

DBSCAN

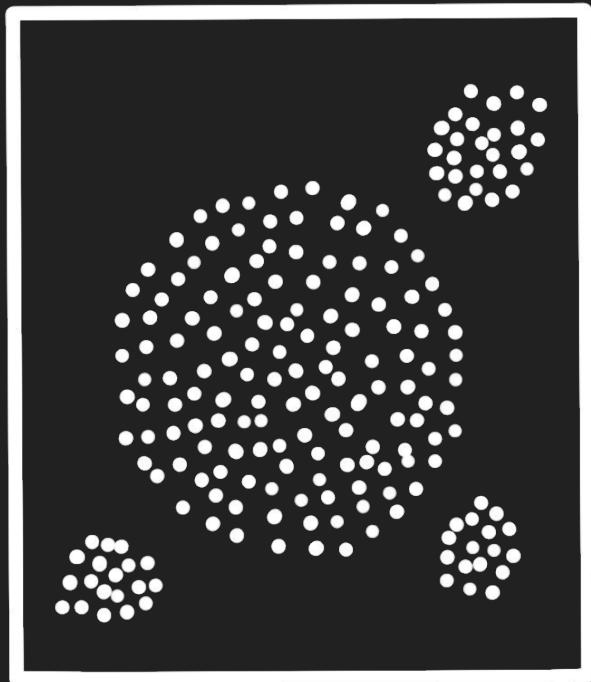
Chelsea Parlett-Pelleriti

DBSCAN

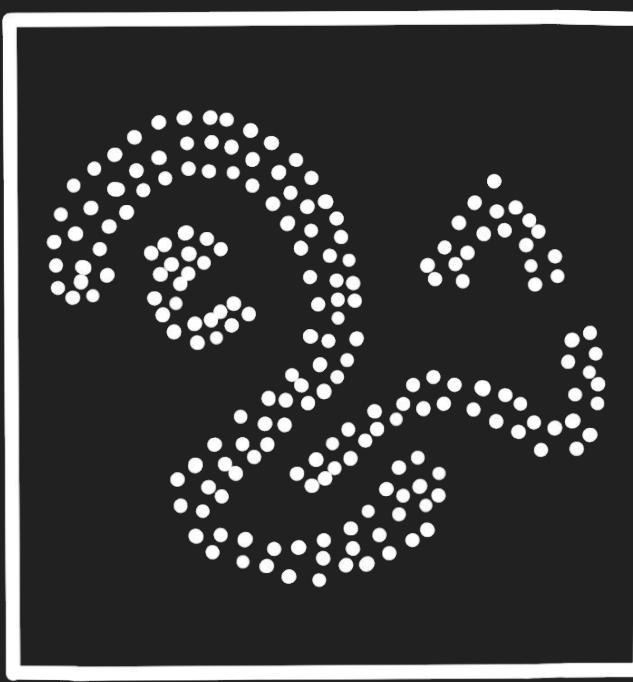
Density Based Spatial Clustering of Applications with Noise

- Distance Metric
- Epsilon (eps)
- Minimum Points (minpts)

