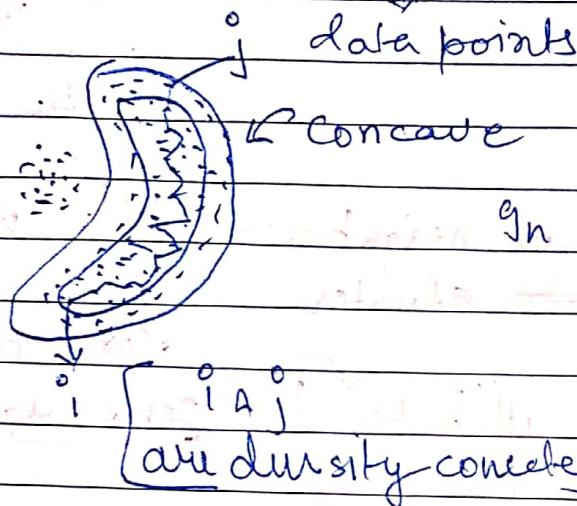


# Density Based Clustering

DBSCAN: Density based Spatial Clustering  
Application and Noise



No way to Classify Outlier.

→ every data is defined in the following way

- 1) Core
- 2) border
- 3) outlier

2 parameters

- 1)  $M_p$ : points in neighbour pts
- 2)  $\epsilon$ : neighbour radius



# of pts in radius  $\geq M_p$

Classify it as core points.

border

If no. of points  $\leq M_p$  - Then is border

Necessary cond<sup>n</sup>

border

Density.

Sufficient cond<sup>n</sup>: reachable from any core points and  $\# \text{pts} \leq M_p$ .

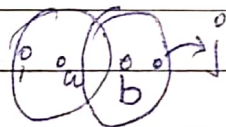
Outlier

Outlier: pts which are not core or border are outlier

for parameter choosing

$\epsilon$  should as small as possible  
 $M_p$  should be as large as possible.  
 for high Density.

Density connected



$i$  &  $j$  are not density reachable to each other

but  $i, j$  are density reachable to some core pts so  $i, j$  are density connected

Density connected is transitive property

$i \rightarrow j$      $j \rightarrow k$

so  $i \rightarrow k$  are density connected



## Clustering Process

→ Randomly choose a data pts which is not part of any cluster or outlier

→ Find out if data pts is core pts.

→ ~~Ad Start~~

if pt is core pt:

~~pt~~ start new cluster

check for neighbours of new cluster.

Keep on doing until u don't have any pt in neighbour that is not classified as core or border (part of cluster)

## ~~Hierarchical~~

Extra

### ADVANTAGES

→ Doesn't require a-priori specification of number of cluster.

→

→ Able to identify noise data while clustering

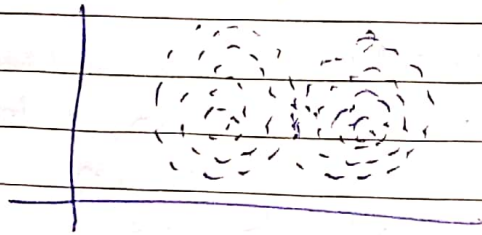
→ find arbitrarily size and arbitrarily shaped cluster

### ~~Disadvantage~~



## DBSCAN Disadvantage.

- 1) fails in case of varying density clusters
- 2) fails in case of neck type of Data.



DBSCAN fails to identify 2 clusters

- 3) Doesn't work well in case of high Dimensional data.

Q

## Nearest- SNN (Shared Neighbour)

SNN Tried clip the on the disadvantages to DBSCAN

→ For higher dimension SNN introduced some new notation for distance & density

Major points of SNN.

Distance

- euclidean distance is not good choice of higher dimension
- use different similarity measure in terms of KNN.



Page No. \_\_\_\_\_  
Date \_\_\_\_\_

→ Define Similarities in terms of new distance.

### Jarvis-Patrick algorithm

Step 1: SNN sparsification.

- Construct an SNN Graph from data matrix
- if  $p \neq q$  have each other in kNN list
- then create link between them

Step 2: Weighting

Weight the links with  $\text{sim}(p, q) = \frac{1}{|NN(p) \cup NN(q)|}$   
Where  $NN(p) \neq NN(q)$  are k neighbours of  $p \neq q$  resp

Step 3: Filtering

- then filter the edges
- remove all edges with weight less than some threshold

Step 4: clusters.

lets all connected component be cluster.

### Density

→ In euclidian space, density is no. of points per unit volume.

→ but as dimension increases volume decrease rapidly  $\Rightarrow$  density tends 0 for same  $M_p$ .

