Note ⇒ Farid Sir Coming Back
From America.

Page-72

## ⊞ Clustering:

Approach of grouping data.

⊞ In Clustering we try to group the number of rows.

⊞ N সংখ্যক instance কে k সংখ্যক cluster এ ভাগ করা হয়।

⊞ আমার কাছে যদি 30 টি instances থাকে তবলে 30 টা cluster তৈরি হবে, 30 এর বেশি cluster তৈরি হবে, সর্বনিম্ন ২ টা cluster তৈরি হবে, আর কোন cluster ফাঁকা থাকবে না।

⊞ Clustering হল Un-supervised learning.

$$X = \{ x_1, x_2, x_3, \ldots, x_N \}$$

এখানে X হল unlabeled data এর সেট।

$c_i \neq \emptyset$ , $i = 1, \ldots, k_i$ ॥ মানে cluster সংখ্যা শূন্য হবে না।

$\bigcup_{i=1}^{k} c_i = X$ ॥ মানে এই সবগুলো cluster কে merge করি।

$c_i \cap c_j = \emptyset$ , $i \neq j$ , $i, j = 1, \ldots, k_i$ ॥ i এবং j হল disimilar. অর্থাৎ ২ টা instance একটা cluster এ থাকবে।

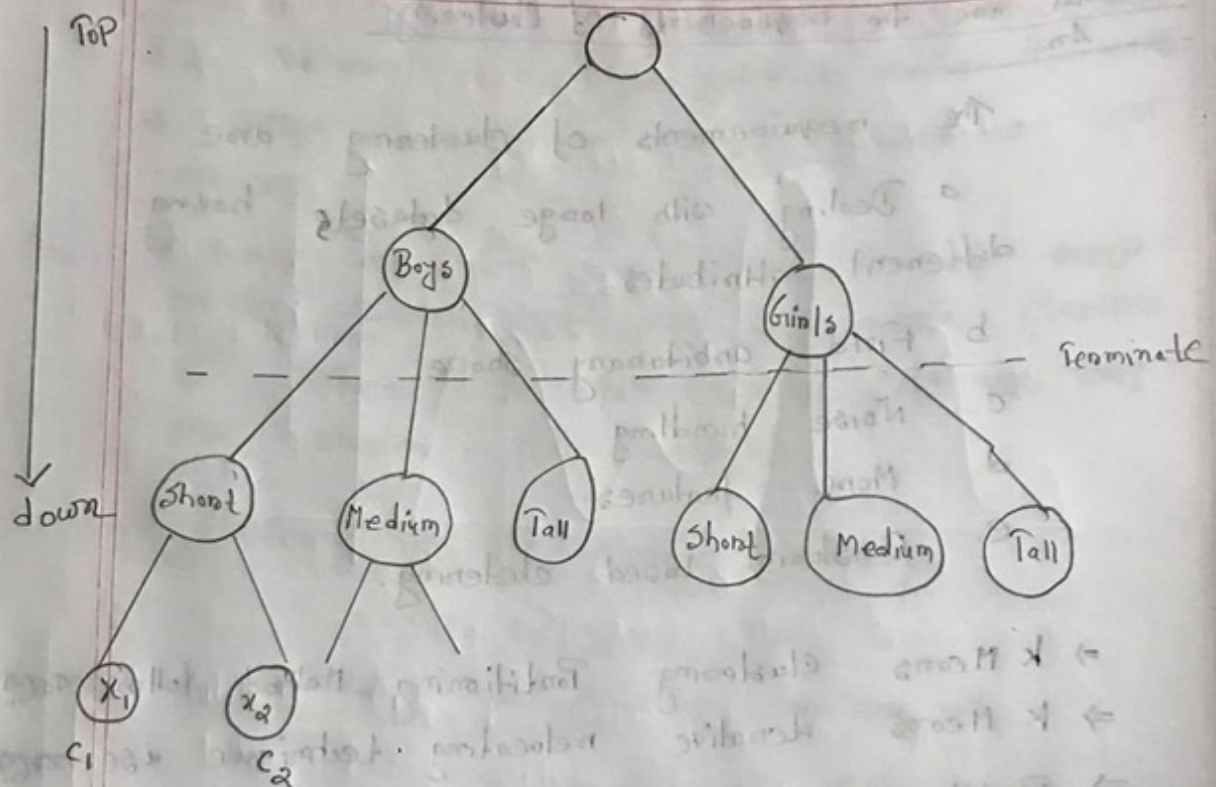☘ What are the requirements of Clustering

