

# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

Conducted by:

iHUB Dhyesha Shukla, Roorkee  
and

Ritvij Bharat Private Limited (RBPL)

# DBSCAN

Presented by:

Shreyas Shukla

## DBSCAN

Conducted by:

Density-based spatial clustering of applications with noise is a powerful technique which can be used for clustering and outlier detection.

Presented by:  
Shreyas Shukla

11 Sep 2023 - 20 Oct 2023

- Intuition of DBSCAN
- DBSCAN vs. K-Means Clustering
- DBSCAN Hyperparameters Theory
- DBSCAN Hyperparameters Coding

Presented by:  
Shreyas Shukla

An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

Conducted by:

iHUB Divya Sampark, IIT Roorkee

and

Ritvij Bharat Private Limited (RBPL)

## Theory and Intuition

Presented by:

Shreyas Shukla

DBSCAN stands for **Density-based spatial clustering of applications with noise.**

iHUB Divya Sampark, IIT Roorkee  
and  
Ritvij Bharat Private Limited (RBPL)

Presented by:  
Shreyas Shukla

## Some Questions:

- How does DBSCAN work?
- Advantages and disadvantages of DBSCAN?
- How does it deal with outliers and noise?

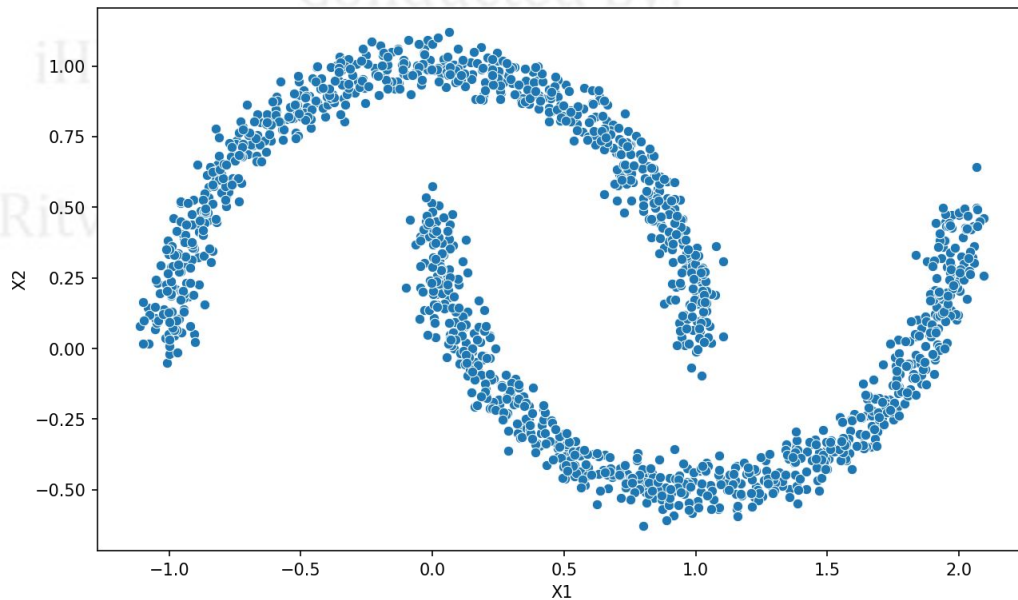
Presented by:  
Shreyas Shukla

## Key Ideas

- DBSCAN focuses on using **density** of points as its main factor for assigning cluster labels.
- This creates the ability to find cluster segmentations that other algorithms have difficulty with.

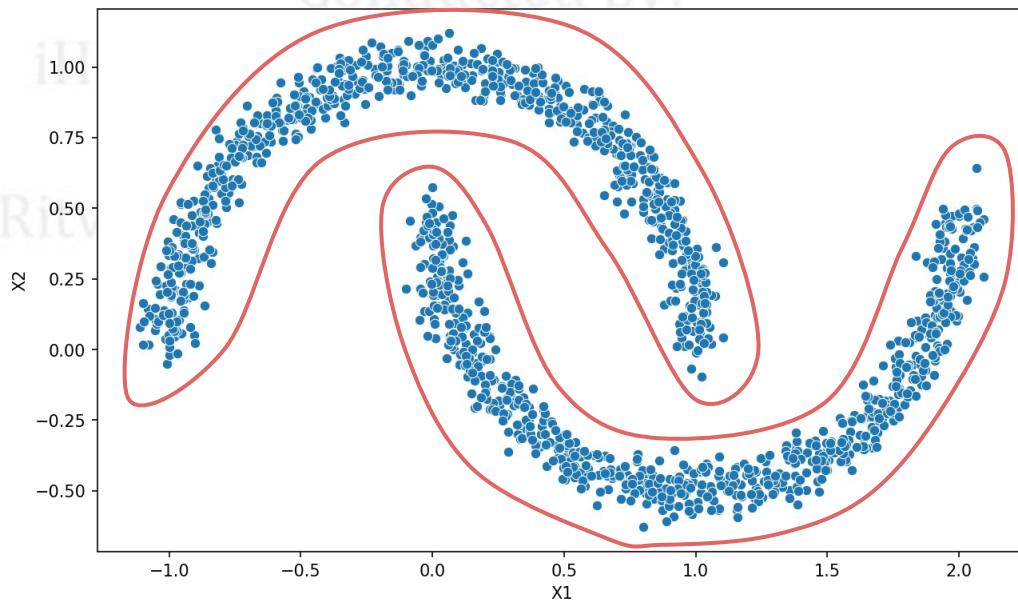
Consider the following data set:

Conducted by:

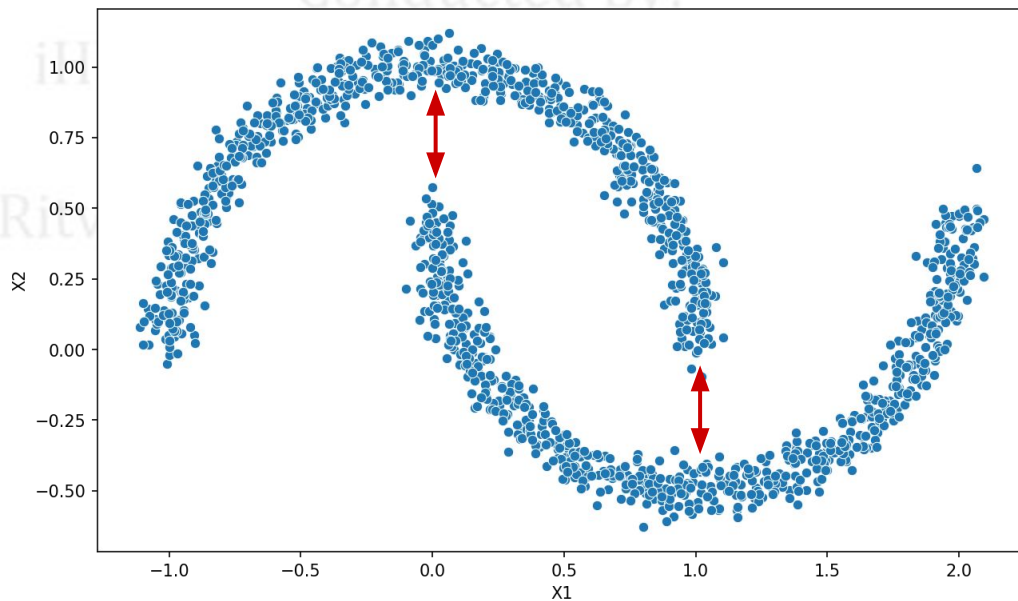




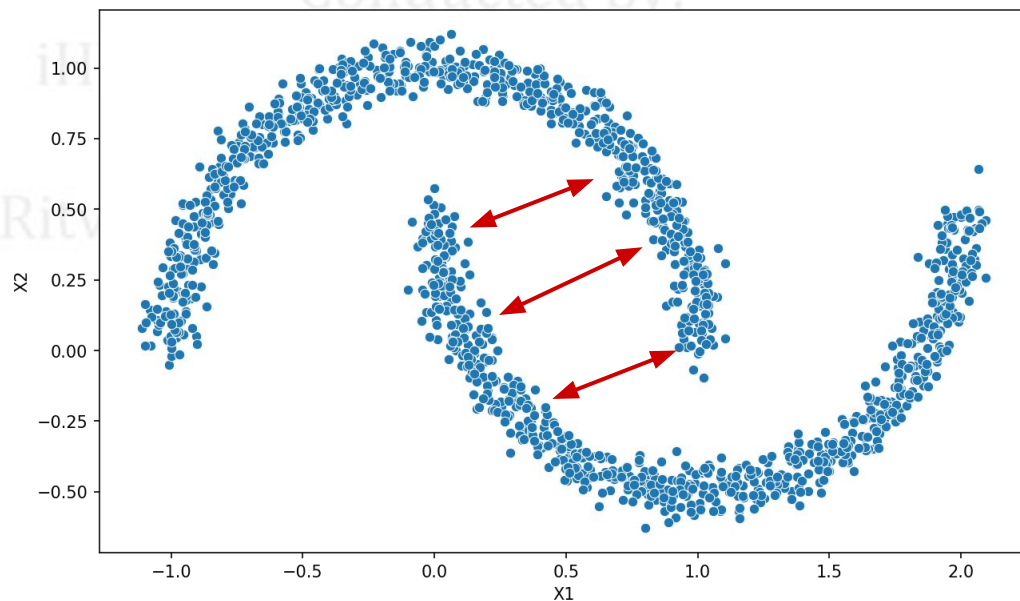
Clearly two “moon” shaped clusters:



But distance based clustering has issues:

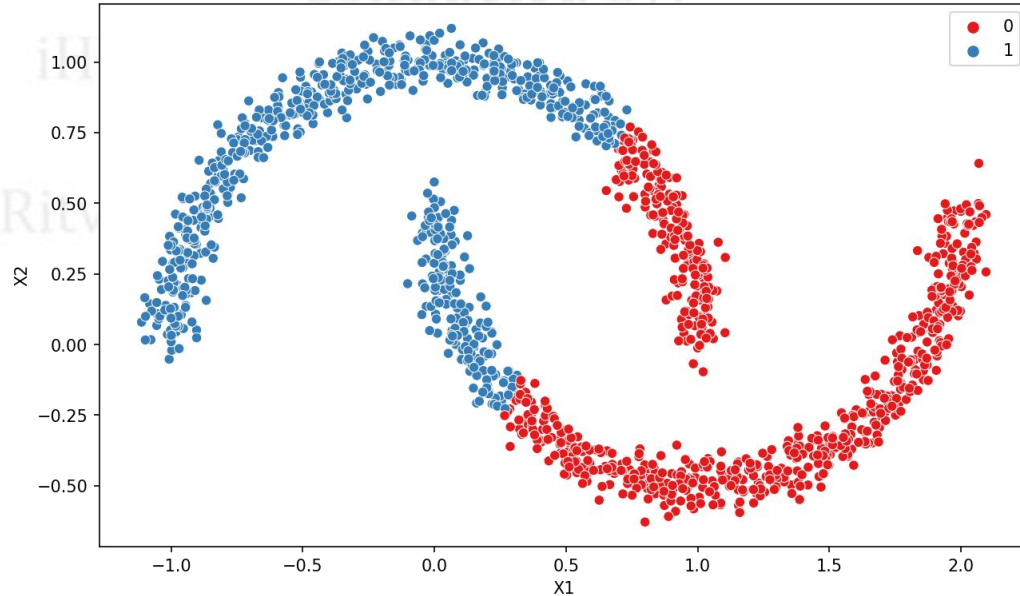


But distance based clustering has issues:



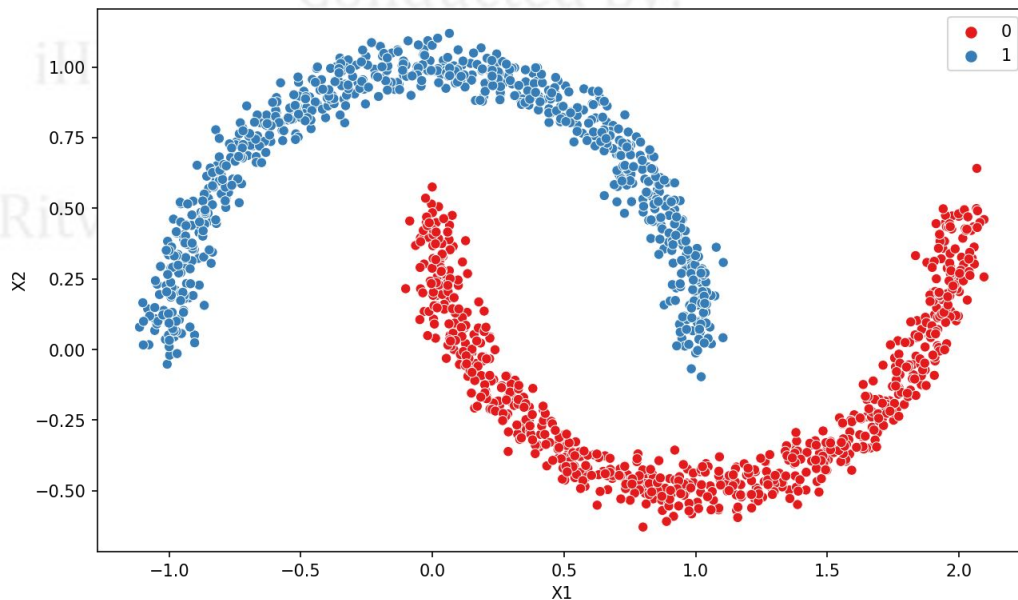
## Results of K-Means:

Conducted by:



## Results of DBSCAN:

Conducted by:



DBSCAN iterates through points and uses two key hyperparameters (epsilon and minimum number of points) to assign cluster labels.

Unlike K-Means, it focuses on density as the main factor for cluster assignment of points.

Presented by  
Shreyas Shukla

## Key Hyperparameters:

- Epsilon:
  - Distance extended from a point.
- Minimum Number of Points:
  - Minimum number of points in an epsilon distance.

Presented by:  
Shreyas Shukla

## DBSCAN Point Types:

- Core
- Border
- Outlier

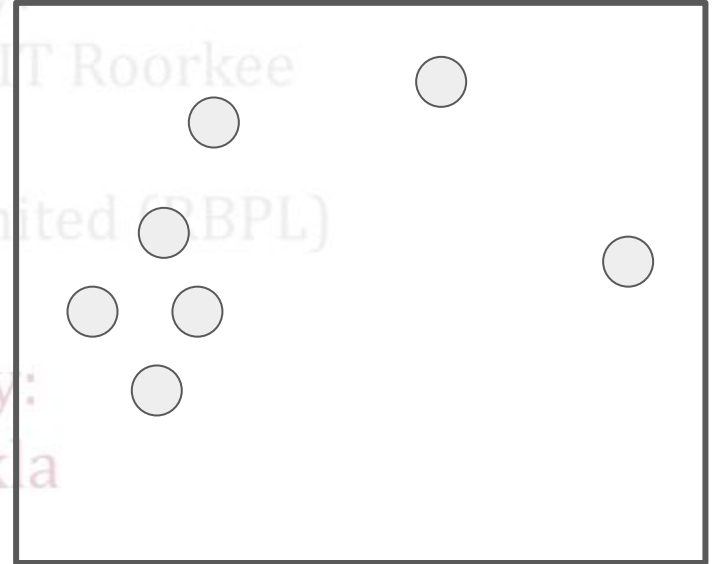
Conducted by:  
HUB Divya Sampark, IIT Roorkee  
and  
Kilvij Bharat Private Limited (RBPL)

Presented by:  
Shreyas Shukla



## DBSCAN Point Types:

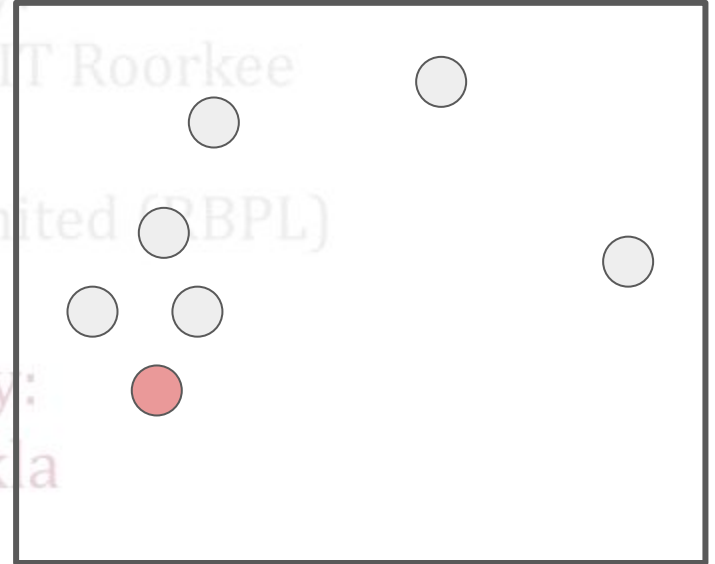
- Core
- Border
- Outlier



Presented by:  
Shreyas Shukla

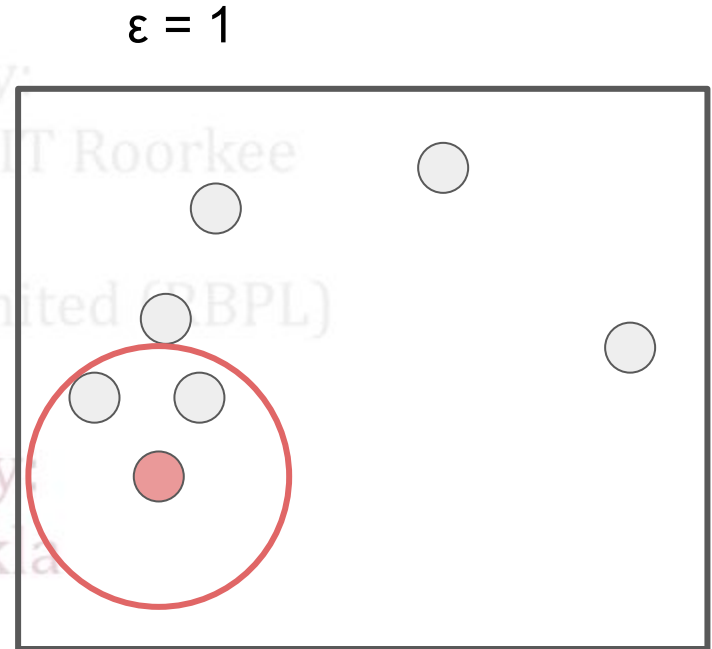
## DBSCAN Point Types:

- Core



## DBSCAN Point Types:

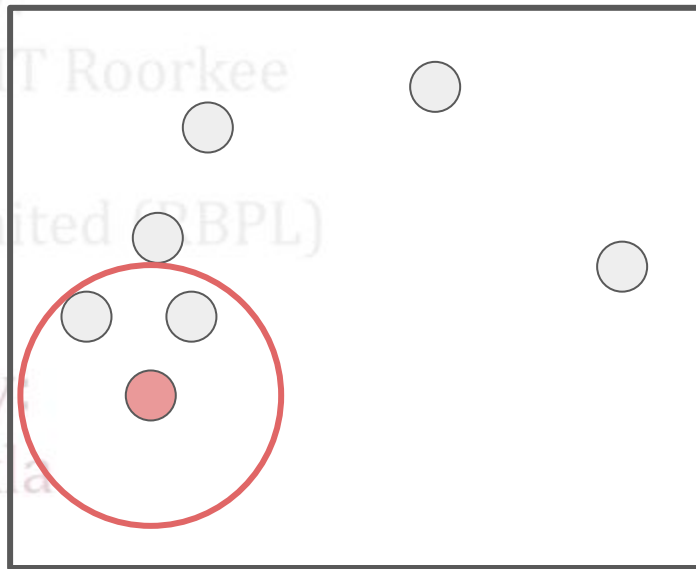
- Core



## DBSCAN Point Types:

- Core

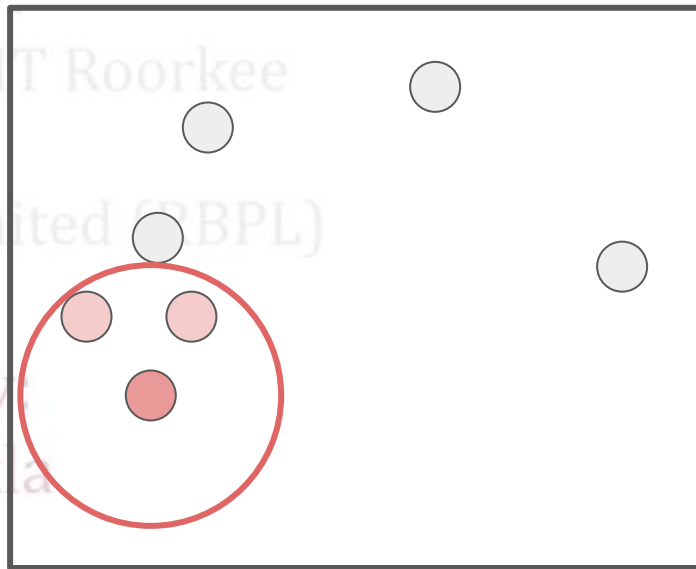
$\epsilon = 1$  and Min Points = 2



## DBSCAN Point Types:

- Core

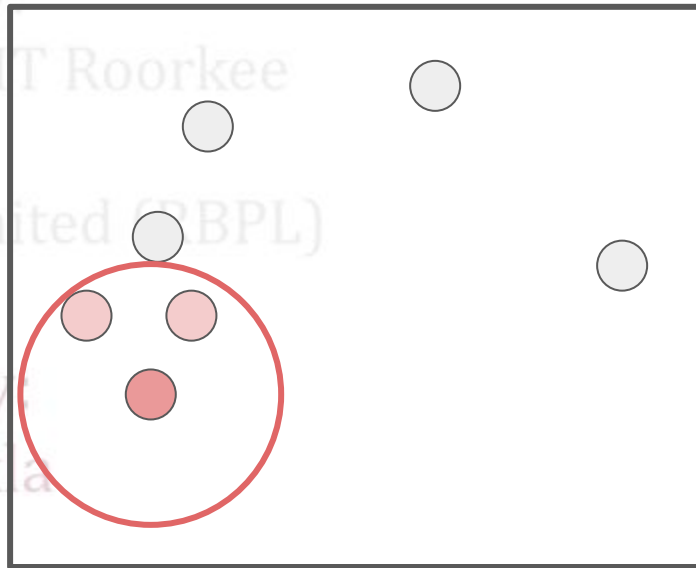
$\epsilon = 1$  and Min Points = 2



## DBSCAN Point Types:

- Core:
  - Point with min. points in epsilon range.

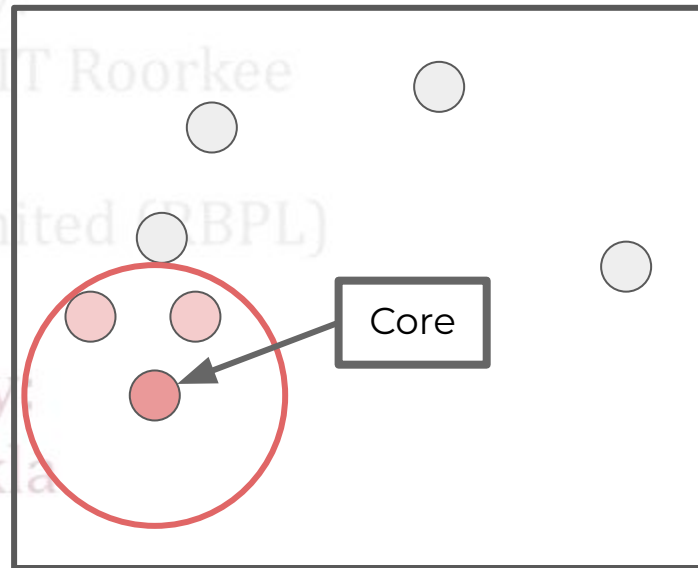
$\epsilon = 1$  and Min Points = 2



## DBSCAN Point Types:

- Core:
  - Point with min. points in epsilon range.

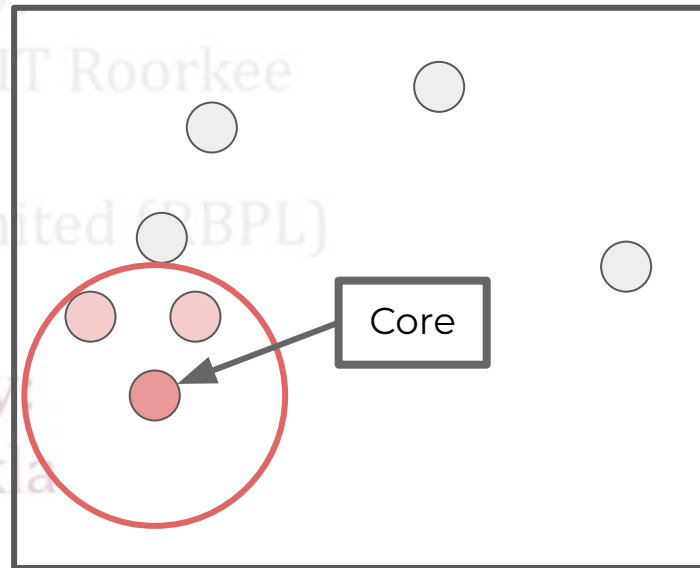
$\epsilon = 1$  and Min Points = 2



## DBSCAN Point Types:

- Core:
  - Point with min. points in epsilon range (including itself).

