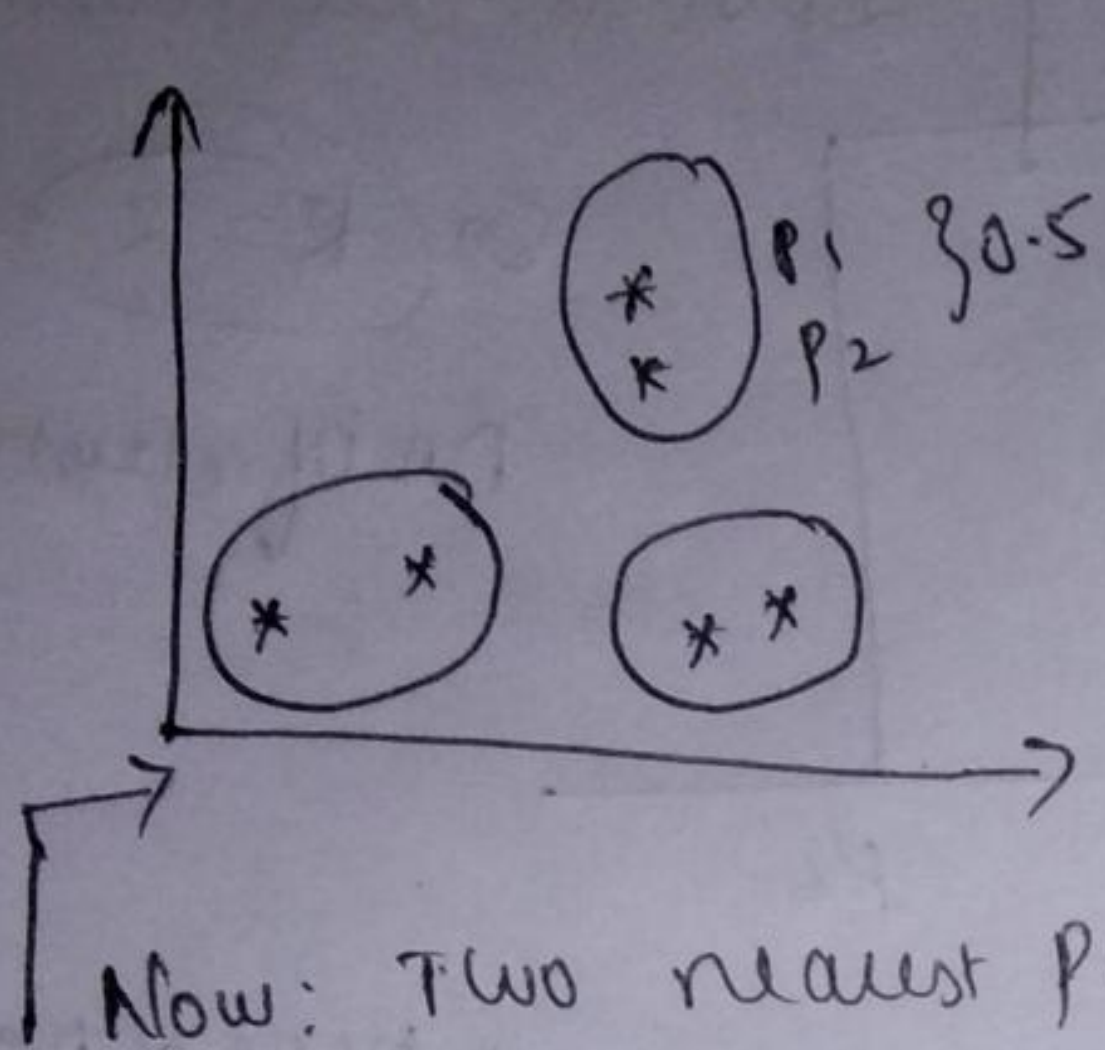


Hierarchical clustering

⇒ unsupervised

⇒ like K-means but technically different



In hierarchical cluster firstly considering all points are different cluster.

