



Data Mining Unit 4 Handwritten

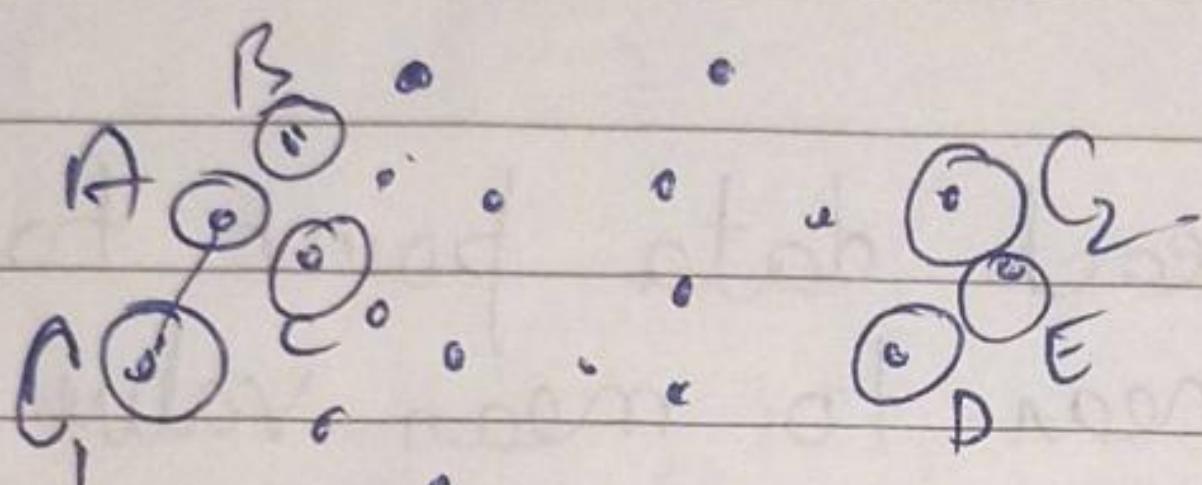
Data Mining And Analytics (SRM Institute of Science and Technology)

Unit - 4

Clustering → Clustering is the task of grouping a set of objects in such a way that objects in same group are more similar to each other than those in other group.

① K-means

→ K - no. of clusters.



| C ₁ | C ₂ |
|----------------|----------------|
| A | D |
| B | E |
| C | F |
| | |

Suppose we want to have 2 cluster. k=2.

Step 1 → Select 2 random point as centroid C₁, C₂.
Now measure distance of all point from C₁ and C₂.

$$\begin{array}{ll}
 C_1 & C_2 \\
 \hline
 A & D \\
 B & E \\
 C & F \\
 \hline
 \frac{A+B+C}{3} & , \quad \frac{D+E}{3}
 \end{array}$$

→ Mean of above will become centroid.

* Distance is measured as Euclidean distance.

Algorithm k-means :-

Input → k - Number of clusters.

Data Set (D)

Output → Set with k clusters.

- Algo
- ① Select k point randomly from the data set D .
 - ② Repeat
 - a. Assign each data point to the cluster that is near to mean value (centroid) of the cluster.
 - b. Update the centroid value of the cluster with the mean value of members.
 3. Until there is no change in cluster.

Note → Do it 3-4 times in exam

Q. Apply k-means clustering to the following data-points.

| | x_1 | x_2 |
|---|-------|-------|
| A | 5 | 1 |
| B | 4 | 2 |
| C | 3 | 3 |
| D | 1 | 4 |
| E | 6 | 7 |

where $k = 2$.

Sol → Let A represent C_1 and B represent C_2 .

| | Distance from C_1 | Distance from C_2 | Cluster |
|---|---------------------------------------|---------------------------------------|---------|
| A | $\sqrt{(5-5)^2 + (1-1)^2} = 0$ | $\sqrt{(5-4)^2 + (2-1)^2} = \sqrt{2}$ | C_1 |
| B | $\sqrt{(5-4)^2 + (2-1)^2} = \sqrt{2}$ | $\sqrt{(4-4)^2 + (2-2)^2} = 0$ | C_2 |
| C | $2\sqrt{2}$ | $\sqrt{2}$ | C_2 |
| D | 5 | $\sqrt{13}$ | C_2 |
| E | $\sqrt{37}$ | $\sqrt{29}$ | C_2 |

New, $C_1 = A$ [\because only 1 (A) has C_1]

$$C_2 = \frac{B + C + D + E}{4} \leftarrow \text{point form.} \quad (\text{mean of } A, B, C, D)$$

$$C_2 = \left(\frac{4+3+1+6}{4}, \frac{2+3+4+7}{4} \right)$$

$$\Rightarrow C_2 = (3.5, 4).$$

| | Distance Square from C_1 | Distance Square from C_2 | Cluster |
|---|----------------------------|-------------------------------|---------|
| A | $(5-3.5)^2 + (1-4)^2 = 0$ | $(5-3.5)^2 + (1-4)^2 = 11.25$ | C_1 |
| B | $\sqrt{2}$ | 4.15 | C_1 |
| C | $2\sqrt{2}$ | 1.25 | C_2 |
| D | 5 | 6.25 | C_2 |
| E | $\sqrt{37}$ | 25.25 | C_2 |

$$\text{New } C_1 \rightarrow \left(\frac{A+B}{2} \right)$$

$$\rightarrow \left[\left(\frac{5+4}{2} \right), \left(\frac{1+2}{2} \right) \right]$$

$$= (4.5, 1.5)$$

$$\text{New } C_2 \rightarrow \left(\frac{C+D+E}{3} \right)$$

$$\rightarrow \left[\left(\frac{1+6+3}{3} \right), \left(\frac{3+4+7}{3} \right) \right]$$

$$C_2 \rightarrow (3.3, 4.7)$$

| | Distance Square from C_1 | Distance Square from C_2 | Clusters |
|---|----------------------------|----------------------------|----------|
| A | 1.25 | 16.58 | C_1 |
| B | 0.5 | 7.78 | C_1 |
| C | 4.5 | 2.98 | C_2 |
| D | 18.5 | 5.78 | C_2 |
| E | 32.5 | 12.58 | C_2 |

Since Now, Cluster 1 does not change.

