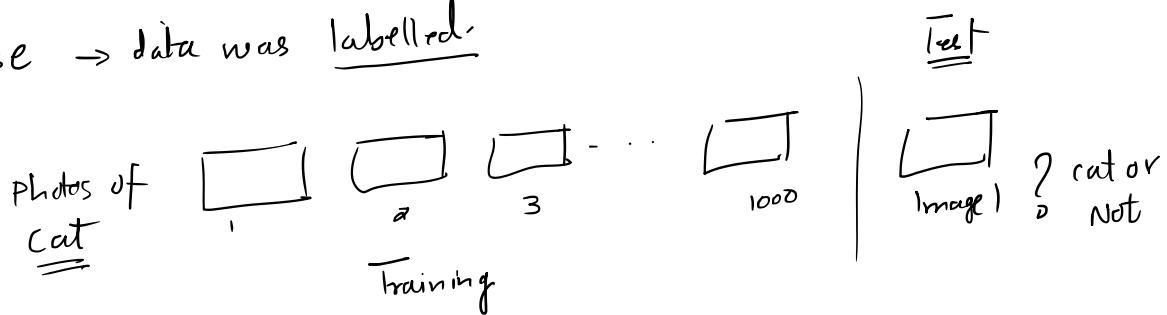


Module 5

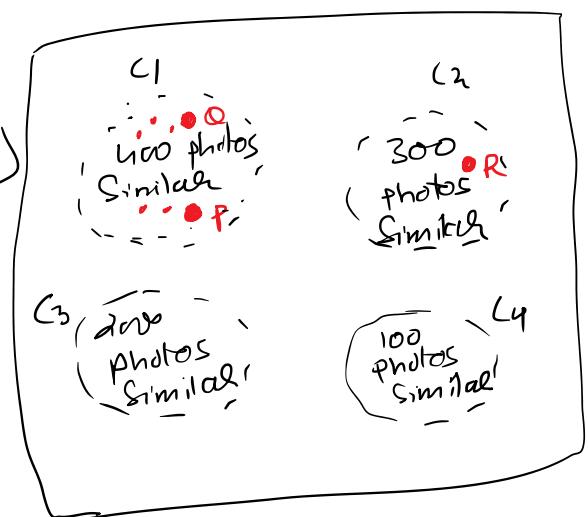
Clustering →

In earlier case → data was labelled.



Supervised learning

Now consider → 1000 photos (not labelled)



Clustering is the task of dividing the population or data points into no of groups such that the data points in the same group are more similar to other data points in same group than those in other group.

* Aim is to Segregate groups with Similar traits/ features and assign them into clusters.

Distance Metric : It is a parameter that tells how close two points are

Given element a, b, c in a set, a distance metric is defined as a function with following properties:

(1) Non-negativity $\rightarrow d(a, b) \geq 0$

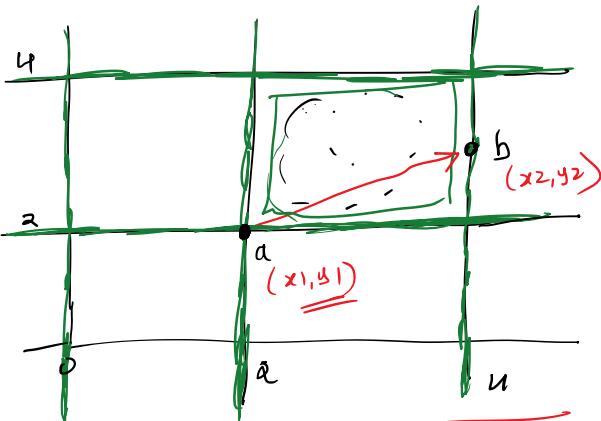
(2) Symmetry $\rightarrow d(a, b) = d(b, a)$

