

## \* Clustering

WORLD STAR™  
Page: 1  
it is not a L  
{clustering is Data Mining  
Algorithm}

→ clustering (group)

Enhance versioning  
→ K-means

→ K-means and K-means++

→ Mean shift algo

→ Hierarchical algo

→ DB Scan algo

→ Mini batch algo

we are going to create a group

Q) How to grouping the data?

→ Based on Distance.

→ Between Data point we are going to calculate distance

clustering → Grouping of data based on Distance → Distance can be many type (Distance formula)

## \* K-Means Algo

① Concept of the Centroid

② Distance

(i) Intra (within cluster)

(ii) Inter (Between cluster)

Evaluation of  $\rightarrow$  WCSS → elbow plot

of cluster

	$x_1$	$x_2$
Data point	Height	Weight

1 180 60

2 170 50

3 160 40

4 165 30

5 185 45

6 160 65

7 167 75

8 190 80

9 195 90

10 200 100

K-Means ++ (No. of clusters)  $\times$  (mean or Average)

conversion of mean

~~K-Means~~  $\rightarrow$  ~~Mean~~

It is all about centroid

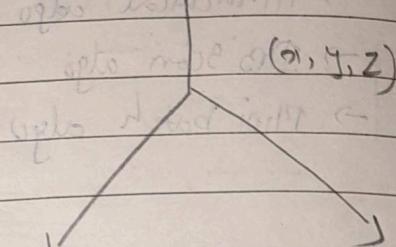
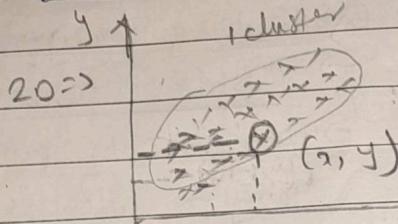
Centroid  $\rightarrow$  (0.8 - 0.01)  $\times$  (0.01)

→ for initializing the centroid  
the method is random method

① Calculate a centroid (data point)

→ Minimum = 1 cluster we can create

→ Maximum = No. of data points (No. of rows)



Minimum two  $\Rightarrow K = \text{no. of centroid} \because K=2$

1<sup>st</sup> centroid

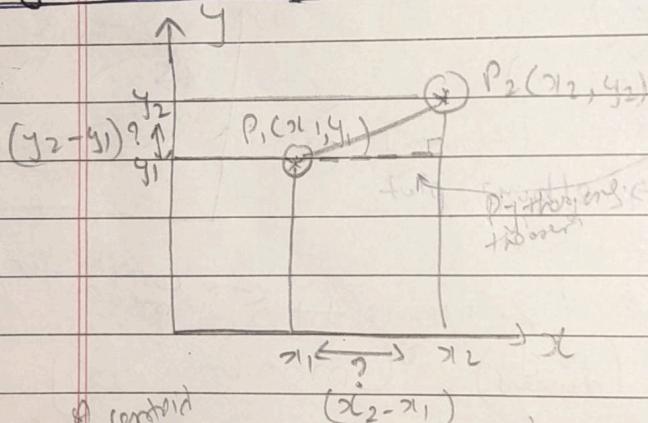
$C_1 (180, 60)$   
This is my 1<sup>st</sup> cluster

2<sup>nd</sup> centroid

$C_2 (170, 50)$  2 point

This is my 2<sup>nd</sup> cluster (1)

## Proof of Euclidean Distance



$$H^2 = B^2 + P^2$$

$$H^2 = (x_2 - x_1)^2 + (y_2 - y_1)^2$$

$$H = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

→  $C_1 (180, 60)$        $C_2 (170, 50)$       (from  $C_1$  and  $C_2$  we will calculate difference for the 3<sup>rd</sup> point)