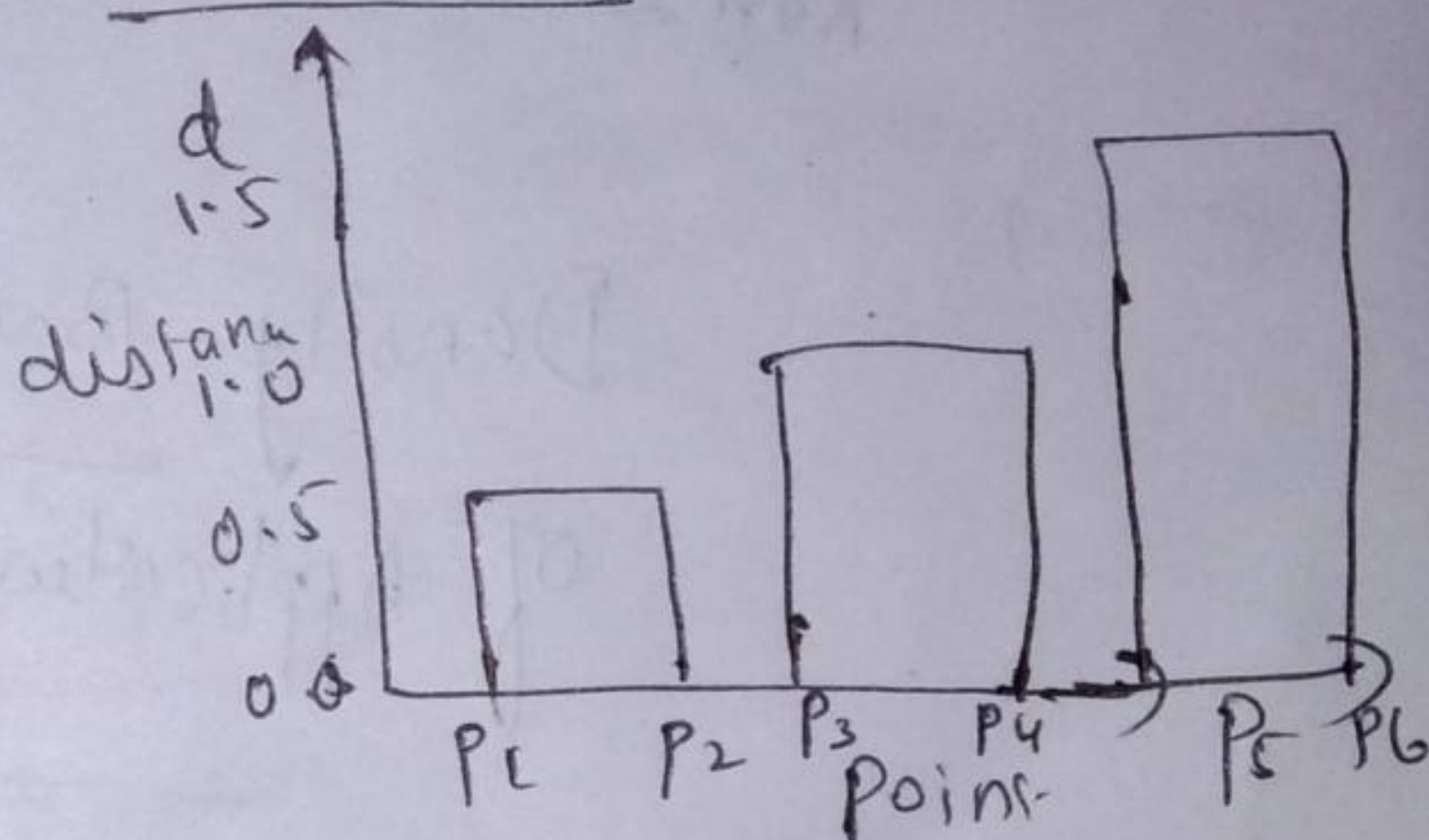
Now: Two nearest point are considered as same cluster.

