$\forall C_i \exists$ rep point to summarize the cluster

commonly mean/center

e.g. for $C_i \ni |C_i| = n_i$

$$\text{Good! } \vec{m}_i = \frac{1}{n_i} \sum_{\vec{x}_j \in C_i} \vec{x}_j \quad \text{Bad!}$$



How to get started:

1. opt criterion

2. generate all k -partitions of points

3. pick the partition w/ highest (lowest) score

How many ways to create k -partitions of n parts

$O\left(\frac{k^n}{k!}\right)$ -- even for tiny n too much work

K-means algo

1. pick opt criterion sum of squared error

$$SSE(C) = \sum_{i=1}^k \sum_{\substack{x_j \in C_i}} \|x_j - \mu_i\|^2$$

Goal: find $C^* = \underset{C}{\operatorname{argmin}} \left\{ SSE(C) \right\}$



clustering that minimizes SSE

Idea of k-means

- greedy
- iterative

Note: we can get stuck in a local opt

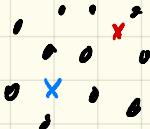
High level

1. init each μ_i
 (typically, randomly sample from range of each dim
 or bounding box of data if range unknown)

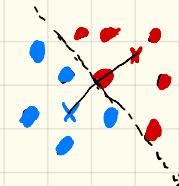
2. at each iteration

- a. cluster assignment

assign each point to closest center



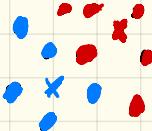
assign each point to closest center



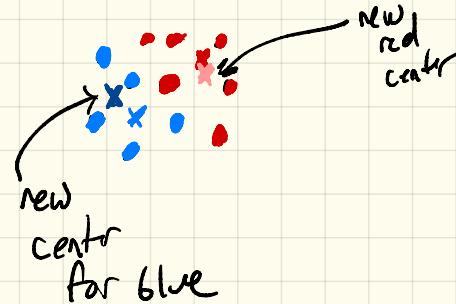
assign each point to cluster

$$C_{j^*} \text{ w/ } j^* = \arg \min_{i=1}^k \{ \|x_j - \mu_i\|^2 \}$$

- b. centroid update



update center to mean of points in cluster



3. continue until centers don't move too much

e.g. μ_i^t is center for cluster i at iteration t

$\epsilon > 0$ (user specified)

Stop when $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\| < \epsilon$

Algorithm 13.1: K-means Algorithm

K-MEANS (\mathbf{D}, k, ϵ):

- 1 $t = 0$
 - 2 Randomly initialize k centroids: $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$
 - 3 **repeat**
 - 4 $t \leftarrow t + 1$
 - 5 $C_i \leftarrow \emptyset$ for all $i = 1, \dots, k$
 // Cluster Assignment Step
 - 6 **foreach** $\mathbf{x}_j \in \mathbf{D}$ **do**
 - 7 $i^* \leftarrow \arg \min_i \left\{ \|\mathbf{x}_j - \mu_i^{t-1}\|^2 \right\}$
 - 8 $C_{i^*} \leftarrow C_{i^*} \cup \{\mathbf{x}_j\}$ // Assign \mathbf{x}_j to closest centroid
 // Centroid Update Step
 - 9 **foreach** $i = 1, \dots, k$ **do**
 - 10 $\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$
 - 11 **until** $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$
-

Things to note:

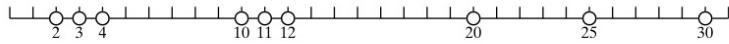
- usually run the algo a few times w/ diff initial centers report the clustering w/ min SSE (helps avoid really bad initial guesses)



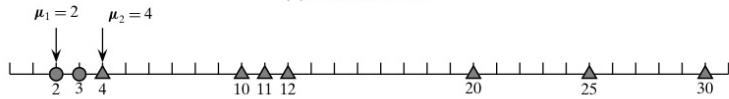
- Clusters are convex



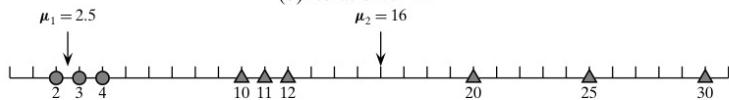
← each color is not convex which is a problem for k-means (come back to later)