3<sup>rd</sup> point  $\rightarrow (160, 40)$

①  $C_1 \text{ to } 3 \quad \sqrt{(160 - 180)^2 + (40 - 60)^2}$

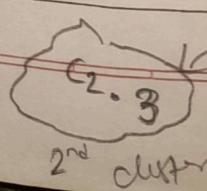
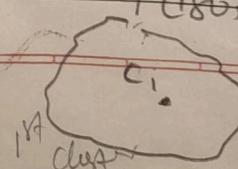
$$= 28.28$$

(distance is v.v. high)

$C_2 \text{ to } 3 \quad \sqrt{(160 - 170)^2 + (40 - 50)^2}$

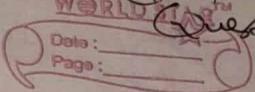
$$= 14.14$$

(less)  $\therefore C_2$  the 3<sup>rd</sup> point is v.v. near



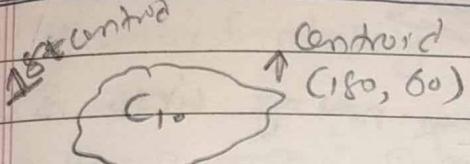
New point going over here. So, we have to update this cluster

# How we are maintaining centroid in k-means? (Interview Question)



Date: \_\_\_\_\_  
Page: \_\_\_\_\_

- 1) Initialize a centroid (Randomly)
- 2) Distance
- 3) fill the cluster
- 4) Update the centroid ← (mean conversion of centroid)



→ New point going here. C1 3rd point  
so we have to update this centroid

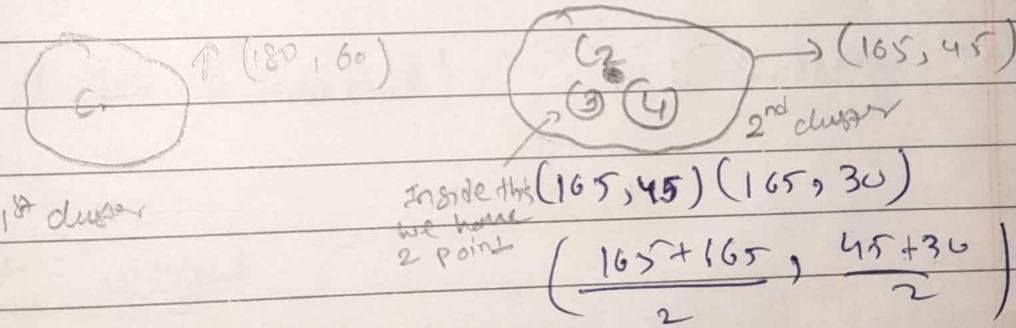
$$(170, 50) \quad (160, 40)$$

$$\left( \frac{170+160}{2}, \frac{50+40}{2} \right) \text{ AVG}$$

How centroid would be the average of the base centroid  
 1<sup>st</sup> centroid      2<sup>nd</sup> centroid      (165, 45) ← updated centroid  
 C<sub>1</sub>                  C<sub>2</sub>  
 (180, 60)            (165, 45)      4<sup>th</sup> (165, 30)

$$② C_1 \rightarrow 4 \sqrt{(180 - 165)^2 + (60 - 30)^2} \approx 33.5$$

$$C_2 \rightarrow 4 \sqrt{(165 - 115)^2 + (45 - 30)^2} = 15$$



Inside this (165, 45) (165, 30)

