

- ③ DBSCAN Intro
- ④ Key Concepts in DBSCAN
 - ① Min-points
 - ② Eps (Epsilon)
 - ③ Core-points
 - ④ Border Point
 - ⑤ Noise Point
- ⑤ DBSCAN Algo
- ⑥ Key Hyperparameters
- ⑦ Code Implementation
- ⑧ Introduction to Anomaly/Novelty/Outlier Detection
- ⑨ Distribution Based Anomaly Detection
 - ③ RANSAC (Random Sample Consensus)
 - ④ Elliptic Envelope
- ⑩ Sklearn Implementation

DBSCAN

- Density-Based Spatial Clustering of Application with Noise
- DBSCAN handles outlier effectively
- Scales well with Big Datasets
- No need to define n -clusters

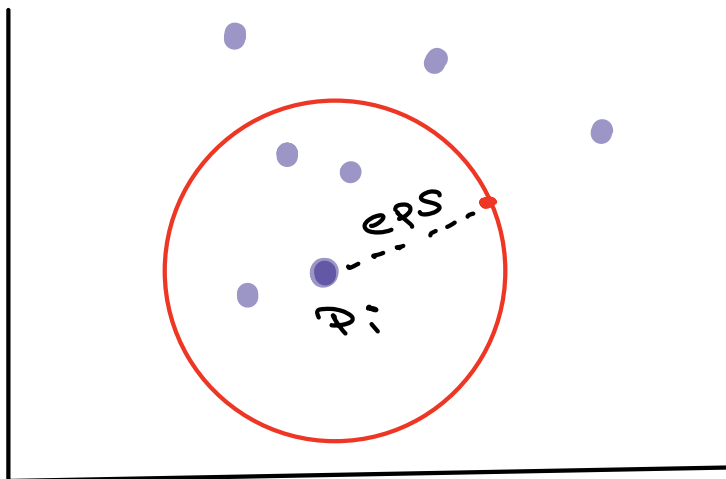
Core Concepts

- ① Min points
- ② ϵ (eps)
- ③ Core-point
- ④ Border-point
- ⑤ Noise-point

density @ P_i = # of points

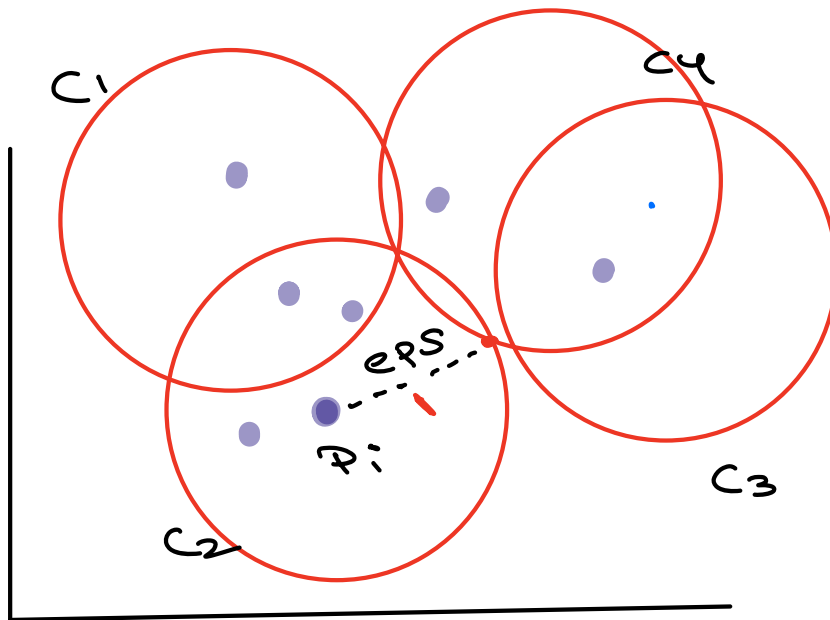
that are within

some radius ϵ



- 4 (including P_i)
- 3 (Excluding Self)