this points are specify dendrogram

$P_1 \rightarrow P_2$ (dist: 0.5)

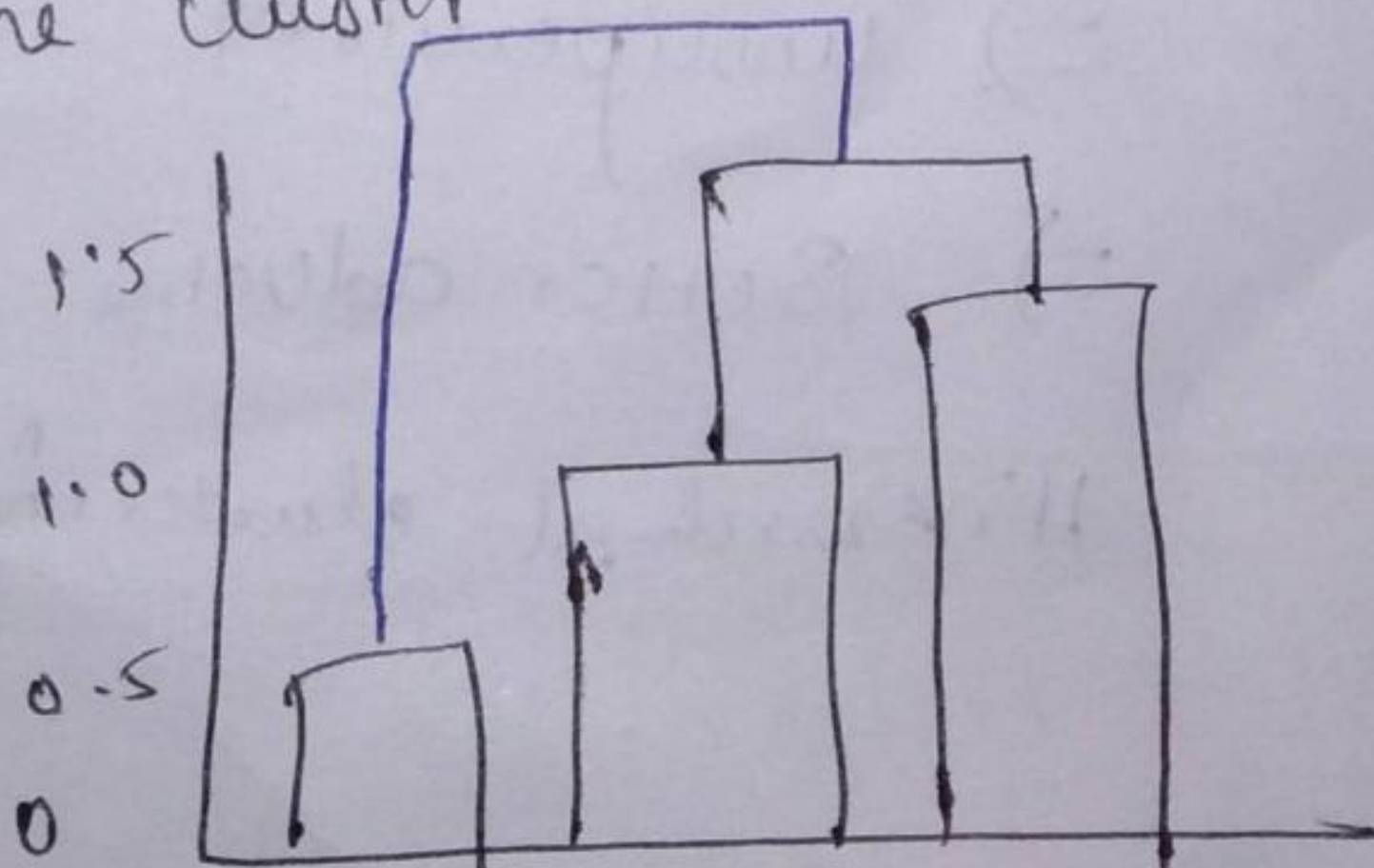
$P_3 \rightarrow P_4$ (1.0)