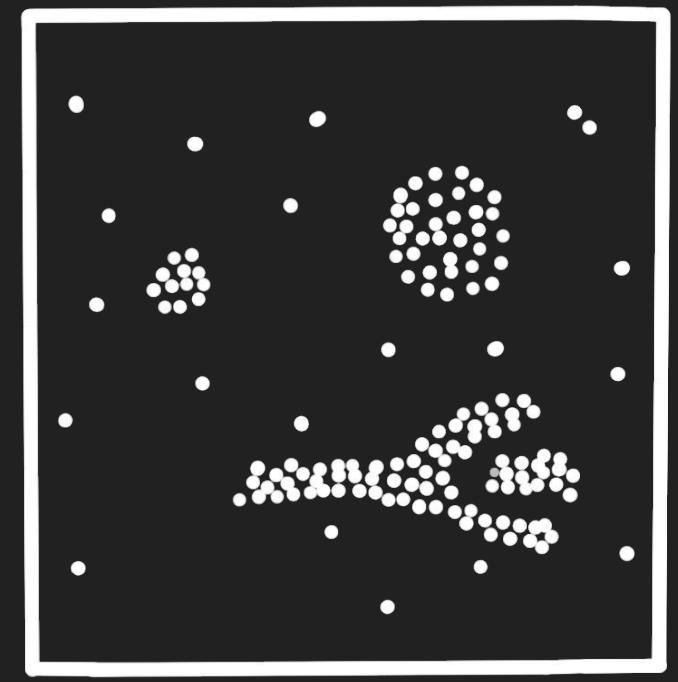
Benefits of DBSCAN



database 1



database 2

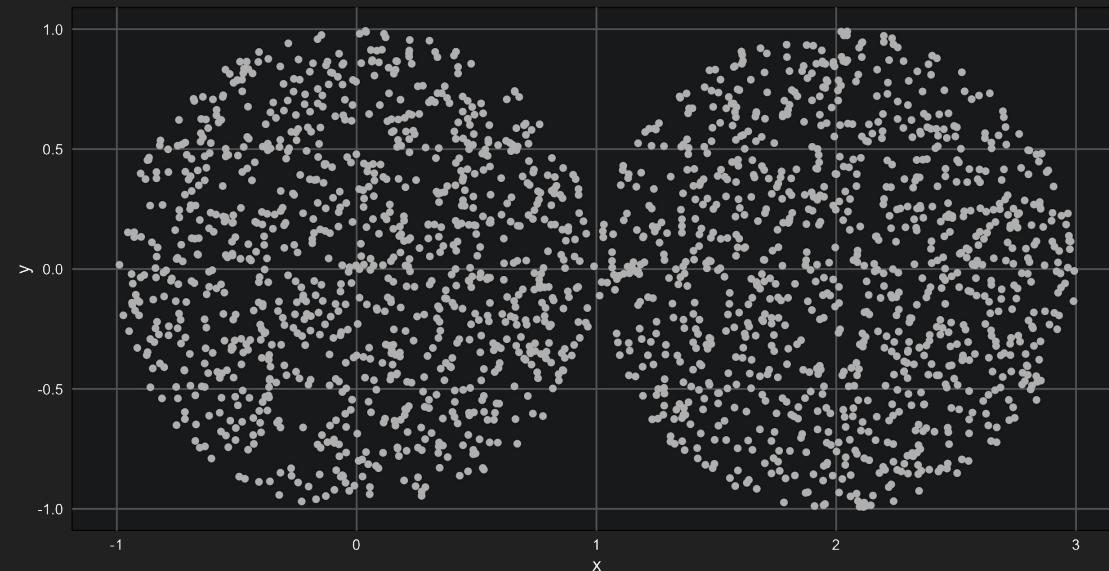


database 3

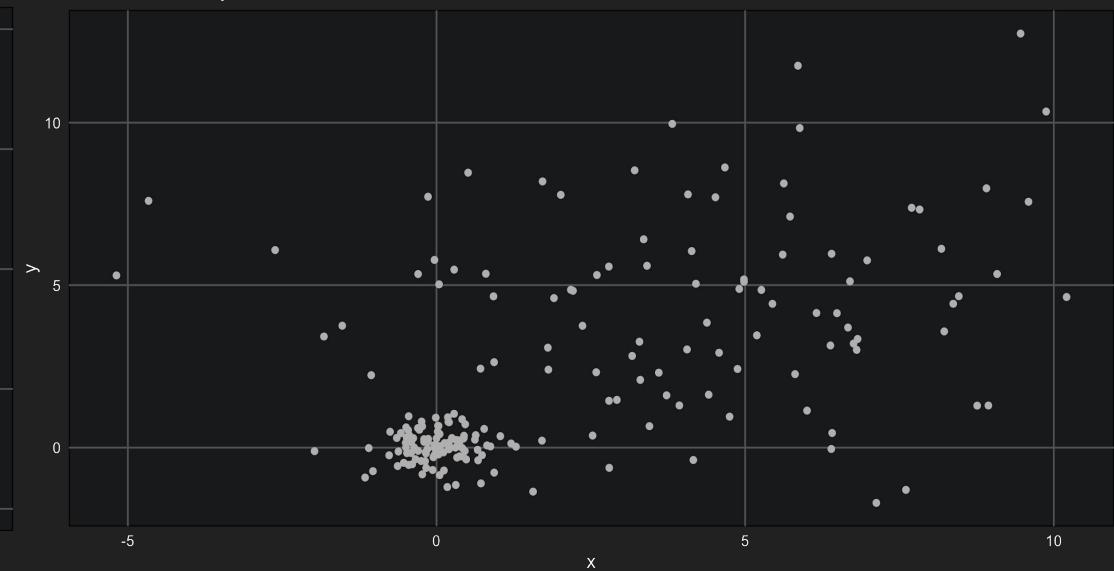
Disadvantages of DBSCAN