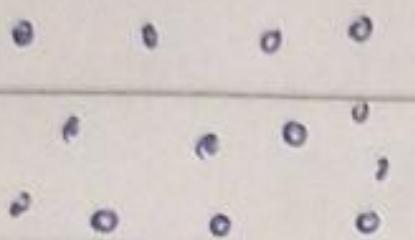
So, we stop.

So, Cluster 1 $\rightarrow A, B$
 Cluster 2 $\rightarrow C, D, E$

Disadvantage of k-mean Clustering :-

1. The k-means algorithm is sensitive to outliers!

O ← outlier

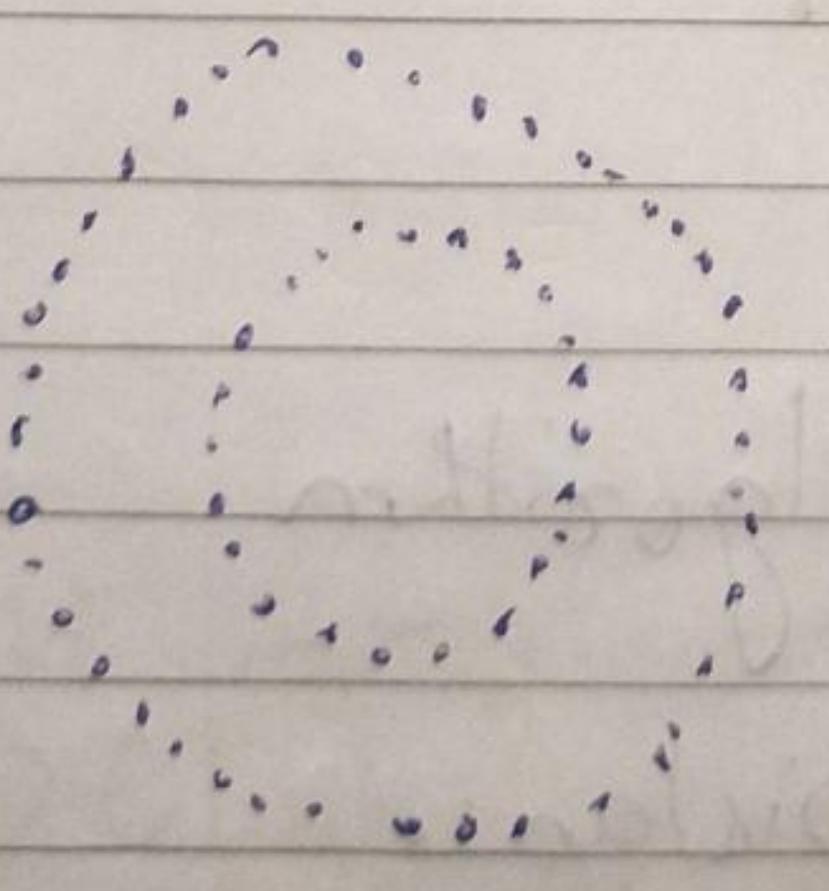


Cluster is impacted by outlier in this.

2. In this when we find new ~~cluster~~^{centroid} point, we might get a decimal point which may not be any actual point in cluster.

So to overcome this we used new algo which is based on median instead of mean.

3.



for any cluster in any shape k-mean is not able to identify shape or define cluster in form of that shape. For this we used Density Based cluster.

K-Medoide :-

| | x | y | Distance from C, (4, 5) | Distance from (5, 8) |
|-----|---|--------------------|----------------------------|-------------------------|
| 1. | 8 | 7 | 6 | 2 |
| 2. | 3 | 7 | 3 | 7 |
| 3. | 4 | 9 | 4 | 8 |
| 4. | 9 | 6 | 6 | 2 |
| 5. | 8 | 5 | 4 | 6 |
| 6. | 5 | 8 → C ₂ | — | — |
| 7. | 7 | 3 | 5 | 3 |
| 8. | 8 | 4 | 5 | 1 |
| 9. | 7 | 5 | 3 | 1 |
| 10. | 4 | 5 → C ₁ | — | — |

C₁ → (2, 3, 5) , C₂ → (1, 4, 7, 8, 9).

$$\text{Cost} = (3+4+4) + (2+2+3+1+1) \\ = 20.$$

Now, in next step, we choose any other point than C₂ as C_{2'} & find cost.

Let say (7, 5) as C_{2'}.

Then Cost = 22

As 22 > 20

So, undo swap.
& follow this step till end.

Algorithm k-medoids :-

Input :- K :- no. of clusters.

D :- data set containing n objects

Output :- A set of k -clusters.

1. Randomly chose k object in D as initial representative object or seeds :-
2. Repeat
3. assign each remaining object to the cluster with nearest representative object;
4. randomly select a non-representative object Orandom.
5. compute the total cost, S , of swapping representative object O_j with Orandom.
6. if $S < 0$, then swap Orandom, O_j to form new set of k representative object.
7. Until no change.

