

DBSCAN

* Density Based Spatial Clustering of Application With Noise

* Why?

→ K-Mean specify no of clusters
DBSCAN automatically does this

→ DBSCAN can handle outlier automatically

→ DBSCAN can handle irregular shapes.

* Developed by Martin Ester and Peter Kriegl in 1996.

↳ Limitation of K-Means
↳ fix.

General

Same as K-Means
(Different in mathematical approach)

Mathematical

<u>Point</u>	<u>X_1</u>	<u>X_2</u>
A	1	1
B	2	2
C	2	3
D	8	8
E	8.5	8.5
F	7.5	8
G	25	80
H	24	81

Two factors

① Epsilon (ϵ)

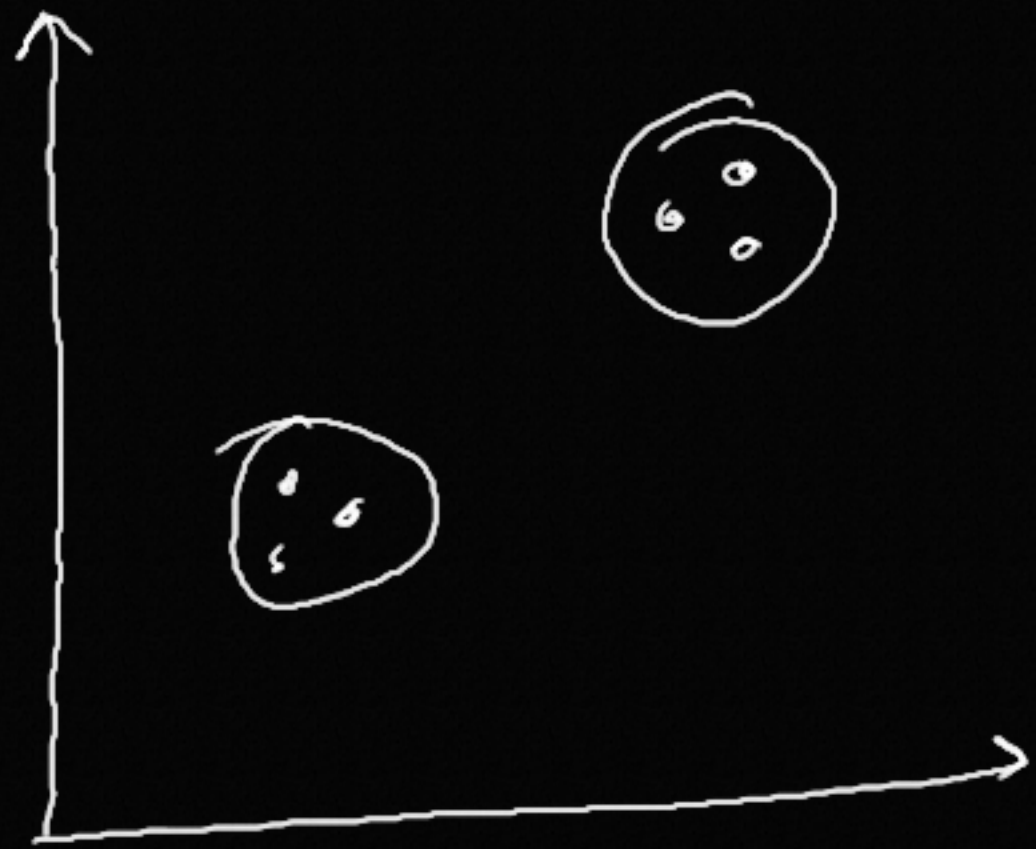
we set $\epsilon = 2$

$A \leftrightarrow B = 1.8 \rightarrow$ Neighbour

$A \leftrightarrow C = 1.5 \rightarrow$ "

$A \leftrightarrow D = 4 \rightarrow$ Not
neighbour

② MinPts → We set MinPts = 3 =
Min points reqd to make
a cluster



Classifying Criteria

- ① Core point \Rightarrow Point that satisfy MinPts condition
- ② Border point \Rightarrow Not a core point but is within ϵ distance of core point
- ③ Noise point \Rightarrow Neither core point nor border points

