# (12) DBSCAN (Density Based Spatial Clustering of Applications with Noice.)

→ It is an unsupervised machine learning algorithm. This algorithm defines clusters as continuous regions of high density.

## Key words :-

### ⊙ EPSILON :- (EPS)
This is one distance till which we look for the neighbouring points.

### ⊙ min-points :-
The min no. of points specified by the user.

### ⊙ core points :-
If the no. of points inside the epsilon radius of a point is greater than or equal the min-points then it is called a __core point__.

### ⊙ Border points :-
If the no. of points inside the epsilon radius of a point is less than the min points and it lies within the eps radius region of a core point, then it is called border point.

### ⊙ noise :-
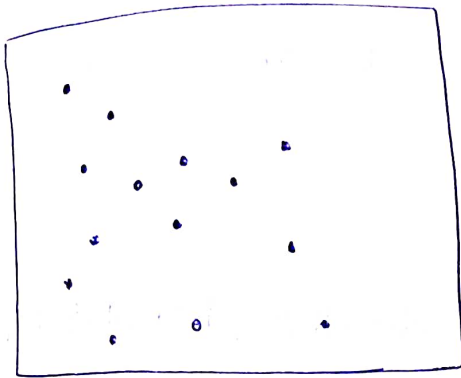A point which is neither core nor a border point is a noise point.

# Algorithm steps

① The algorithm starts with a random point in the dataset which has not been visited yet and its neighbouring points are identified based on the eps value.

② If the point contains ≥ points than min points then the cluster formation starts and this point becomes a _core point_; else its considered as _Noise_!

→ The point to note here is that a point initially classified as noise can later become a border point if its in the eps radius of a core point.

③ If the point is a _core point_, then all its neighbours become a part of cluster. If the points in the neighbourhood turn out to be core points then their neighbours are also part of cluster.

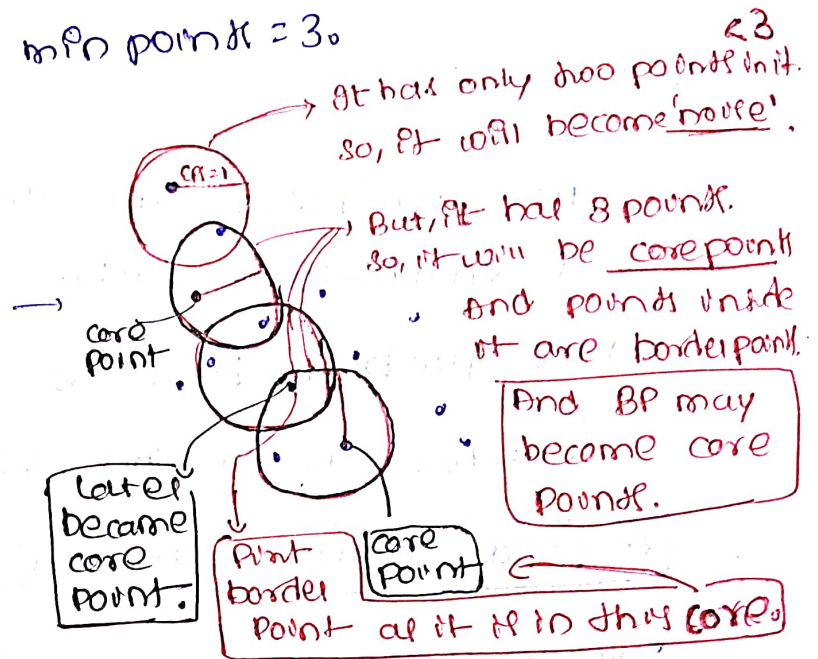④ Repeat the steps above until all points are classified into different cluster or noise.

this Algo works well if all clusters are dense enough, and they are well seperated by low dense regions.

<u>Eg:</u>

Let's say eps=1, min point = 3.



**It has only two points in it. so, It will become 'noise'.** <3

**But, It has 8 points. so, It will be core point And points inside it are border point.**

**And BP may become core points.**

Point border Point as it is in this core.

core point

Later became core point.

Point border Point | core Point

data.

→ In short, DBSCAN is a very simple yet powerful Algorithm, capable of identifying any noof clusters of any shape, it is robust to outliers, and it has just two hyper parameters. [ep & min samples]

→ However, if the density varies significantly across the clusters, it can be impossible for it to capture all the clusters properly. moreover its computational complexity is roughly $O(m \log m)$, making it pretty close to linear with regards to the no. of instances.

————————————→//