$\epsilon = 1$  and Min Points = 3

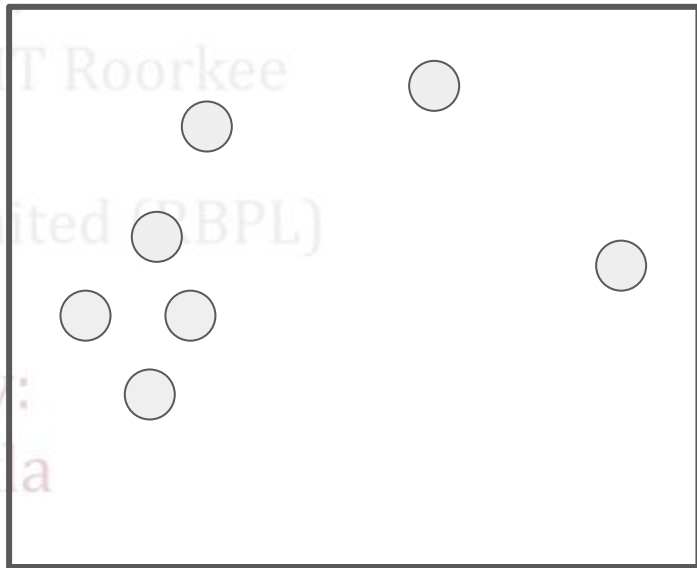




## DBSCAN Point Types:

- Border:
  - In epsilon range of core point, but does not contain min. number of points.

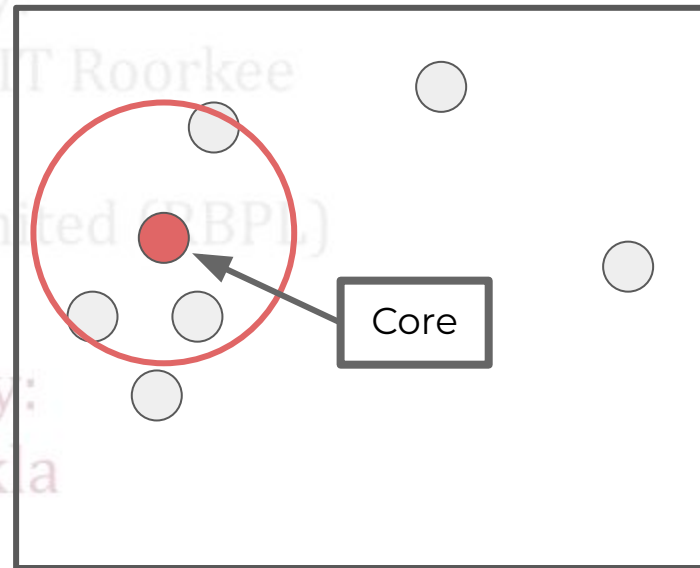
$\epsilon = 1$  and Min Points = 3



## DBSCAN Point Types:

- Border:
  - In epsilon range of core point, but does not contain min. number of points.

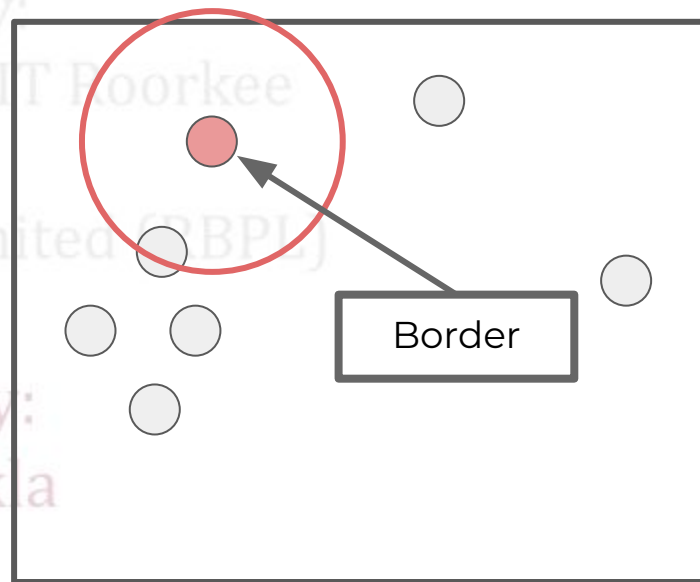
$\epsilon = 1$  and Min Points = 3



## DBSCAN Point Types:

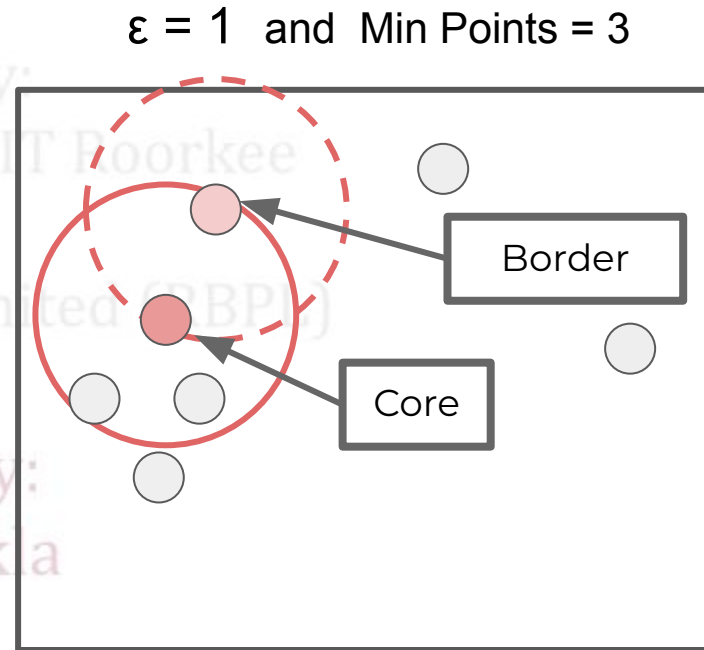
- Border:
  - In epsilon range of core point, but does not contain min. number of points.