(3) Triangle Inequality $d(a, c) \leq d(a, b) + d(b, c)$. If $a, b, c \in S$

Different Distance Metric

(1) Euclidean Distance \rightarrow Most commonly used

\rightarrow When data is dense or continuous, this is the best proximity measure.



$$\text{Euclidean Distance (vector)} = d(a, b) = \sqrt{\sum_i (a_i - b_i)^2}$$

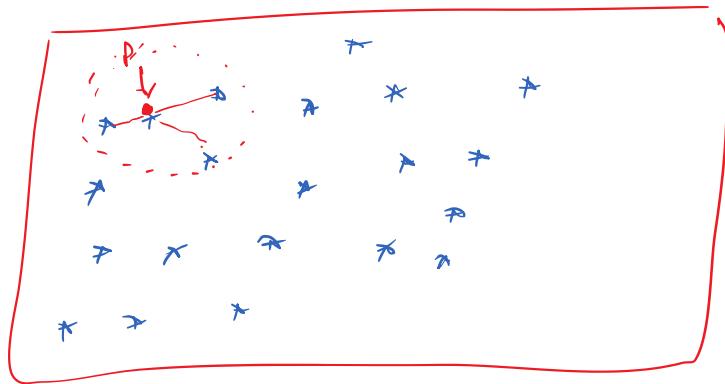
$$= \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

If for a point $P = (p_1, p_2, \dots, p_n)$ and $q = (q_1, q_2, \dots, q_n)$.

$$d(p, q) = d(q, p) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

$$= \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

$$= \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

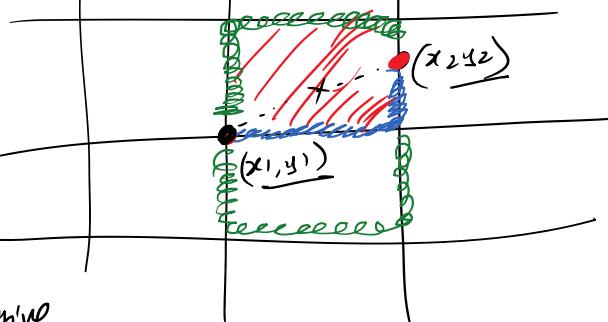


Manhattan Distance

* If we are trying to find the distance betⁿ two buildings that are separated by several blocks. We can only walk through sidewalk. that are parallel and cannot walk diagonally through buildings

In such case
we cannot use
Euclidean
distance.

Here Euclidean
distance will not give
realistic estimate for the distance



Manhattan distance -

$$= |x_1 - x_2| + |y_1 - y_2|$$

* Manhattan Distance is a metric in which the distance between two points is the sum of absolute difference of their Cartesian coordinates.

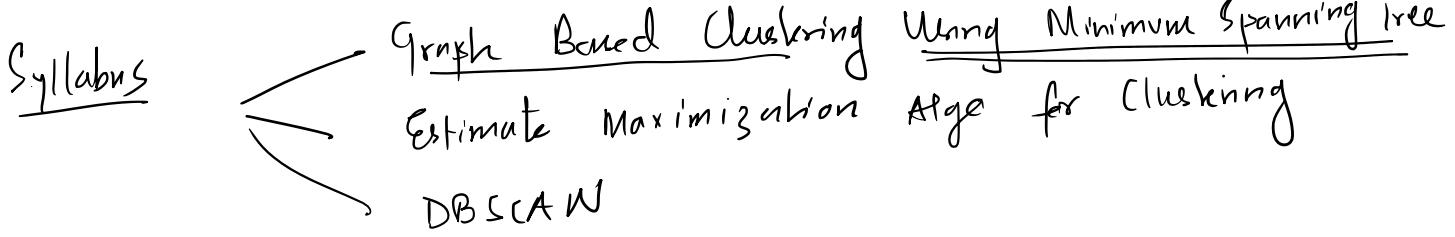
* Total sum of difference betⁿ the x coordinate & y coordinate.