$P_5 \rightarrow P_6$ (1.5)



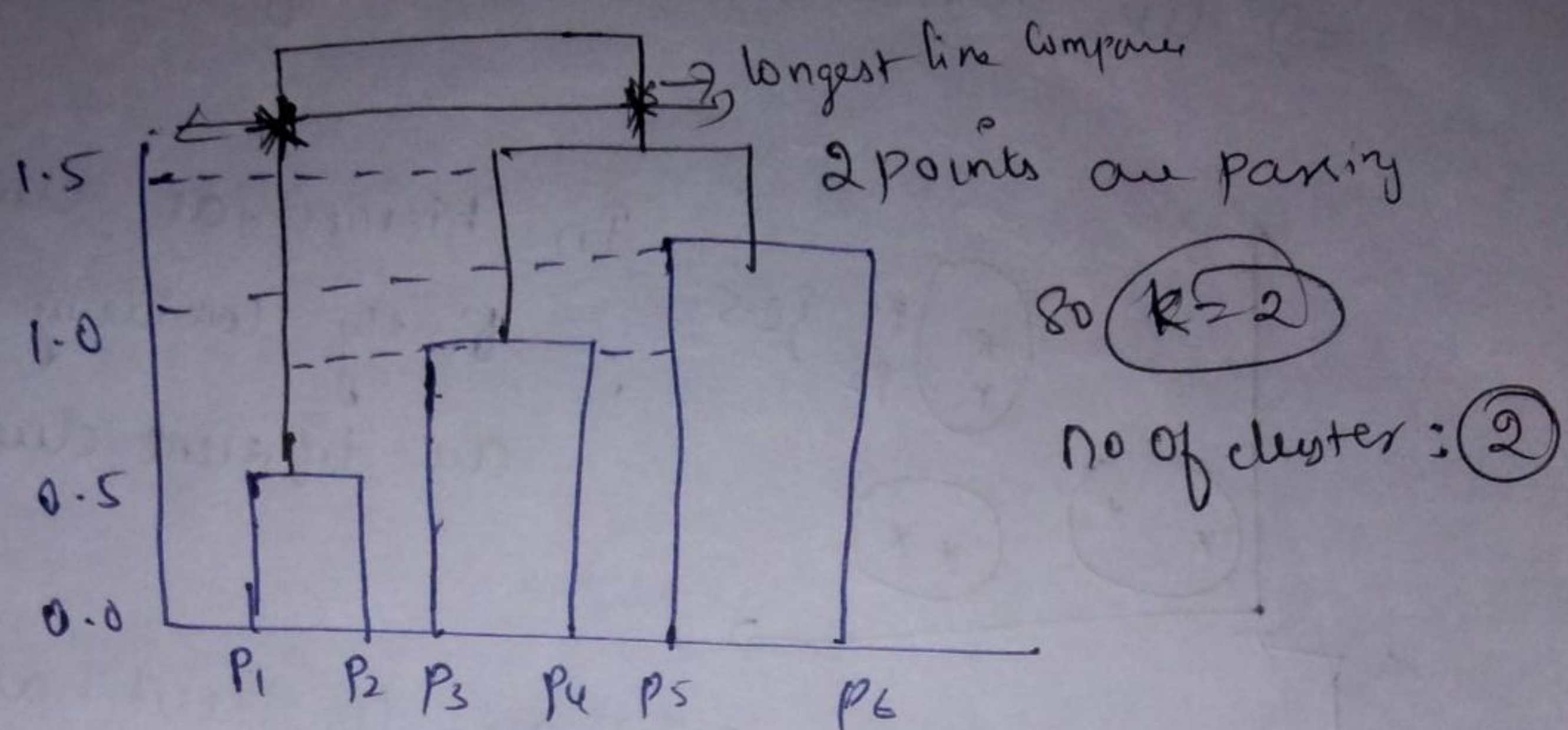
After that which one is next neighbour (nearest) in ~~that~~ that new cluster $\{P_3, P_4, P_5, P_6\}$ is nearest