Working

$$\underline{\underline{\text{MinPts}}} = 3$$

$$\underline{\underline{\epsilon}} = 2$$

① Start with A

Calculate all points distance with A

$A \leftrightarrow H$

$A \leftrightarrow A, A \rightarrow B, A \leftrightarrow C, A \leftrightarrow D, \dots$

Using euclidean distance formula

Check which distance is less than ϵ

Neighbour points $\Rightarrow A, B \Rightarrow \underline{\underline{2}}$ (Not satisfied)

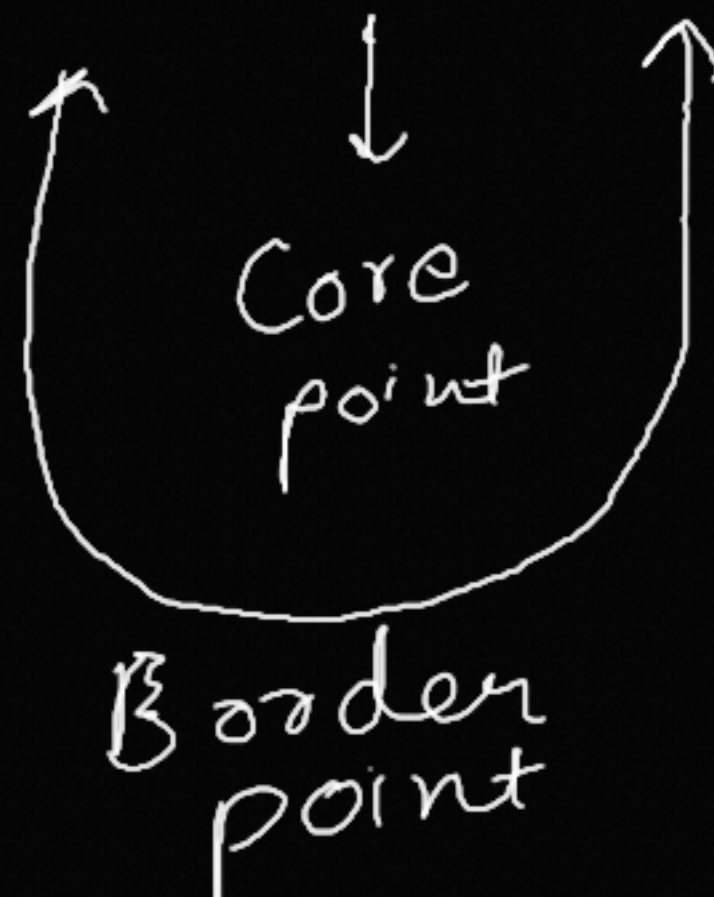
It is not a core point

② Now point B.

Neighbour points = A, B, C \Rightarrow 3 (satisfied) condition

(B is a core point)

Cluster 1 = { A, B, C }



③ Point D.
Neighbours = D, E, F \Rightarrow 3 (satisfied condition)

(D is a core point)
Cluster 2 = $\left\{ \underset{\substack{\downarrow \\ \text{core}}}{D}, \underset{\substack{\downarrow \\ \text{border point}}}{E, F} \right\}$

④ Remaining points = $\{G, H\} = \underline{\underline{2}} \rightarrow$ Never ever satisfy condition

(Noise point)

final clusters

cluster 2

Noise

cluster 1
 $\{A, B, C\}$

$\{D, E, F\}$

$\{G, H\}$

Model
trained

Prediction

New point = $(7.9, 7.9)$

Calculate all distances again

Neighbours = $\{D, E, F\}$
($k=2$)

Check majority of point belong to which cluster.

Output \Rightarrow Cluster 2

When to Use

- * Irregular cluster shape
- * Data has noise / outlier
- * No of cluster is unknown
- * Cluster are of varying sizes.

When Not to Use

- * Dataset is very large & high-dimensional (Computational Cost ↑)