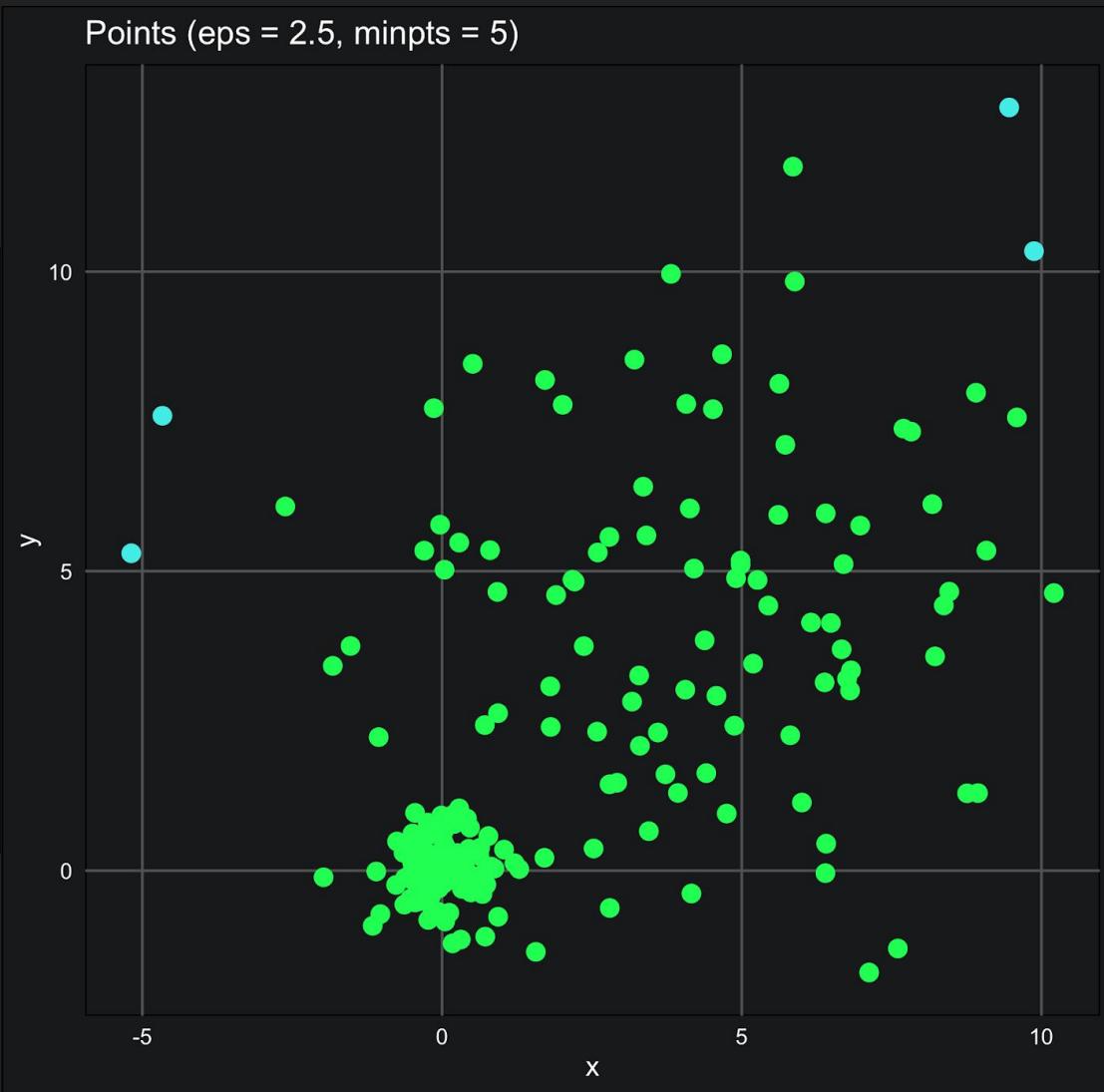
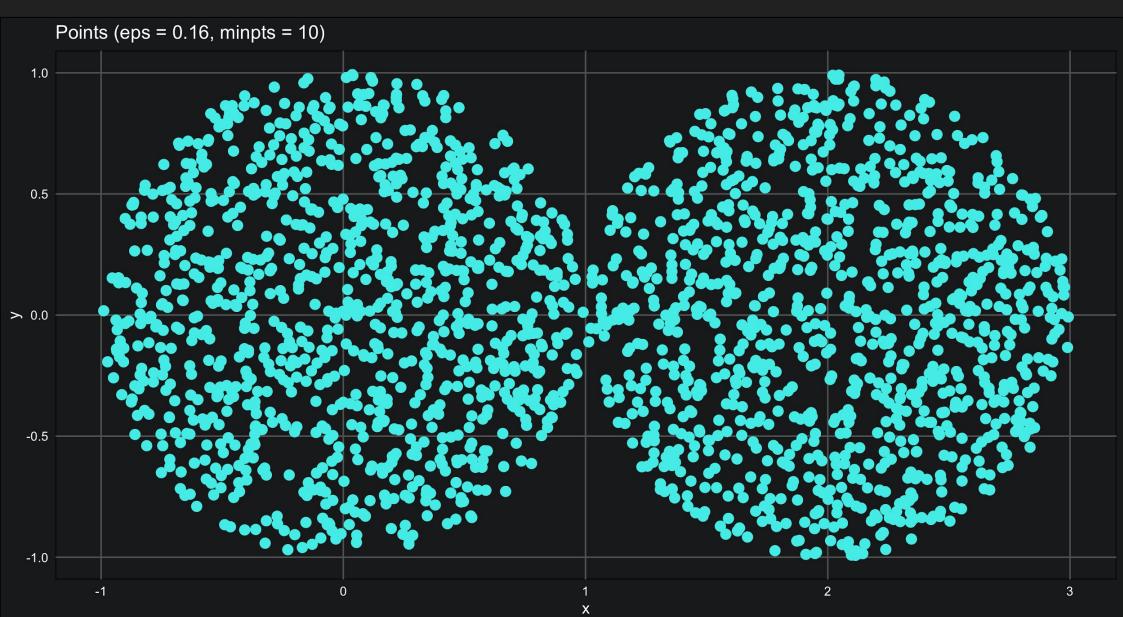
- Can be less effective in High Dimensional Data
- Not great with overlapping/touching clusters
- Suboptimal when clusters have different densities

