

AND MANY  
OUTLIER  
NOVELTY  
DETECTION

# Outliers

Q Why do outliers exist?

\* Human error

\* Sensor error / faulty machine

\* "Unusual" data

→ Anomaly  $\Rightarrow$  Not Normal

→ Novelty  $\Rightarrow$  New, Never happened before.

eg. \* Fraud detection -

\* Car Mileage EV, hybrid.

C	Mileage	Range
	10-20 km	600
	EV	300
	hybrid	1000

Q what are some ways to detect outlier?

A

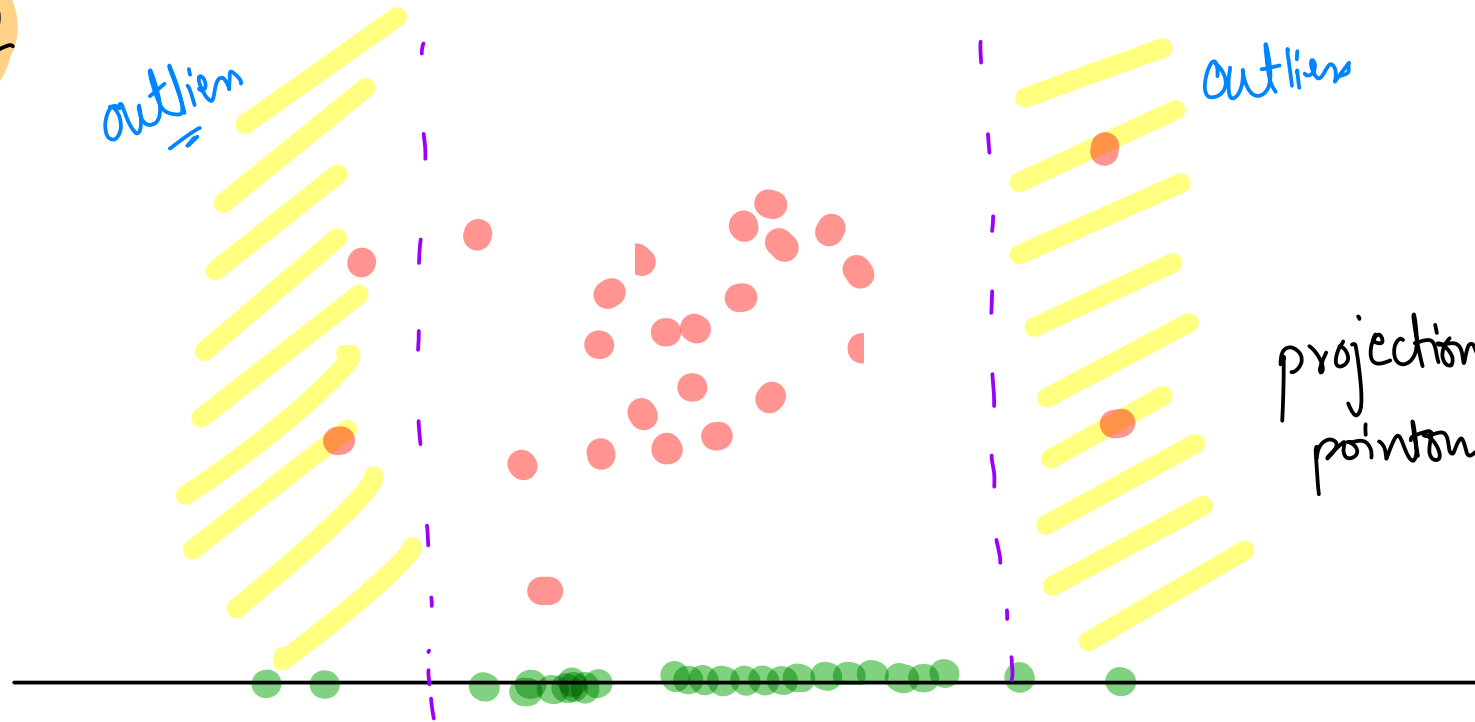
IQR, KNN → DBSCAN

11

outlier

outliers

