

# 3/26/2024: K-Means Clustering

Machine Learning: Algorithms that learn from data

Supervised Learning

Training dataset:

$$D = \{ (x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \}$$

input to model  $\uparrow$   $\uparrow$  output (model predicts this)

Goal: Learn mapping from  $x$  to  $y$

Unsupervised Learning

Training dataset only contain  $x$ 's

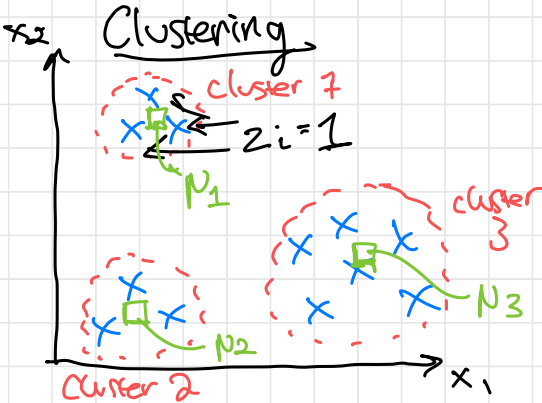
$$D = \{ x^{(1)}, \dots, x^{(n)} \}$$

No "correct output"

Goal: Learn what structure is present in data

- 1) Group/subpopulation/cluster structure
- 2) low-dimensional structure
- 3) Similarity/Relationships between words (word2vec)

Clustering



Input: Dataset =  $\{ x^{(1)}, \dots, x^{(n)} \}$  and # of clusters  $K$

Output: An assignment

$z_1, \dots, z_n$  where each  $z_i \in \{1, 2, \dots, K\}$  denotes the cluster assigned to example  $x^{(i)}$

K-means Clustering Algorithm

Idea #1: Write down a loss function to define "badness" of assignment  $z_1, \dots, z_n$

Idea #2: Add more parameters to help define loss func.  
Learn the "centroid" of each cluster

$\mu_1, \dots, \mu_k$  where

$\mu_j \in \mathbb{R}^d$  is "center of mass" of cluster  $j$

Payoff: Loss of assignment & choice of centroids  
is how far each  $x^{(i)}$  is from its  
assigned centroid

Loss:

$$L(z_{1:n}, \mu_{1:k}) = \sum_{i=1}^n \|x^{(i)} - \mu_{z_i}\|^2$$

$\uparrow$   
=  $[z_1, \dots, z_n]$

centroid of cluster assigned to  $x^{(i)}$

cluster ID assigned to  $x^{(i)}$

"Reconstruction Error":

If we only knew  
assignments  $z_{1:n}$  &  
means  $\mu_{1:k}$ ,

how far are we from  
reconstructing  $\{x^{(1)}, \dots, x^{(n)}\}$ ?

Goal: Minimize  $L(z_{1:n}, \mu_{1:k})$  with respect to  
 $z_{1:n}, \mu_{1:k}$

Can't do gradient descent because  $z_i$ 's are  
discrete i.e. only be one cluster or another,  
no "in-between"

Solution: Alternating Minimization

① Start with random choice of  $\mu_1, \dots, \mu_k$

Alternate until **Convergence** ← when  $z_{1:n}$ 's &  $N_{1:k}$ 's stop changing.

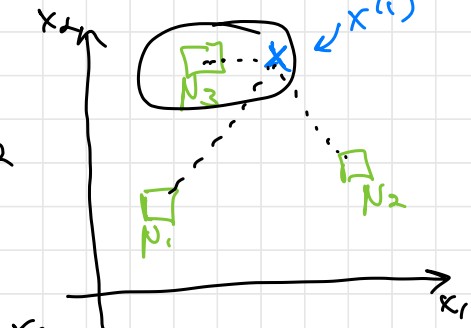
② Choose  $z_{1:n}$  to minimize  $L$  given current  $N_{1:k}$

③ Choose  $N_{1:k}$  to minimize  $L$  given current  $z_{1:n}$

**Step 1** Choose each  $N_j$  to be random distinct  $x^{(i)}$

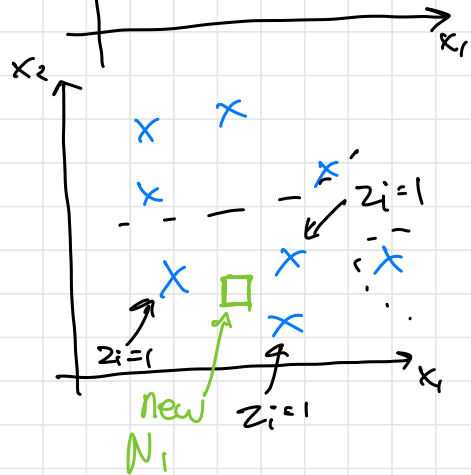
**Step 2** Minimizing  $L$  w.r.t.  $z_{1:n}$

For each  $i=1, \dots, n$   
 set  $z_i = \operatorname{argmin}_{j=1, \dots, k} \|x^{(i)} - N_j\|^2$



**Step 3** Minimizing  $L$  w.r.t.  $N_{1:k}$   
 Fortunately:  $N_j$  should be mean of all points where  $z_i = j$

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^n \|x^{(i)} - N_{z_i}\|^2 \\ = \quad & \sum_{j=1}^k \sum_{i: z_i=j} \|x^{(i)} - N_j\|^2 \end{aligned}$$



