

Clustering: Partitioning Methods, Hierarchical Methods, Density-Based Methods, Grid-Based Methods.

Partitioning → ~~K-Means~~
→ K-Medoids
→ CLARA

Properties

① Scalable

② High dimensionality

③ Diff data

④ Unstructured data

Partitioning \rightarrow K-Means

16, 16, 17, 20, 21, 21, 22, 23

K = 2

Centroid (C_1) = 16

Centroid (C_2) = 21

$$S_1 = [16, 16, 17] \rightarrow 16.33$$

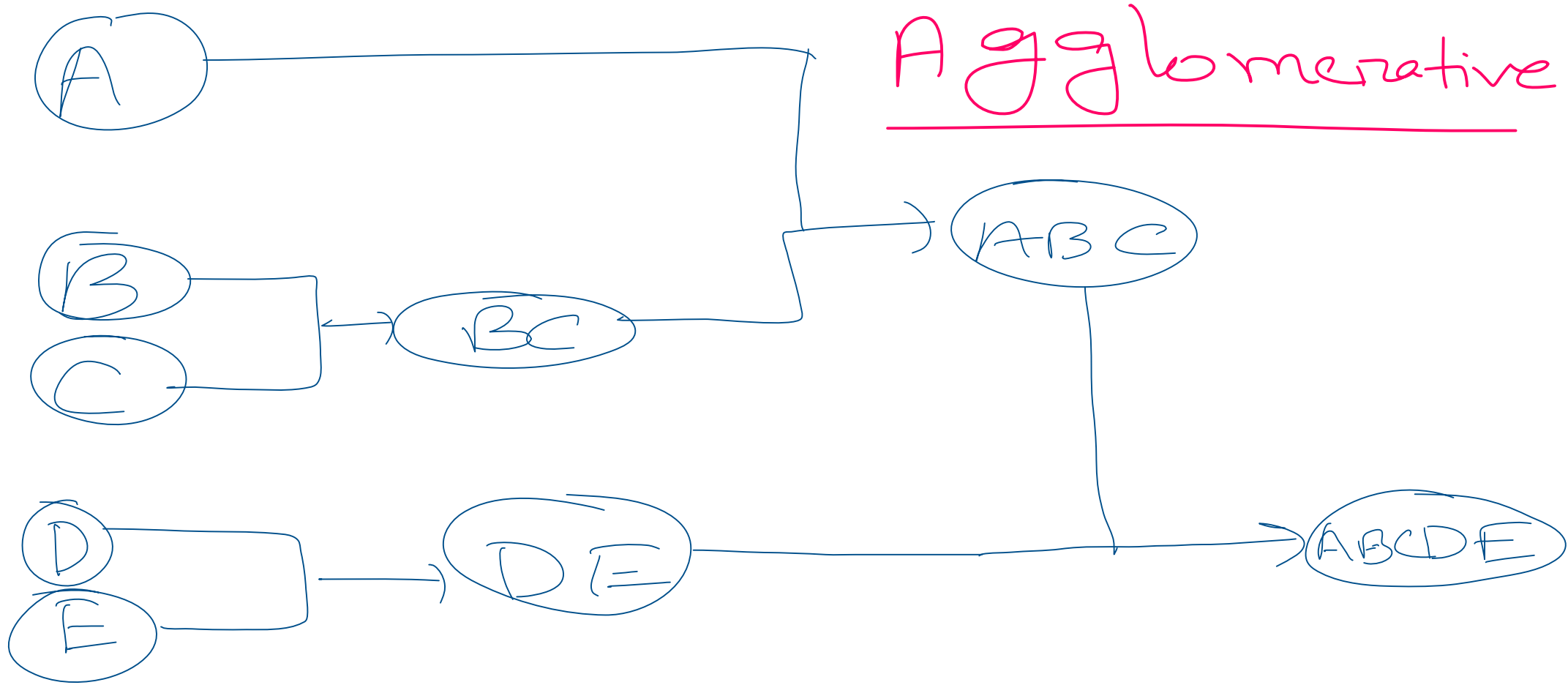
$$S_2 = [20, 21, 21, 22, 23] \rightarrow 21.4$$