So that is same cluster



Then $P_1, P_2, P_3, P_4, P_5, P_6 \rightarrow$ same cluster

Okay now how to select k -values
(no of cluster) that we want



we need to find longest vertical line of the horizontal

Density Based spatial clustering
of Application with Noise

DBSCAN:

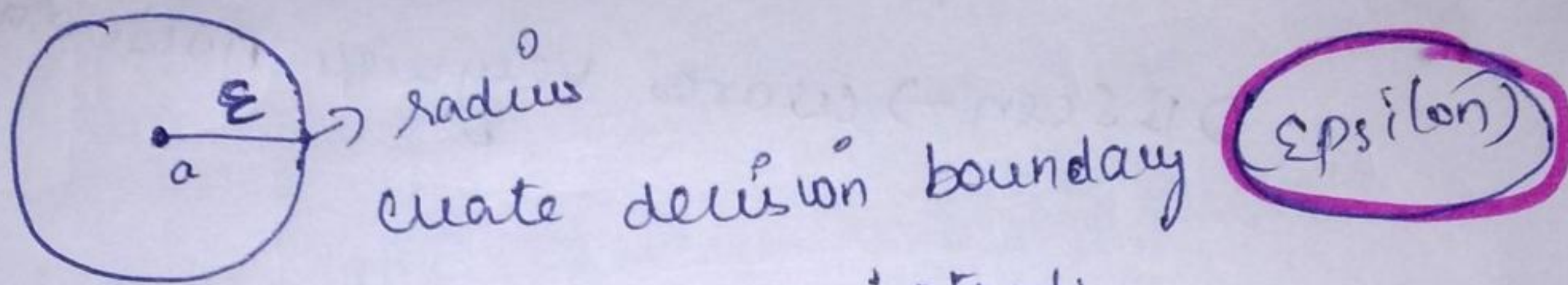
⇒ unsupervised

⇒ Better advance of k -means,

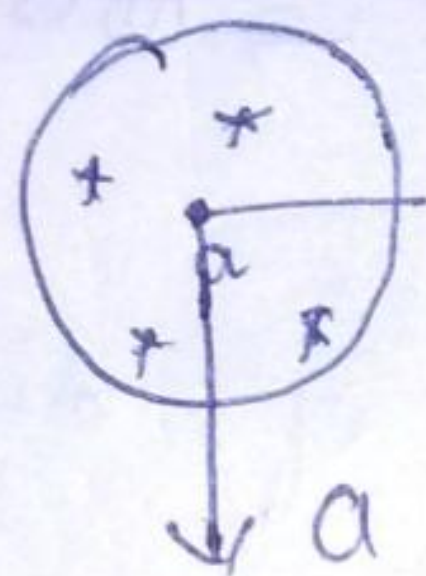
Hierarchical clustering.

4 important concepts:

- ϵ psilon
 - Core point
 - Border points
 - Noise
- } create cluster.



Now assume $minpts = 4$.



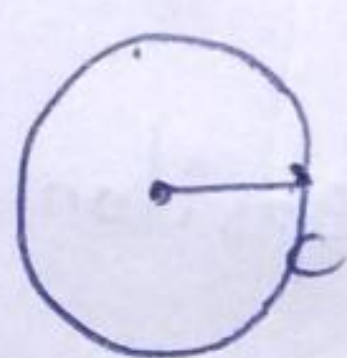
midpoints present in boundary

core point

when core point achieve

at least 4 (or) more points are
in decision boundary then decision boundary
create.