Touching/Overlapping Clusters



Different Density Clusters



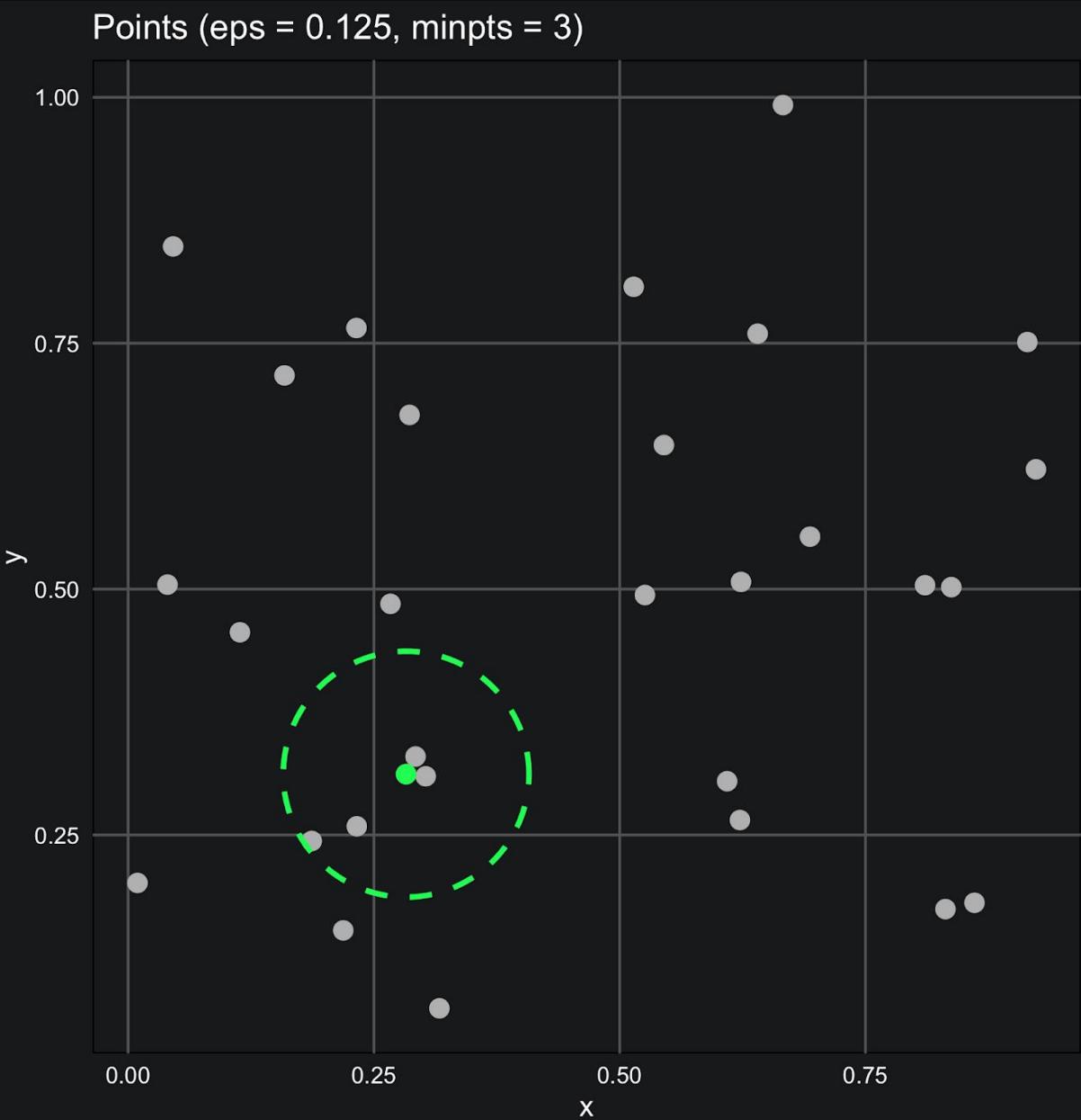


DBSCAN

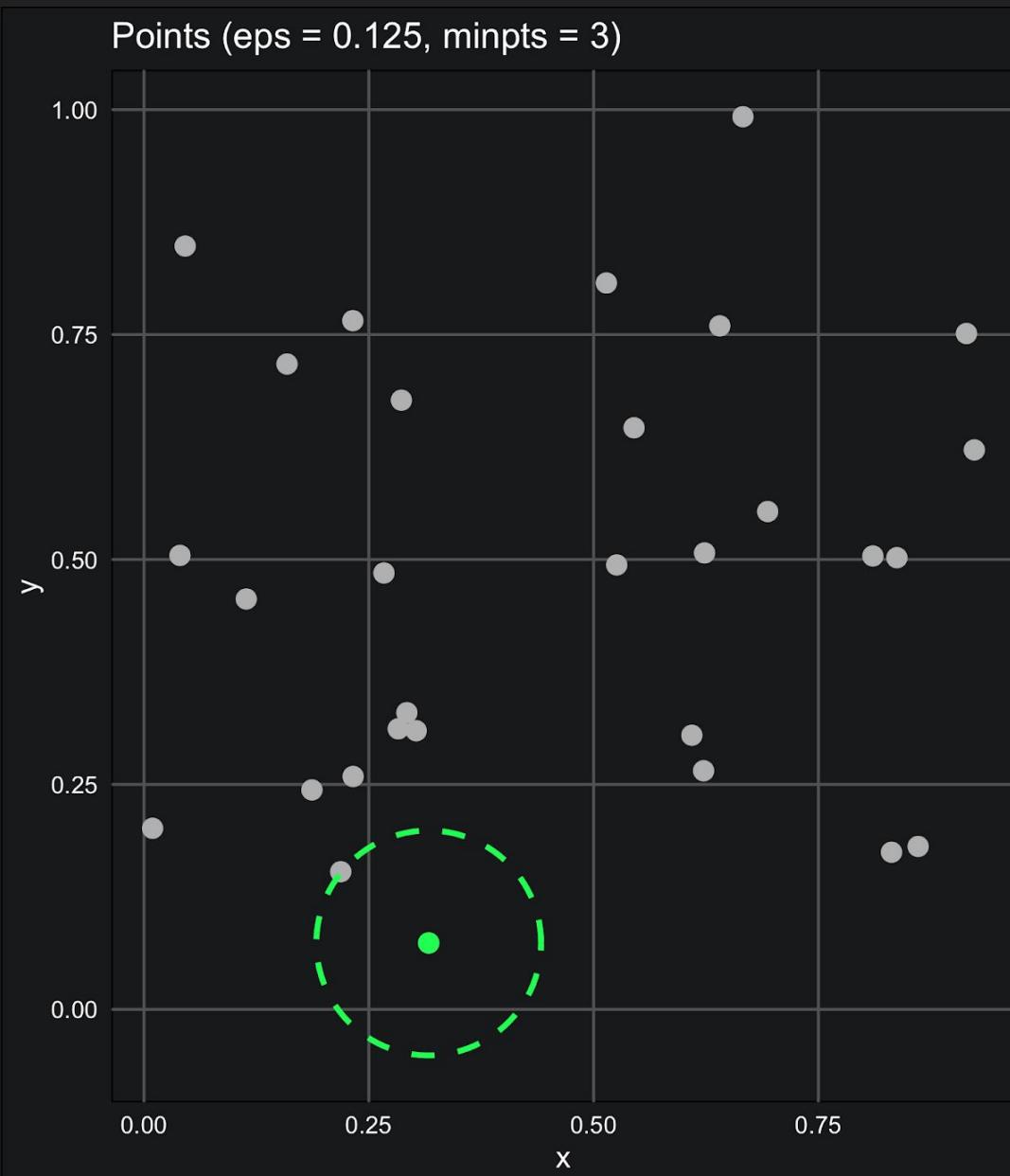
Density Based Spatial Clustering of Applications with Noise

- Distance Metric
- Epsilon (eps)
- Minimum Points (minpts)