For particular index  $j$ , minimize w.r.t.  $N_j$

$\Leftrightarrow$  minimize  $\sum_{i: z_i=j} \|x^{(i)} - N_j\|^2$  Take gradient & set to 0

$$\nabla_{N_j} \sum_{i: z_i=j} \|x^{(i)} - N_j\|^2 = \sum_{i: z_i=j} 2(x^{(i)} - N_j) \cdot (-1) = 0$$

$$\sum_{i: z_i=j} x^{(i)} = \underbrace{\left| \{i: z_i=j\} \right|}_{\substack{\text{\# of points assigned} \\ \text{to cluster } j}} \cdot \mu_j$$

$$\mu_j = \frac{1}{\left| \{i: z_i=j\} \right|} \cdot \sum_{i: z_i=j} x^{(i)} = \text{mean of all } x^{(i)} \text{ where } z_i=j$$

No guarantee of finding optimal solution

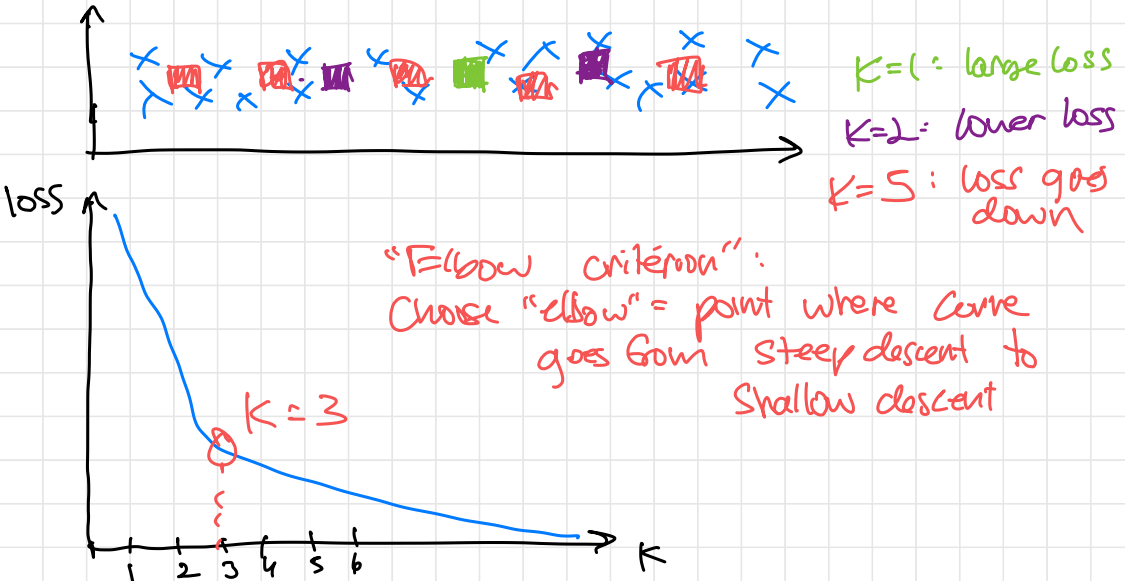
- Algorithm runs until it finds local optimum
- Random initialization affects final result

How to choose  $K$ ?

$K$  This is a hyperparameter

Wrong way: Choose  $K$  to minimize loss on dev set

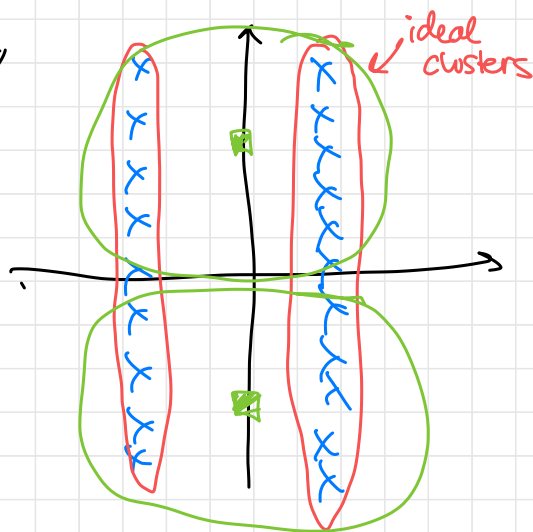
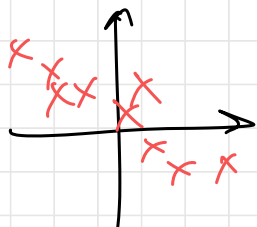
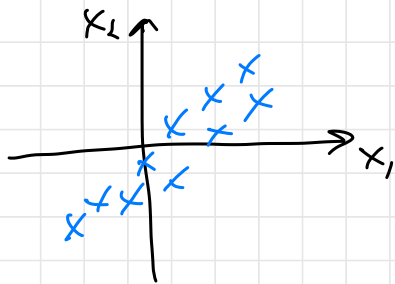
Why not? Larger  $K$  always makes loss lower



K-means uses Euclidean distance, so it assumes that clusters are spherical

Goal: new algorithm that can learn location and shape of clusters

= Covariance



$x_1$  &  $x_2$  are positively correlated  $\Leftrightarrow$   
positive covariance

negative correlation  $\Leftrightarrow$   
negative covariance