

# 3/21/2022 : K-Means Clustering

## Machine Learning

### Supervised Learning

Training Dataset

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

Input to model  
Correct output

Goal: Learn a function from  $x \rightarrow y$

### Unsupervised Learning

Dataset only contains  $x$ 's

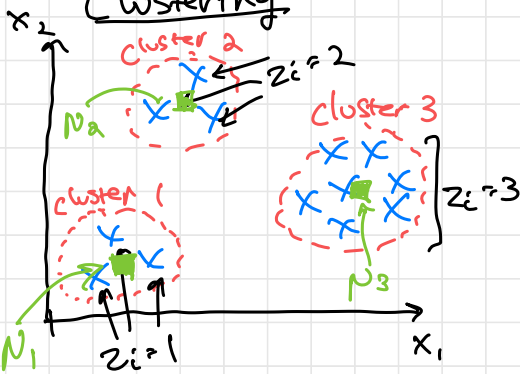
$$D = \{x^{(1)}, \dots, x^{(n)}\}$$

No input  $\rightarrow$  output mapping  
(learnable from dataset)

Goal: Learn about structure of dataset

- ① Clusters
- ② Sequential structure
- ③ Subspace structure (low-dimensional structure)
- ④ Similarity or relationships structure

## Clustering



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

Assume  $K$  clusters  $1, \dots, K$

Goal of clustering:

Produce an assignment  $z_1, \dots, z_n$  where  $z_i \in \{1, \dots, K\}$  and  $z_i$  denotes cluster assigned to  $x^{(i)}$

Need loss function that measures how bad an assignment is

Today: K-Means clustering

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

Loss function more formally:

$$L(\underbrace{Z_{1:n}}_{\text{assignments}}, \underbrace{N_{1:k}}_{\text{centroids}}) = \sum_{i=1}^n \|x^{(i)} - N_{z_i}\|^2$$

cluster ID assigned to Example  $i$   
centroid of cluster for example  $i$

"Reconstruction Error"

If we replaced each  $x^{(i)}$  with its cluster centroid, how wrong is that?

Can't directly do gradient descent -  $z_i$ 's are discrete

Strategy: Alternating minimization

① start with random choice of  $N_1, \dots, N_k$

Alternate until convergence:

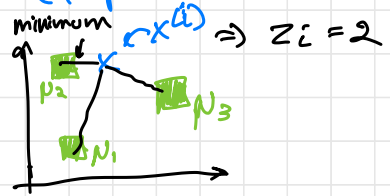
② Choose  $Z_{1:n}$  to minimize  $L$  given choice of  $N_{1:k}$

③ Choose  $N_{1:k}$  to minimize  $L$  given choice of  $Z_{1:n}$

Step ①: Choose each  $N_j$  to be a random example in dataset

Step ②: Minimizing w.r.t.  $Z_{1:n}$

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



Step ③: minimizing w.r.t.  $N_{1:k}$

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

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

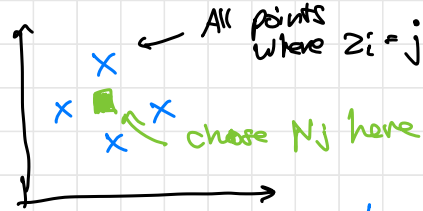
Now consider each  $j$  independently

For  $j=1$ :  $\nabla_{N_1} L(z_{1:n}, N_{1:k}) = \nabla_{N_1} \sum_{i: z_i=1} \|x^{(i)} - N_1\|^2$

$$= \sum_{i: z_i=1} 2(x^{(i)} - N_1) \cdot (-1) = 0$$

Average of all points in cluster 1

$$N_1 = \frac{1}{|\{i: z_i=1\}|} \sum_{i: z_i=1} x^{(i)}$$

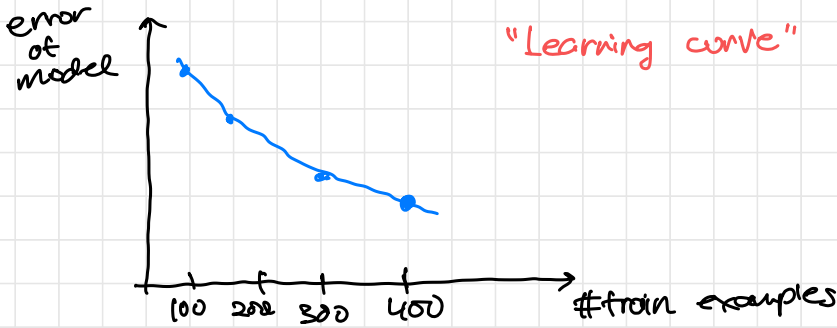


Note: Eventually we will converge  
At every step  $L$  decreases (or stays same)

This is not guaranteed to find global optimum

### Announcements

- Midterm graded
- Progress reports due Thurs
- HW 3 due 4/11
- "How much training data?"