STOP

7

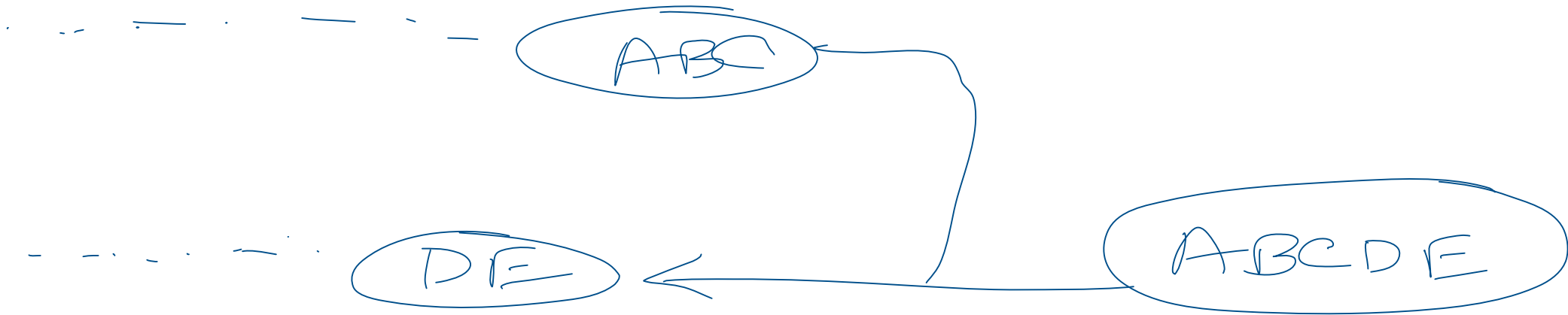
Hierarchical

- Tree of clusters
- Bottom-up



Agglomerative

Divisive



Adv

- ① Diff size, density
- ② Missing & noisy
- ③ Relationship

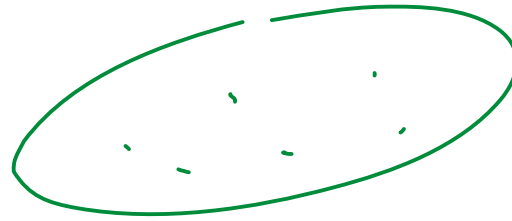
Disadv

- ① Computational cost $\uparrow\uparrow$
- ② Criterion \rightarrow #clusters
stopping

DBSCAN



DB-1



DB-2

1000

1

Queries

— large diff in density

Density-based Spatial Clustering of Applications with Noise

DBSCAN

① ϵ \rightarrow neighborhood

dist, $d[A, B] \leq \epsilon$

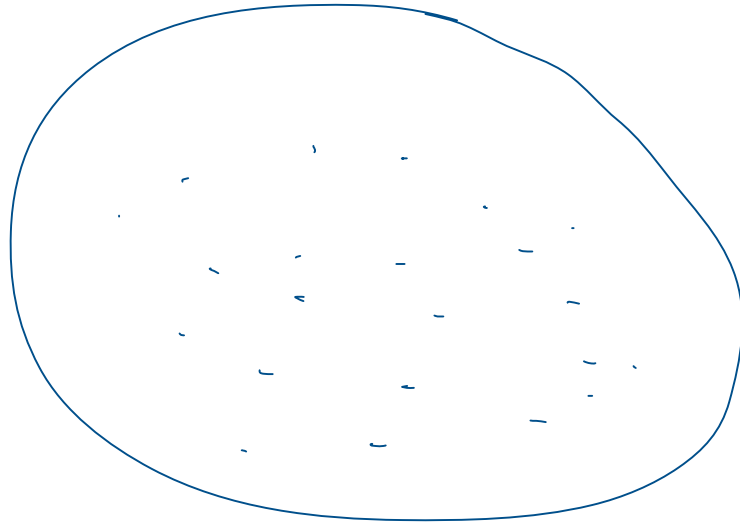
neighbors

- \angle - distance graph

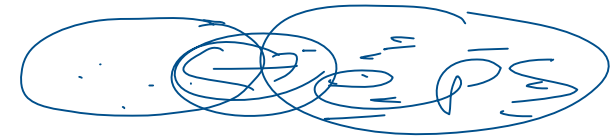
$$\underline{\underline{cp=2}}$$

② MinPts \rightarrow min number
of datapoints
(neighbor)