③ Min-pts \rightarrow No. of points allowed within ϵ radius for to be considered Dense regions



min-pts \ni 3

ϵ \ni 1

$C_1 \ni 3$

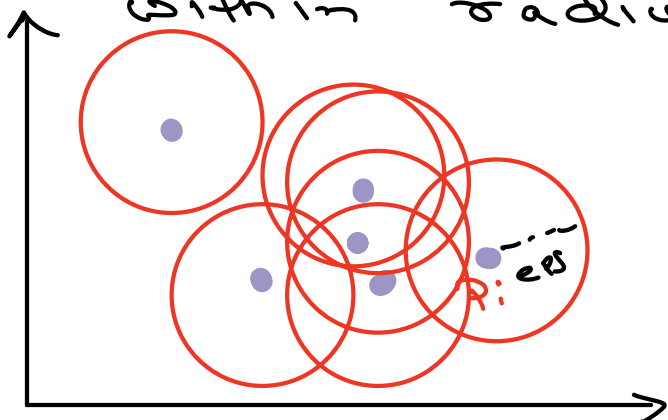
$C_2 \ni 4$

$C_3 \ni 1$

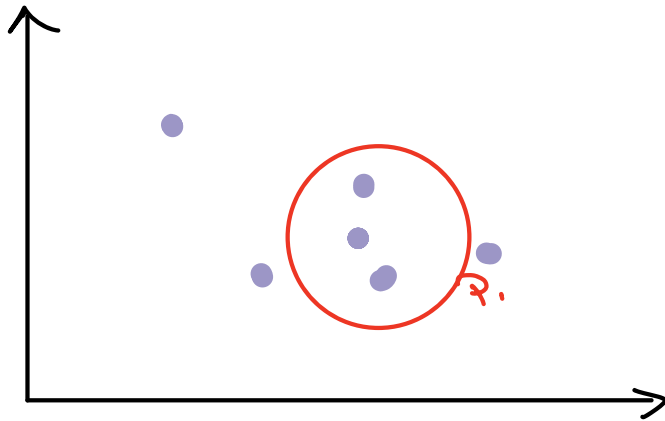
$C_4 \ni 2$

Core-points

A point which has $\#p \geq \text{min-point}$ within radius of ϵ



min-pts \ni 2



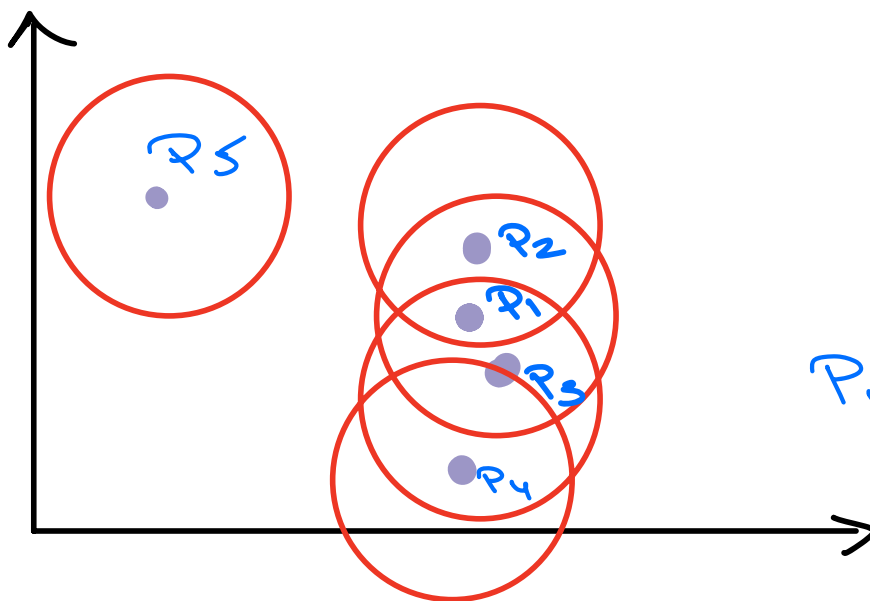
min-pts ≥ 2

D.P \ni min-pts $\geq n$

C.P $\ni n-1$

Border-point

- ① A point which is not a Core point
- ② It lies in neighbourhood of point Q which is a Core-point