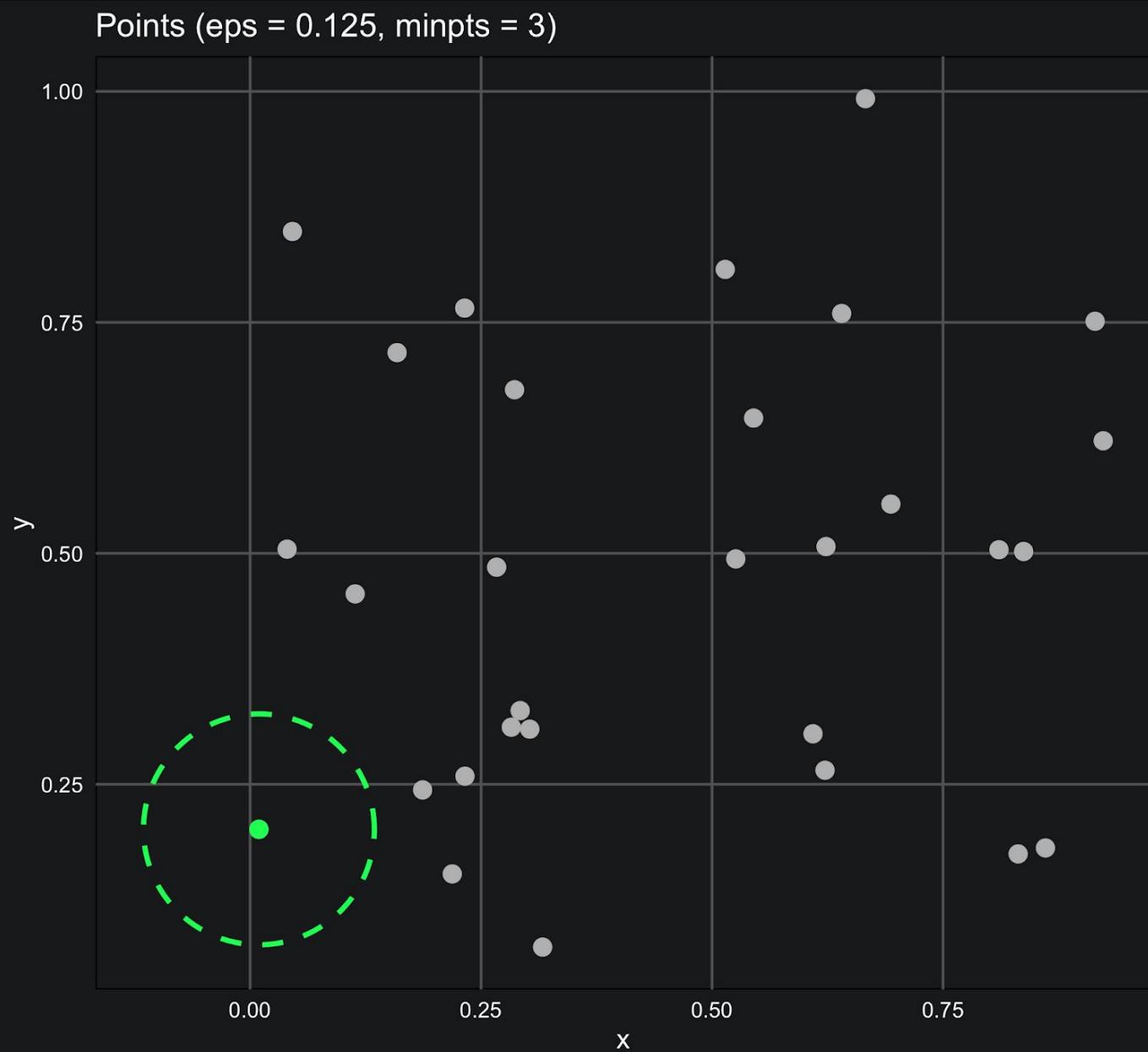
Core Point: p is a core point if it has at least $minpts$ neighbors within eps distance of itself



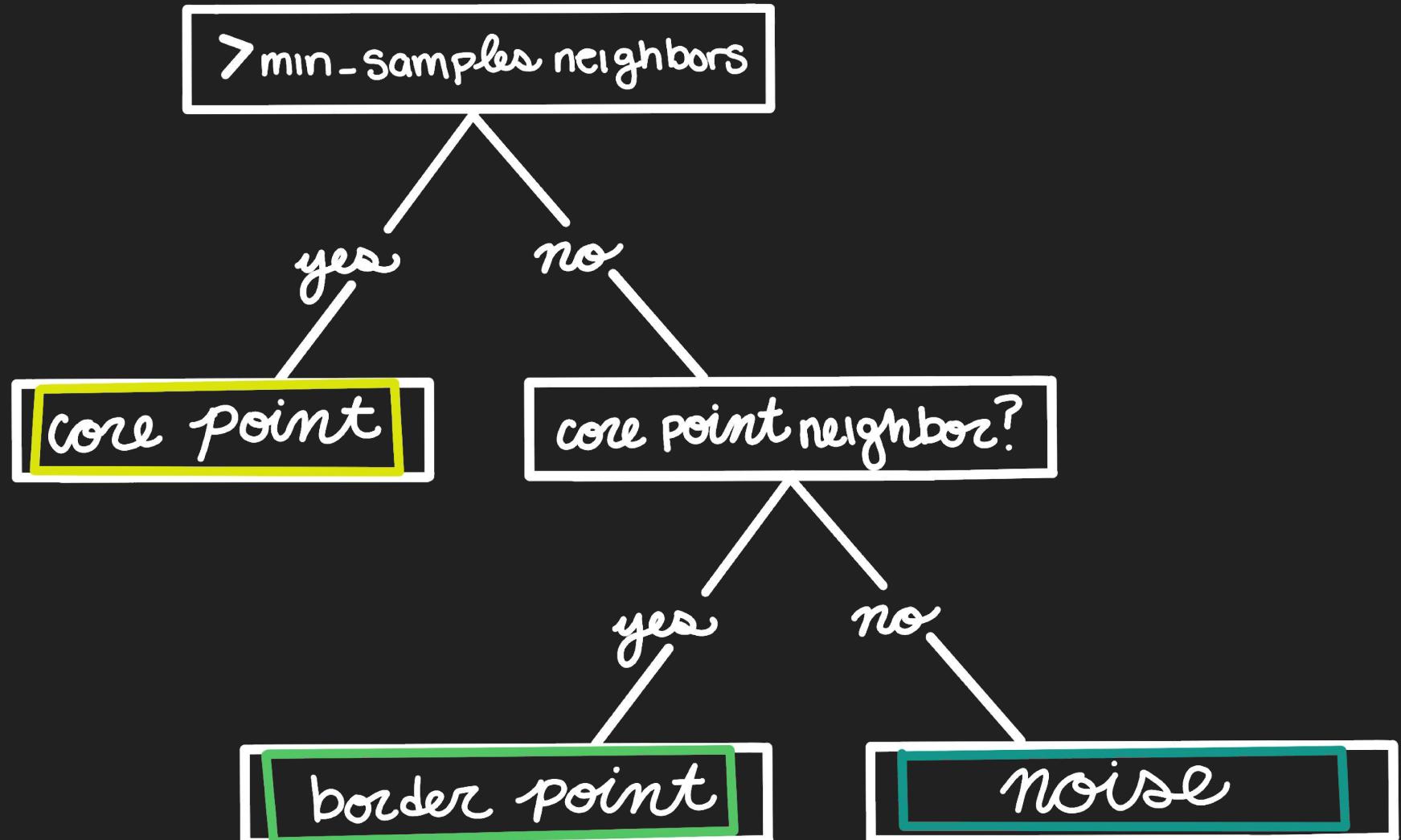
Border Point*: \mathbf{p} is a border point if it DOES NOT have at least $minpts$ neighbors within eps distance of itself, but is a neighbor of a core point