$$\geq 3$$



- dataset \uparrow MinPts \uparrow

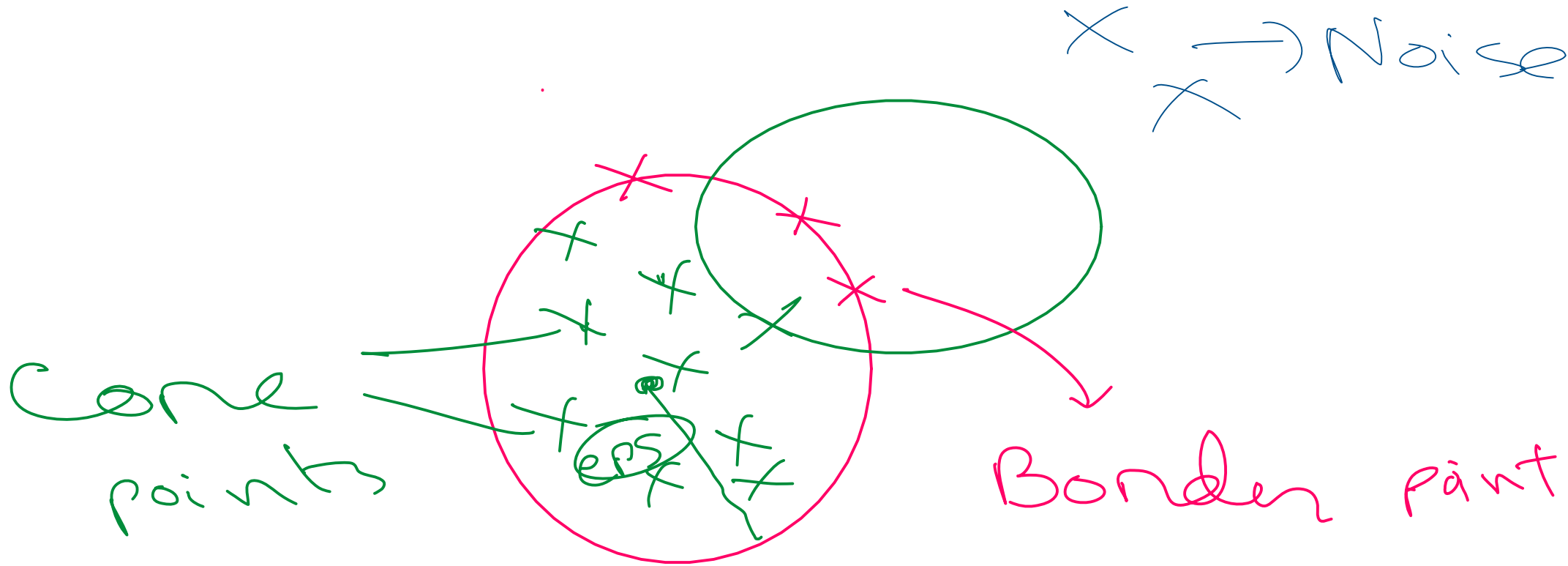


3-types of data-points

① Core \rightarrow $>$ Minpoints

② Border \rightarrow $<$ Minpts & in the
neigh of core point

③ Noise/Outlier \rightarrow Neither
Core / Border



Steps

- ① Find all neighbor points within ϵ and identify core pts
 - ② For each core pts, if not assigned to a cluster \rightarrow new cluster
 - ③ Recursively find all density - connected points.
- Assign the 1 sample at 100

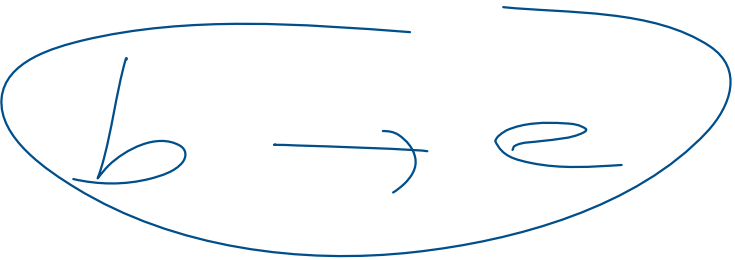
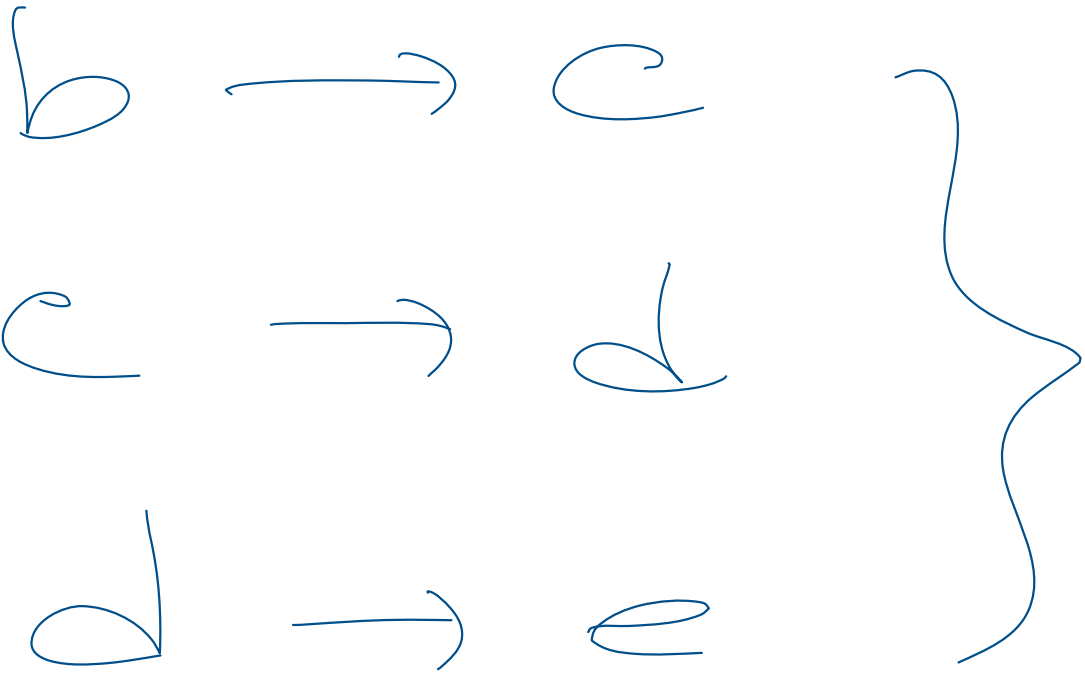
Assign them to same cluster
as the core point.