min-pts ≥ 2

$P_4, P_2 \ni$ B.P

$P_1, P_3 \ni$ C.P

Noise-points

① A point which is neither Core nor Border

Ex: P_5 is Noise point

Density Edge/Connection

① If Edge connecting two Core points such that $\text{len}(\text{edge}) \leq \epsilon_{\text{pts}}$



Density Edges \Rightarrow Edge-PQ

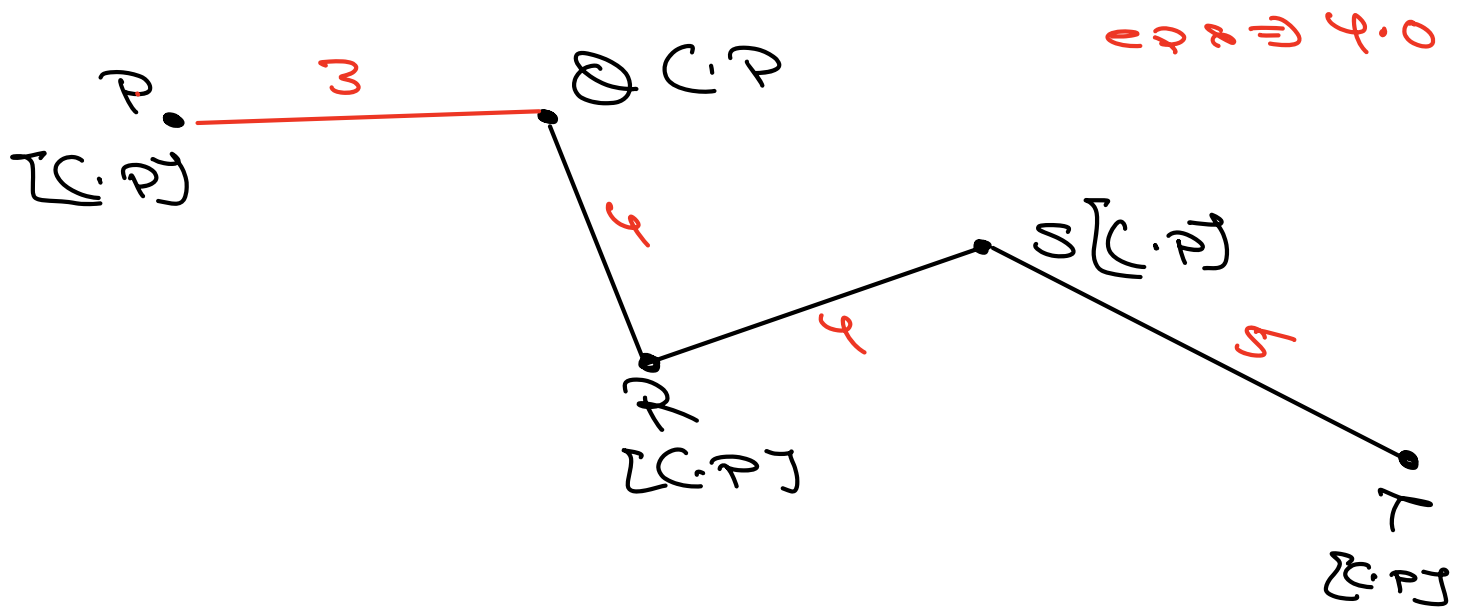
Density Connected points

Two points P and Q are called

D.C.P.,

① if Both P and Q are Core points

② if there exist Density Edges connecting point P and Q



$P \rightarrow S : D.C.P$

$P \rightarrow T : No \ D.C.P$

Q For two points to be D.C.P
can the Dist. between them
be greater than $C \cdot 2$?

