# Deg-7    Cluster Analysis

It is a collection of data objects.. simmilar within the same group. dissimmilar objects in other groups.
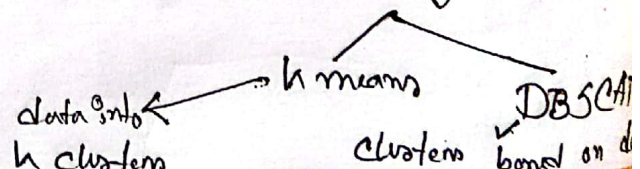
## Cluster analysis (Clustering)

The process of identifying these similarities and grouping similar data objects into clusters based on thier characterstics.

# Clustering is a form of unsupervised learning because there are no predefined classes for the data points.

## Typical Application

Insight: understanding data distribution.

Preprocessing: preparing data for other Algorithm.

data into ⟵ k means        DBSCAN
k clusters              Clusters based on d

# Clustering as preprocessing Tool.

**Summarization** : Groups similar data ~~togth~~ together, making it easier to work with large datasets.

- helps regression, PCA, Classification by revealing patterns within groups.

**Compression** : Image processing $\rightarrow$ vector Quantization.
to reduce img size

**K-nearest neighbors** : ~~Its helps~~ clusters helps quickly finding. K-nearest neighbors by focusing nearby cluster.

Clustering can help identify outliers.

## #Good clustering

A good clustering method Creates clusters with high intra-clan similarity and low inter similarity. The qual of clustering depends on the similarity measure and the implementation of Algorithm and its ability to discover hidden pattern in the data revealing relationships.

# Measure Quality of Clustering

Quality of Clustering accessed using dissimilarity and similarity metric.

↳ measured through distance function.
denoted $(d.J)$

* Distance function depends on the

Interval-Scaled → temp, Boolean → yes/no, Categorical → Colors, ordinal data → ranking;

* Weights depends on Application

III The Quality of clustering is typically evaluated using quality function that measure how well clustering achieves the desired grouping.

Euclidean distance $= d(g_1, g_2) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$

Manhattan distance $= d(g_1, g_2) = \sum_{i=1}^{n} |(x_i - y_i)|$

Minkowski " $= d(g_1, g_2) = \sqrt[m]{\sum_{i=1}^{n} (x_i - y_i)^m}$

Partitioning Approach: Creating various partitions of the Data and evaluating them using Criterion

methods → K means, K - medoids.

Hierarchical approach: Its creating tree like structure by either merging Clusters or dividing them.

methods - Dinna, Agnus

Density based on connectivity function.

methods - DBSACN

Grid based on multiple lvl granularity → Em, Som

model hypothesizes fits it to the data.

Link based on use relationship or links between objects → Link clas.

# Partitioning Algo

$$E = \sum_{i=1}^{k} \sum_{p \in c_i} (p - c_i)^2$$

K-means → Each cluster → represents by center of cluster

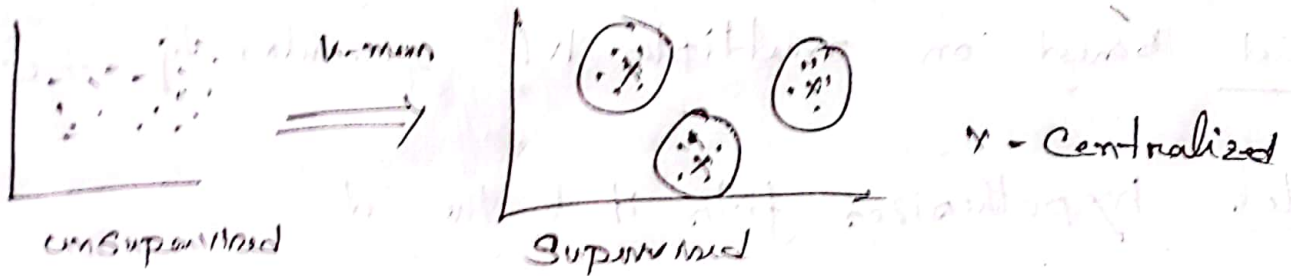K-medoids → Each " " " one of the objects
in the cluster

## K-means (unsupervised)
↳ center of Cluster.

K → no- of pre-defined Cluster

$V = 2$ अर्थात् 2 Clusters.



unsupervised          Supervised          → Centralized

* Elbow and Silhouette method used for.