Noise: p is noise if it DOES NOT have at least $minpts$ neighbors within eps distance of itself, and IS NOT a neighbor of a core point



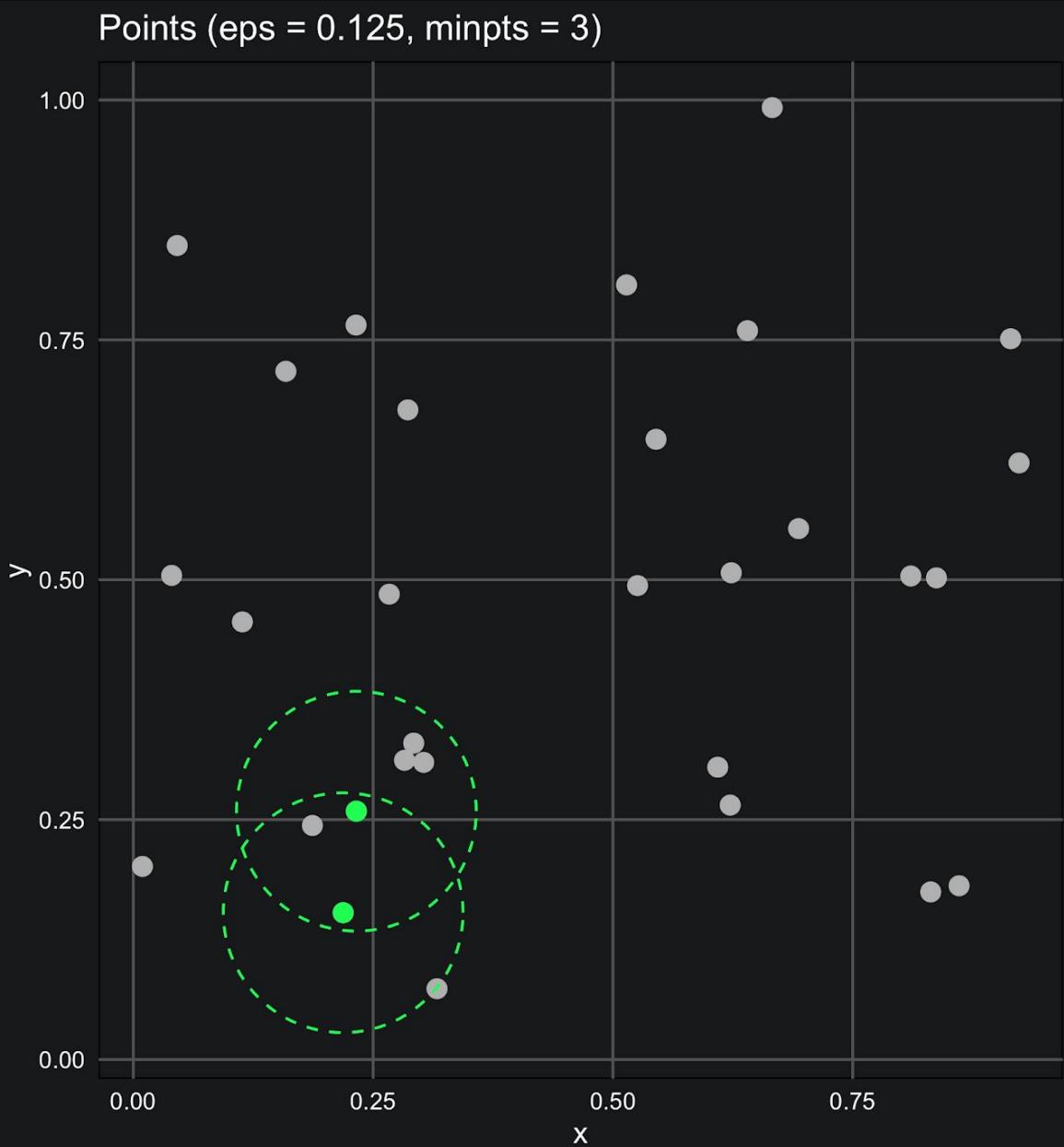
DBSCAN



Definitions

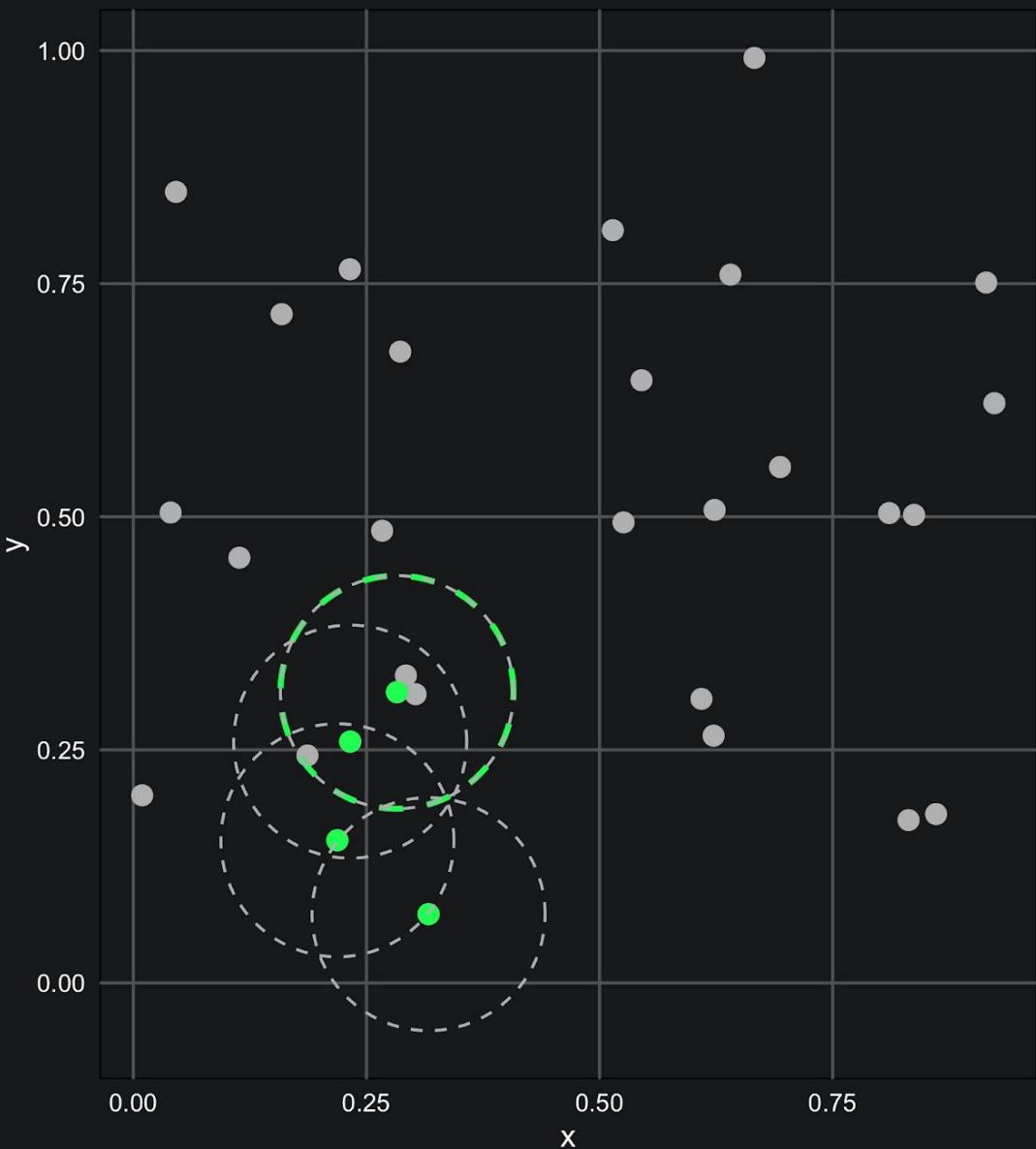
- Directly density reachable: \mathbf{p} is directly density reachable from core-point \mathbf{q} if it is in the neighborhood of \mathbf{q}
 - Density reachable: \mathbf{p} is directly reachable from \mathbf{q} if there are a chain of points that are directly density reachable from \mathbf{q} to \mathbf{p}