$\epsilon = 1$  and Min Points = 3



## DBSCAN Point Types:

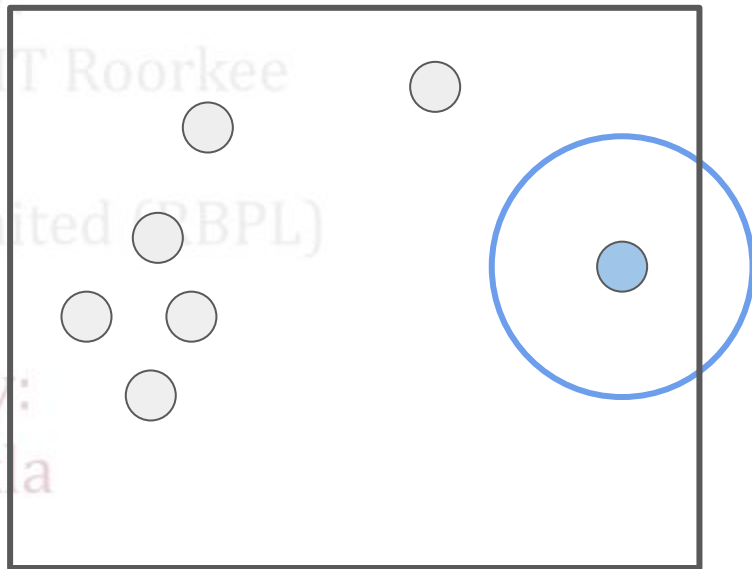
- Border:
  - In epsilon range of core point, but does not contain min. number of points.



## DBSCAN Point Types:

- Outlier:
  - Can not be “reached” by points in a cluster assignment.

$\epsilon = 1$  and Min Points = 3



Let's review the actual process of DBSCAN for  
assigning clusters.

Conducted by:  
iHUB Divya Sampark, IIT Roorkee  
and  
Ritvij Bharat Private Limited (RBPL)

Presented by:  
Shreyas Shukla

## DBSCAN Procedure:

- Pick a random point not yet assigned.
- Determine the point type.
- Once a **core** point has been found, add all directly reachable points to the same cluster as core.
- Repeat until all points have been assigned to a cluster or as an outlier.

Presented by:

# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

Conducted by:  
iHUB Divya Sampark, IIT Roorkee  
and  
Ritvi Bharat Private Limited (RBPL)

## Coding Example on Data Sets

Presented by:  
Shreyas Shukla



# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

Conducted by:

iHUB Divya Sampark, IIT Roorkee

and

Ritvij Bharat Private Limited (RBPL)

## Key Hyperparameters

Presented by:

Shreyas Shukla

## Two key hyperparameters for DBSCAN:

- Epsilon:
  - Distance extended from a point to search for Min. Number of Points.
- Min. Number of Points:
  - Min. Number of Points within Epsilon distance to be a core point.

Adjusting these hyperparameters have two main outcomes:

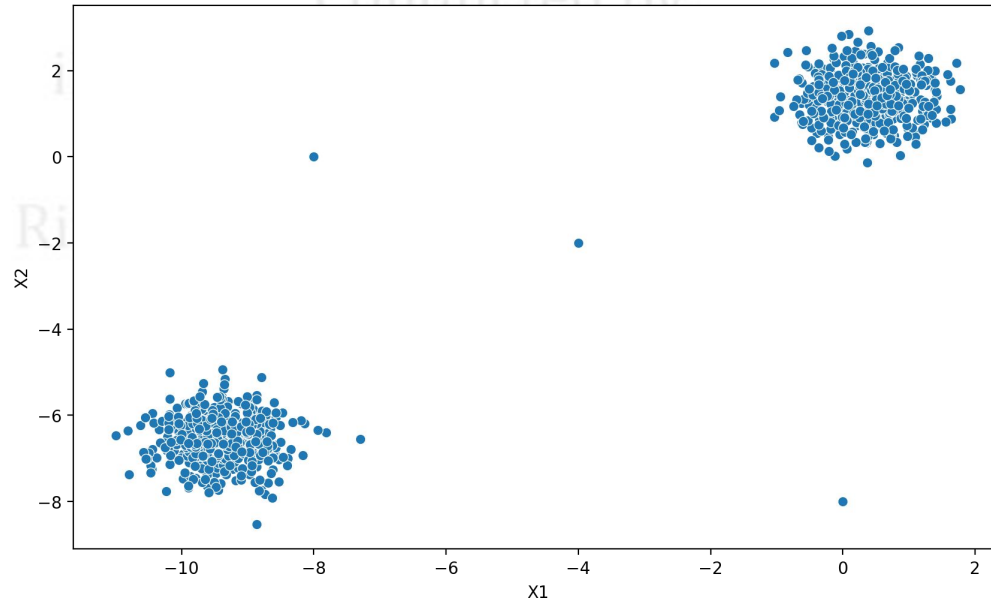
- Changing number of clusters.
- Changing what is an outlier point.

Presented by:  
Shreyas Shukla

# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

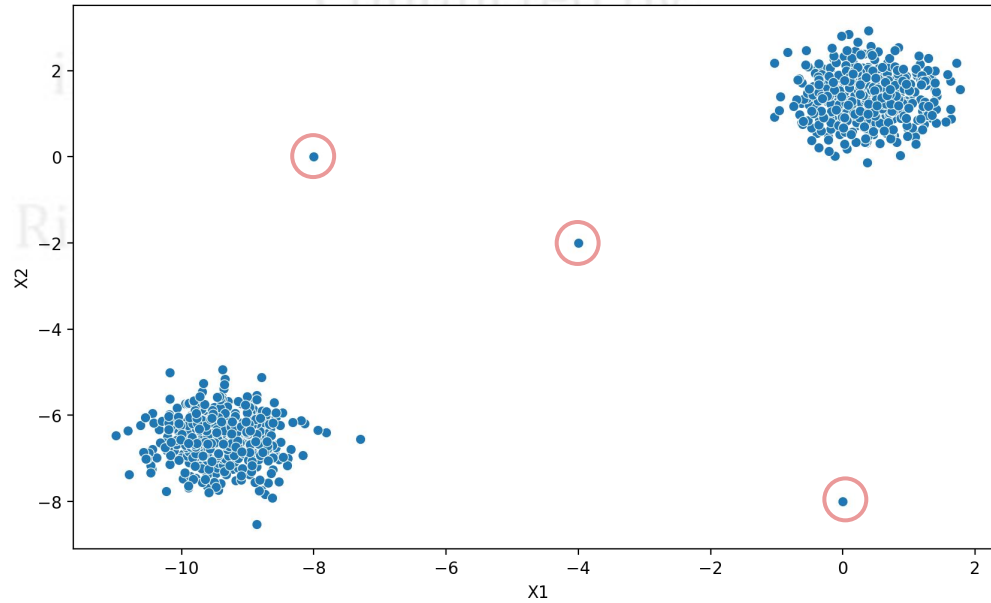
Conducted by:



# An Introduction to Machine Learning with Python Programming

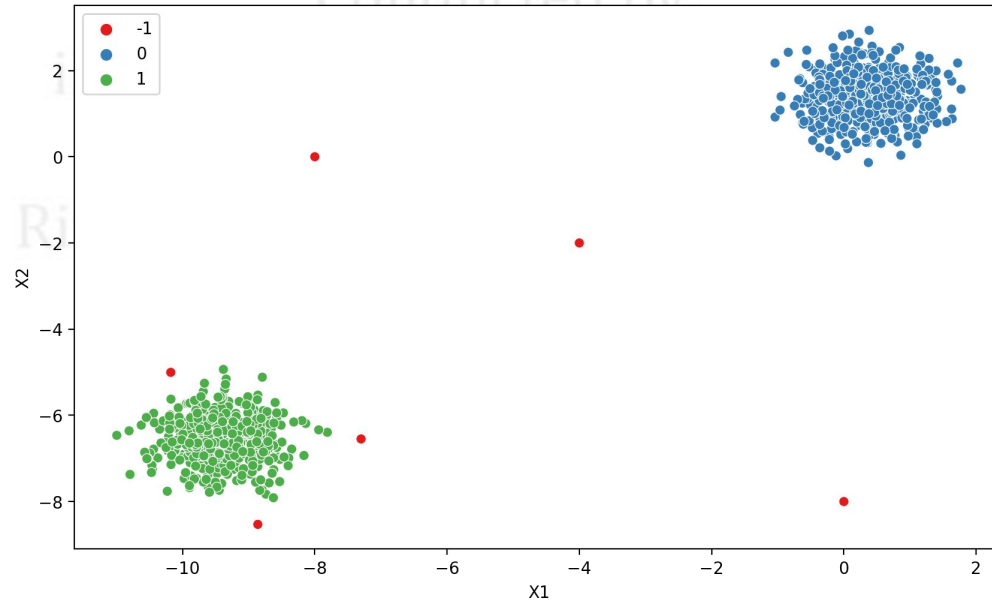
11 Sep 2023 - 20 Oct 2023

Conducted by:



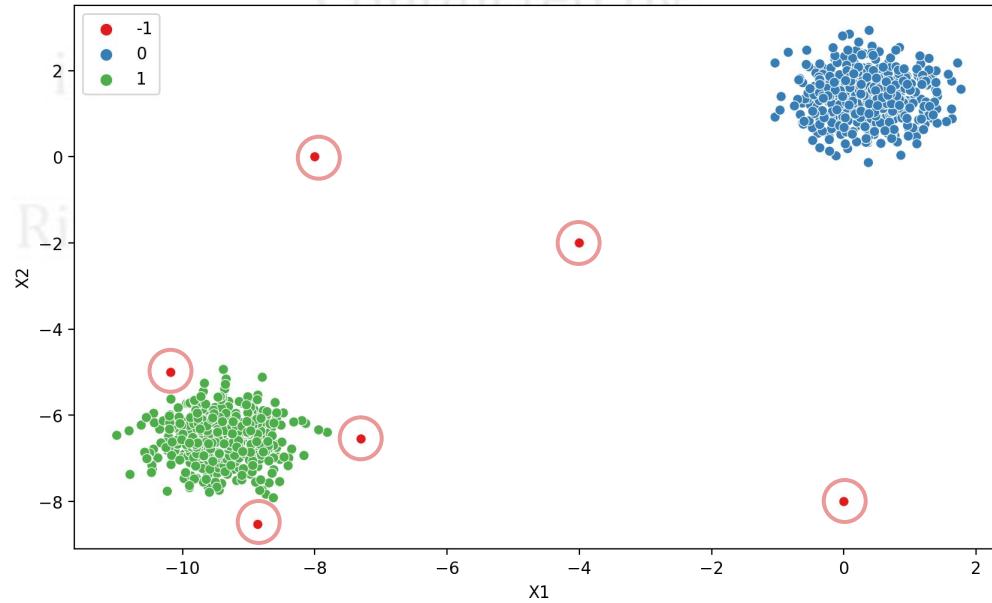
# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023



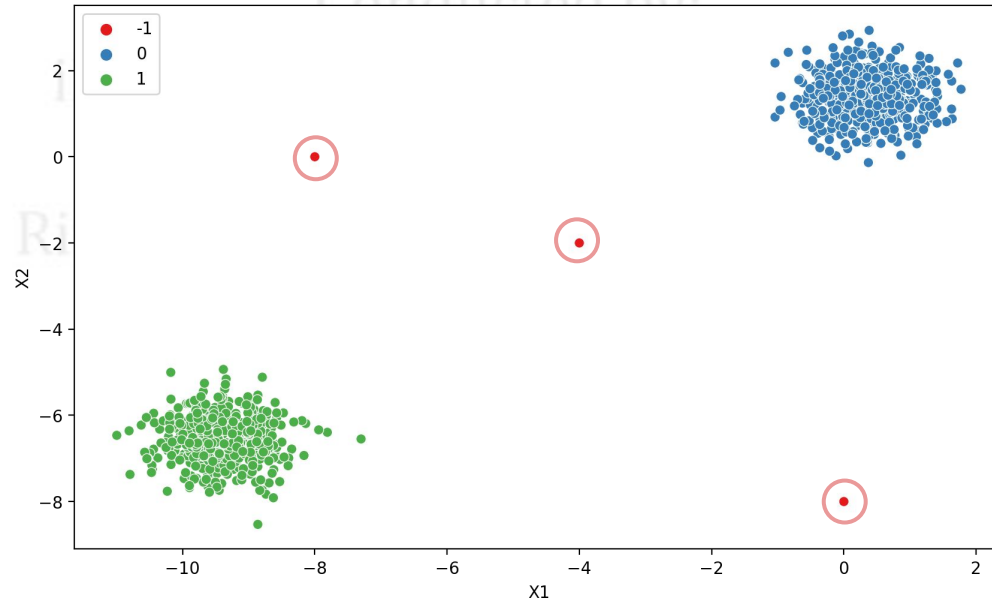
# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023



# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023





## Epsilon Intuition:

- Increasing epsilon allows more points to be **core** points which also results in more **border** points and less outlier points.
- Imagine a huge epsilon, all points would be within the neighborhood and classified as the same cluster!
- Decreasing epsilon causes more points not to be in range of each other, creating more unique clusters.
- Imagine a tiny epsilon, the range would not extend far out enough to come into contact with any other points!

Presented by  
Shreyas Shukla

Methods for finding an epsilon value:

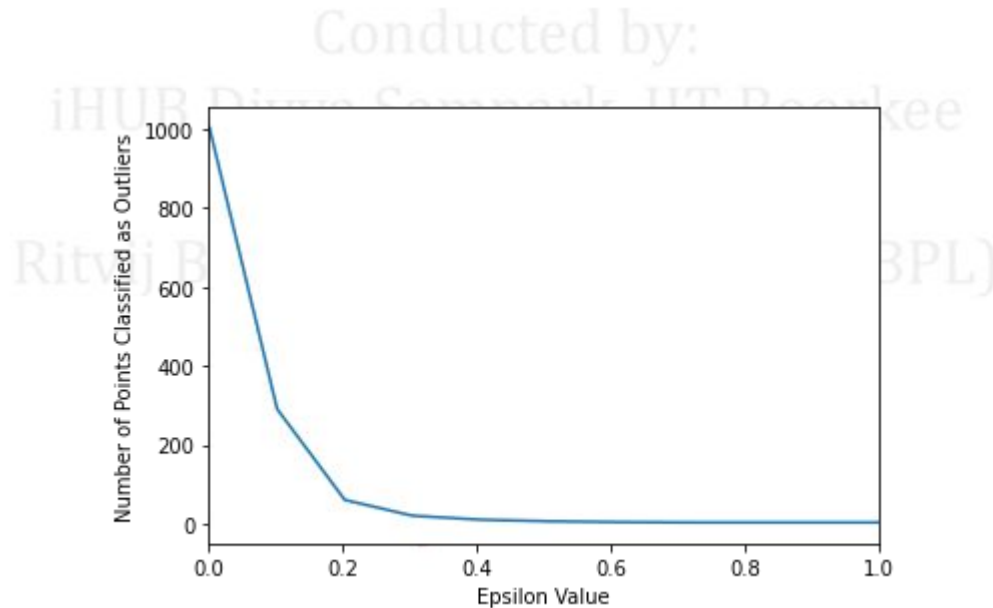
Run multiple DBSCAN models varying epsilon and measure:

- Number of Clusters
- Number of Outliers
- Percentage of Outliers
- 

Extremely dependent on the particular data set and domain space.

Requires user to have some expectation or intuition about number of clusters and relative percentage of outliers.

## Plot “elbow/knee” diagram comparing epsilon values:



## Minimum Number of Samples/Points:

- Number of samples in a neighborhood for a point to be considered as a **core** point (including the point itself).

Presented by:  
Shreyas Shukla

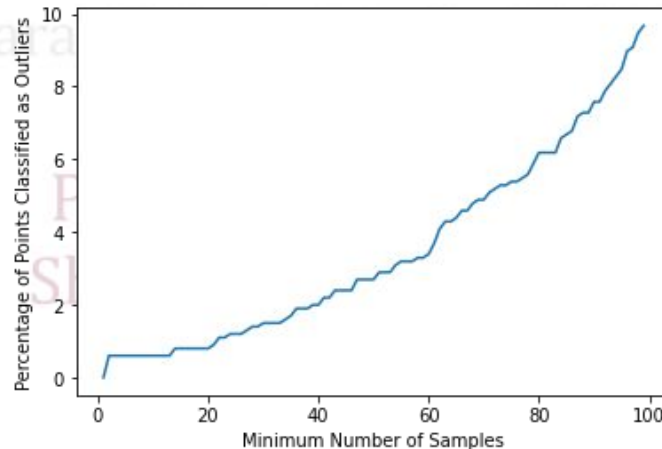
## Min. Number of Samples Intuition:

- Increasing to a larger number of samples needed to be considered a core point, causes more points to be considered unique outliers.
- Imagine if min. number of samples was close to total number of points available, then very likely all points would become outliers.

Presented by:  
Shreyas Shukla

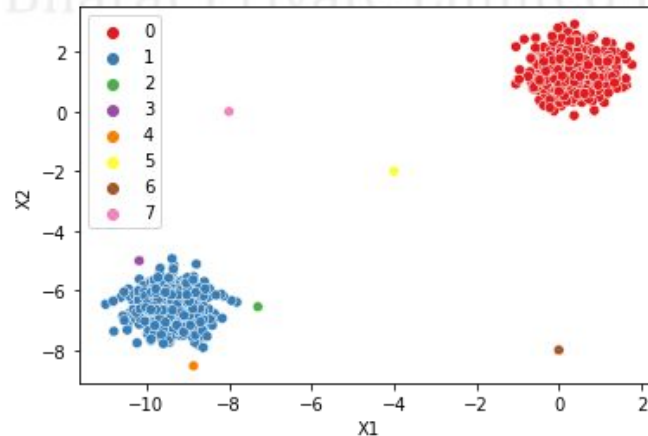
## Choosing Min. Number of Samples:

- Test multiple potential values and chart against number of outliers labeled.



## Min. Number of Samples Note:

- Useful to increase to create potential new small clusters, instead of complete outliers.



# An Introduction to Machine Learning with Python Programming

11 Sep 2023 - 20 Oct 2023

Let's continue by exploring hyperparameters with code and data examples!

Conducted by:  
iHUB Divya Sampark, IIT Roorkee  
and  
Ritvij Bharat Private Limited (RBPL)

Presented by:  
Shreyas Shukla