* Manhattan distance metric is also known as Manhattan length, Rectilinear distance, L1 distance, L1 Norm, City block distance, taxi cab metric.,

* If the data points has n features -

Manhattan Distance

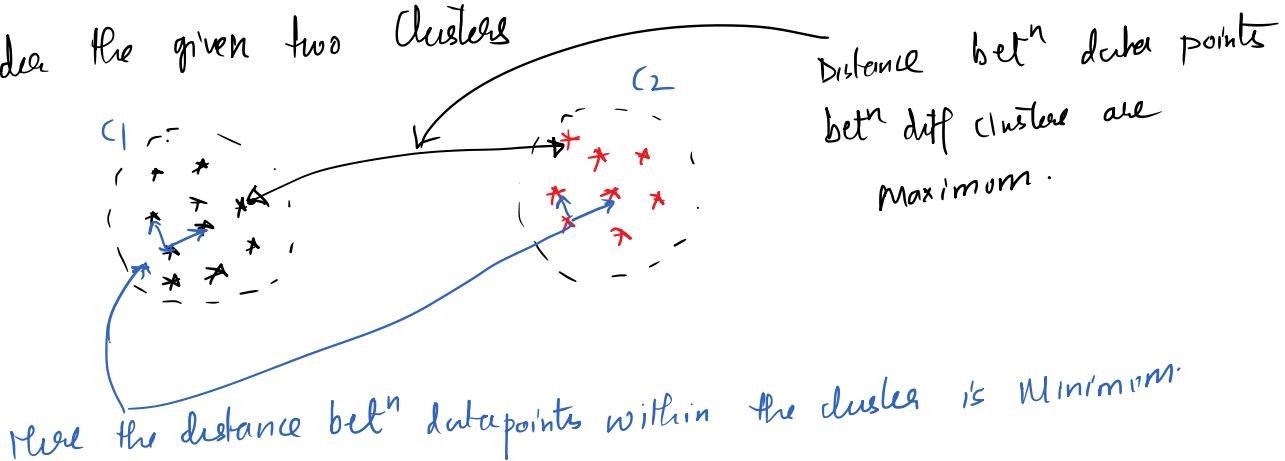
$$d(a, b) = |a_1 - b_1| + |a_2 - b_2| + |a_3 - b_3| + \dots + |a_n - b_n|$$
$$= \sum_{i=1}^n |a_i - b_i|$$



① Graph Based Clustering →

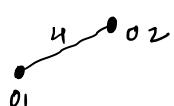
- (1) First Consider How to Construct Graph of given data.
- k Neighbourhood graph ✓
- ↳ ε Neighbourhood Graph ✓

Consider the given two clusters

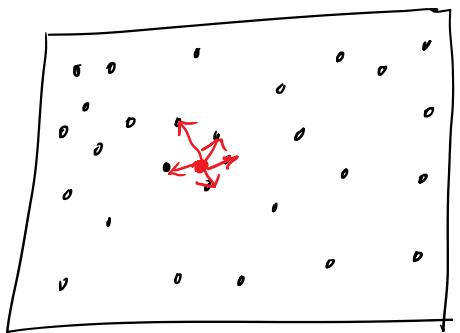


* In Graph Based Clustering the objects are represented as nodes in graph.

* The weight of edge/branch betⁿ the objects
defines the distance betⁿ the objects



① k Neighbourhood Approach



* To find distance of a point with all other points will be computationally challenging.

* we will consider near~~to~~ k points

Let k=5.

So we will calculate distance of
nearest 5 data points and will be

$n = 1000$

So we will calculate distance !
 nearest 5 data points are will be
 labelled as Nearest Neighbours of node

Challenge → choosing K value is challenging task
 in graph construction.

(2) Epsilon Neighbourhood Graph → Objects within a radius of epsilon from
 a given object are considered nearest
neighbours.



Here P_1 & P_2 are neighbours of P_1 .

Challenge → Selecting a proper epsilon value.