## New Concept

if KNN is close, then region is most likely to be high density

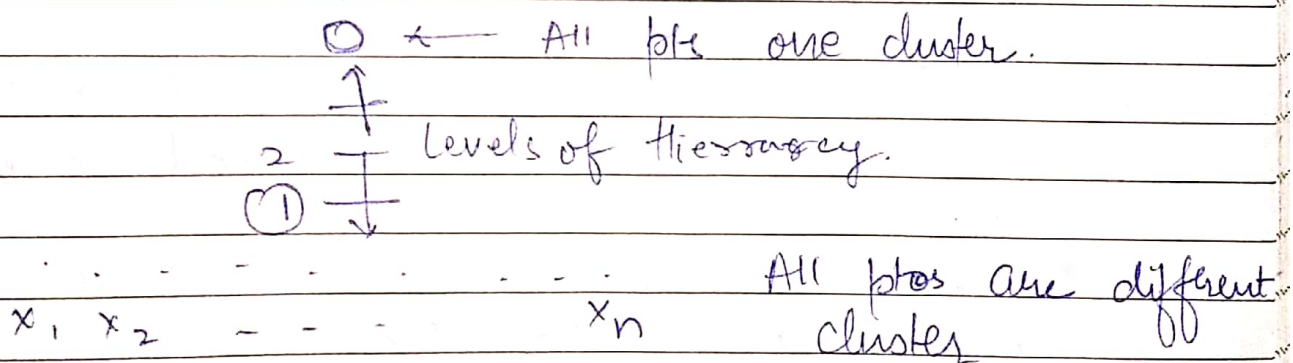
→ So distance to Kth neighbour gives a measure of density of a pt

~~more detailed measure of density~~

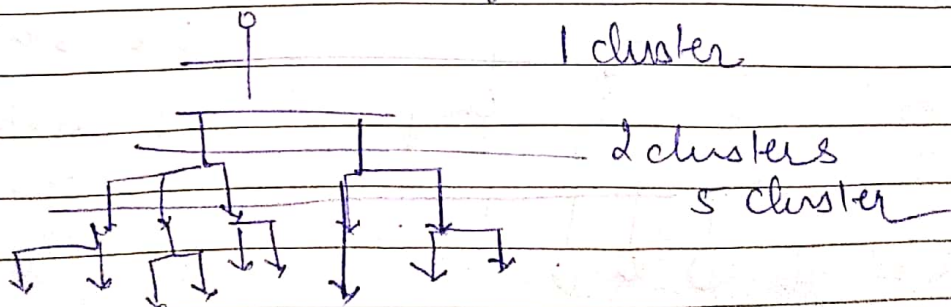
## Hierarchical Clustering

- 1) Bottom up / Agglomerative
- 2) top-down / Divisive

levels are defined



Tree structure / Dendrogram.





Eg 500 cities name of nation

○ → All belong to one nation

○ ○ ○ → National Capital

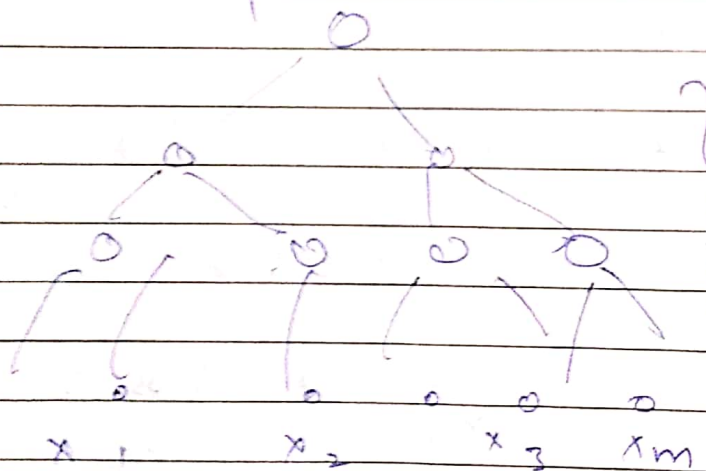
○ ○ ○ ○ ○ → State Capital

○ ○ ○ ○ ○ ○ → districts

○ ○ ○ ○ ○ → 500 city name

~~Speech data having multiple speaker.~~  
~~also eg. of hic.~~

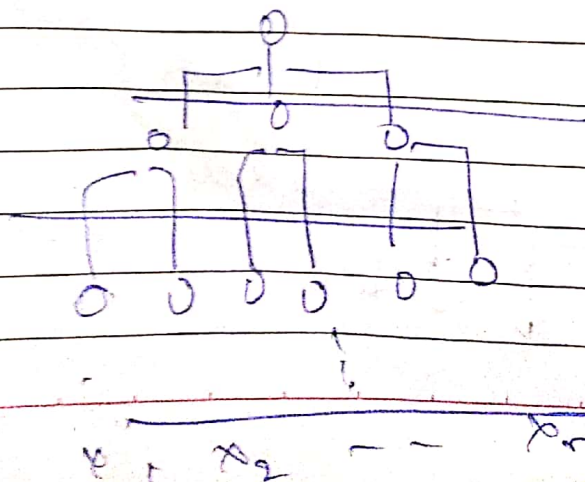
Bottom up.