we have 2 point

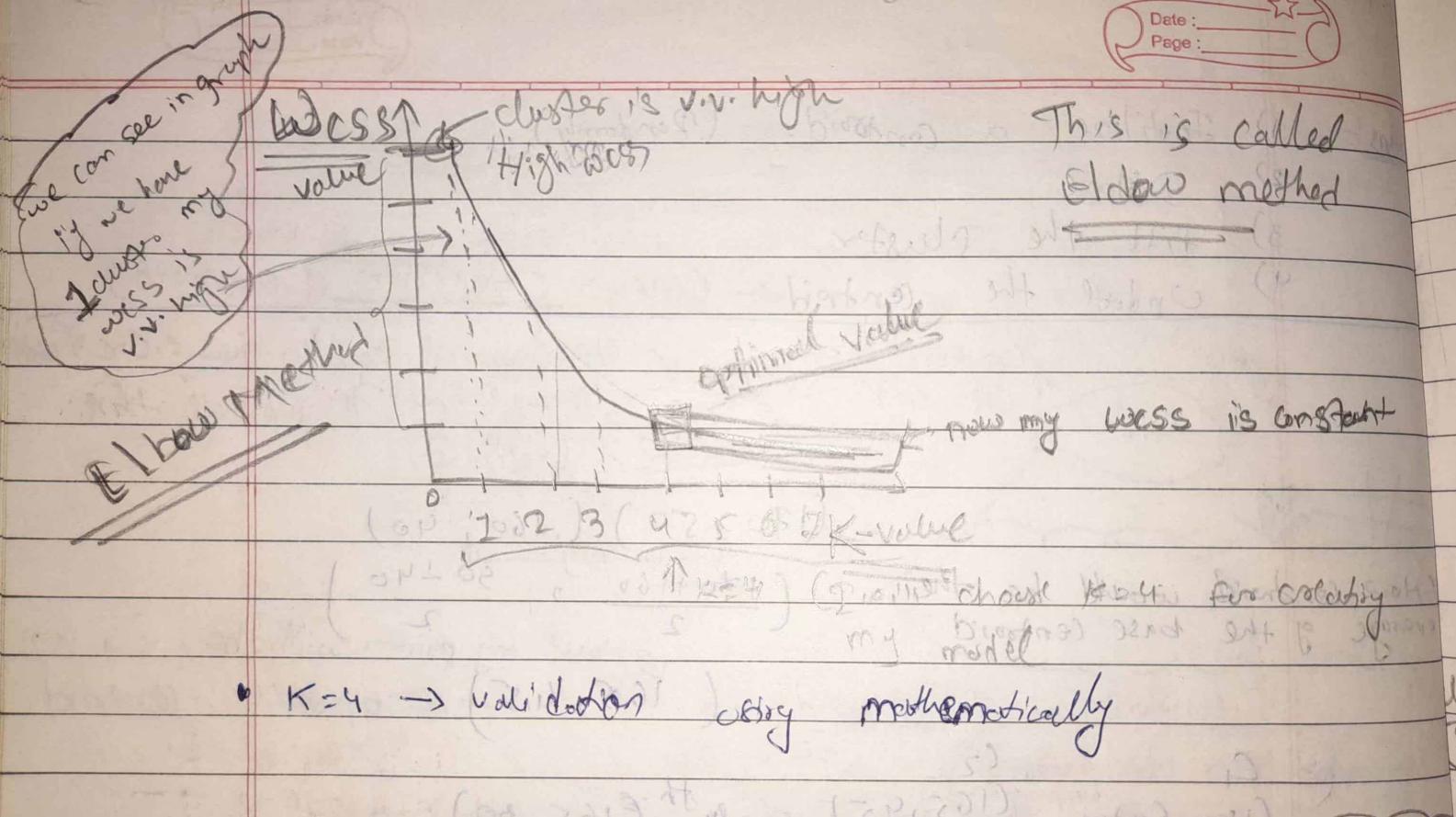
$$\left( \frac{165+165}{2}, \frac{45+30}{2} \right)$$

$$C_1 \in (180, 60) \quad C_2 \in (165, 37.5) \quad \underbrace{(165, 37.5)}_{\text{new centroid value of } C_2}$$

(Q) If my Intra-Distance is less, so my cluster is Good or Bad?

⇒ Good

Dense cluster ⇒ less Intra distance (within cluster distance) it is good cluster  
 • Greater or Big or more inter-distance is bad (Sparse cluster)



## \* Hierarchical Clustering

Hierarchy of cluster  $\rightarrow$  grouping of data  $\rightarrow$  no. of cluster = no. of data point

Hierarchical clustering  $\rightarrow$  Agg. (Bottom up)  $\rightarrow$  Div. (Top down)

(Initially we are going to consider single cluster)

$\Rightarrow$  Dendrogram : Diagram which shows its representing cluster

Agglomeration

single cluster

TOP

Entire this is my single cluster

Dir

Close to each other forming cluster

TOP

Agg.