Algo ↑

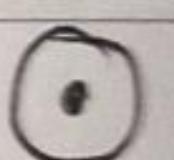
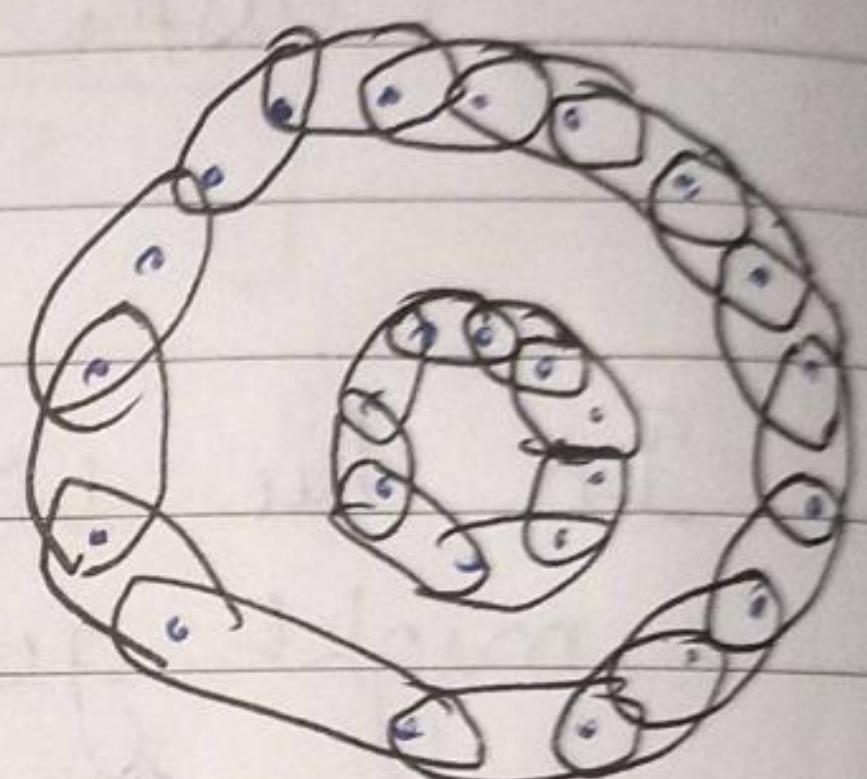
* K-medoid is more robust than K-mean in presence of noise or outliers because a medoid is less influenced by outlier or other extreme values than mean. However processing of K-medoid is more costly than K-mean.

Distance in clustering Algorithm :-

- 1) Minimum Distance
- 2) Maximum Distance
- 3) Mean (Average) Distance
- 4) Average Distance

Nearest Neighbour :-

2-D ~~Circle~~ (Circle)
3-D ~~Circle~~ (Sphere).



K-Nearest Neighbour (KNN)

↳ Supervised Algorithm

| | x_1 | x_2 | y | Distance |
|---|-------|-------|----------------|----------|
| A | 1 | 2 | c ₁ | 1 |
| B | 2 | 3 | c ₁ | 1 |
| C | 1 | 4 | c ₂ | 1 |
| D | 1 | 5 | c ₂ | 4 |

(1, 3) → ? → to which cluster it will go

Density Based Clustering :- Continue growing the cluster as density of neighbourhood exceeds threshold.

DBSCAN Algorithm :-

1. Mark ALL OBJECTS as UNVISITED.
2. do
3. Randomly Select some point P
4. mark it as visited
5. If E -neighbourhood of P is atleast minpts.
Create a cluster C
Let N points are objects of E -neighbourhood of P
for each point P' from N
if P' is unvisited
mark P' as visited
if E -neighbourhood of P' has atleast minpts.
add these point to N.
if P' is not a part of cluster then add it to cluster C.

end for

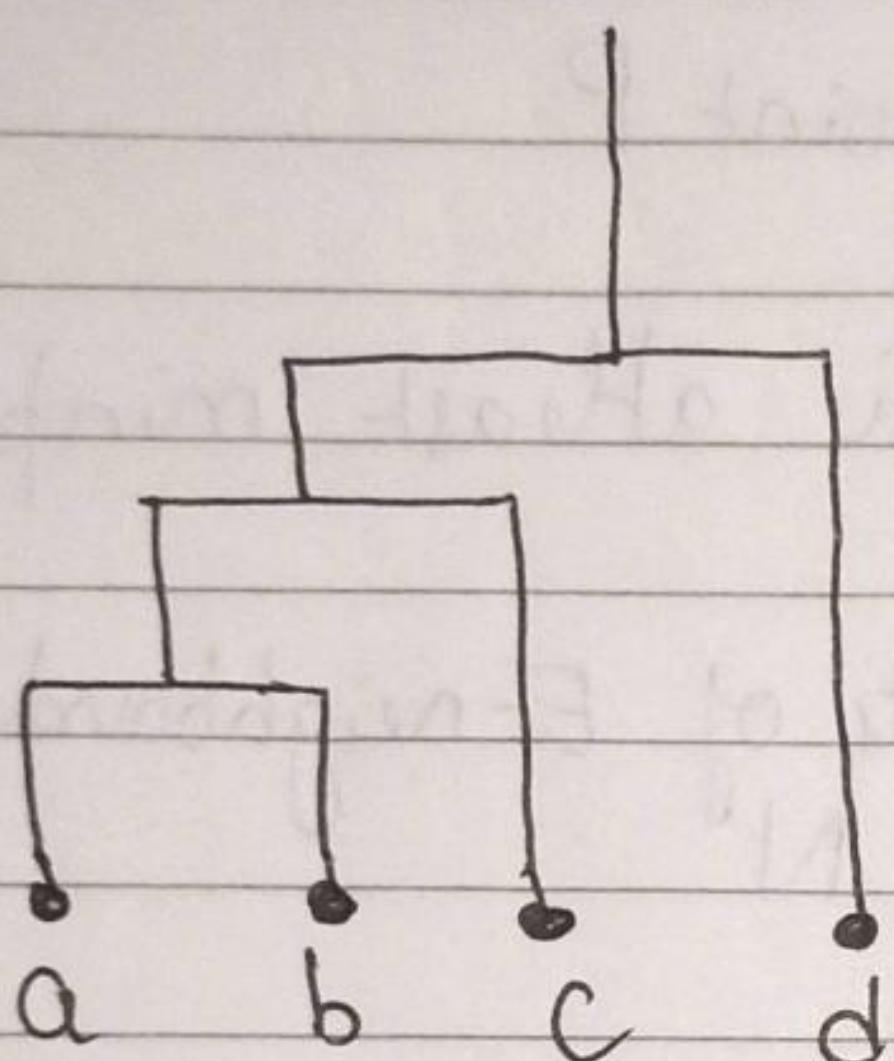
Output C

else mark P as noise

until all points are visited.