* Graph Based Clustering Using Minimum Spanning Tree →

Alg

1. Determine MST of given graph. ✓
2. Delete branch iteratively. ✓
3. Each connected component is a cluster ✓

Different strategies to Delete Branches.

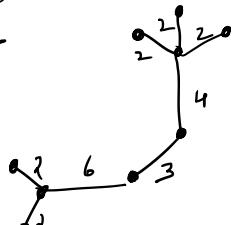
- ✓ (1) Delete Branch with Max Weight → ✓
- ✓ (2) Delete Inconsistent Branch → ✓
- ✓ (3) Delete by Analysis of weight.

(1) Deletion of Branch with Max Weight

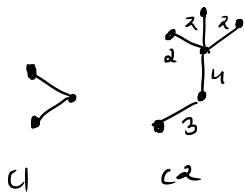
..... also results to clusters by deleting the branch with

- (1) Deletion of Branches
- In each step create two clusters by deleting the branch with max weight
 - Repeat until the desired no of clusters is reached
no of cluster desired is 3

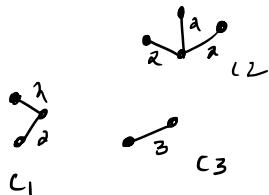
Consider MST



Step 1 \Rightarrow delete edge with weight 6



Step 2 \Rightarrow delete edge with weight 4.

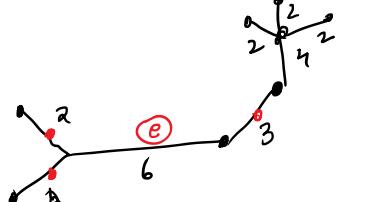


3 clusters

(2) Delete Inconsistent Branches-

- * A branch is inconsistent if the corresponding weight (\underline{de}) is much larger than the reference value \bar{de}
- * The reference value \bar{de} can be defined by the average weight of all the branches adjacent to edge e

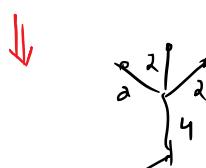
Consider MST

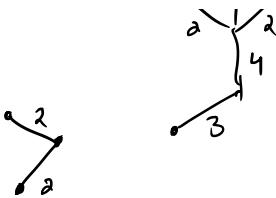


$$\text{Here } \bar{de} = \frac{2+2+3}{3} = \frac{7}{3} = 2.\underline{3}$$

Here $de = 6 > \bar{de}$

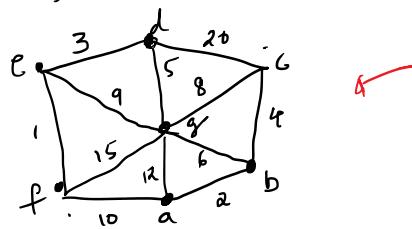
So delete the edge e ,



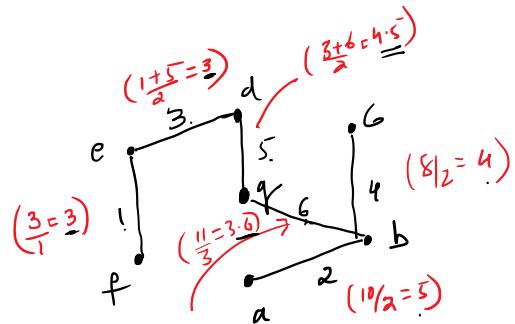


2 cluster

Q) Construct MST and provide clustering of graph Using Inconsistent branches.



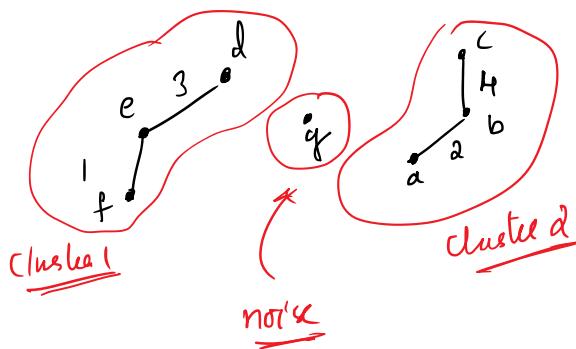
Step 1) MST of the Graph



Step 2) For deletion we are using Inconsistent branch approach.