Ans:

The requirements of clustering are

a. Dealing with large datasets having different attributes.

b. Find arbitrary shape

c. Noise handling

d. More features.

e. Constraint based clustering.

⇒ k Means Clustering Partitioning Method follow করে,

⇒ k Means iterative relocating technique use করে,

⇒ Partition based clustering এর বিগ সমস্যা কম এই বড় ডাইমেনশন এর ডাটা handle করতে পারে না,

Hierarchical Clustering:

এটা ২ Approach এ কাজ করে,

i. Top-down

ii. Bottom-up

TOP

down →

Boys

Girls

Short   Medium   Tall   Short   Medium   Tall

$X_1$   $X_2$

$C_1$   $C_2$

— — — — — — — — Terminate

ক) Bottom Up Approach কে Similarity Approach বলা হয়। Bottom UP এ Similar cluster গুলোকে Merge করা হয়ে শেষে গিয়ে ১টি Cluster এ পরিণত হয়।

খ) Hierarchical Clustering কে যেকোন ১টি level পর্যন্ত terminate করানো যায়।

গ) Density Based Clustering: কোন একটি Space G dataPoints কোথায় কিরকম আছে সেটা density based clustering দিয়ে measure করা হয়।

প্র What is Similarity Measure? Write the similarity rules.

Ans:

Based on Column value; similarity is measured.

The rules are

a. Each cluster must contain at least one instance.

b. Each instance must belong to exactly one cluster.

প্র Normal k Means randomly center Pick করবে, এর ফলে Global Optima তে algorithm stuck হয়ে -যায়।

প্র k Means ++ algorithm একটা ম্যাথমেটিক্স follow করে Center Pick করে, Min, Max, Average value গুলন Center হিসেবে Pick করে।
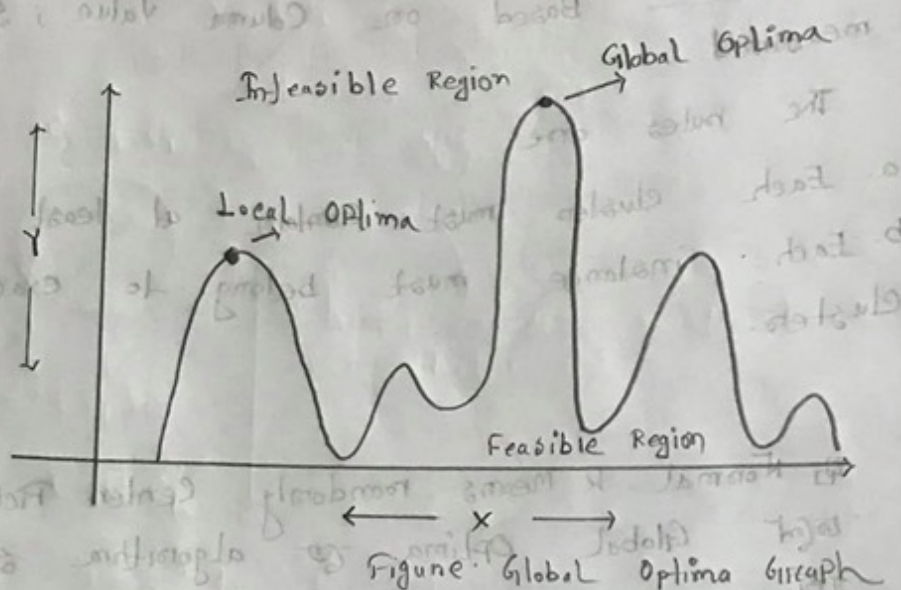
প্র Define Iterative Relocation

Ans:

It's a technique to improve the quality of the initial cluster. It implies that cluster membership is changed to find the local optima.