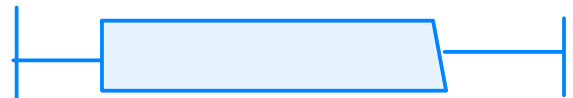
projection of point on axis



Boxplot

outlier

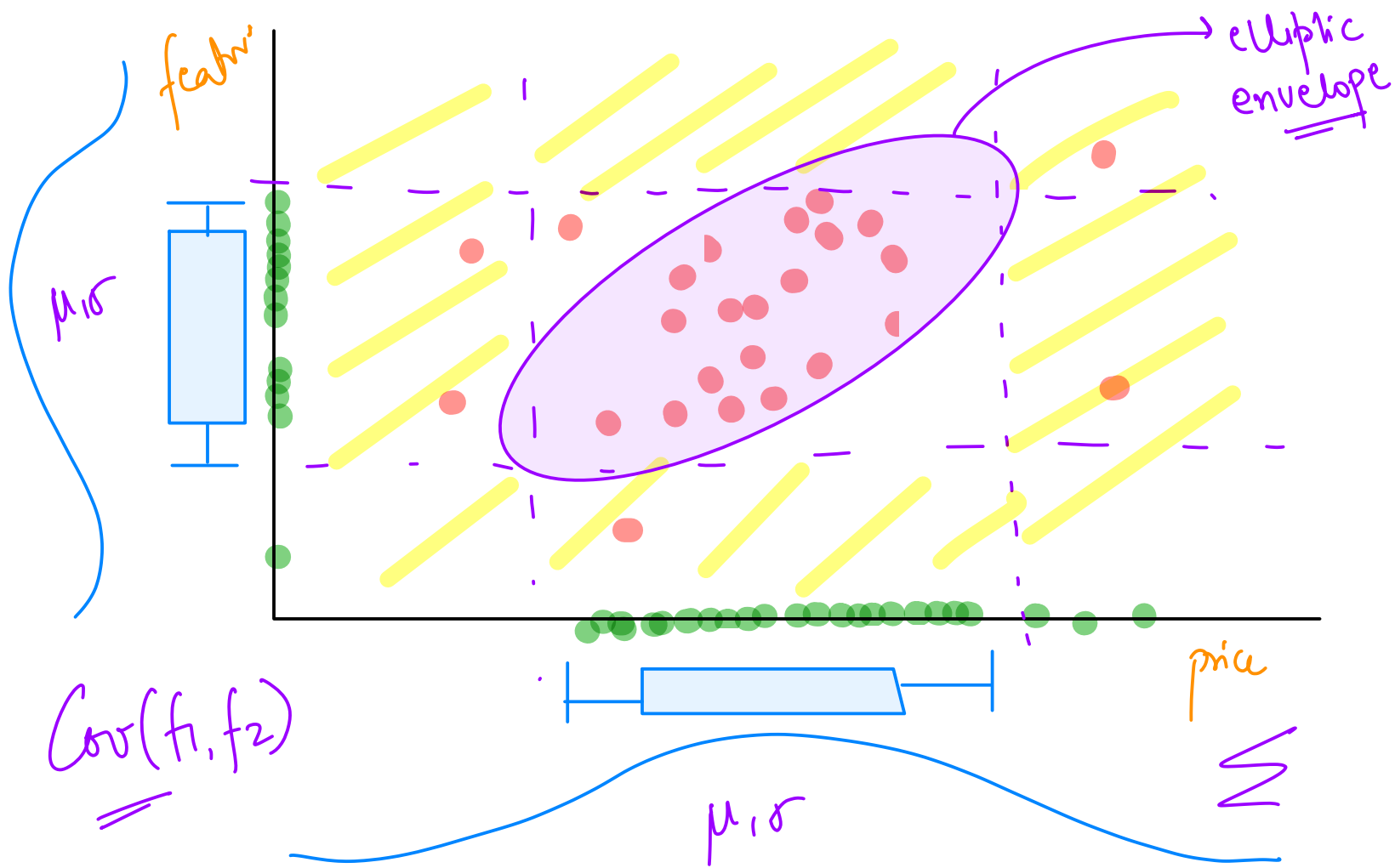
2%



outliers

96% inliers

2% outlier



① IQR (BOXPLOT)

② Multi Variate Gaussians  
↳ Elliptical Envelope..

learn  
estimate  $\{\mu, \sigma, \Sigma\} \Rightarrow$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 & \sigma_1 \sigma_3 \\ \sigma_1 \sigma_2 & \sigma_2^2 & \sigma_2 \sigma_3 \\ \sigma_1 \sigma_3 & \sigma_2 \sigma_3 & \sigma_3^2 \end{bmatrix}$$

