

|          |     |
|----------|-----|
| PAGE No. |     |
| DATE     | / / |

## DBSCAN (Density-Based spatial clustering of Application with Noise)

invented by Ester, Kriegel, Sander, Xu in 1996.

\* here Noise means outliers.

\* when to use DBSCAN

- when data is large.
- when data is not in linear fashion.
- Not Sensitive to Noise.

\* in K-means it has some drawbacks

- 1) we have to specify the no. of K.
- 2) K-means does not perform well on finding non-convex / non-spherical shapes of clusters.
- 3) K-means is sensitive to noise data.

## \* Hierarchical clustering drawback.

- Not Suitable for big dataset.
- high computational complexity
- Need clusters (linkage) that affects the clustering result.
- sensitive to noise

## \* DBSCAN ~~intuition~~

It discards clusters of non spherical shape the DBSCAN clustering method can represent cluster of arbitrary shape and to handle noise.

\* DBSCAN is like human intuition it means DBSCAN will find out clusters just like humans way.

based on there 2 parameters  
DBSCAN will create clusters -

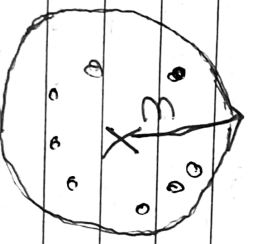
DBSCAN parameters

① epsilon distance. ( $\epsilon$ ) (user - define)

$$N_{\epsilon}(x) = B_{\epsilon}(x, \epsilon) = \{y | d(x, y) \leq \epsilon\}$$

② epsilon distance should be bigger than 0 that because it radius.

it defines the size and borders of each neighbourhood.

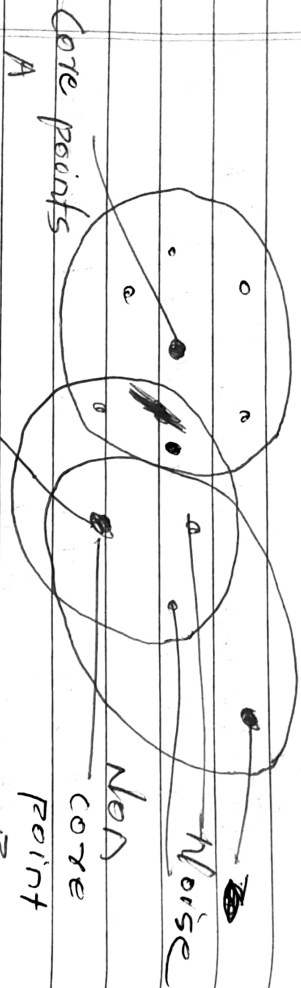


② minPTS (minimum points) (user - define)  
should be at least 3 points in minPTS.  
its a density threshold for example  
lets say minPTS = 5 so in  $\epsilon$   
neighbourhood include at least 5.  
then it will consider cluster otherwise not -

### 3 Border Point.

To understand this Border Point let take an example.

$E \rightarrow 45$  minpts  $\rightarrow 4$



(it is close to the core points)

In this case core A has more than 4 points so it core point but B is not a core but it doesn't have more than 5 point but it closer to the Neighborhood of core point so it called Border point but it count as noise. The noise will count when we don't have core point further to make cluster as you can see the diagram this is also called Directly Density Reachable

in general  $\epsilon$  should be chosen as small as possible.

\* metrics for measuring DBSCAN's performance.

for performance.

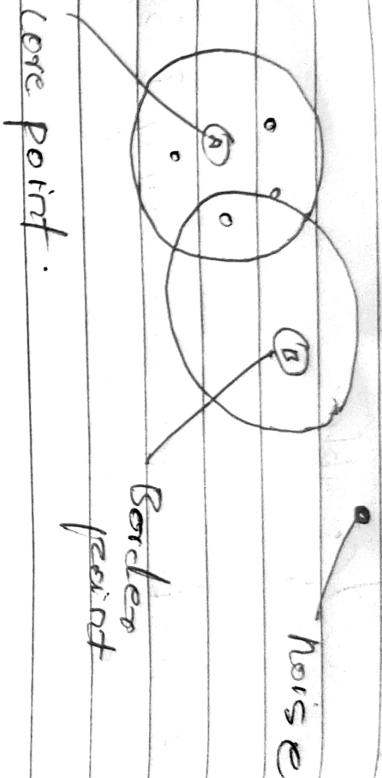
K-means  $\rightarrow$  inertia

DBSCAN  $\rightarrow$  Silhouette Score.

/

-1 (worst)  
1 (best)  
0 (overlapping)  
clusters.

Note - when two core point meets  
two each then it's going to merge to  
the cluster.



steps in DBSCAN algorithms.

- ① Find the neighbor points.
- ② if core point is not assigned to a cluster create a new cluster.
- ③ these who are non core points but the neighbor points belongs to the core points then it will take as a cluster but that non core points will not further be extended to make epsilon distance

④ iterate the process.

Pros

- ① ~~used~~ identify randomly shapes clusters.
- ② doesn't need literature
- ③ Handle noise

cons

- ① varying densities are problematic
- ② input of DBSCAN may be difficult to determine.
- ③ computational complexity  $\rightarrow$  when the dimensionality is high