

# **SNS COLLEGE OF TECHNOLOGY**

**Coimbatore-35 An Autonomous Institution** 

Accredited by NBA – AICTE and Accredited by NAAC – UGC with 'A++' Grade Approved by AICTE, New Delhi & Affiliated to Anna University, Chennai

# **DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING**

### **23AMB201 - MACHINE LEARNING**

**II YEAR IV SEM** 

## **UNIT IV – UNSUPERVISED LEARNING ALGORITHM**

**TOPIC 22 – Clustering – K-Means** 

K-Means/ML/S.Rajarajeswari/AP/AIML

Redesigning Common Mind & Business Towards Excellence







Build an Entrepreneurial Mindset Through Our Design Thinking FrameWork



## Clustering



K-Means/ML/S.Rajarajeswari/AP/AIML





• **Clustering** is the <u>classification</u> of objects into different groups, or more precisely, the <u>partitioning</u> of a <u>data set</u> into <u>subsets</u> (clusters), so that the data in each subset (ideally) share some common trait - often according to some defined distance measure.

• Applications:

- Market Segmentation 1.
- 2. Statistical data analysis
- 3. Social network analysis
- Image segmentation 4.
- 5. Amazon and Netflix





- **1. Hierarchical algorithms**: Find successive clusters **1.Agglomerative** ("bottom-up"): Begins with each element as a separate cluster and merge them into successively larger clusters. **2.Divisive** ("top-down"): Begins with the whole set and proceed to divide it into successively smaller clusters.
- **2. Partitional clustering:** Partitional algorithms determine all clusters at once. They include:

### K-means and derivatives

Fuzzy *c*-means clustering

- **3. Density based clustering**
- 4. Fuzzy clustering





# **Common Distance measures**

- *Distance measure* will determine how the *similarity* of two elements is calculated and it will influence the shape of the clusters.
  - They include:
- 1. The <u>Euclidean distance</u> (also called 2-norm distance) is given by:

2. The <u>Manhattan distance</u> (also called taxicab norm or 1-norm) is given by:



$$d(\mathbf{p},\mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i-p_i)^2}$$

$$\mathbf{d}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{m} |\mathbf{x}_i - \mathbf{y}_i|$$



 $d(x, y) = \max_{1 \le i \le v} |x_i - y_i|$ 3. The <u>maximum norm</u> is given by:

4. <u>Hamming distance</u> (sometimes edit distance) measures the minimum number of substitutions required to change one member into another.

$$D_H = \sum_{i=1}^k \left| x_i \right| -$$











- K-Means Clustering is an Unsupervised Machine Learning algorithm which groups the unlabeled dataset into different clusters.
- The **k-means algorithm** is an algorithm to <u>cluster n</u> objects based on attributes into k <u>partitions</u>, where k < n.
- K-means clustering is a technique used to organize data into groups based on their similarity. • For example online store uses K-Means to group customers based on purchase frequency and
- spending creating segments like:
  - Budget Shoppers
  - Frequent Buyers
  - Big Spenders for personalised marketing.









# **K-MEANS CLUSTERING-work flow**



K-Means/ML/S.Rajarajeswari/AP/AIML



- Form K clusters by assigning all points to the closest centroid. Recompute the centroid of each cluster.





Dataset

# **How to Apply K-Means Clustering Algorithm?**



K-Means/ML/S.Rajarajeswari/AP/AIML





### Recompute the centroids of newly formed clusters







- 1. Centroids of newly formed clusters do not change
- 2. Points remain in the same cluster
- 3. Maximum number of iterations is reached

K-Means/ML/S.Rajarajeswari/AP/AIML







Data set {2, 4, 10, 12, 3, 20, 30, 11, 25}

Iteration 1

M1, M2 are the two randomly selected centroids/means where

M1=4, M2=11

and the initial clusters are

 $C1=\{4\}, C2=\{11\}$ 

Calculate the Euclidean distance as

$$d(\mathbf{p},\mathbf{q})=\sqrt{\sum_{i=1}^n(q_i-p_i)^2}$$

K-Means/ML/S.Rajarajeswari/AP/AIML



| Datapoint | D1 | D2 | Cluster |
|-----------|----|----|---------|
| 2         | 2  | 9  | C1      |
| 4         | 0  | 7  | C1      |
| 10        | 6  | 1  | C2      |
| 12        | 8  | 1  | C2      |
| 3         | 1  | 8  | C1      |
| 20        | 16 | 9  | C2      |
| 30        | 26 | 19 | C2      |
| 11        | 7  | 0  | C2      |
| 25        | 21 | 14 | C2      |

541 Septem (W)

Therefore

C1= {2, 4, 3}

C2= {10, 12, 20, 30, 11, 25}





Data set {2, 4, 10, 12, 3, 20, 30, 11, 25}

$$d(\mathbf{p},\mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i-p_i)^2}$$

| Datapoint | D1 | D2 | Cluster |
|-----------|----|----|---------|
| 2         | 2  | 9  | C1      |
| 4         | 0  | 7  | C1      |
| 10        | 6  | 1  | C2      |
| 12        | 8  | 1  | C2      |
| 3         | 1  | 8  | C1      |
| 20        | 16 | 9  | C2      |
| 30        | 26 | 19 | C2      |
| 11        | 7  | 0  | C2      |
| 25        | 21 | 14 | C2      |