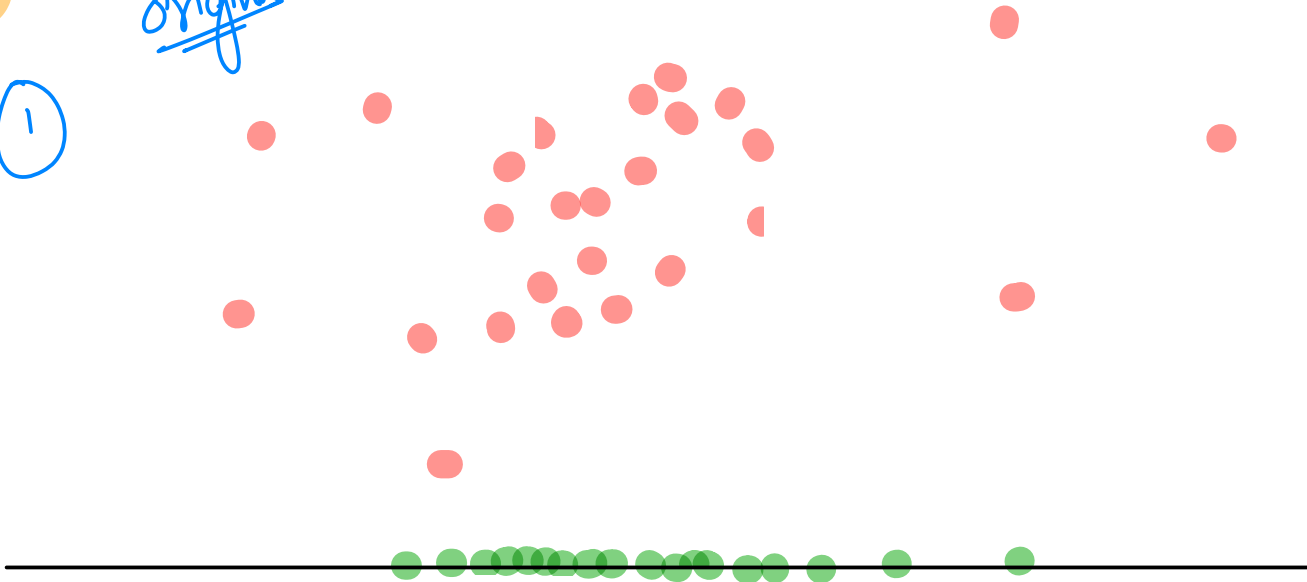
RANSAC

Random  
Sample  
Consensus

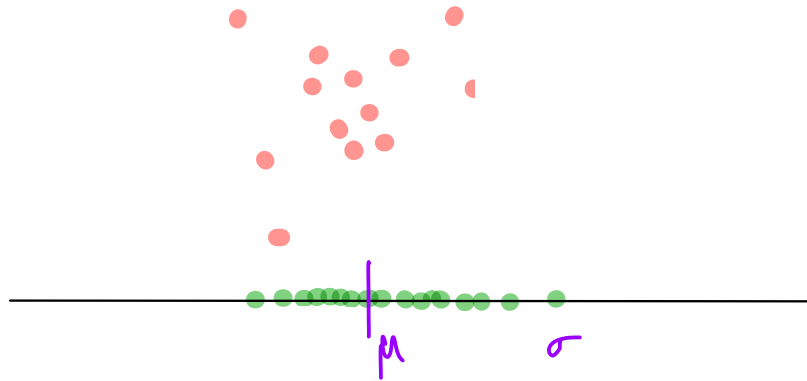
①

original

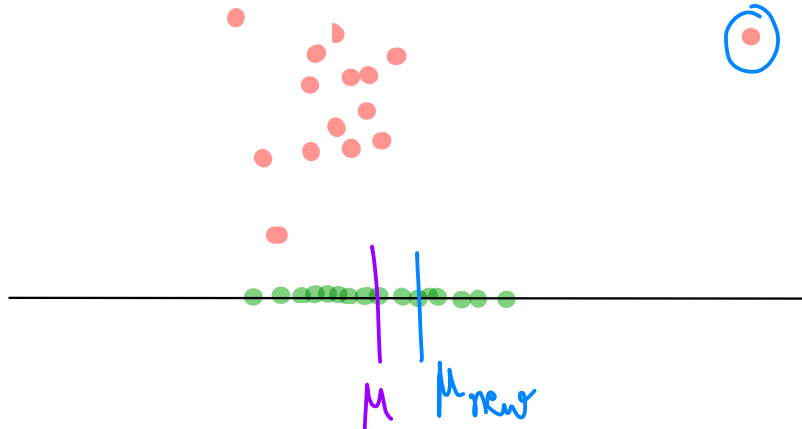
Robust



2



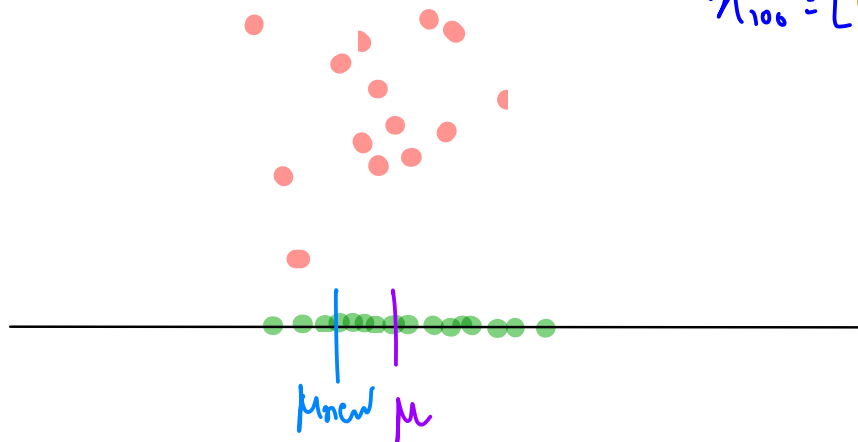
3



4



3

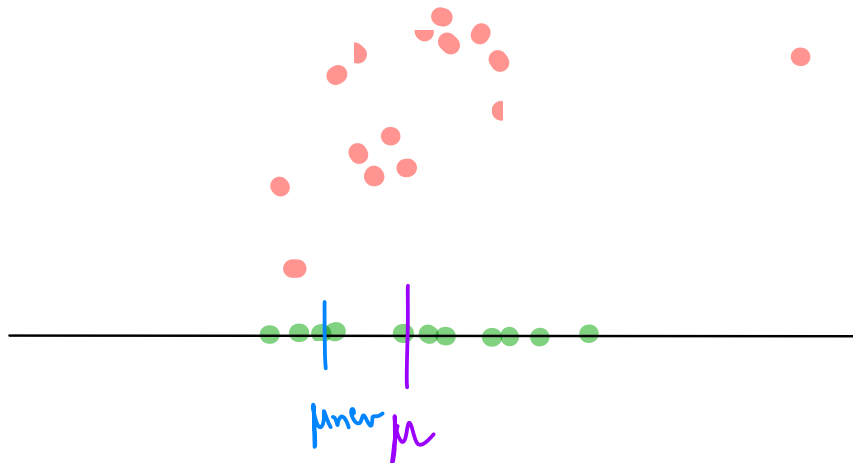