Yes

DBSCAN ALGO

Step 1: Annotate Every point as C.P,
B.P or N.P $B.F \propto N^2$

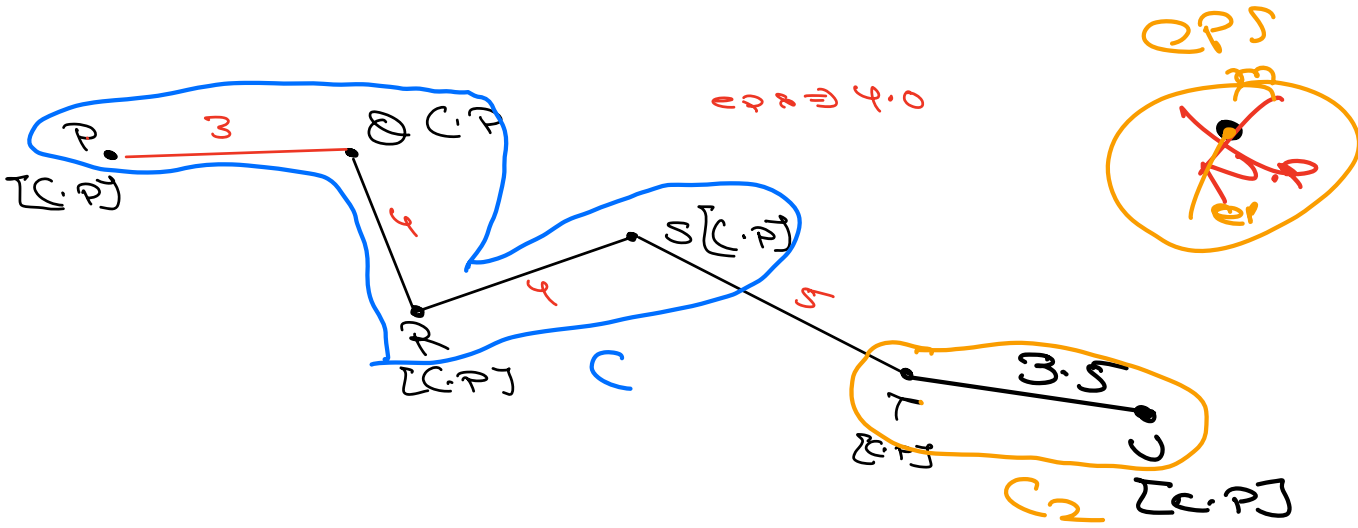
⑤ Range Query (D, x_i , E_{ps})
T.C. N^2
 $N \log N$