Agglomerative Approach

- In this objects are merged to form a cluster.
- It uses bottom up approach.



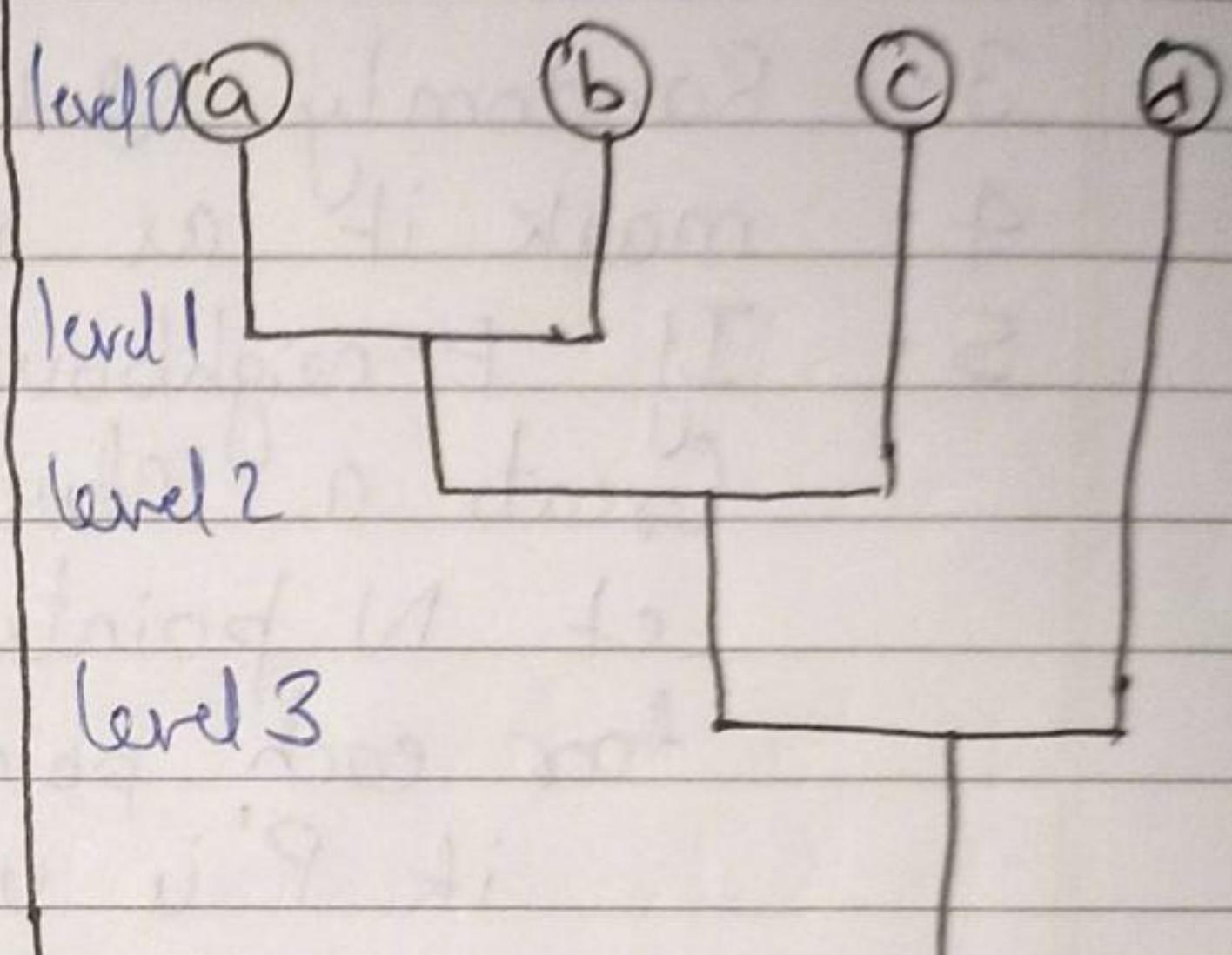
→ Represented using Dendrogram.

→ Uses closed group objects to merge so euclidean distance is calculated between various points to decide the merging point.

*Note → On viewing from top of the divisive approach diagram, it looks like Agglomerative clustering and from bottom it looks like divisive approach.

Divisive Approach

- In this clusters are split into objects.
- It uses top down approach.



→ Represented using Dendrogram.

BIRCH (Balanced Iterative Reduced Clustering Using Hierarchical Hierarchy) :-

K-means usually have time complexity of $O(n^2)$. So to improvise then a new method BIRCH was introduced which has the time complexity of $O(n)$. It provides scalability.

* BIRCH is designed for clustering a large amount of numerical data by integration of hierarchical clustering and other clustering methods such as iterative partitioning.

It overcomes two difficulties of agglomerative clustering methods :-
i) Scalability
ii) the inability to undo what was done in previous step.