Here edge with $5 >$ reference value (4.5) so delete edge 5

Also edge with $6 >$ reference value (3.6) so delete edge 6 .



② Estimation Maximization Algo

→ This Algo is for maximum likelihood Estimation in presence of latent variable (missing)

→ The most discussed application of EM Algorithm is used for Clustering with Mixture Model. (where data points have missing value)

* E-step → Estimate are expectations for machine & then classify the data into some classes.

* M-step → Whatever we estimated should now be maximized.

* Iteratively repeat E-step and M-step until the values converges (no difference)

Imagine There are 2 coins A & B (Biased)
 Here One is likely to get Head more no of times and other is likely to get tail more no of tail

* You pick one at random and toss it
 Which one was it

Sol: Repeat this 5 time.

↳ * pick a coin randomly (known whether A or B)

- * Toss 10 times
- * Read No of Head & Tail

Consider the result

	Coin A	Coin B
First Round	5H, 5T	
Second Round	9H, 1T	

Consider the

First Round	5H, 5T
Second Round	9H, 1T
Third Round	8H, 2T
Fourth Round	4H, 6T
Fifth Round	7H, 3T
Total	(24H 6T) (9H, 11T)

Probability of getting Head with Coin A = $\frac{24}{30} = 0.80$

" " " " " Coin B = $\frac{9}{20} = 0.45$

Coin A yields head 80% of time and coin B 45% of time.

* What if only results are given.

& we need to guess the percentage of heads that each coin yields & Part A

Also need to guess which coin was picked at each round of toss. Part B

First Round	5H, 5T	?
Second Round	9H, 1T	?
Third Round	8H, 2T	?
Fourth Round	4H, 6T	?
Fifth Round	7H, 3T	?

Assume \rightarrow the probability of Head for $\left\{ \begin{array}{l} \text{Coin A} = 60\% \Rightarrow P_A = 0.6 \\ \text{Coin B} = 50\% \Rightarrow P_B = 0.5 \end{array} \right\}$ Randomly

Initial

First Round \Rightarrow

~~5H, 5T.~~

Compute the likelihood that it was coin A & coin B Using

Binomial Distribution i.e.

$$\boxed{P(k) = \Omega^k (1-\Omega)^{n-k}}$$

$n = \text{no of toss}$
 $k = \text{successful (head)}$

\Rightarrow is probability on n trials with k success

Likelihood of A = $P_A(h)^k (1 - P_A(h))^{n-k}$

for $n=10, k=5$

$$= (0.6)^5 \times (1 - 0.6)^{10-5} = \boxed{0.0007962624}$$

Likelihood of B = $P_B(h)^k (1 - P_B(h))^{n-k}$

$n=10, k=5$

$$= (0.5)^5 \times (1 - 0.5)^{10-5} = \boxed{0.0009765625}$$

Normalize by Using $\frac{A}{A+B}$ [To convert likelihood into Probability]

for coin A = $\frac{0.0007962624}{(0.0007962624 + 0.0009765625)}$

$$= 0.45\bar{4}$$

coin B = $\frac{0.0009765625}{(0.0007962624 + 0.0009765625)}$

$$= 0.55\bar{5}$$

Estimating Likely no of head & tails for coin A & B for First Round.