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

2<sup>nd</sup> May

5 clusters to 1 cluster

\* Agglomerative  $\rightarrow$  Bottom to Top  $\rightarrow$  5 cluster  $\rightarrow$  1 cluster  
 Divisive  $\rightarrow$  Top to Down  $\rightarrow$  1 cluster  $\rightarrow$  5 clusters

(n = 5) without P2

Representation cluster:

Dendrogram

We can calculate distance from

these clusters

with the help of Euclidean Distance (Weight)

Growth the help of point matrix

With the help of point matrix we are going to find out cluster

we have to

calculate a distance from  $P_1$

1 point another point

Randomly

Rating values

$P_1 \quad P_2 \quad P_3 \quad P_4 \quad P_5$

$O_{P_1, P_1}$

Upper value

will be redundant

$O_{P_2, P_1}$

$O_{P_3, P_1}$

$O_{P_4, P_1}$

$O_{P_5, P_1}$

$O_{P_2, P_2}$

$O_{P_3, P_2}$

$O_{P_4, P_2}$

$O_{P_5, P_2}$

$O_{P_3, P_3}$

$O_{P_4, P_3}$

$O_{P_5, P_3}$

$O_{P_4, P_4}$

$O_{P_5, P_4}$

$O_{P_5, P_5}$

$P_5 \quad 11 \quad 6 \quad 10 \quad 13 \quad 8$

$P_4 \quad 6 \quad 10 \quad 13 \quad 8 \quad 0$

$P_3 \quad 13 \quad 8 \quad 0 \quad 11 \quad 6$

$P_2 \quad 8 \quad 0 \quad 11 \quad 6 \quad 13$

$P_1 \quad 6 \quad 13 \quad 11 \quad 8 \quad 0$

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

$P \rightarrow \text{Points}$

Step 1  $\rightarrow$  least distance in  $P_i$ -clm  $\rightarrow$  2

Diagonal will be zero

$(P_5, (P_1, P_3))$

$(P_1, P_3)$

$P_1 \quad P_3$

$P_5 \quad P_2 \quad P_4$

intensity  $\rightarrow$  high + low and low + high = high intensity

Step 2  $\Rightarrow$  Now, we are going to reduce this cluster  $(P_1, P_3)$

$(P_1, P_3) \quad P_2 \quad P_4 \quad P_5$

$(P_1, P_3) \quad 0 \quad \dots \quad \dots$

$P_2 \quad 7 \quad \dots \quad \dots$

$P_4 \quad 6 \quad \dots \quad \dots$

$P_5 \quad 3 \quad 5 \quad 6 \quad \dots$

Minimum 10 8 0

$(P_5) \leftrightarrow (P_4)$  } we can easily calculate this point?

But,  $(P_5) \leftrightarrow (P_1, P_3) = ?$  How we can calculate distance

$$\frac{d(P_5, (P_1, P_3))}{d((P_5, P_1), (P_5, P_3))} \leftarrow \begin{cases} \text{Single linkage Method} \\ \text{Method based on breadth} \end{cases}$$

$$[d(1, 3)] \leftarrow 3 \text{ is minimum distance}$$

$$d(P_2, (P_1, P_3))$$

$$d((P_2, P_1), (P_2, P_3))$$

$$[d(9, 7)] \leftarrow 7 \text{ is minimum, removes as redundant}$$

$$d(P_4, (P_1, P_3))$$

$$d((P_4, P_1), (P_4, P_3))$$

$$[d(6, 9)] \leftarrow 6 \text{ is minimum}$$

Now, least distance in  $(P_1, P_3)$  - cm  $\rightarrow [3]$

\* DBSCAN clustering [density-based] (partial clustering) of  
 (application with noise) Noise

→ It is a soft clustering method to incorporate noisy data outliers based on density.