 - Density connected: \mathbf{p} and \mathbf{q} are density connected if they are both density reachable from a third point, \mathbf{o}
-
- Cluster: choose core point \mathbf{q} , a cluster \mathbf{C} contains all points density reachable by \mathbf{q}
 - Noise: any point not in a cluster

Directly density reachable: \mathbf{p} is directly density reachable from core-point \mathbf{q} if it is in the neighborhood of \mathbf{q}



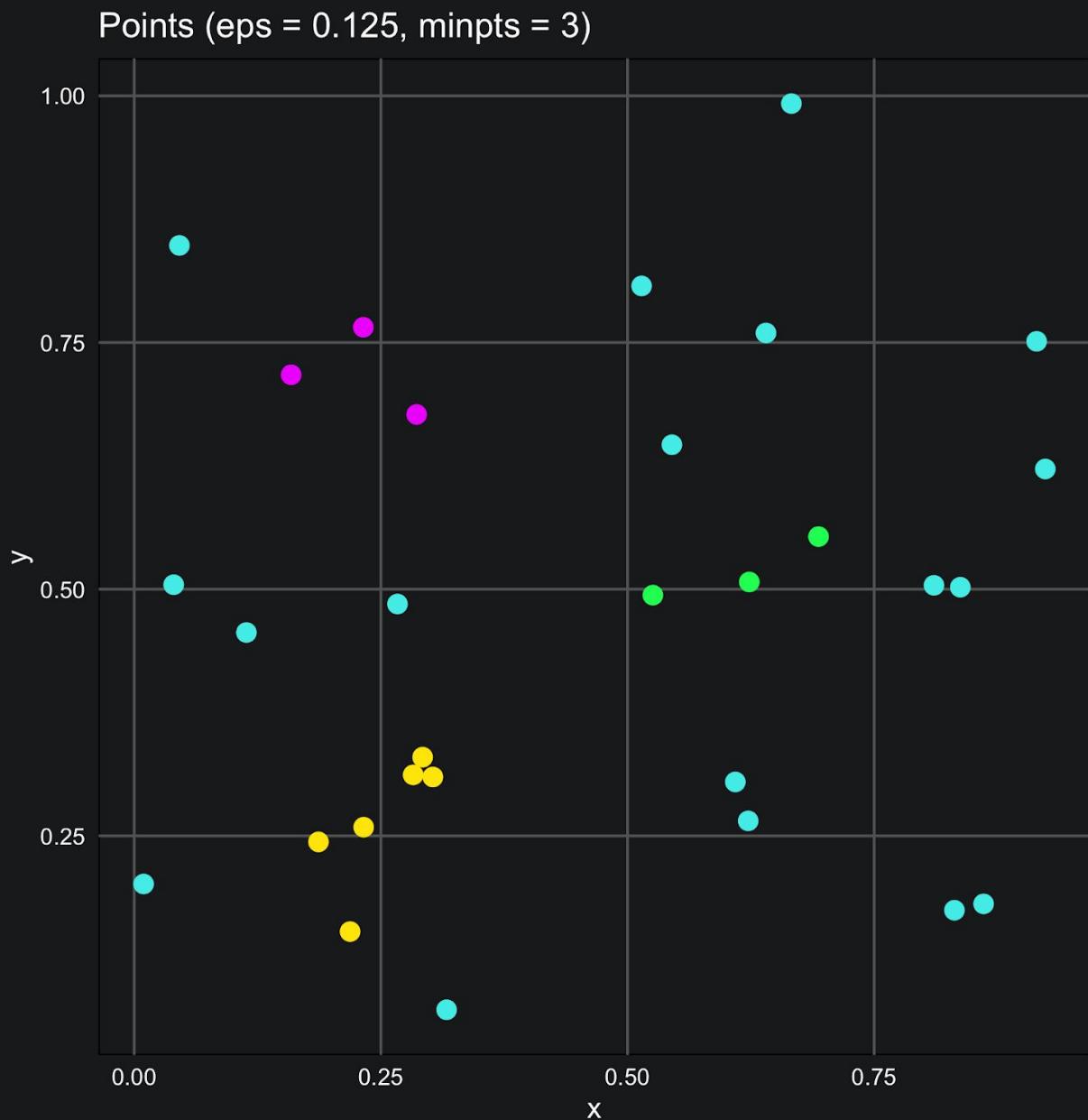
Density reachable: **p** is density reachable from **q** if there are a chain of points that are directly density reachable from **q** to **p**

