

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Unsupervised Learning



Table of contents

- 1. K-means clustering**
- 2. Project Part- 2**

What is unsupervised learning?

- Given a training set of n unlabeled samples - $x(i)$

==

- What can we learn from the samples?

- Estimate the overall distribution of the data without knowing their label.



- Figure out the groupings of the samples (if any). — Clustering (K-means)

- Identify some features that may be more important than others.

Feature selection

Finding clusters

- How to represent the clusters?
 - Centroid (mean of all samples)
- Which cluster a sample should be assigned to (e.g., membership)?
 - Some similarity measure like Euclidean dist.
- What similarity measure to use?
 - Different measures
 - Euclidean dist.
 - Cosine similarity.

Clustering Objective function

- The sum-of-squared-error criterion/cost-

$$J_e = \sum_{i=1}^C \sum_{x \in D_i} \|x - m_i\|^2$$

(centroid of
ith cluster.)

The diagram shows two clusters of data points. The first cluster, labeled m_1 , has several points with arrows pointing towards its central centroid. The second cluster, labeled m_2 , also has a central centroid and surrounding data points. Arrows indicate the distance from each point to its respective centroid.

$$\underline{\underline{m_i = \frac{1}{n_i} \sum_{x \in D_i} x}}$$

Reduce the intra-cluster distance
increase the inter-cluster distance.

K-Means Clustering

Given: n samples, a number k .

No. of
clusters .

Begin

initialize $\underline{\mu_1, \mu_2, \dots, \mu_k}$ (randomly selected)

do classify n samples according to
nearest μ_i \Rightarrow Membership

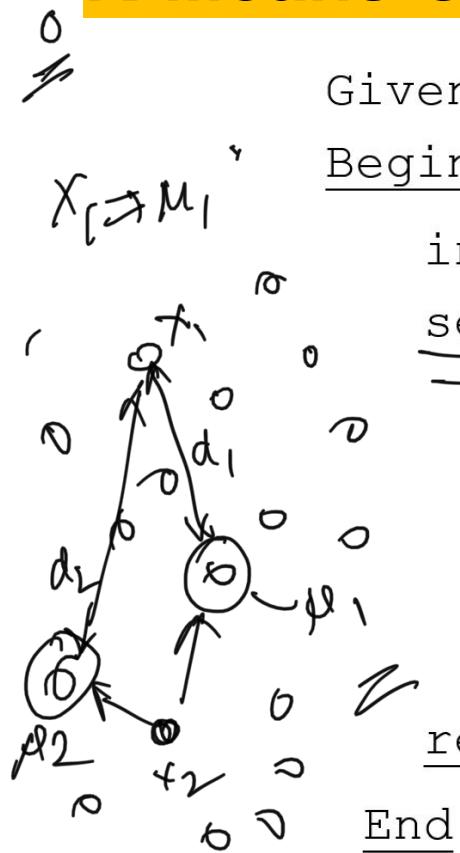
recompute μ_i

Assignment

until no change in μ_i

return $\mu_1, \mu_2, \dots, \mu_k$

End



K-Means Clustering - Some scenarios

- What happens if the distance of a point is same from more than one centroid?

Randomly assigned one of the clusters.

- Hand assignment

- What happens in case of outliers?

Sensitive to outliers.

K-Means Clustering - Some scenarios

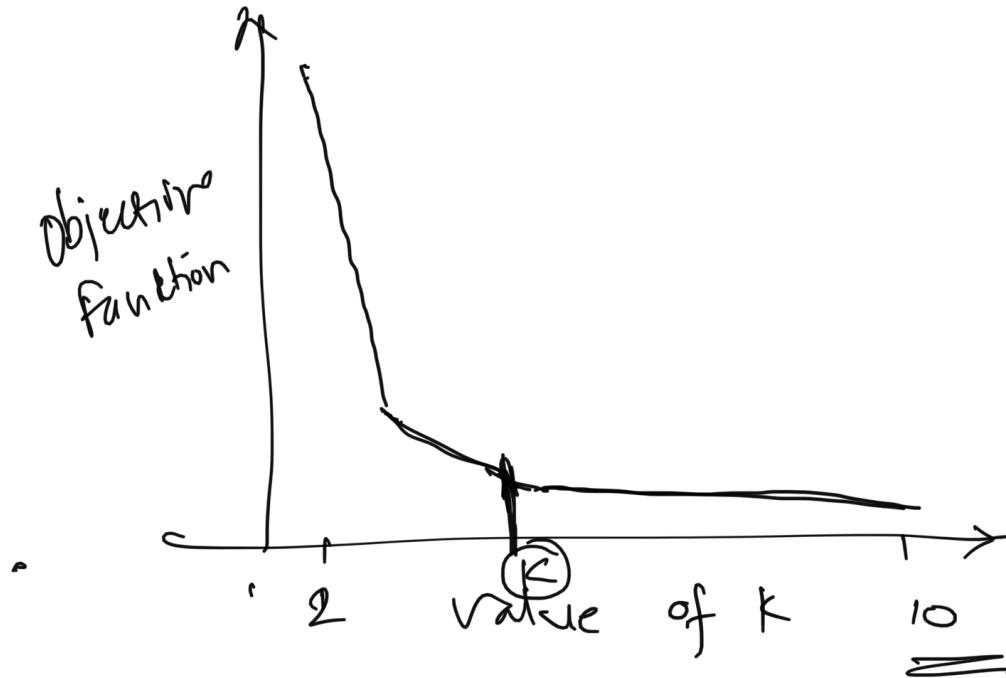
- Different initializations?
 - Might lead to sub-optimal results.
- Different number of clusters than k?