* BIRCH introduces two concepts :-
i) Clustering factor
ii) clustering factor tree

* Clustering factor \rightarrow is a vector (3-D) summarising information about cluster. $\langle n, LS, SS \rangle \leftarrow CF$

no. of points in cluster Linear sum of n point Square sum of data points

Q. Given three points $(2, 5)$, $(3, 2)$, & $(4, 3)$.
Compute the clustering factor.

Sol →

$$n = 3.$$

$$\begin{aligned}\text{Linear Sum} &= \sum_{i=1}^n x_i \\ &= (2+3+4), (5+2+3) \\ &= (9, 10).\end{aligned}$$

$$\begin{aligned}\text{Square Sum} &= \sum_{i=1}^n x_i^2 \\ &= 4+9+16, 25+9+4 \\ &= (29, 38).\end{aligned}$$

$$CF_1 = \langle 3, (9, 10), (29, 38) \rangle \quad \underline{\text{Ans}}$$

How to merge 2-clustering factor?

$$CF_2 = \langle 3, (35, 36), (412, 440) \rangle$$

$$\begin{aligned}CF_3 &= CF_1 + CF_2 \\ &= \langle 3, (9, 10), (29, 38) \rangle + \langle 3, (35, 36), (412, 440) \rangle \\ &= \langle 3, (44, 46), (446, 478) \rangle. \quad \underline{\text{Ans}}\end{aligned}$$

How to find Centroid of C.F :-

$$* \text{Centroid} = \frac{\sum x_i}{n} = \frac{LS}{n}$$

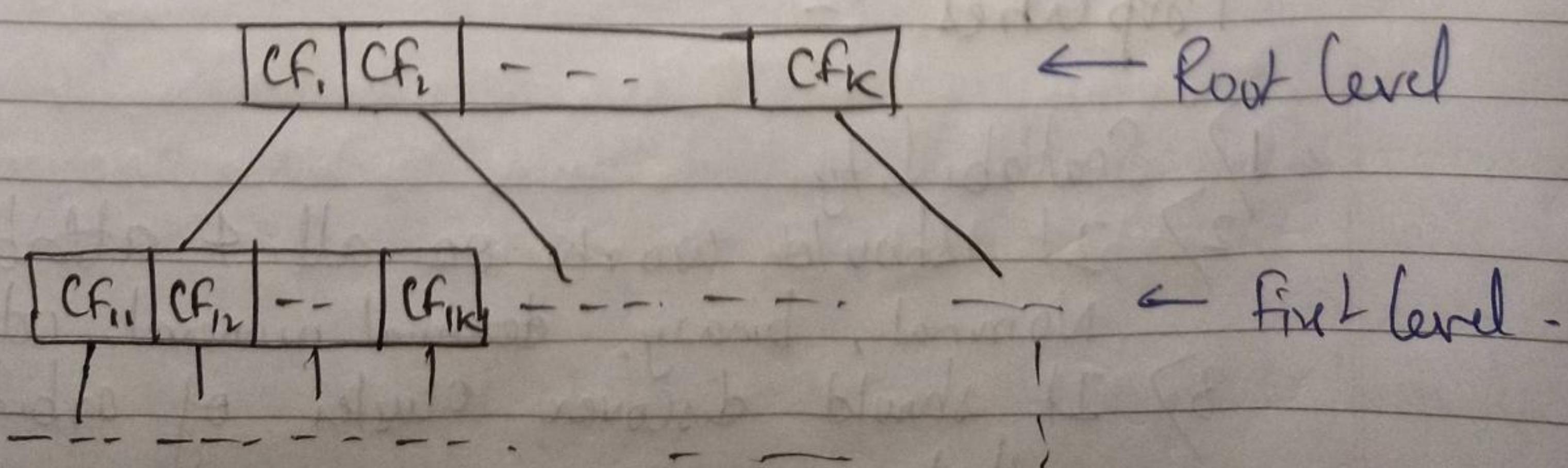
$$= \left(\frac{9}{3}, \frac{16}{3} \right) \text{ in above examp.}$$

$$* \text{Radius} = \sqrt{\frac{\sum (x_i - x_0)^2}{n}} = \sqrt{\frac{n \cdot SS - 2LS^2 + n \cdot LS}{n^2}}$$

$$* \text{Diameter} = \sqrt{\frac{2 \cdot n \cdot SS - 2 \cdot LS^2}{n(n-1)}}$$

* CF Tree :- CF Tree is a height balanced tree that stores the clustering factor for a hierarchical clustering.

Example :-



- * A CF Tree has 2 parameters :
 - i) Balancing Factor :- Specifies maximum no. (Branching Factor) of children per non-leaf node.
 - ii) Threshold → specifies maximum diameter of subset subclusters stored at the leaf node of the tree.

These two parameters influence the size of resulting tree.

BIRCH Algorithm works in 2 phases :-

- 1) Phase 1 :- Scan the database to build an initial in-memory CF Tree with multilevel compression of data.
- 2) Phase 2 :- Apply clustering algorithm to cluster the leaf node of the CF Tree, which removes sparse clusters as outliers and ~~remove group~~ merge large dense clusters into larger ones.

Properties :-

- 1) Scalability
- 2) It should work on all 4 attributes (i.e Nominal, binary, decimal numerical, ordinal).
- 3) It should discover cluster of arbitrary shape.

- 4) It should be incremental.
- 5) It should work on high dimensional data.