## ⟶ Define Global Optima Problem Graph

Ans:



Figure: Global Optima Graph

Global Optima represents a very best Solution. Local optima is better than it's immediate neighbors.

## ⟶ Define Multimodal

Ans. Local optima is also known as multimodal

## ⟶ k-means clustering numerical সংখ্যা নিয়ে কাজ করে।

## ⟶ Similarity based clustering categorical value নিয়ে কাজ করে।

## ⟶ Sub-space clustering; ensemble clustering এর একটা Part.

৫. k- Means Clustering এ Cluster কয়টা হবে, সেই অগ্র্যাডা বলে দিতে হয়, Similarity based Clustering এ Cluster এর সংখ্যা বলা লাগে না, Similarity based Clustering automatically Cluster তৈরী করে ফেলে।

৬। Similarity Based Clustering threshold-value use করে। Threshold-value দিয়ে আমরা কতগুলো feature -এর উপর base করে Clustering করতে চাই। অর্থাৎ আমার কাছে 4 টা feature / Column আছেন আমি 1/4 ; 2/4 ; 3/4 ; 4/4 threshold আকারে Clustering করতে পারি। 4/4 threshold হল tight - Coupling situation.

**A. Write Similarity Based Clustering Algorithm**

Ans:

Input: $X = \{x_1, x_2, x_3, \ldots, x_N\}$ // A set of Unlabelled instances

Output: A set of clusters, $C = \{C_1, C_2, \ldots, C_k\}$

Method:

1. $C = \phi$;

2. $k = 1$

3. $C_k = \{x_1\}$;   Dataset -এর প্রথম row/instance টা — দিয়ে cluster তৈরি করে

4. $C = C \cup C_k$

5. for $i = 2$ to $N$ do

6. for $l = 1$ to $k$ do

7. find the $l$th cluster center $x_l \in C_l$ to maximize the similarity measure $sim(x_i, x_l)$;

8. end for

9. if $sim(x_i, x_l) \geq$ threshold_value then

10. $C_l = C_l \cup x_i$

11. else

12. $k = k + 1$

13. $C_k = \{x_i\}$;

14. $C = C \cup C_{ki}$

15. end if

16. end for

Answer

Let

$t = 75\% = \frac{75}{100} = 0.75$ // Which means at least 3 features should be similar.
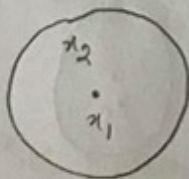
Step 1:

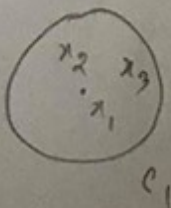$C_1$ is the cluster $x_1$ would be the center.

$x_1 \rightarrow C_1$



$C_1$

Step 2:

$Sim(x_1, x_2)$ // $Sim(x_1, x_2) \geq t$ : true

Center of cluster $(C_1)$ → Instance



$C_1$

Step 3:

$Sim(x_1, x_3)$ // $Sim(x_1, x_3) \geq$ : true



$C_1$

Step 4:

$$Sim(x_1, x_4) \text{ // } Sim(x_1, x_4) \geq t : false$$

So $x_4$ would become Center of cluster 2

$$x_4 \to C_2$$



$C_1$

$C_2$

Step 5:

$$Sim(x_1, x_5) \text{ // } Sim(x_1, x_5) \geq t : false$$
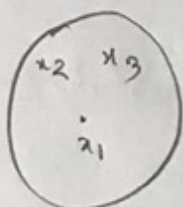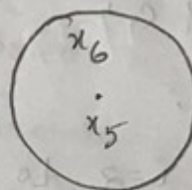$$Sim(x_4, x_5) \text{ // } Sim(x_4, x_5) \geq t : false$$

So, $x_5$ would become Center of Cluster 3

$$x_5 \to C_3$$



$C_1$

$C_2$

$C_3$

Step 6:

$sim(x_1, x_6) \parallel sim(x_1, x_6) \geq t$ : false

$sim(x_4, x_6) \parallel sim(x_4, x_6) \geq t$ : false

$sim(x_5, x_6) \parallel sim(x_5, x_6) \geq t$ : true



$C_1$        $C_2$        $C_3$

Step 7:

$sim(x_1, x_7) \parallel sim(x_1, x_7) \geq t$ : false

$sim(x_4, x_7) \parallel sim(x_4, x_7) \geq t$ : false

$sim(x_5, x_7) \parallel sim(x_5, x_7) \geq t$ : false

$$x_7 \rightarrow C_4$$

So, $x_7$ would become clu· center of cluster 4.