$$x_{100} = \begin{bmatrix} 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ | & | & | & | & | & | & | & | & | & | & | & | \end{bmatrix}$$

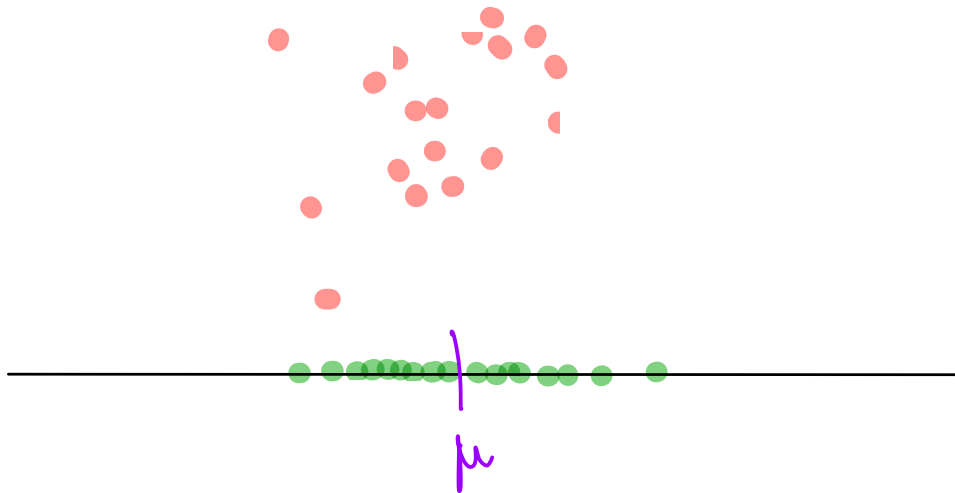
present not outlier . 1  
 present outlier . 2  
 absent . 0