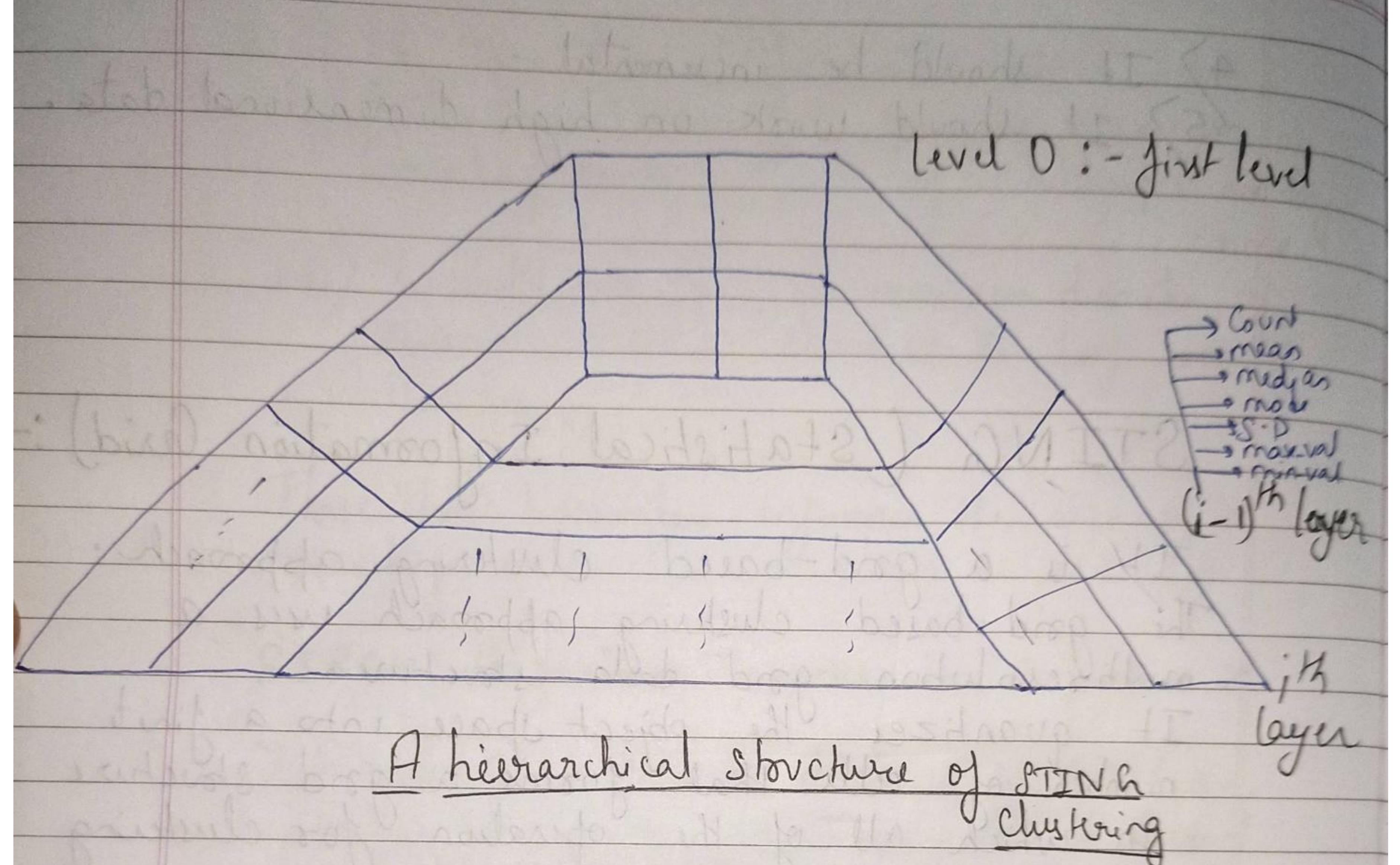
STING (Statistical Information Grid) :-

It is a grid-based clustering approach. The grid-based clustering approach uses a multiresolution grid data structure. It quantizes the object space into a finite number of cells that forms a grid structure on which all of the operations for clustering are performed.

The main advantage of the approach is its fast processing time, which is typically independent of the number of data objects, yet dependent on only the number of cells in each dimension in the quantized space.

STING:- Sting is a grid based ~~multiresolut~~ multiresolution clustering technique in which the spatial area is divided into rectangular cells.

There are usually several levels of each rectangular cells corresponding to different levels of resolution, and these cells form a hierarchical structure, each cell at higher level is partitioned to form a number of cells at the next lower level.



Statistical parameters of higher level cells can be easily computed from parameters of lower-level cells. These parameters include the following :-

- the attribute independent parameters, (count)
- the attribute dependent parameters, (mean, S.D, min, max)
- the type of distribution that the attribute value in the cell follows, such as normal, uniform, exponential or none.

Advantages of STING :-

STING offers several advantages :-

- i) the grid-based computation is query independent, because the statistical information stored in each cell represent the summary information of data in grid cell, independent of query.
- ii) the grid structure facilitates parallel processing and incremental updating.
- iii) the med method's efficiency is a major advantage. time complexity of generating clusters is $O(n)$.