Step 2: Remove All Noise-points

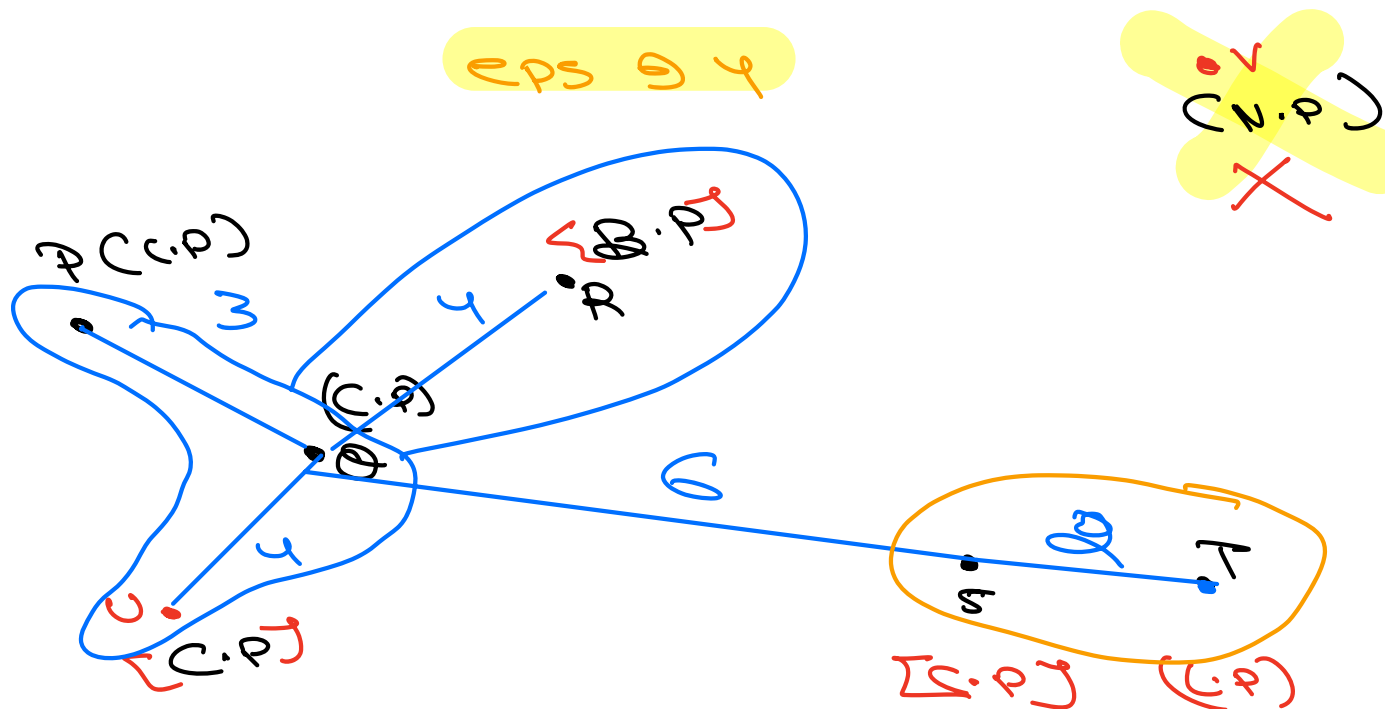
Step 3: For each Coe-point 'i.e. not part of any Cluster

② Create a new Cluster with Point p

5 Add all the density connected points to this Cluster



Step 4: Assign each Border-point to its nearest Cluster



C1 = PQRU

C2 = ST

100, 5
6, 10

Hyperparameter

min-pts \Rightarrow Rule of thumb \Rightarrow dimensions
 \rightarrow min-pts $\geq d+1$
 \rightarrow Typically $d \times d$

if D.S is very Noisy \Rightarrow Higher Value of min-pts

② $\in \mathbb{R}$:

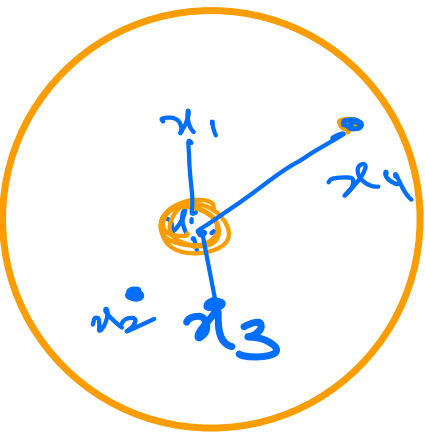
With a fix value of min-pts

$x_1 \rightarrow d_1$ distance of
 $x_2 \rightarrow d_2$ farthest

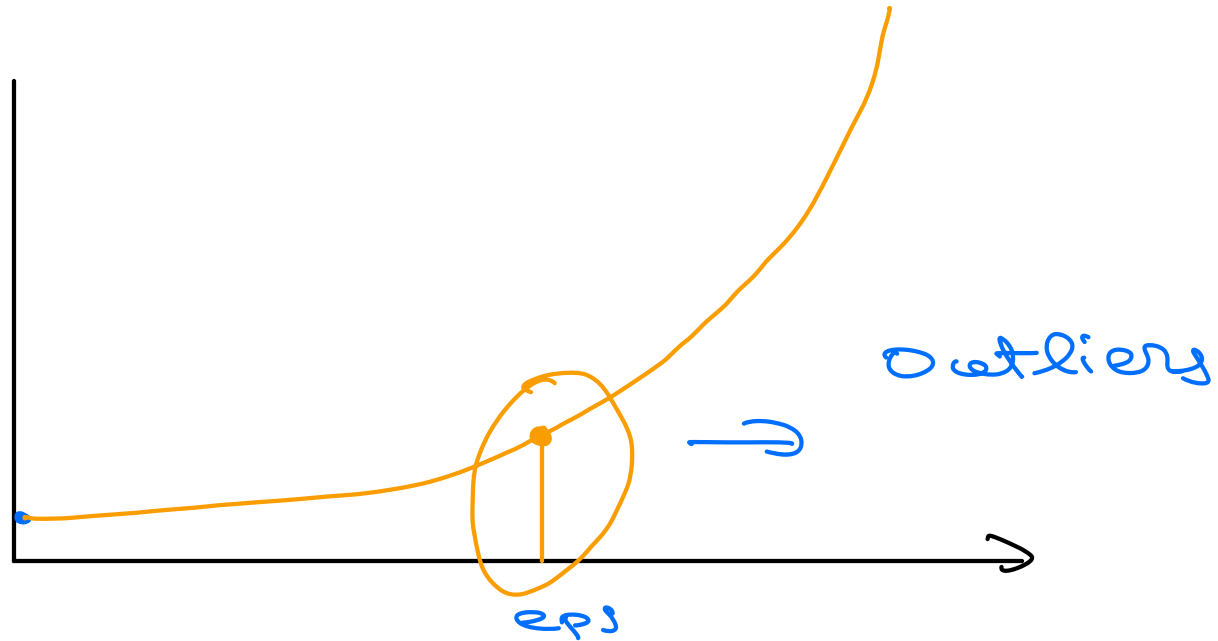
B.P or C.P

$x_n \rightarrow d_n$

•
Cig min-pts = 4
then we take
 d_i as 4th high



③ Sort all distance and plot them



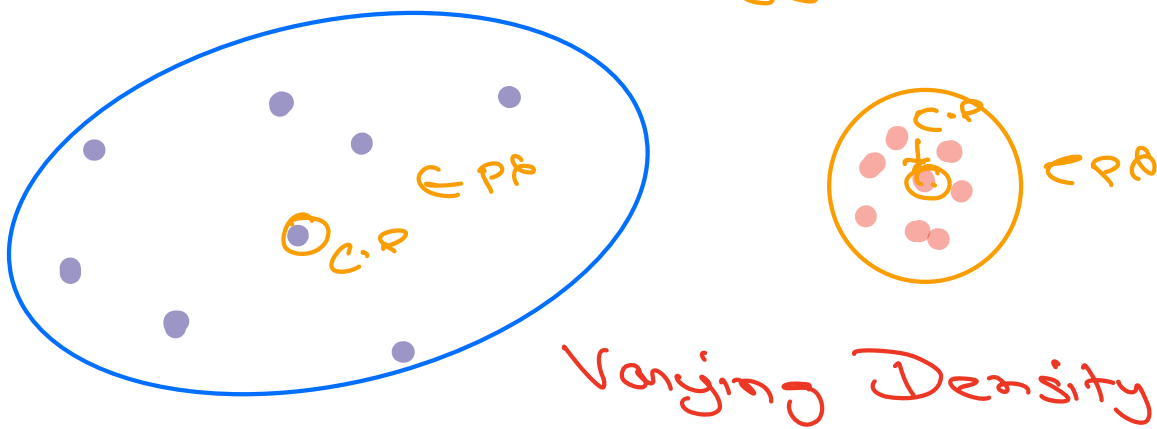
euclidean Distance

Advantages

- ① Resistant to Noise
- ② Can Handle any shape and Size of Cluster
- ③ Does n't require k -clusters a priori
- ④ Only two parameters $\Rightarrow \epsilon, \minPts$
- ⑤ Speed : Thanks to DB Community $N \log N$