• DBSCAN :- 1) Epsilon Distance

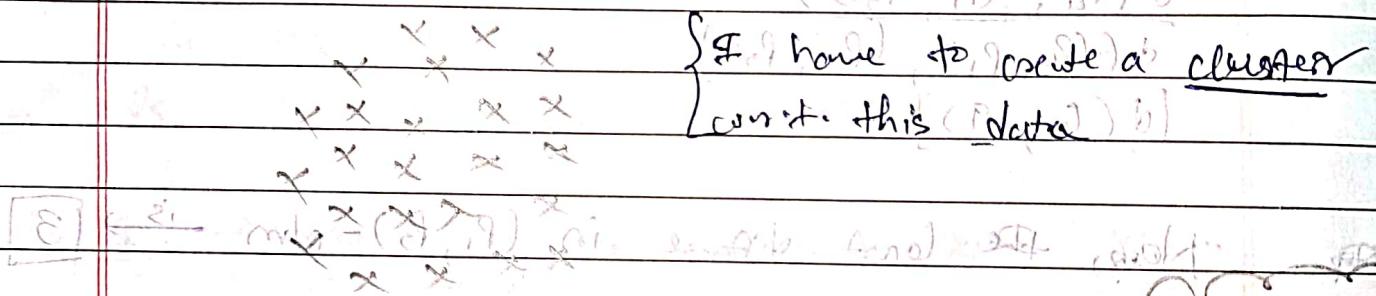
2) No. of points

3) Core point

4) Boundary point

5) Noise Point (i.e. outliers)

Explanation → whenever we are talking about clustering, these data has been given. (Drawing some randomly data point over here)



③ How we can create a cluster?

→ Over here we have some requirements

① Epsilon distance ?

② No. of points → These both are Hyperparameters tuning (we can select different epsilon values and no. of points)

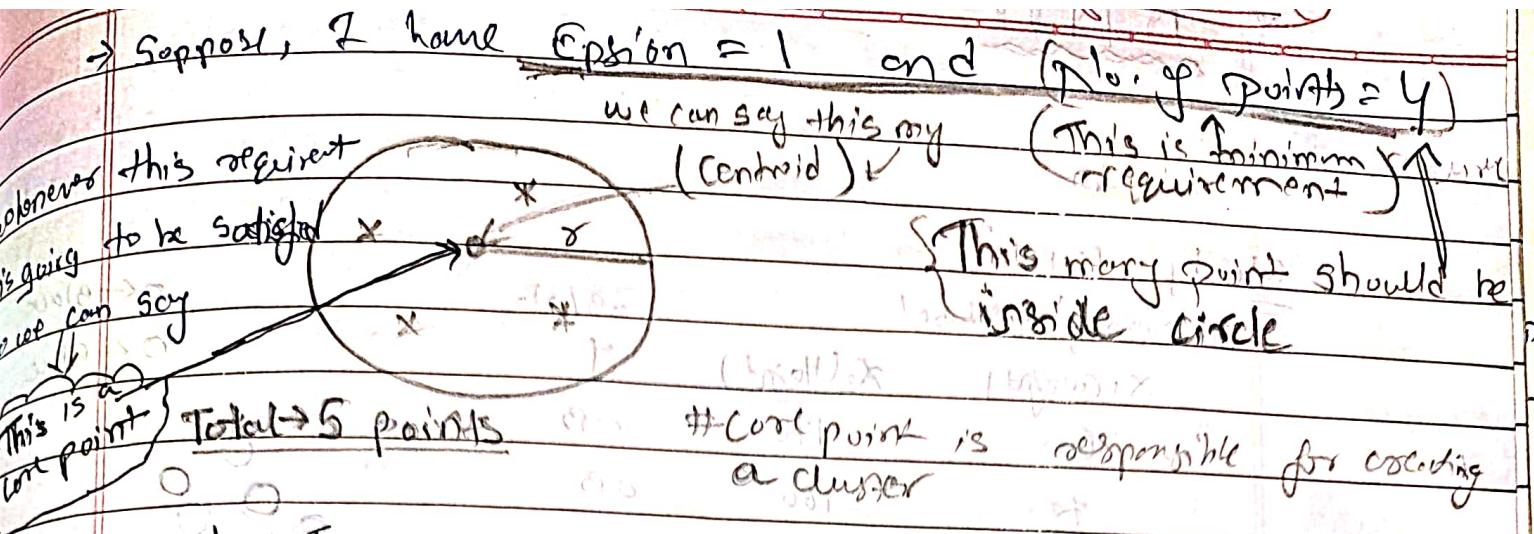
→ we can select Epsilon distance → {1, 2, 3, 4, 5...n}