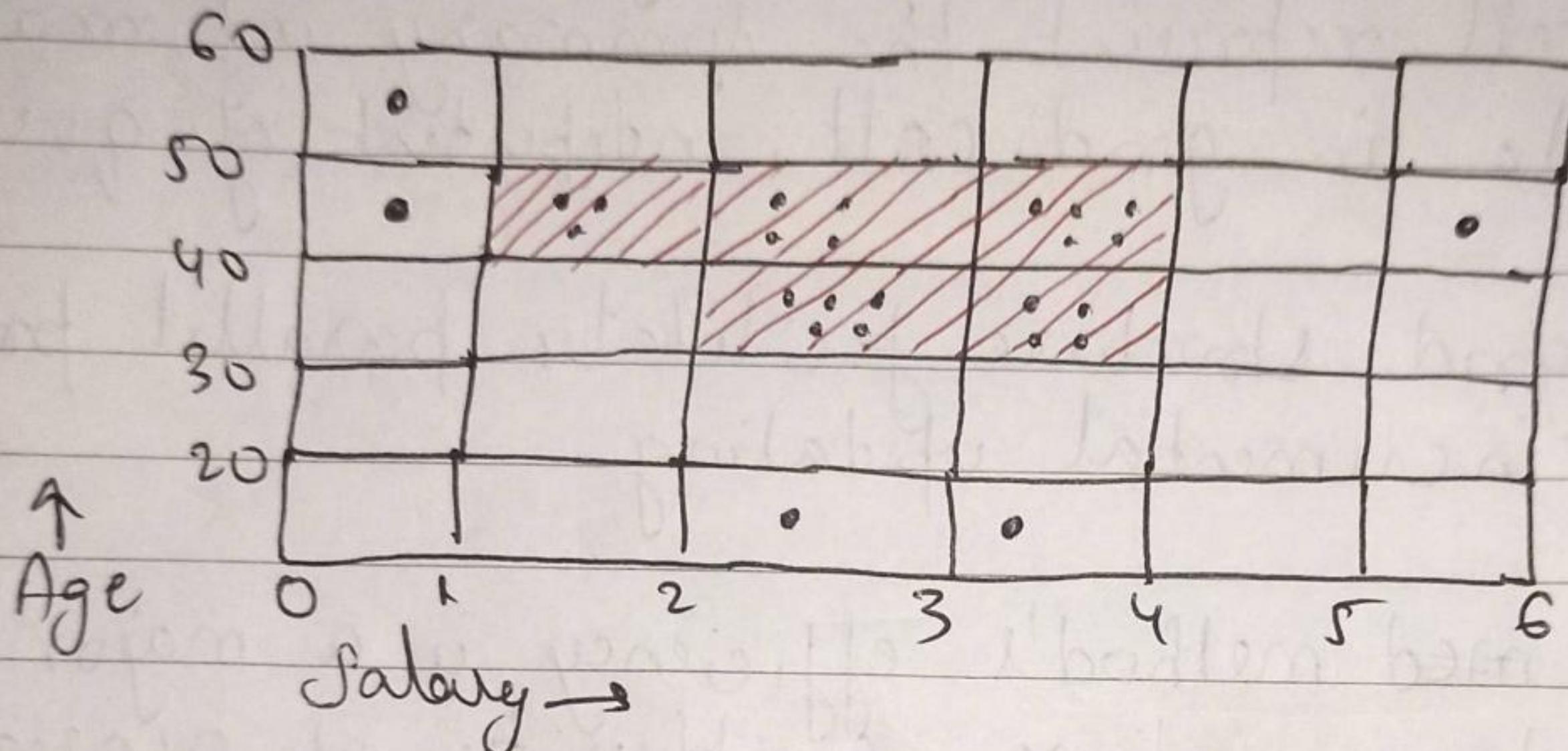
CLIQUE (clustering in Quest) :-

- 1. It is a grid based clustering algorithm
- 2. It is density based clustering algorithm
- 3. It uses Apriori Property.

CLIQUE was the first algorithm proposed for dimension-growth subspace clustering in high dimensional space.

Cells

Since it is based on density Based, there will be a threshold in there?



Threshold value define whether the grid is dense or sparse.

If Threshold = 3.

Then,

| | |
|------------|----------------------|
| $T \geq 3$ | \rightarrow dense |
| $T < 3$ | \rightarrow sparse |

Steps :-

- 1) Identify the subspace with clustering
- 2) Identify the cluster
- 3) Generate the minimum descriptor for clustering

Advantage :-

- 1) It is faster than K-means • Moreover, we need not have to find neighbours that we need to find DBSCAN, so it is faster than DBSCAN and K-means both.
- 2) It can find any shape of cluster
- 3) It is one of the simplest method.

Disadvantage :-

- 1) Not suitable for high dimension because
- 2) It is hard to visualize when dimension increases

Evaluation of Clustering :-

① Assessing the clustering technique :-

Whether they access the non-random structure or not?

Hopkin's Statistics :- used to identify whether data set follow any pattern or not?

→ Given a dataset D as a sample of random variable σ , determine how far away σ is from being uniformly distributed in data space.

→ Sample n points p_1, \dots, p_n uniformly from D .
 For each p_i find its nearest neighbour in D .

$$u_i = \min_{\substack{v \in D \\ v \neq p_i}} \text{dist}(p_i, v) \quad \text{where } v \in D.$$

\uparrow
 Population
 \uparrow
 Set of neighbours

→ Sample n points q_1, \dots, q_n uniformly from D ,
 for each q_i finds its nearest neighbour in $D - \{q_j\}$:-
 $y_i = \min_{\substack{v \in D \\ v \neq q_i}} \text{dist}(q_i, v)$ where $v \in D$ & $v \neq q_i$

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

→ If Uniformly distributed, $\sum x_i = \sum y_i$

$$\Rightarrow H = \frac{0.5}{\downarrow}$$

Then cluster is uniformly distributed.