Iteration 1

|     | _ |        |
|-----|---|--------|
| New | C | luster |

Therefore

C1= {2, 4, 3}

C2= {10, 12, 20, 30, 11, 25}

| Datapoint | D1 | D2 | Cluster |
|-----------|----|----|---------|
| 2         | 1  | 16 | C1      |
| 4         | 1  | 14 | C1      |
| 3         | 0  | 15 | C1      |
| 10        | 7  | 8  | C1      |
| 12        | 9  | 6  | C2      |
| 20        | 17 | 2  | C2      |
| 30        | 27 | 12 | C2      |
| 11        | 8  | 7  | C2      |
| 25        | 22 | 7  | C2      |

'S

M1 = (2+3+4)/3 = 3

Iteration 2

M2= (10+12+20+30+11+25)/6= 18

New Clusters C1= {2, 3, 4, 10} C2= {12, 20, 30, 11, 25}



Data set {2, 4, 10, 12, 3, 20, 30, 11, 25}

$$d(\mathbf{p},\mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i-p_i)^2}$$

| Datapoint | D1 | D2 | Cluster |
|-----------|----|----|---------|
| 2         | 1  | 16 | C1      |
| 4         | 1  | 14 | C1      |
| 3         | 0  | 15 | C1      |
| 10        | 7  | 8  | C1      |
| 12        | 9  | 6  | C2      |
| 20        | 17 | 2  | C2      |
| 30        | 27 | 12 | C2      |
| 11        | 8  | 7  | C2      |
| 25        | 22 | 7  | C2      |

Iteration 2

New Clusters C1= {2, 3, 4, 10} C2= {12, 20, 30, 11, 25}



**New Clusters** 

M1 = (2+3+4+10)/4 = 4.75M2= (12+20+30+11+25)/5= 19.6

K-Means/ML/S.Rajarajeswari/AP/AIML



| Datapoint | D1    | D2   | Cluster |
|-----------|-------|------|---------|
| 2         | 2.75  | 17.6 | C1      |
| 4         | 0.75  | 15.6 | C1      |
| 3         | 1.75  | 16.6 | C1      |
| 10        | 5.25  | 9.6  | C1      |
| 12        | 7.25  | 7.6  | C1      |
| 20        | 15.25 | 0.4  | C2.     |
| 30        | 25.25 | 10.4 | C2      |
| 11        | 6.25  | 8.6  | C1      |
| 25        | 20.25 | 5.4  | C2      |

Iteration 3

New Clusters C1= {2, 3, 4, 10, 12, 11} C2= {20, 30, 25}



Data set {2, 4, 10, 12, 3, 20, 30, 11, 25}

$$d(\mathbf{p},\mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i-p_i)^2}$$

| Datapoint | D1    | D2   | Cluster |
|-----------|-------|------|---------|
| 2         | 2.75  | 17.6 | C1      |
| 4         | 0.75  | 15.6 | C1      |
| 3         | 1.75  | 16.6 | C1      |
| 10        | 5.25  | 9.6  | C1      |
| 12        | 7.25  | 7.6  | C1      |
| 20        | 15.25 | 0.4  | C2      |
| 30        | 25.25 | 10.4 | C2      |
| 11        | 6.25  | 8.6  | C1      |
| 25        | 20.25 | 5.4  | C2      |

Iteration 3

New Clusters C1= {2, 3, 4, 10, 12, 11}  $C2=\{20, 30, 25\}$ 

| Datapoint | D1 | D2 | Cluster |
|-----------|----|----|---------|
| 2         | 5  | 23 | C1      |
| 4         | 3  | 21 | C1      |
| 3         | 4  | 22 | C1      |
| 10        | 3  | 15 | C1      |
| 12        | 5  | 13 | C1      |
| 11        | 4  | 14 | C1      |
| 20        | 13 | 5  | C2      |
| 30        | 23 | 5  | C2      |
| 25        | 18 | 0  | C2      |

**New Clusters** 

M1=(2+3+4+10+12+11)/6=7

M2 = (20+30+25)/3 = 25

K-Means/ML/S.Rajarajeswari/AP/AIML



Iteration 4

New Clusters C1= {2, 3, 4, 10, 12, 11} C2= {20, 30, 25}

No Change between Iteration 3 and 4





- Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin, —Learning from Data, AML Book Publishers, 2012.
- P. Flach, —Machine Learning: The art and science of algorithms that make sense of data<sup>I</sup>, Cambridge University Press, 2012.
- W3school.com
- https://discourse.opengenus.org/t/using-id3-algorithm-to-builda-decision-tree-to-predict-the-weather/3343



15/15