for coin A = $0.45 \times 5 = 2.2 H$
 $= 0.45 \times 5 = 2.2 T$

for coin B = $0.55 \times 5 = 2.8 H$
 $= 0.55 \times 5 = 2.8 T$

Consider Round 2 (9H, 1T)

$$\text{Likelihood of } A = 0.6^9 (1-0.6)^{10-9} = 0.0040310784$$

$$\text{Likelihood of } B = 0.5^9 (1-0.5)^{10-9} = 0.0009765625$$

Normalize A & B

$$P(A) = 0.80$$

$$P(B) = 0.20$$

Estimate likely no of Head & Tail for coin A & B in Round 2

$$\begin{array}{l} \text{No of H for coin A} = 0.80 \times 9 = 7.2H \\ \text{T " " A} = 0.8 \times 1 = 0.8T \end{array} \quad \left. \right\}$$

$$\begin{array}{l} \text{No of H for coin B} = 0.2 \times 9 = 1.8H \\ \text{T " " B} = 0.2 \times 1 = 0.2T \end{array} \quad \underline{\underline{}}$$

After 5 Rounds		Coin A	Coin B
Round 1	2.2H, 2.2T	2.8H, 2.8T.	
Round 2	7.2H, 0.8T	1.8H, 0.2T.	
Round 3	5.9H 1.5T	2.1H 0.5T	
Round 4	1.4H 2.1T	2.6H 3.9T.	
Round 5	4.5H 1.9T	2.5H 1.1T	
Total	21.3H 8.6T	11.7H 8.4T.	

Probability of getting Head

$$Q_A = \frac{21.3}{(21.3+8.6)} = \underline{\underline{0.71}}$$

$$Q_B = \frac{11.7}{(11.7+8.4)} = \underline{\underline{0.58}}$$

After 1 step of E & M $\begin{cases} Q_A = 0.71 \\ Q_B = 0.58 \end{cases}$ ✓

n ... until New Value of Q_A & Q_B is more similar

Repeat until New Value of $\underline{\theta_A \& \theta_B}$ is very small
to Equalise.

In This Example

The Value Converge to

$$\begin{cases} \theta_A = 0.8 \\ \theta_B = 0.52 \end{cases} \quad \checkmark \quad \checkmark$$

Part 1 Answered } Percentage of Head Coin A yields = 80%
B yields = 52%

$$\begin{aligned} \theta_A &= 0.8 \\ \theta_B &= 0.52 \end{aligned}$$

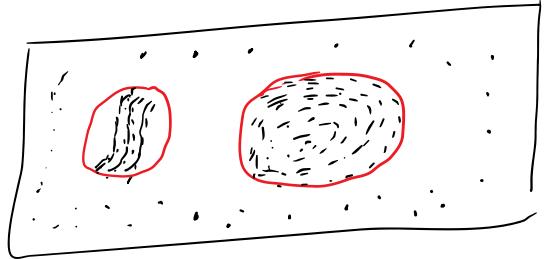
Now Consider the Result →

	H	T	Which Coin?
Round 1	5	5	$P(H) = \frac{5}{10} = 0.5$ = Coin B
Round 2	9	1	$P(H) = \frac{9}{10} = 0.9$ = Coin A
Round 3	8	2	$P(H) = \frac{8}{10} = 0.8$ = Coin A
Round 4	4	6	$P(H) = \frac{4}{10} = 0.4$ = Coin B
Round 5	7	3	$P(H) = \frac{7}{10} = 0.7$ = Coin A

Go to Part 2

DBSCAN (Density Based Spatial Clustering of Application with Noise)

- * In Dense Region there is possibility of cluster than in Sparse Region
- * Density based approach is better if the cluster is of arbitrary shape.



Key Features

To understand DBSCAN Algo we need to understand

ϵ Maximum radius of neighbourhood.

Minpts: Minimum no of points in an ϵ neighbourhood of that point.



$$\epsilon \text{ neighbourhood of } q = \{4\}$$

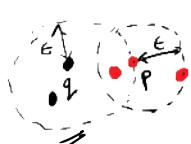
$$\epsilon \text{ ————— ————— } - p = \{3\}$$

Density of q is high

Density of p is low

Terminologies > Given ϵ and Minpts then categorize the objects into 3 group

① Core Point, if for a point the neighbourhood has the points $\geq \text{Minpts}$. Then it is Core point.