K value detection

# k means Example

$A_1(2,10)$  $A_2(2,5)$  $A_3(8,4)$  $B_1 = 5,8$  $B_2 = 7,5$  $B_3 = 6,4$

$C_1 = 1,2,$  $C_2 = 4,9$

$A_1$  $B_1$  $C_1$ → Center    [using Euclidean distance)

$\Rightarrow$ Euclidean distance $d(P_1, P_2) = \sqrt{(n_2 - n_1)^2 + (y_2 - y_1)^2}$

$\sqrt{(2-2)^2 + (10-10)^2} = 0$ same item 0

| Data | | | Center | | | Cluster | New Cluster |
|------|------|------|------|------|------|------|------|
| | $n_x$ | $y_2$ | $A_1$ | $B_1$ | $C_1$ | | |
| | | | 2  10 | 5  8 | 1  2 | | |
| $A_1$ | 2 | 10 | 0 | 8.61 | 8.06 | ↳8 big → smaller বা বড়? | |
| $A_2$ | 2 | 5 | 5 | 4.24 | 3.16 | 1 / 3 | $A_1 = 1$ এবং আছে আর 5টে |
| $A_3$ | 8 | 4 | 8.49 | 5 | 7.28 | 2 | $A_1$ ① only |
| | 5 | 8 | 3.61 | 0 | 7.21 | 2 | $n_3 = 2, y_1 = 10$ |
| | 7 | 5 | 7.07 | 3.61 | 6.71 | 2 | $B_1 = 2.5$ বর |
| | 6 | 4 | 7.21 | 5.39 | 7.62 | 2 | $\frac{2+8+7+4+6}{6}$ |
| | 1 | 2 | 8.06 | 4.12 | 0.0 | 3 | $n_5 = \frac{8+5+7+4+6}{5}$ |
| | 4 | 9 | 8.66 | 7.21 | 5.39 | 2 | $= 6$ |
| | | | 2.24 | 1.41 | 7.24 | 3 | ৫ই for $B_1$ |
| | | | | | 7.62 | 2 | $\frac{30}{5} = 6$ |

$B_1 = (6, 6)$

$C_1 = (1.5, 3.5)$

Same way

new center :

$A_1 = (2, 10)$

$B_1 = (6, 6)$

$C_1 = (1.5, 3.5)$

— এর Same way -e আবার cluster গুলো Same আসবে?

— যদি আমাদের কাছে এমন change আসতো iteration

— হলে পাবে new center কে তৎ?

<u>Step 2</u>

|  | $A_1^{(2,10)}$ | $A_2^{(2,5)}$ | $A_3^{(8,4)}$ | $B_1^{(5,8)}$ | $B_2^{(7,5)}$ | $B_3^{(6,4)}$ | $C_1^{(1,2)}$ | $C_2^{(4,9)}$ | Center |
|---|---|---|---|---|---|---|---|---|---|
| $A \rightarrow$ | 0 | 5 | 8 | 3 | 7 | 7 | 8 | 2 | $A_1 = 2, 10$ |
| $B \rightarrow$ | 5 | 4 | 2 | 2 | 1 | 2 | 6 | 3 | $B_1 = 6, 6$ |
| $C \rightarrow$ | 6 | 1 | 6 | 5 | 5 | 9 | 1 | 6 | $C_1 = 1.5, 3.5$ |

— এখানে যে value smallest সেটা

| 1 | 3 | 2 | 2 | 2 | 2 | 3 | 1 | ← যে(1) দিলাম |

next Center

$A_1 = 1$ টা আছে আরে যদি $= \frac{2+4}{2} = 3 , \frac{19}{2} = 9.5 \Rightarrow (3, 9.5)$

$B_1 = (6.5, 5.25)$

$C_1 = (1.5, 3.5)$

① ← যে(1) দিলাম

② টা, so change হও (সেটা
আগে কম
হয়েছে)

$$D_3 = \begin{bmatrix} 1.2 & 4.61 & 7.43 & 2.5 & 6. & 6 & 7 \\ 6 & 4 & 4 & 9. & 1 & 3 & 64' \\ 6 & 1 & 6 & 5 & 5 & 4 & 1 & 6 \end{bmatrix}$$

|   | 3 | 2 | 1@ | 2 | 2 | 3 | - |

(12) mil,
এর Next Iteration

New coordinates

$A_1 = 3.67, 3$

$B_1 = (7, 4.63)$

$C_1 = (1.5, 3.5)$.