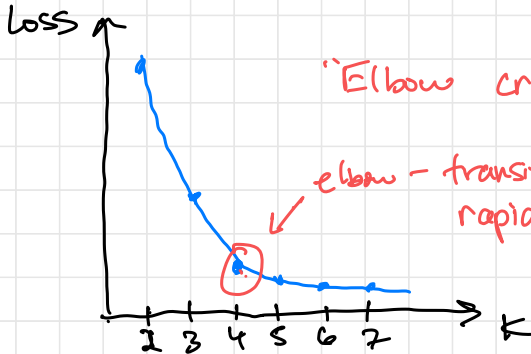


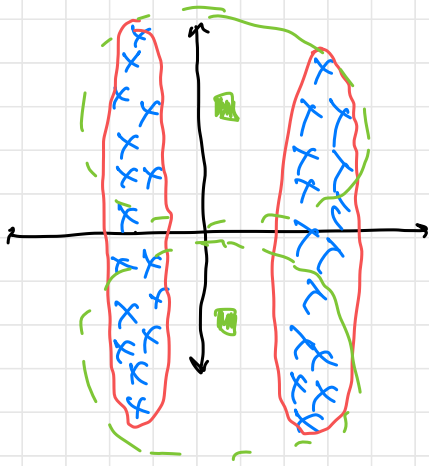
How do you choose  $K$ ?

Wrong Answer: Choose based on dev set



Larger  $k$  essentially always decreases loss

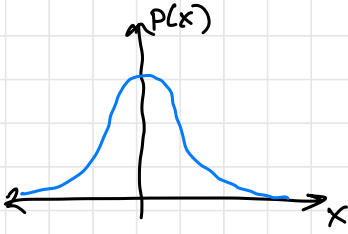




↳ Means is looking for spherical clusters  
(because it uses Euclidean distance)

Need new algorithm that can learn  
both location AND shape  
of clusters

Plan: Describe clusters as multivariate Gaussian distribution

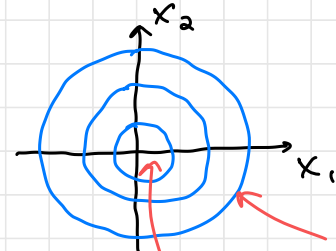


Standard univariate Gaussian

$$(N=0, \sigma^2=1)$$

↑  
mean

↑  
variance



Standard multivariate Gaussian

$$(N = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix})$$

↑  
mean

↑  
covariance matrix

Covariance Matrix

$$\Sigma = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_1, x_2) & \text{Var}(x_2) \end{pmatrix}$$

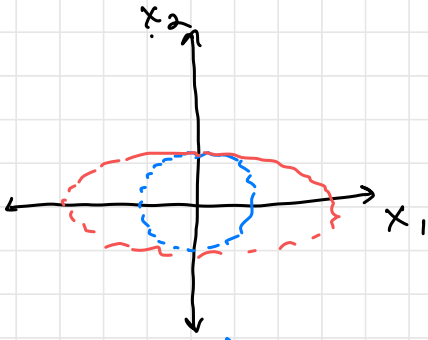
$$\text{Var}(x_1) = \mathbb{E}[(x_1 - \mathbb{E}[x_1])^2]$$

$$\text{Cov}(x_1, x_2) = \mathbb{E}[(x_1 - \mathbb{E}[x_1])(x_2 - \mathbb{E}[x_2])]$$

$$\text{Correlation}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)} \cdot \sqrt{\text{Var}(X_2)}}$$

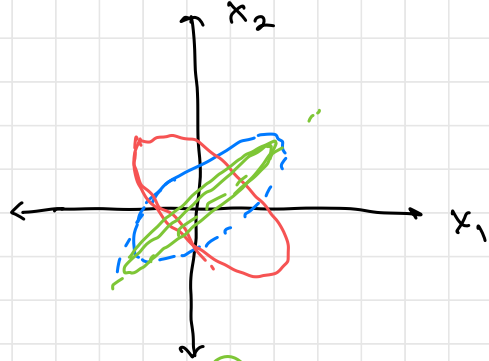
$\text{Cov} > 0 \Leftrightarrow$  positively correlated

$\text{Cov} < 0 \Leftrightarrow$  negatively correlated



$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 1 & -1/2 \\ -1/2 & 1 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}$$