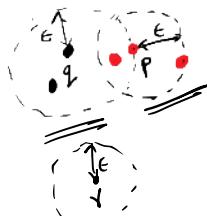


$$\text{Minpts} = 4$$

No of points in neighbourhood of $q = 4 \rightarrow$
q is core point

So q is core point

② Border Point \rightarrow If for a point the neighbourhood has points $< \text{Minpt}$ but it should be neighbours of core point



$$\text{Minpt} = 4$$

No of points in neighbourhood of p = 3
 $< \text{Minpt}$

& p is neighbour of q (core point)

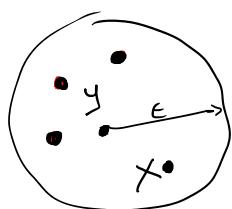
So 'p' is border point

③ Outlier \rightarrow If a point is neither a core point nor a border point, it is called outlier.

"x is outlier"

* Directly Density Reachable \rightarrow A point X is directly density reachable from point Y w.r.t. epsilon & Minpt if

1. X belongs to neighbourhood of Y ($\text{dist}(X,Y) \leq \epsilon$)
2. Y is core point -



$$\text{Minpt} = 4$$

① as Y a core point \rightarrow

no of point in neighbourhood of Y = 5 $> \text{Minpt}$
 $\therefore Y$ is core point

② X belongs to neighbourhood of Y.

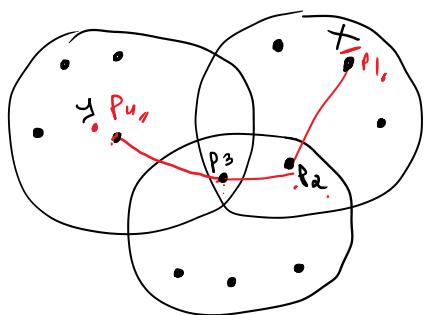
$\therefore X$ is Directly Density Reachable from Y.

② Density Reachable \rightarrow

A point X is density reachable from Point Y w.r.t. ϵ & Minpt

A point X is densely reachable from Point Y wrt ϵ & M_{inpt}
 if there is a chain of points $P_1, P_2 \dots P_n$ and $P_1 = X$
 & $P_n = Y$ such that P_{i+1} is directly density reachable from P_i

$M_{\text{inpt}} = 4$ here X & Y are core points



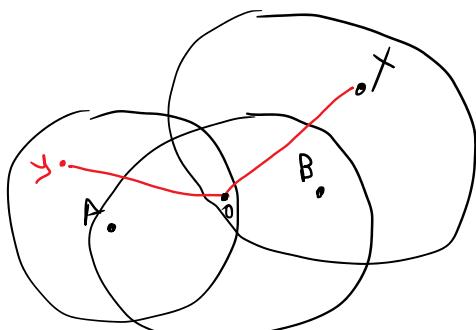
Here P_2 is directly density reachable from X

P_3 is directly density reachable from P_2

\nexists is directly density reachable from P_3

$\therefore X$ is Density Reachable from Y

Density Connected \rightarrow A point X is density connected from point Y wrt ϵ and M_{inpt} if there exists a point O such that both X & Y are density reachable from O wrt ϵ & M_{inpt}



Y & O are direct density reachable
 X & O are direct density reach

thus X & Y are density connected

Direct $\xrightarrow{\text{within } \epsilon}$ Direct Density Reachable.
 Via Some chain \rightarrow "

2 points are densely reachable through common pt \rightarrow Dense
Connected

Algo

1. Arbitrary Select a Point P

2. Retain all points density reachable from P wst ϵ & Min Point

3. If P is core point then cluster is formed

4. If P is border point then DBSCAN visits next point of dataset

5. Continue the process until all the points are processed.