

10/26/2023 : K-Means Clustering

Machine Learning

Supervised Learning

Training Dataset:

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

input to model \uparrow \uparrow desired output

Goal: Learn function mapping x to y

Unsupervised Learning

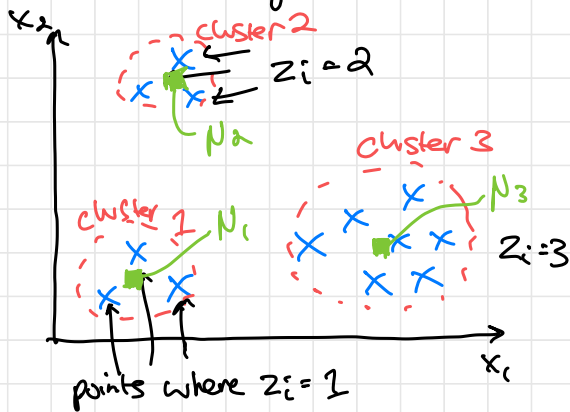
Training Dataset only contains x 's

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

No "correct output" per example
Goal: Learn what structure is present in the data

- 1) Sub-population / group / clusters
- 2) low-dimensional structure (subspace structure)
- 3) Similarity / relationship structure

Clustering



Input:

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

Given: # of clusters in data K

Output: An assignment z_1, z_2, \dots, z_n where $z_i \in \{1, \dots, K\}$ denotes cluster we assign to $x^{(i)}$

Today: K-Means Clustering Algorithm

Idea 1: Write down a loss function to define what is a good assignment z_1, \dots, z_n

Idea 2: Each cluster has a centroid N_j for $j=1, \dots, K$
Loss = how far is $x^{(i)}$ from centroid it was assigned to

Loss function (formally):

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

$= (z_1, \dots, z_n)$

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

Centroid for the assigned cluster of $x^{(i)}$

"Reconstruction Error" Loss
If we only knew
the cluster assignments $z_{1:n}$
and cluster means $N_{1:k}$

How well could we reconstruct
 $\{x^{(1)}, \dots, x^{(n)}\}$?

Goal: minimize $L(z_{1:n}, N_{1:k})$ with respect to $z_{1:n}, N_{1:k}$
Note: Can't do gradient descent because z_i 's are discrete,
Can't take a derivative w.r.t z_i

Strategy: Alternating Minimization

- ① Start with a random choice of N_1, \dots, N_k
Alternate until **convergence** ← when $z_{1:n}$'s and $N_{1:k}$'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 $z_{1:n}$

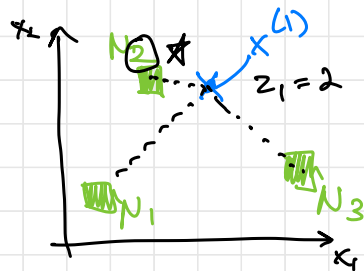
Step 1: Choose each N_j to be a different random $x^{(i)}$
from the dataset

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

For each $i = 1, \dots, n$:

$$\text{set } z_i = \underset{j=1, \dots, k}{\operatorname{argmin}} \|x^{(i)} - N_j\|^2$$

English: choose cluster whose
 N_j is closest



Step 3: Minimizing L w.r.t. $N_{1:k}$

Intuitively: N_j should be mean of all points where $z_i = j$

$$\text{minimize } \sum_{i=1}^n \|x^{(i)} - N_{z_i}\|^2$$

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

For particular value of j :

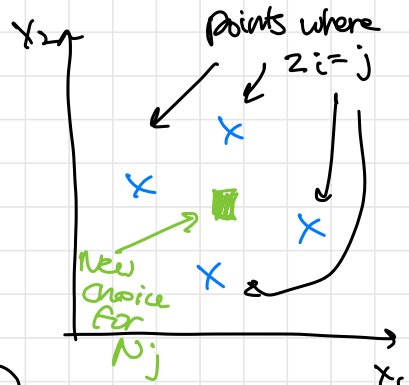
$$\text{minimize } \sum_{i: z_i=j} \|x^{(i)} - N_j\|^2$$

Take gradient w.r.t. N_j , set = 0

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

$$\sum_{i: z_i=j} x^{(i)} = |\{i: z_i=j\}| \cdot N_j$$

$$\Rightarrow N_j = \frac{1}{|\{i: z_i=j\}|} \sum_{i: z_i=j} x^{(i)}$$



Average of points in cluster j

Note: Eventually we will converge why? At every step L gets better

No guarantee of converging to **optimal solution**
Random initialization influences final clustering

↳ i.e. global minimum of L

How do you choose k ?

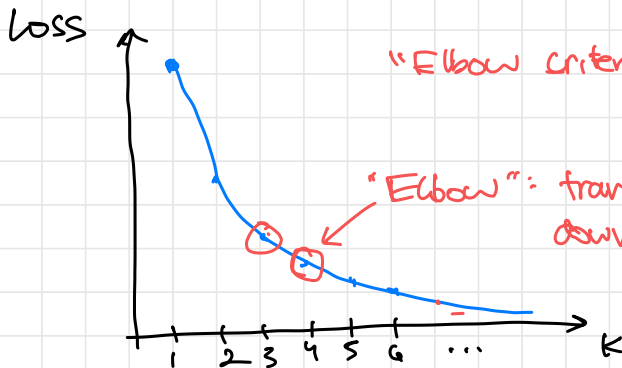
← this is a hyperparameter

Wrong answer: Fit on a dev set

Why? larger k always makes loss go down



$k=1$: loss large
 $k=2$: loss goes down
 $k=6$: loss even lower



"Elbow criterion"

"Elbow": transition between loss going down rapidly vs. loss going down slowly

Choose k at elbow

k -means looks for spherical clusters (because it uses Euclidean distance)

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