$$x_{100} = \begin{bmatrix} 1 & 2 & 2 & 0 & 0 \end{bmatrix}$$

4



5

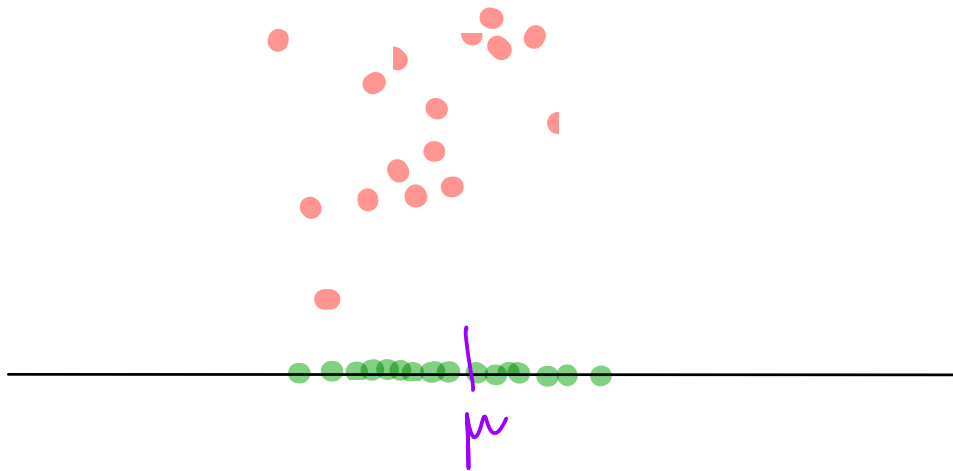


100 times

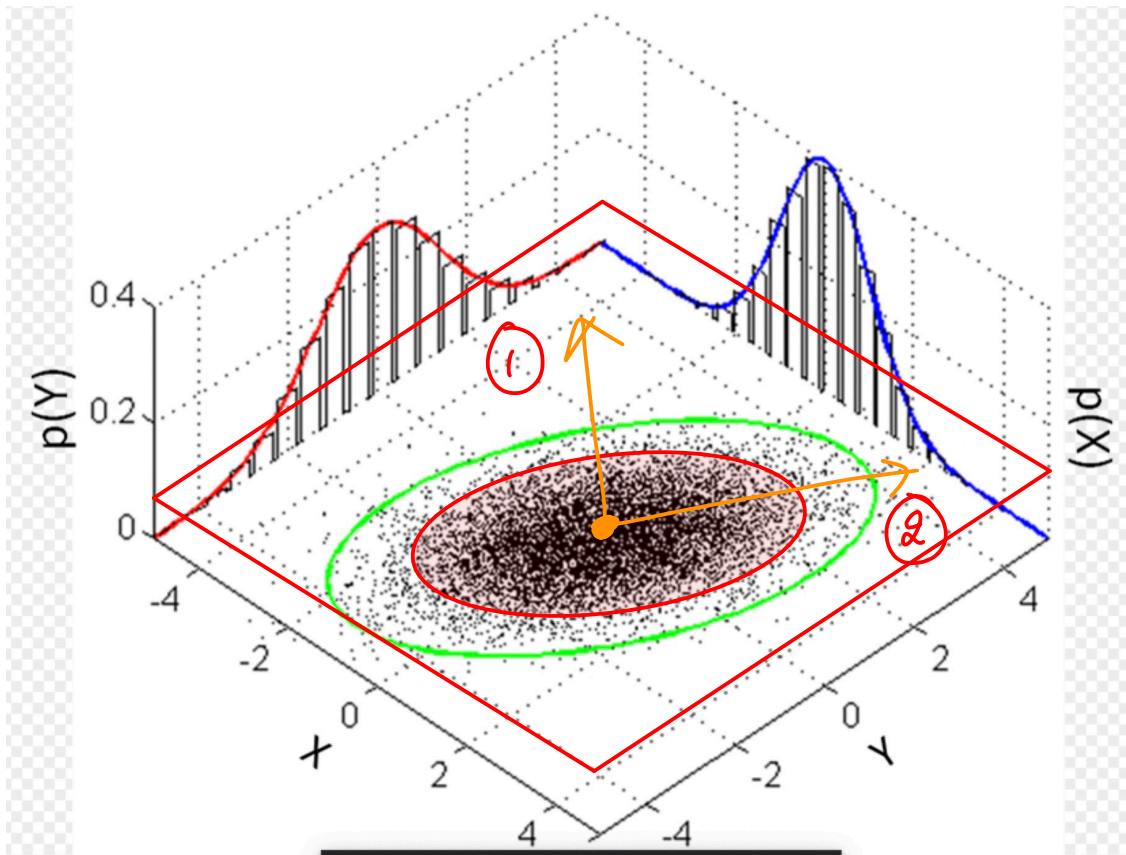
96%  $\mu$

10%  $\mu + \delta$   
 $\mu - \delta$