$C_1$        $C_2$        $C_3$        $C_4$

**Q.** Write Nearest Neighbor Clustering Algorithm

**Ans**

Input:
$$D = \{x_1, x_2, x_3, \ldots, x_n\} \quad // \text{A set of}$$
instances.

A // Adjacency matrix showing distance between instances

Output:
A set of C clusters.

Method.

1. $C_1 = \{x_1\};$

2. $C = \{C_1\};$

3. $k = 1;$

4. for $i = 2$ to $n$ do

5. find $x_m$ in some cluster $C_m$ in C so that dis $(x_i, x_m)$ is the smallest;

6. if dis $(x_i, x_m) \leq t$, threshold_value then

7. $C_m = C_m \cup X;$

8. else

9. $k = k + 1;$

10. $C_k = \{x_i\};$

11. $C = C \cup C_k;$

12. end if

13. end for

☘ Similarity Based Clustering ∪ Nested for loop ব্যবহার করা হয়,

☘ Nearest Neighbor Clustering ও Single for loop ব্যবহার করা হয়।

☘ Write the advantage of Similarity Based Clustering (sem)

__Ans:__

❀ The instances in (sem) can self-organise local optimal clusters & volumes without using cluster validity functions.

☘ Look at the given table.

| Item | A | B | C | D | E |
|------|---|---|---|---|---|
| A | 0 | 1 | 2 | 2 | 3 |
| B |   | 0 | 2 | 4 | 3 |
| C |   |   | 0 | 1 | 5 |
| D |   |   |   | 0 | 3 |
| E |   |   |   |   | 0 |

Given 5 items with the distance between them. Cluster them using Nearest Neighbor algorithm with a threshold t = 2

Solve :    Given

threshold - value ; $t = 2$

Step 1 :

A : $C_1 = \{A\}$

Step 2 :

B : dis $(B, A) = 1 \leq t$ ; So $C_1 = \{A, B\}$

Step 3 :

C : dis $(C, A) = 2 \leq t$

C : dis $(C, B) = 2 \leq t$

So,        $C_1 = \{A, B, C\}$

Step 4 :

D :    dis $(D, A) = 2 \leq t$

D :    dis $(D, B) = 4 \leq t$    it's false

D :    dis $(D, C) = 1 \leq t$

So,    $C_1 = \{A, B, C, D\}$

Step 5:

E: $dis(E,A) = 3 \le t$ it's false

E: $dis(E,B) = 3 \le t$ it's false

E: $dis(E,C) = 5 \le t$ it's false

E: $dis(E,D) = 3 \le t$ it's false

So, new cluster would be formed.

Therefore

$$C_2 = \{E\}$$

Ay:-

---

**Write the advantage of Nearest Neighbor Clustering**

Ans:

No need to know the number of Clusters to discover beforehand. It's different than in k-means & hierarchical.

Nearest Neighbor clustering -এর Complexity depend করে dataset -এর Volume on -অনুসারে instances -আর যায় দেয়,

Nearest Neighbor Clustering -এর Time Complexity $\theta(N^3)$ $O(n^2)$

Space Complexity $\theta(N^2)$ $O(n^2)$

**Write the differences between Classification & Clustering**

Ans

| Classification | Clustering |
|---|---|
| a. known number of classes | a. Unknown number of classes |
| b. Based on training set. | b. No prior knowledge. |
| c. Used to classify future observations. | c. Used to explore data |
| d. Supervised Type | d. Unsupervised Type. |