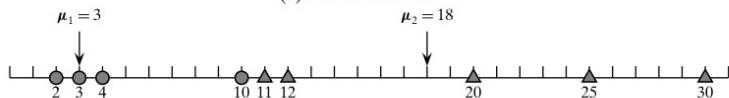
(a) Initial dataset



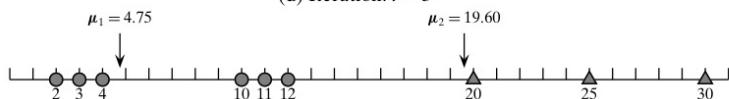
(b) Iteration: $t = 1$



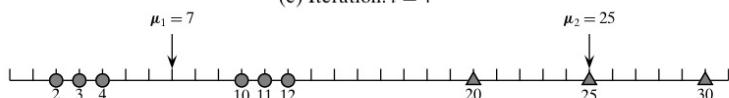
(c) Iteration: $t = 2$



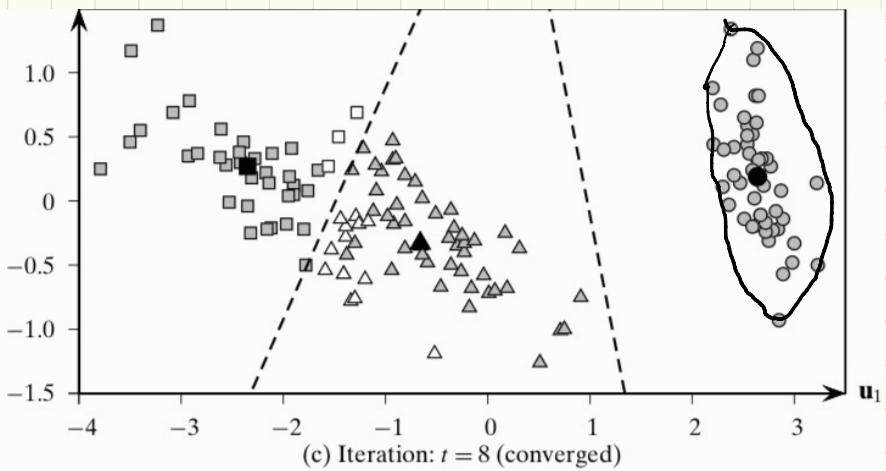
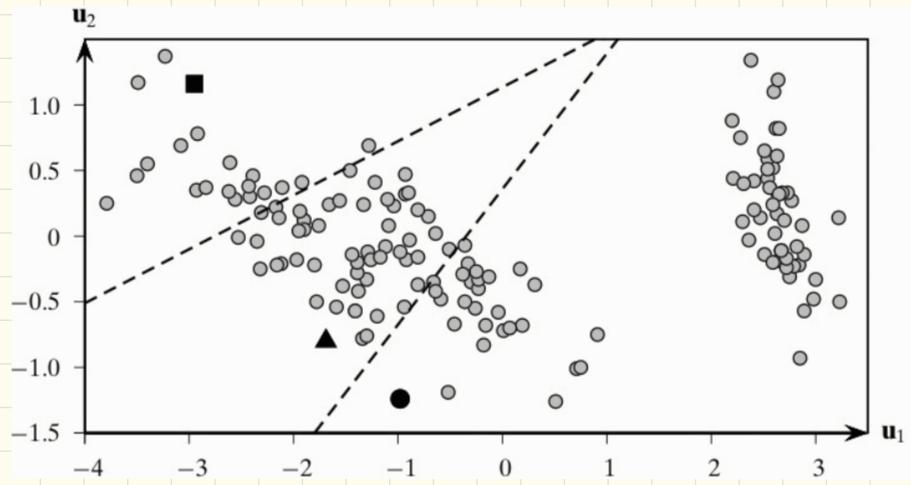
(d) Iteration: $t = 3$



(e) Iteration: $t = 4$



(f) Iteration: $t = 5$ (converged)



Computational Complexity for K-means

- a. Cluster assignment
for each point
for each cluster
compute dist

time proportional to
 $n \cdot k \cdot d$

- b. Recompute Center
for each cluster
recompute center

$d \cdot \# \text{ of points in the cluster}$

$\Rightarrow O(dn)$

Work at each iteration
 $O(nkd + dn) = O(nkd)$

Note
each point
is in exactly
1 cluster

Assume t iteration
total work is $O(tnkd)$

I/O complexity

of DB passes

proposal to

1

- a. Cluster assignment
for each point
- for each cluster
 compute dist

- b. Recompute center
for each cluster
 recompute center

I/O measured in # of
DB passes

each iter $O(1)$ db scans

Total I/O for t iterations is $O(t)$