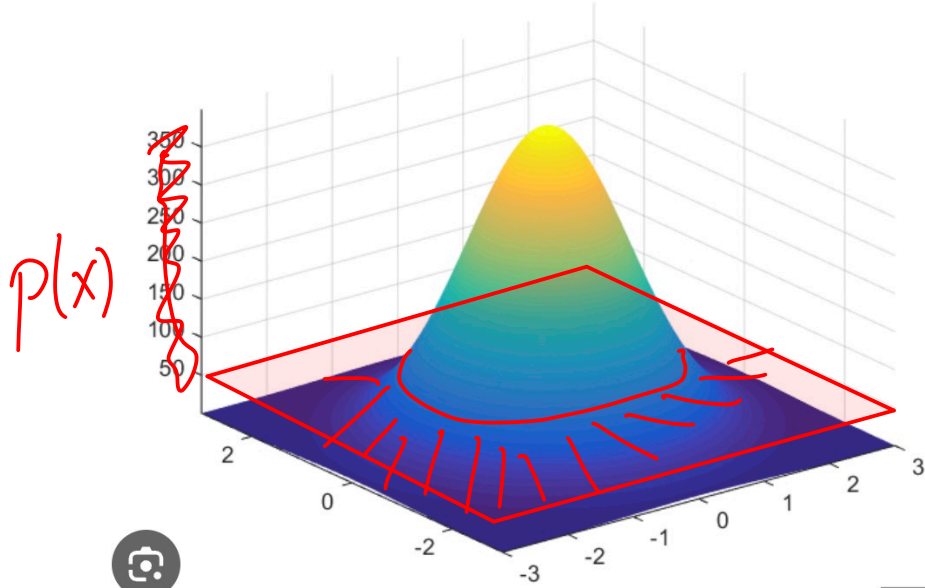
6



Avg  $\Rightarrow$   $\mu$



$$p(x) < 0.2$$



672 x 504



For  $\epsilon$  Radius

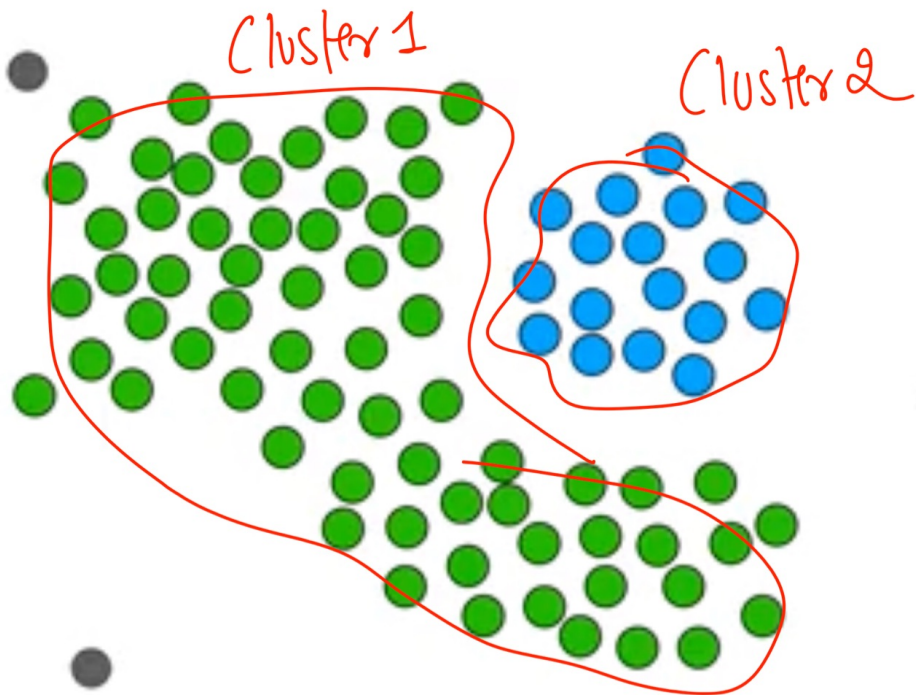
$\Rightarrow$  We will check  
if # of points inside circle  $> n$   
↳ Core point