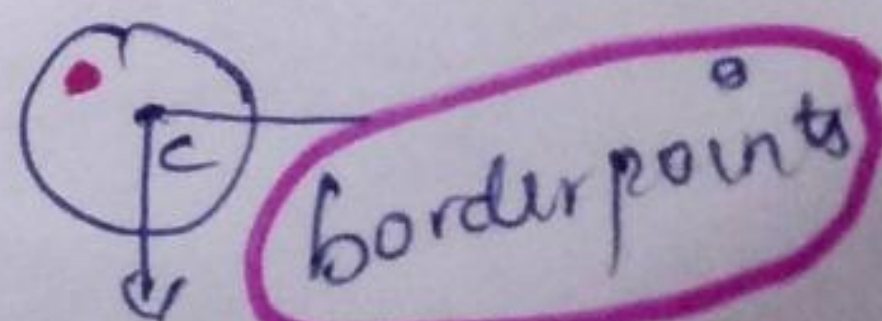
I have point c ($minpts = 4$)



create radius

① create boundaries keep ϵ psilon

② Suppose doesn't satisfy $minpts = 4$
but we have at least one core points (boundary
in boundaries



One point & I want boundary help of
Epsilon:

None of them satisfy

minpoints &
any core points

doesn't meet mid points

noise
(outliers)

DBSCAN \rightarrow works very well noise data

Pros:

work very well \rightarrow noise/outliers

drawback:

dimensionality of difficult to

grouping data together.

too many dimensions dbscan suffer.

doesn't work well when dealing
with cluster of varying densities.

Hyperparameter:

eps = epsilon (radius)

min_samples = mid points

metric = 'Euclidean', 'Manhattan'

Principal Component Analysis:

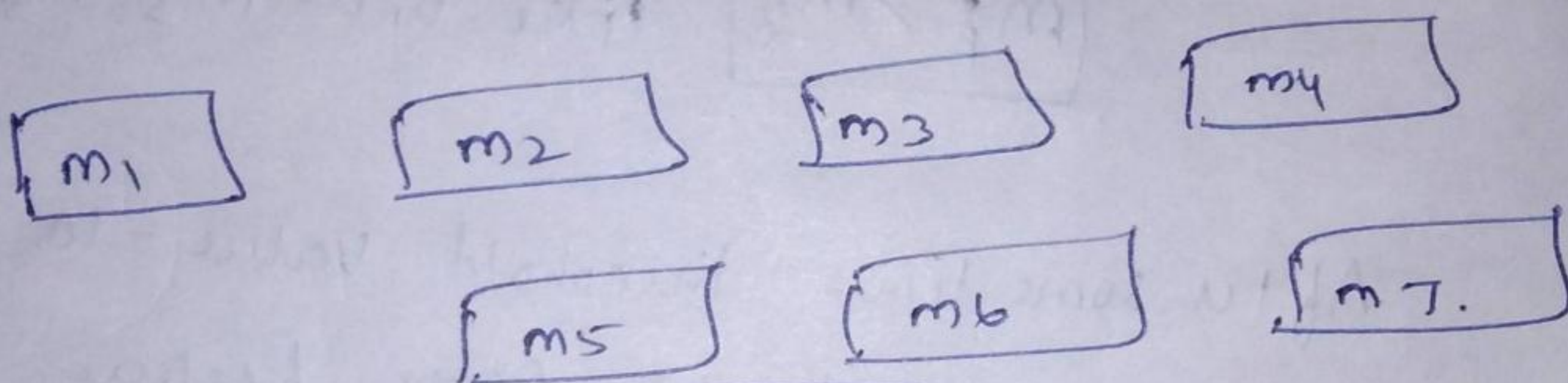
Curse of dimensionality:

↳ dimension (features) (i/p, dep)
↳ attributes

I am going to take dataset that dataset

$m_1 \rightarrow 2$ feature, $m_2 \rightarrow 5$ features, $m_3 \rightarrow 100$ features

$m_4 \rightarrow 1000$, $m_5 \rightarrow 2000$, $m_6 \rightarrow 10000$, $m_7 \rightarrow 20000$



Now I taken House sales price prediction

firstly I take m_1 dataset In m_1 dataset

have 2 features

one is bedroom } i/p
another is size. }

price \rightarrow dependent.