$$D_4 = \begin{array}{cccccccc} A_1 & B_2 & A_3 & B_1 & B_2 & B_3 & C_1 & C_2 \\ \begin{bmatrix} 1 & 4 & 6 & 1 & 5 & 3 & 7 & 0 \\ 7 & 5 & 1 & 4 & 0 & 1 & 6 & 3 \\ 6 & 1 & 6 & 5 & 5 & 4 & 1 & 6 \end{bmatrix} \end{array}$$

| 1 | 3 | 2 | 1 | 2 | 2 | 3 | 1 |

এটি নিজে থেকে আয়তা
তাই আপাতত এই

① → $A_1$ $C_2$, $B_1$

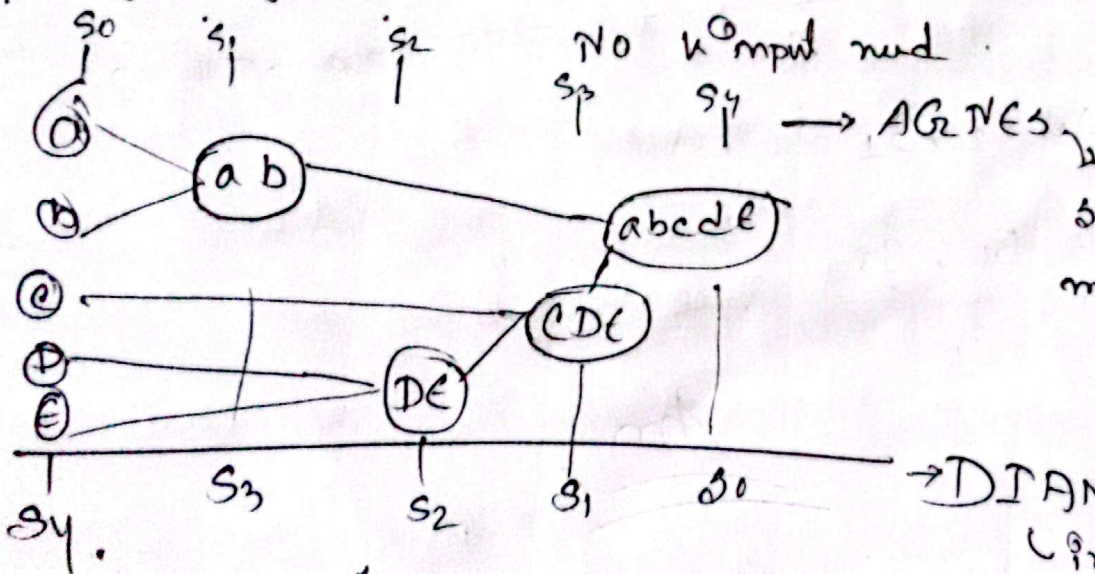② → $A_3$ $B_2$, $B_3$

③ → $A_2$, $C_1$

$Am$

# K modoids Clustering

Its identifier representive object in cluster

PAM - (partitioning around modoids) starts with initial
modoids and iteratively replace one with a non-moda
if it reduces the total distance.

CLARA → improved efficiency in large datasret
     └ PAM samples

# Hierarchical Clustering

Its builds a hierarchy of cluster by either ite..
ratively merging smaller clusters into large one or
spilitting larger clusters into smaller ones.



No k°mput ned

$s_p$   $s_4$ → AGNES

Single Link appru..
marge closest nod..

→ DIANA
  └ inverse of AGNES
Iteratively spilts all data

$s_4$ .

Denodogram

(সবার সম্ভব (আপন
data আছে ) — <u>Density Clustering</u>

It focusing on grouping data points based on local
density. its handle noise, One scan, discover clusters.

<u>Key methods</u>

DBSCAN — can identify noise or outliers.

OPTICS — Extends DBSCAN

— Density based spatial clustering Application. noise.

<u>DBSCAN Parameters</u> :

<u>Eps (C)</u> = minimum distance around a point to
look for other points.

<u>first
এই circle
আলো দিয়া</u>

<u>Minpts</u> = Minimum no of points required within
Eps distance for a point to be considered.