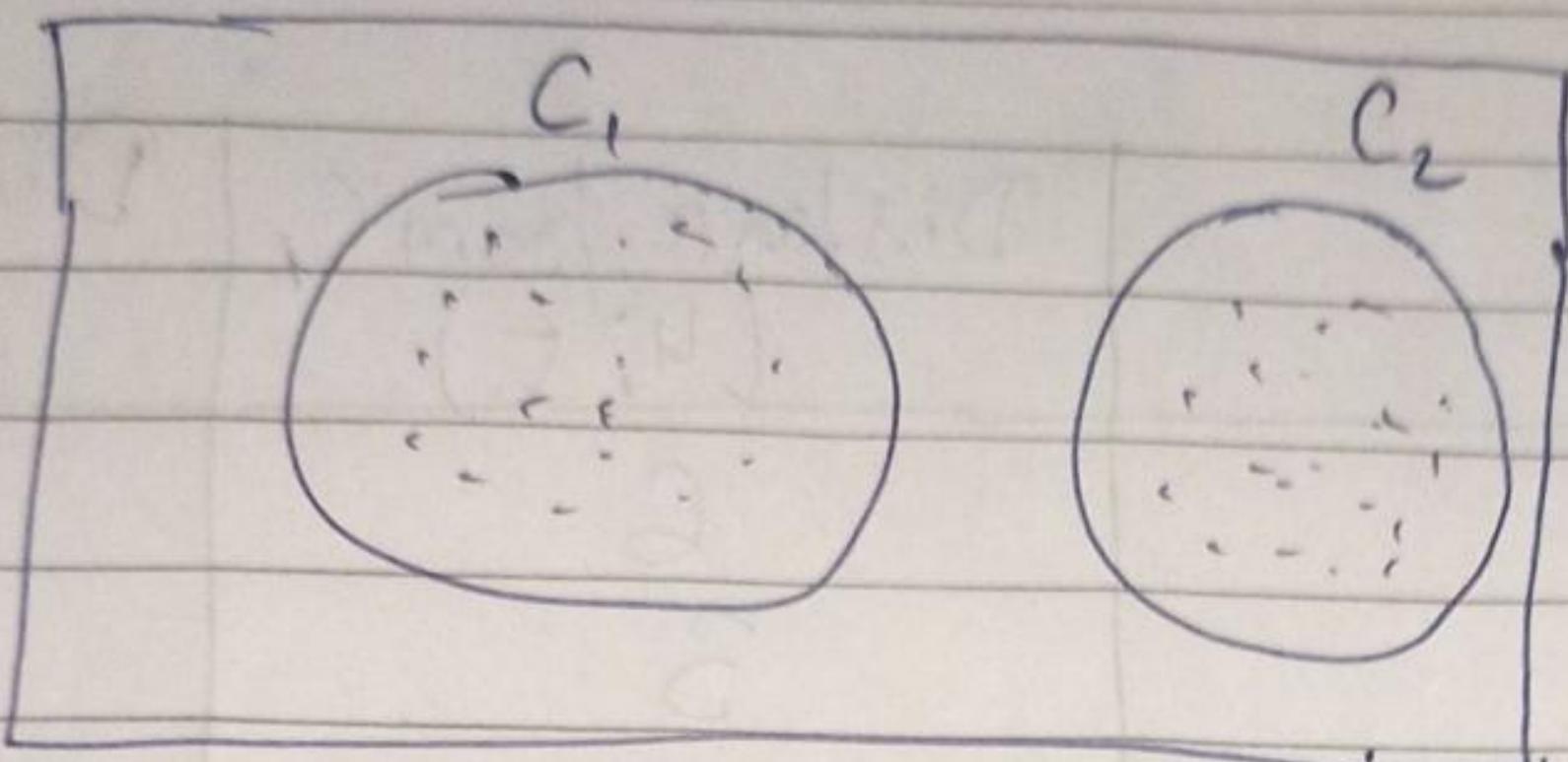
→ If D is highly skewed, H is close to 0.

2) Empirical Method :-

We find No. of clusters in this. For n points, No. of clusters

$k = \sqrt{\frac{n}{2}}$ for dataset of n points.
 each cluster having $\sqrt{2n}$ points.

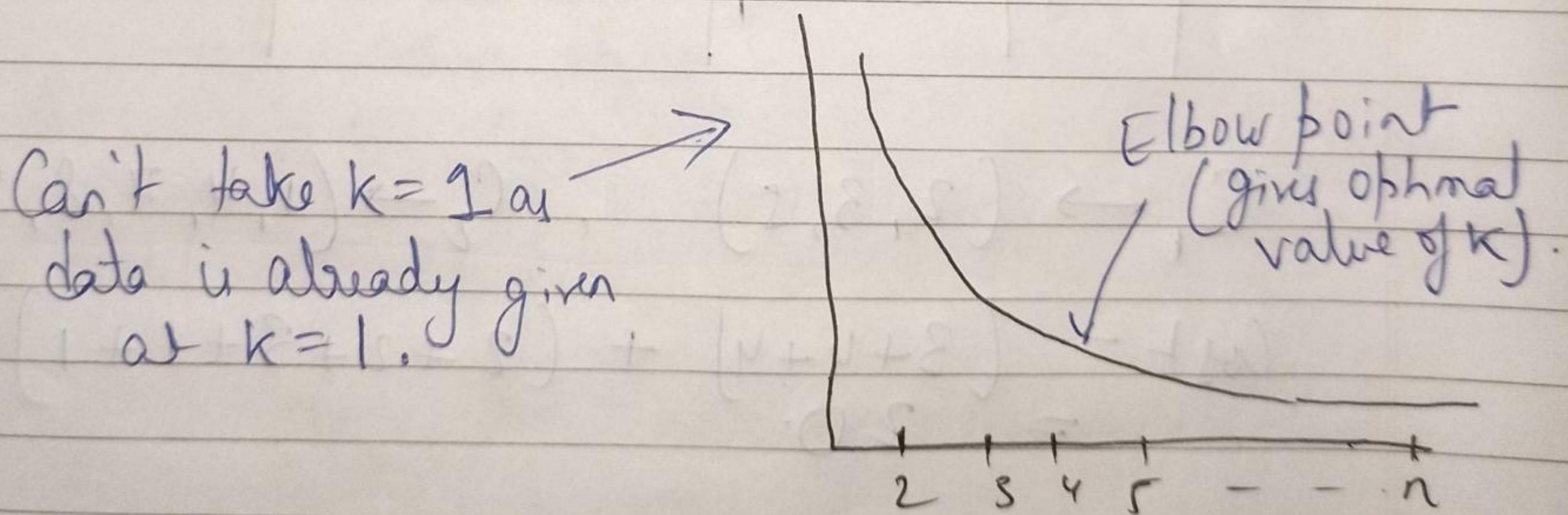
3) Elbow Method :- (intracluster).



→ We take the turning point in the curve of sum of within cluster variance with respect to no. of clusters.

for 2 to n :-

Find distance & draw it on graph



Silhouette Method / Score :-
depend on intra cluster.