Now m_1 accuracy acc = some value
(80% - 85%)

Then take m_2 dataset In m_2 dataset
have 5 features

$m_2.csv$

independent

Size	State	bed	Parking	No of Rooms	Price

dependent (O/P)

Now predict rate use independent (5 features)
So get learn lot of details through 5 features
Compare than 2 features.

So m_2 model accuracy is high

$m_1 > m_2$ like with m_3, m_4, \dots

After some times threshold value = 10.
Assume threshold = 10 After feature reach the
threshold what happens.

$m_4, m_5 \rightarrow$ more than 10 features.
If more than 10 features your feeding
in lot of information in particular so that
model get confused.

machine has confuse model

if it is increasing exponentially

After threshold
accuracy ↓.

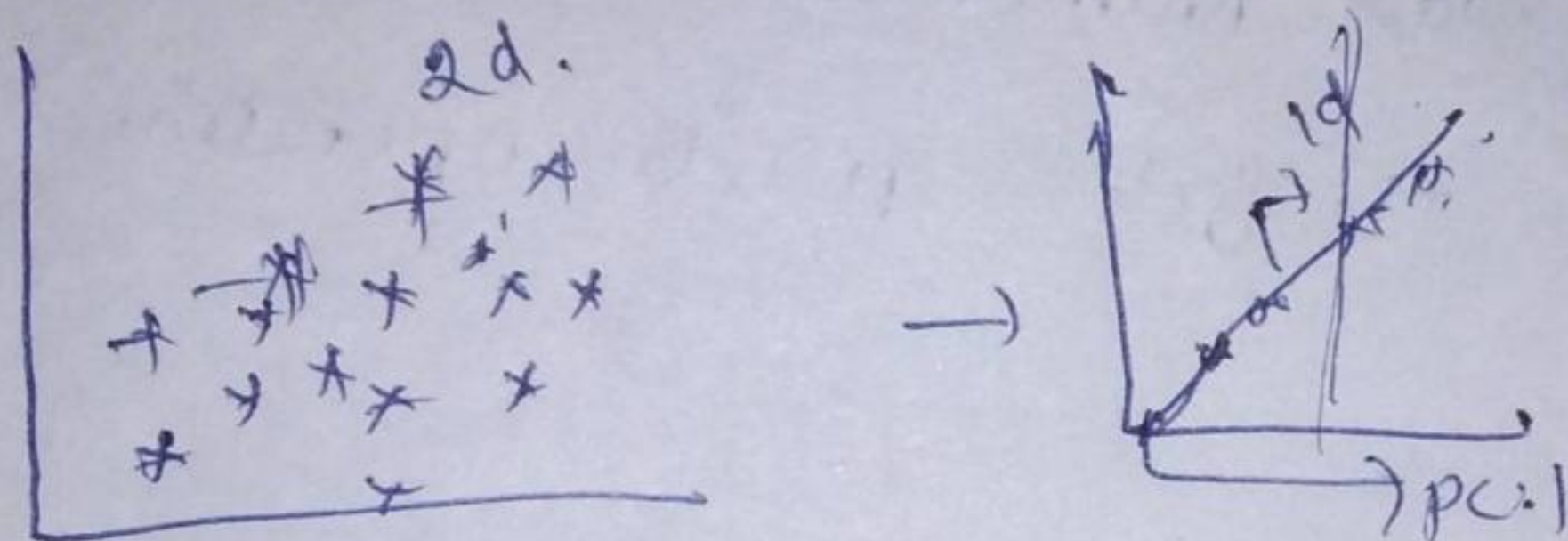
Principal Component Analysis:

Reduce number of dimension

As the number of dimension increase it is a curse.

Because accuracy impact with no. of dimension increase.

2d \rightarrow 1d



- ① best fit line calculate
- ② square error estimate \rightarrow
- ③ Perpendicular (orthogonal) line.