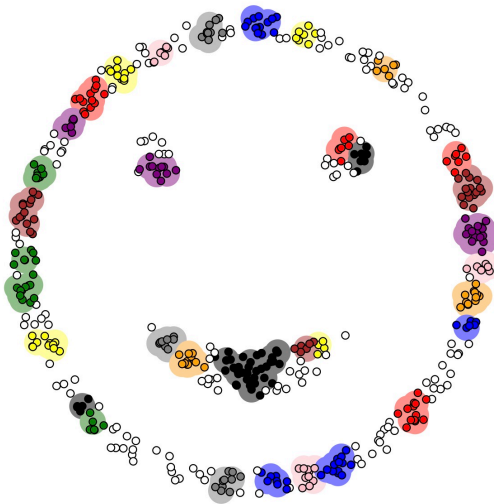
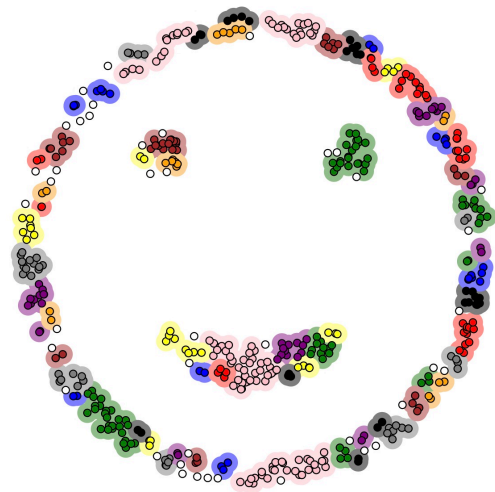
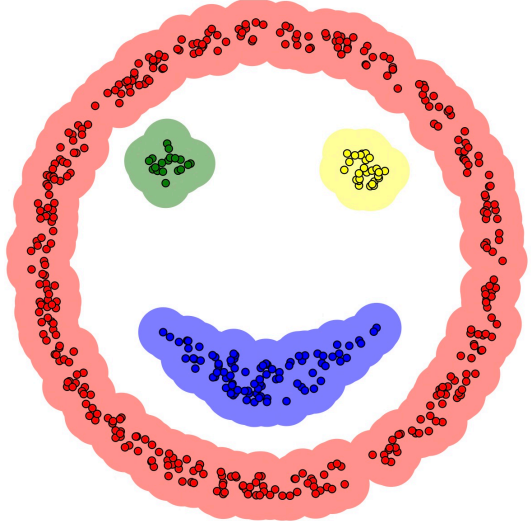
else.  
=

Non Core point

→ Border point  
→ Noise point.

Border point : if a Non Core point lies in the Circle  
of any Core point, it is a Border point  
else: Noise point.







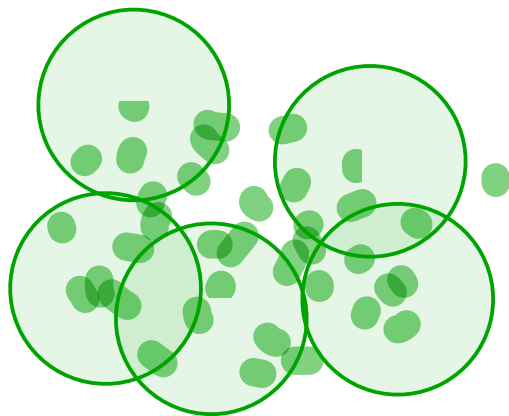
Pros: (1) Finds out outliers & clusters as well  
(2) No need to decide 'K'

Cons: (1) Very Sensitive to hyperparameters.  
(2) Does not work well with sparse data / data with diff densities.