Density - connected points



10

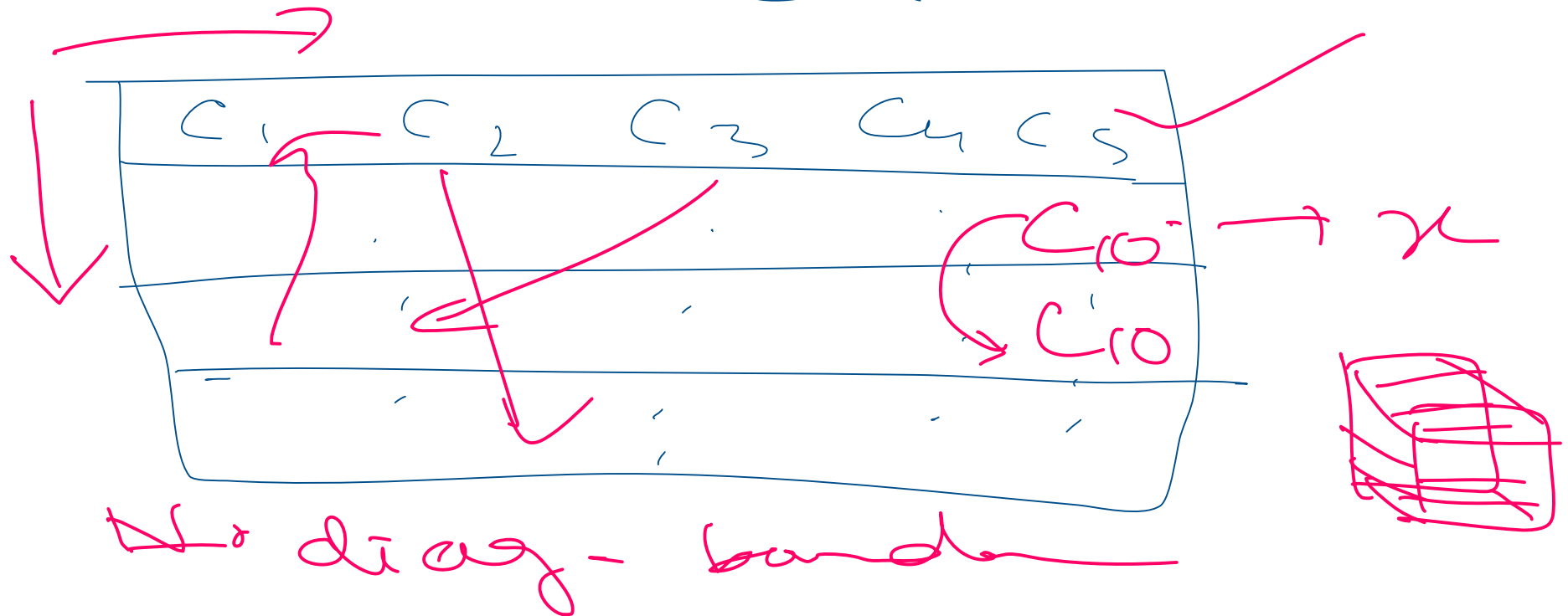




Grid-based → STING

Statistical Infor.

hierarchical



OPTICS → Density based

Ordering Point to

Identify Clustering Struct