( minimum 3 টা circle টা৬ (১) )

2টা আগে ৬টা আছে



outlier
boundary
Core or 2(আ

Core

boundary)
যদি Core ৬।
আগে

Eps=1 cm.
minpts=5 (5 টা
আলা)

# DBSCAN Math

$\epsilon = 1.9$, minpts = 4.

$P_1 = (8,7)$  $P_3 = (5,5)$  $P_5 = (7,3)$  $P_7 = (7,2)$  $P_9 = (3,3)$  $P_{11} = (3,5)$

$P_2 = (4,6)$  $(P_4) = (6,4)$  $P_6 = (6,2)$  $P_8 = (8,4)$  $P_{10} = (2,4)$  $P_{12} = (2,3)$

use euclidian: $\left(\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}\right)$

|        | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ | $P_8$ | $P_9$ | $P_{10}$ | $P_{11}$ | $P_{12}$ |
|--------|------|------|------|------|------|------|------|------|------|------|------|------|
| $P_1$   | 0    |      |      |      |      |      |      |      |      |      |      |      |
| $P_2$   | 1.41 | 0    |      |      |      |      |      |      |      |      |      |      |
| $P_3$   | 2.83 | 1.41 | 0    |      |      |      |      |      |      |      |      |      |
| $P_4$   | 4.24 | 2.83 | 1.4  | 0    |      |      |      |      |      |      |      |      |
| $P_5$   | 5.66 | 4.24 | 2.8  | 1.4  | 0    |      |      |      |      |      |      |      |
| $P_6$   | 5.83 | 4.47 | 3.6  | 2    | 1    | 0    |      |      |      |      |      |      |
| $P_6$   | 6.40 | 5    | 3.1  | 14   | 1    | 1    | 0    |      |      |      |      |      |
| $P_8$   | 5.83 | 4.47 | 2.8  | 2    | 1    | 2    | 2    | 0    |      |      |      |      |
| $P_9$   | 4.00 | 3.16 | 3.1  | 3.9  | 4    | 3    | 4    | 3    | 0    |      |      |      |
| $P_{10}$ | 1.41 | 2    | 2    | 4.4  | 5    | 5    | 6    | 6    | 3    | 0    |      |      |
| $P_{11}$ | 2.00 | 1.4  | 1.4  | 3.16 | 4    | 4.5  | 5    | 5    | 2    | 1    | 0    |      |
| $P_{12}$ | 3.16 | 2.8  | 2.8  | 4.00 | 5    | 4    | 5    | 6    | 1    | 2    | 1    | 0    |

(উত্তরঃ 1.9) Case এর ভিতর horizontal o vertical both ck.

$P_1 \cdot P_2, P_{10}$ ⎮ $P_3 = P_2, P_4$ ⎮ $P_6 = P_5, P_7$ ⎮ $P_9 = P_{12}$ ⎮ $P_{11} = P_3, P_{10}, P_{12}$
$P_2 ; P_1, P_3, P_{11}$ ⎮ $P_4 = P_3, P_5$ ⎮ $P_7 = P_5, P_6$ ⎮ $P_{10} = P_1, P_{11}$ ⎮ $P_{12} = P_9, P_{11}$

| points | status | |
|---|---|---|
| $P_1$ | bordin | $P_3, P_2, P_{10}$ (3)এর _Noise_ |
| $P_2$ | | $P_2 P_1 P_3 P_{11}$ (4) same _Core_ |
| $P_3$ | bordin | Noise |
| $P_4$ | bordin | Noise |
| $P_5$ | | Core |
| $P_6$ | bordin | Noir |
| $P_7$ | bordn | Noise |
| $P_8$ | bordi | 4 |
| $P_9$ | | 9 |
| $P_{10}$ | bordn | 4 |
| $P_{11}$ | | core (4 বা তার) ← 0.b বা সমান |
| $P_{12}$ | bordn | Noise |

$P_9$ ও $P_{12}$ একই
2| core
$P_{12}$ ও

মিন minpoint এর data এর then noise মহান core

① এইটা noise করে
তোমার data core যারা তার noise ও border হবে



only
Noise

# Clustering Application

Biology - Classifies organisms into hierarchical categories

Information Technical - Group document by topic for improve Search and eng.

Land use : Identifies areas with Land use patterns

Marketing : Segment Customers based on behavior

Earthquake Studies : Loc of earthquake plate bound. areas where most are happen

Climate Studies : find patterns of weather and understand climate change