How to choose the k?

- Elbow method.

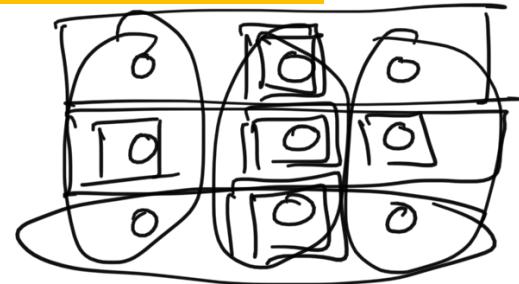
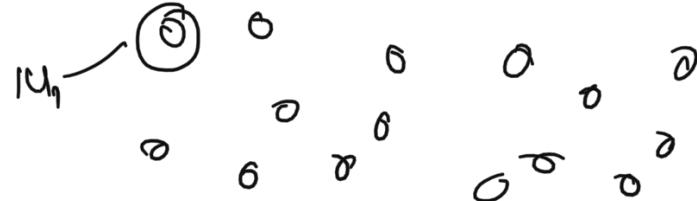
$$\begin{array}{l} 10,000 \quad \underline{k=1} \\ 1-9999 \quad K=n \end{array}$$

$$1 \leq K < n$$



How to deal with sensitivity to initialization problem?

- ① Run the algo for few epochs
= & take the average result
- ② first centroid randomly. → Sensitive to outliers.
- ③ Kmeans++ \Rightarrow Smart Initialization
=



K-Means Clustering - Pros and Cons

Pros

- ① Easy to implement
- ② fast.

Cons

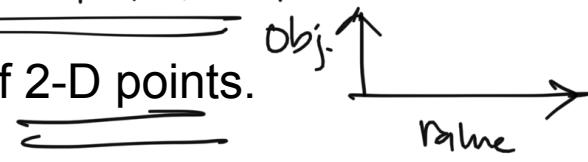
- ① Sensitive to initialization
- ② " " outliers.
- ③ Hard Assignment of data points -
- ④ Choose no. of clusters manually.

Project Part 2 - Unsupervised Learning (K-means)

1st time - $k=2 \Rightarrow$ Obj.
 $\overbrace{\quad\quad\quad}^{k=3}$
 \Leftrightarrow

$k=2, 3, 4, 5, 6, 7, 8, 9, 10$ $k=10$

- Apply k-means algorithm the given dataset of 2-D points.



- Initialization strategies:

$\nearrow 2 \text{ times}$ - $\overbrace{k=2-10}$

- Strategy 1: randomly pick the initial centers from the given samples.
- Strategy 2: pick the first center randomly; for the i -th center ($i > 1$), choose a sample
(among all possible samples) such that the average distance of this chosen one to all previous ($i-1$) centers is maximal.

$\nearrow 1$
 $\nearrow 2 \text{ times}$ - $\overbrace{k=2-10}$

Total - 4 plots

Project Part 2 - Unsupervised Learning (K-means)

- Objective function-

$$\sum_{i=1}^k \sum_{x \in D_i} \|x - \mu_i\|^2 \rightarrow SSE$$

- Number of clusters (k) - from 2-10
- Plot the objective function value vs. the number of clusters k.

Project Part 2 - Deliverables

— Due on 26 March
(Arizona time) 11:59 PM

- Well-commented code
 - Python / MATLAB / Jupyter notebook (.ipynt)
- A report that summarizes the results and includes all the plots.
 - .PDF format .(typed, not handwritten)
- Please submit the code and the report as separate files on Canvas. **Do**

not zip them.

· py
· py
· pdf

Questions?