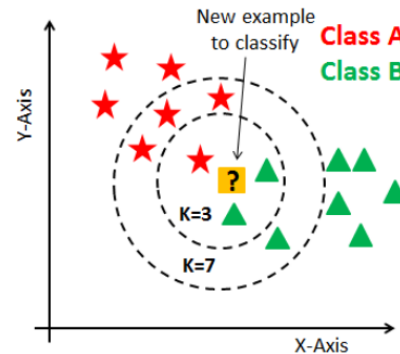# National University of Computer and Emerging Sciences CL
## 461 Artificial Intelligence

*Lab Manual* 11

# KNN & Kmeans

## K Nearest Neighbour

- The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification and regression problems.
- The KNN algorithm assumes that similar things exist in close proximity. In other words, similar things are near to each other.

## KNN Algorithm

1. Load the data
2. Initialize K to your chosen number of neighbors
3. For each example in the data
   a. Calculate the distance between the query example and the current example from the data.
   b. Add the distance and the index of the example to an ordered collection
4. Sort the ordered collection of distances and indices from smallest to largest (in ascending order) by the distances
5. Pick the first K entries from the sorted collection
6. Get the labels of the selected K entries
7. If regression, return the mean of the K labels
8. If classification, return the mode of the K labels

## Distance

$$\text{Euclidean} \qquad \sqrt{\sum_{i=1}^{k}(x_i - y_i)^2}$$

$$\text{Manhattan} \qquad \sum_{i=1}^{k}|x_i - y_i|$$

$$\text{Minkowski} \qquad \left(\sum_{i=1}^{k}(|x_i - y_i|)^q\right)^{1/q}$$

## Example:

$$D = Sqrt[(48-33)^2 + (142000-150000)^2] = 8000.01 \gg Default=Y$$

| Age | Loan | Default | Distance | |
|-----|------|---------|----------|---|
| 25 | $40,000 | N | 102000 | |
| 35 | $60,000 | N | 82000 | |
| 45 | $80,000 | N | 62000 | |
| 20 | $20,000 | N | 122000 | |
| 35 | $120,000 | N | 22000 | 2 |
| 52 | $18,000 | N | 124000 | |
| 23 | $95,000 | Y | 47000 | |
| 40 | $62,000 | Y | 80000 | |
| 60 | $100,000 | Y | 42000 | 3 |
| 48 | $220,000 | Y | 78000 | |
| 33 | $150,000 | Y | 8000 | 1 |
| | | | | |
| **48** | **$142,000** | **?** | | |

Euclidean Distance

$$D = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

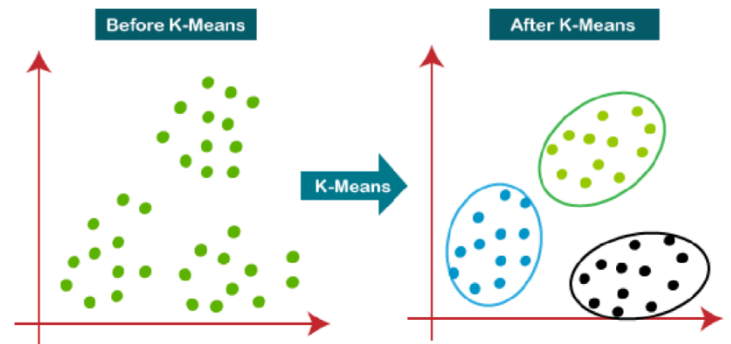## KNN for Text Data

We will use this link to explain in detail.
https://towardsdatascience.com/text-classification-using-k-nearest-neighbors-46fa8a77acc5

# Lab Task:

1. Implement KNN from scratch for the mall_customer dataset (Genre, Age, Annual Income, Spending Score).
   a. As a classification problem where output column is Genre
   b. As a regression problem where output column is Spending Score

# Kmeans

K-Means clustering is a type of unsupervised learning. The main goal of this algorithm to find groups in data and the number of groups is represented by K. It is an iterative procedure where each data point is assigned to one of the K groups based on feature similarity.



# Kmeans Algorithm

1. Specify number of clusters $K$.
2. Initialize centroids by first shuffling the dataset and then randomly selecting $K$ data points for the centroids without replacement.
3. Keep iterating until there is no change to the centroids. i.e assignment of data points to clusters isn't changing.

- Compute the sum of the squared distance between data points and all centroids.
- Assign each data point to the closest cluster (centroid).
- Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

# Example:

Manhattan distance is used to cluster these points into 3 clusters.
Initial centroids are: A1(2, 10), A4(5, 8) and A7(1, 2).

## After 1st Iteration

| Given Points | Distance from center (2, 10) of Cluster-01 | Distance from center (5, 8) of Cluster-02 | Distance from center (1, 2) of Cluster-03 | Point belongs to Cluster |
|---|---|---|---|---|
| A1(2, 10) | 0 | 5 | 9 | C1 |
| A2(2, 5) | 5 | 6 | 4 | C3 |
| A3(8, 4) | 12 | 7 | 9 | C2 |
| A4(5, 8) | 5 | 0 | 10 | C2 |
| A5(7, 5) | 10 | 5 | 9 | C2 |
| A6(6, 4) | 10 | 5 | 7 | C2 |
| A7(1, 2) | 9 | 10 | 0 | C3 |
| A8(4, 9) | 3 | 2 | 10 | C2 |

**After 2nd Iteration**

| Given Points | Distance from center (2, 10) of Cluster-01 | Distance from center (6, 6) of Cluster-02 | Distance from center (1.5, 3.5) of Cluster-03 | Point belongs to Cluster |
|---|---|---|---|---|
| A1(2, 10) | 0 | 8 | 7 | C1 |
| A2(2, 5) | 5 | 5 | 2 | C3 |
| A3(8, 4) | 12 | 4 | 7 | C2 |
| A4(5, 8) | 5 | 3 | 8 | C2 |
| A5(7, 5) | 10 | 2 | 7 | C2 |
| A6(6, 4) | 10 | 2 | 5 | C2 |
| A7(1, 2) | 9 | 9 | 2 | C3 |
| A8(4, 9) | 3 | 5 | 8 | C1 |

# Kmeans for Text Data

We will use this link to explain in detail.
https://towardsdatascience.com/a-friendly-introduction-to-text-clustering-fa996bcefd04

# Lab Task:

2. Implement Kmeans from scratch for the mall_customer dataset (Genre, Age, Annual Income, Spending Score).