Limitation

- ① Very sensitive to choice of Hyperparameter
- ② Can't Handle different Densities

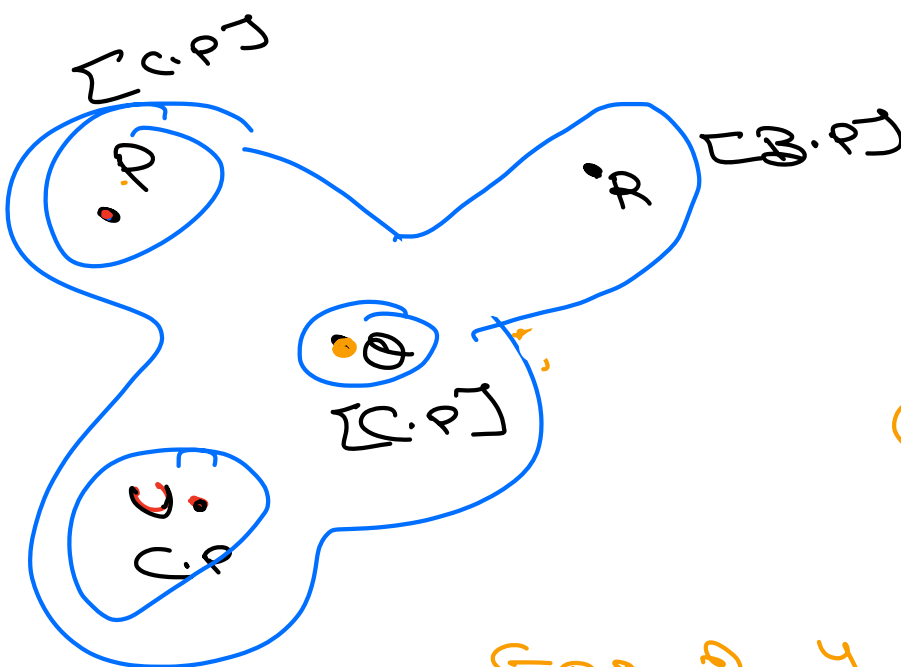


9 Topics Remaining

① DBSCAN Implementation

② Anomaly Detection

③ RANSAC and Elliptic Envelope

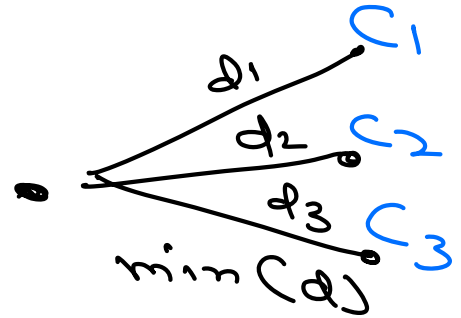


minimize
error

Online Clustering Model

↳ Can predict without retraining

Kmeans \rightarrow Centroid



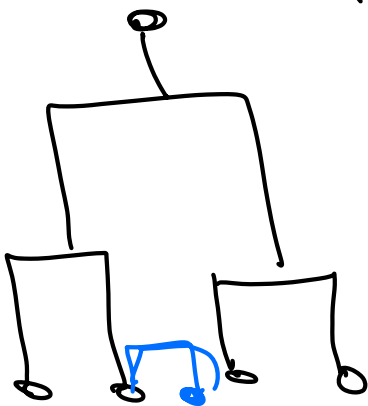
▷ Predict

Offline Clustering Model

▷ Refit the full dataset to perform prediction

▷ Process Full DS + new-point

▷ Predict Method Not Available



Build Dendrogram