→ we can select No. of points → {1, 2, 3, 4, 5...n}

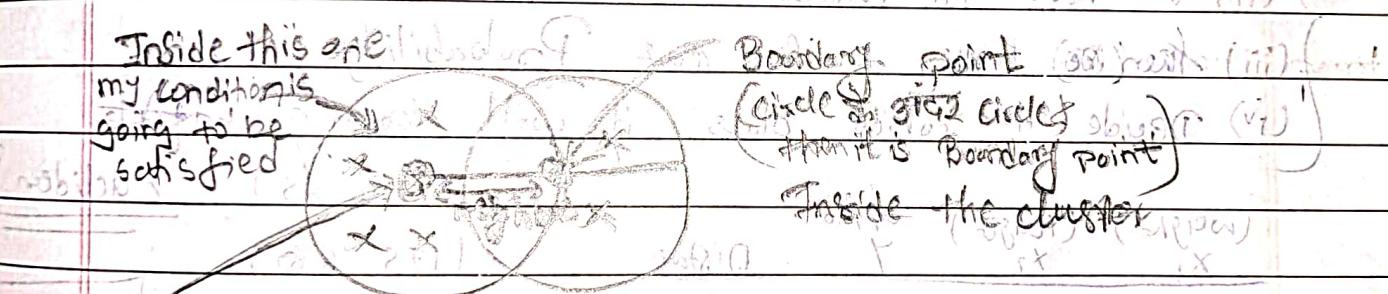
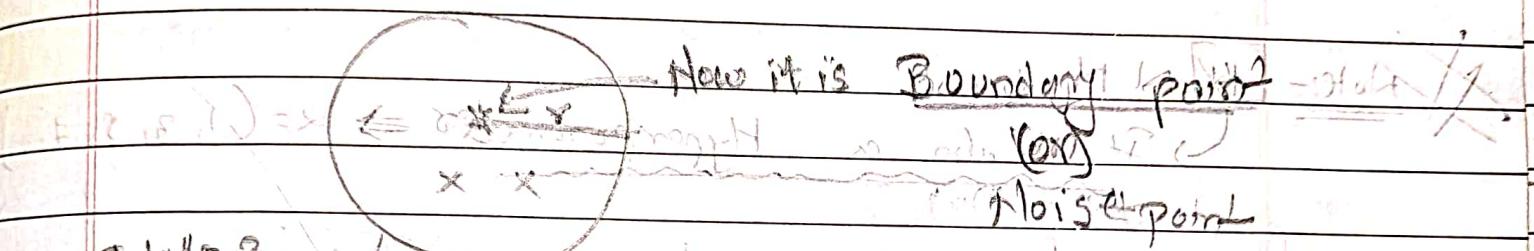
# Epsilon Distance is also called Radius

for Circle

Imp = # It is used to find out the outliers



Core point  $\rightarrow$  If we are taking any random point and we are taking  $\epsilon$  distance (i.e. radius) and we are drawing one circle, and it is fulfilling the requirement



This is clearly visible  
 this is my Core point  
 $(22, 21)$

Condition for Boundary Point :- The point is called as boundary point if that point comes inside the circle where core point is available. (or) If, if the point is neighbour of core point. then we can say it is Boundary point