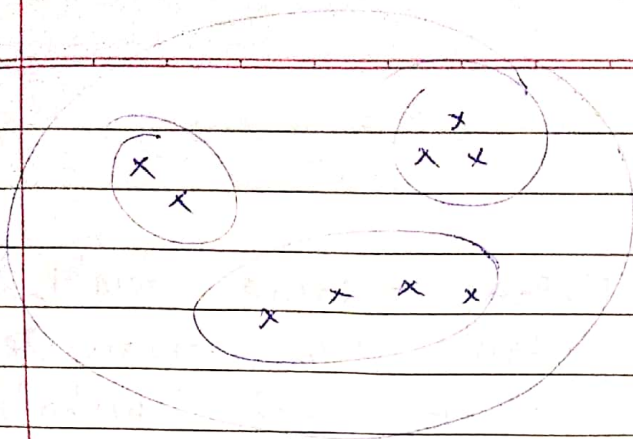
we keep merging until we get one.

We all feature use clusters.

top-down



We start with one & keep on dividing until we have no. of cluster = data set pt.

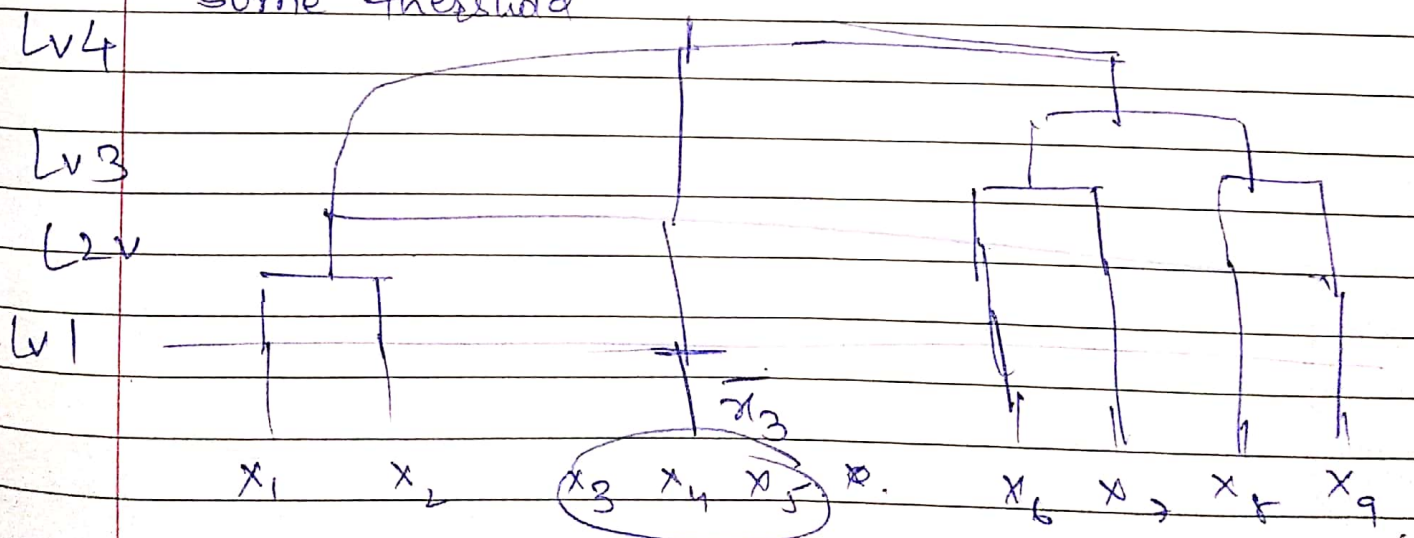


which feature vector are closest of each other

	$x_1$	$x_2$	$x_3$	---	$x_9$
$x_1$	0				
$x_2$		0			
$x_3$			0		
$\vdots$					
$x_9$					0

We have this distance matrix

Now we will find which are closest & less than some threshold



Now we find New distance matrix with  $\bar{x}_3$

by euclidian distance the clustering is done



for top-down

we say they don't belong to same cluster if at some level if distance b/w mean vectors are greater than that level threshold

we can say it is reverse of bottom-up