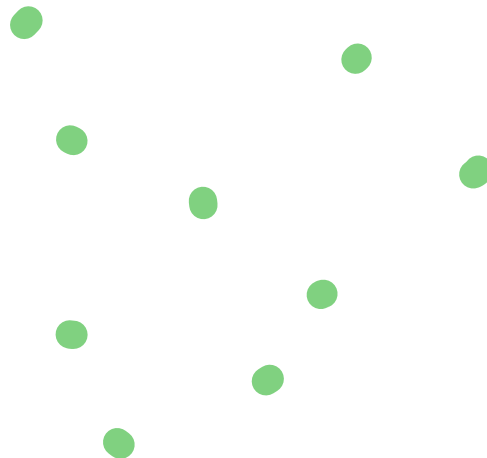
RF  
L 100  
L 120

DBSCAN  
[0.1 → 0.11]

density  $\uparrow$   
Cluster 1

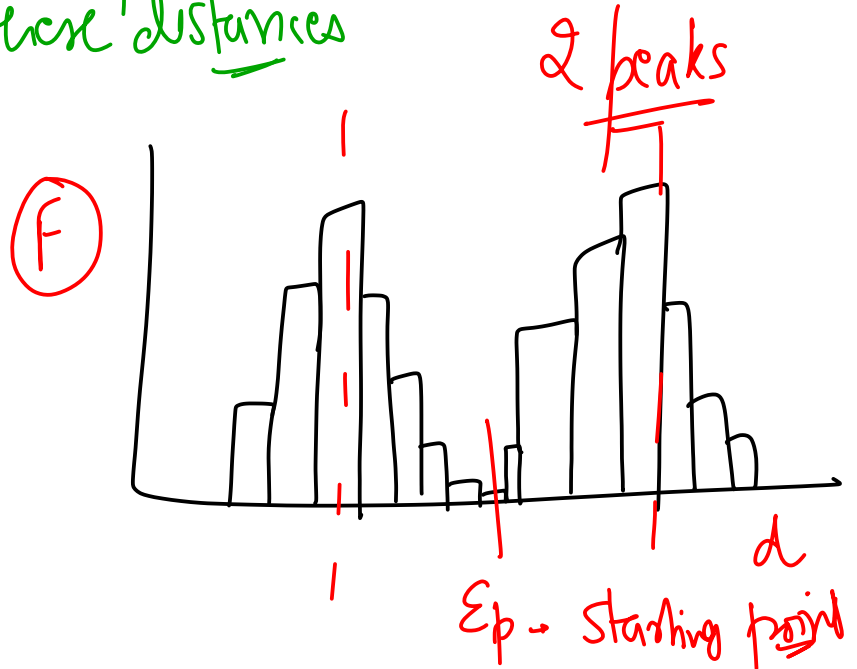


density  $\downarrow$   
Cluster 2



How to estimate good value of  $\epsilon$ ?

- ① Calc. distance for each point.
- ② Plot histogram of these distances



How to estimate min # of points ??

\*  $d$  dimensional

$$\hookrightarrow n \gtrsim d+1$$

\*  $d$ -dimensional

$$\hookrightarrow n \gtrsim d \times 2$$

## Quiz time!

🕒 Quiz Ended!

### What is a core point in DBSCAN?

5 users have participated

- ☐ A A point located at the center of a cluster 0%
- ☐ B A point with the highest density in the dataset 0%
- ☒ C A point that has more than Min Points within the Epsilon radius 100%
- ☐ D A point that is not part of any cluster 0%

🕒 Quiz Ended!

### What is a "density edge" in DBSCAN?

6 users have participated

- ☐ A A line connecting two random points in a dataset 0%
- ☒ B An edge that connects two core points with a distance less than or equal to  $\epsilon$  (Epsilon) 100%
- ☐ C The edge of a dense cluster in the dataset 0%
- ☐ D An edge between two noise points 0%

## Quiz time!

🕒 TIME LEFT: 13 Secs

### How is a border point defined in DBSCAN?

6 users have participated

- ☐ A A point located on the border of the dataset 0%
- ☐ B A point with the highest density in a cluster 0%
- ☐ C A point that is not part of any cluster 0%
- ☒ D A point that is not a core point but is within the Epsilon radius of a core point 100%

### When are two points considered "density connected" in DBSCAN?

6 users have participated

- ☐ A When they are connected by an edge 0%
- ☐ B When they are both noise points in the dataset 0%
- ☒ C When they are core points and there exists a sequence of density edges connecting them 100%
- ☐ D When they are part of the same cluster 0%

[End Quiz Now](#)