Points (eps = 0.125, minpts = 3)

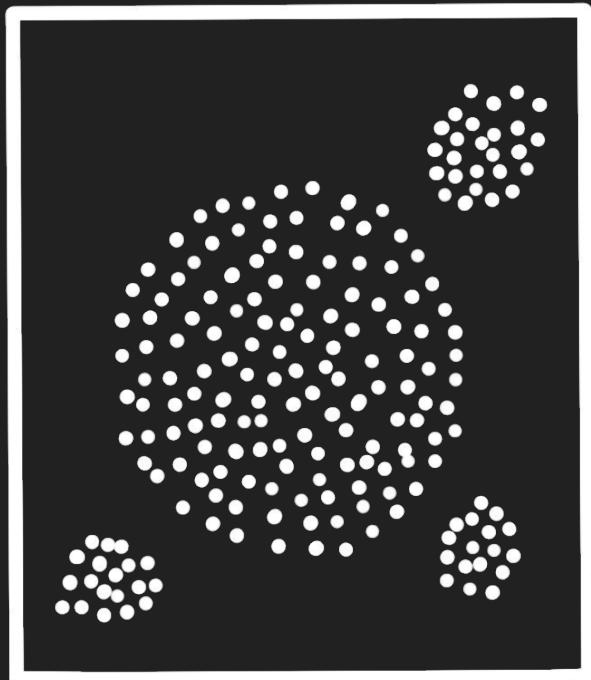


Cluster: choose core point \mathbf{q} , a cluster \mathbf{C} contains all points density reachable by \mathbf{q}

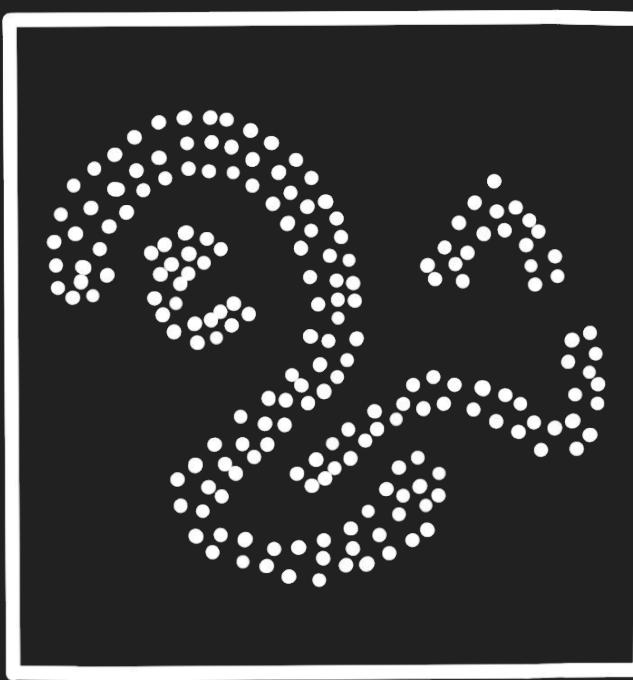
Noise: any point not in a cluster



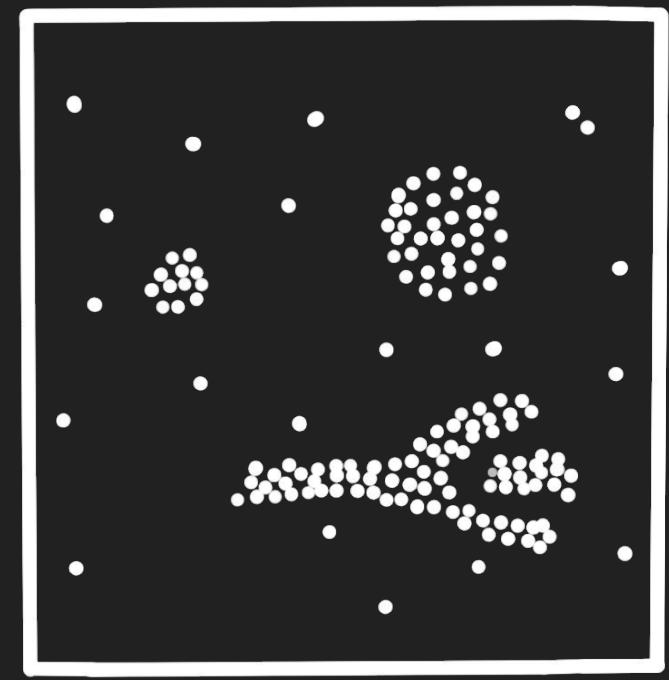
Benefits of DBSCAN



database 1

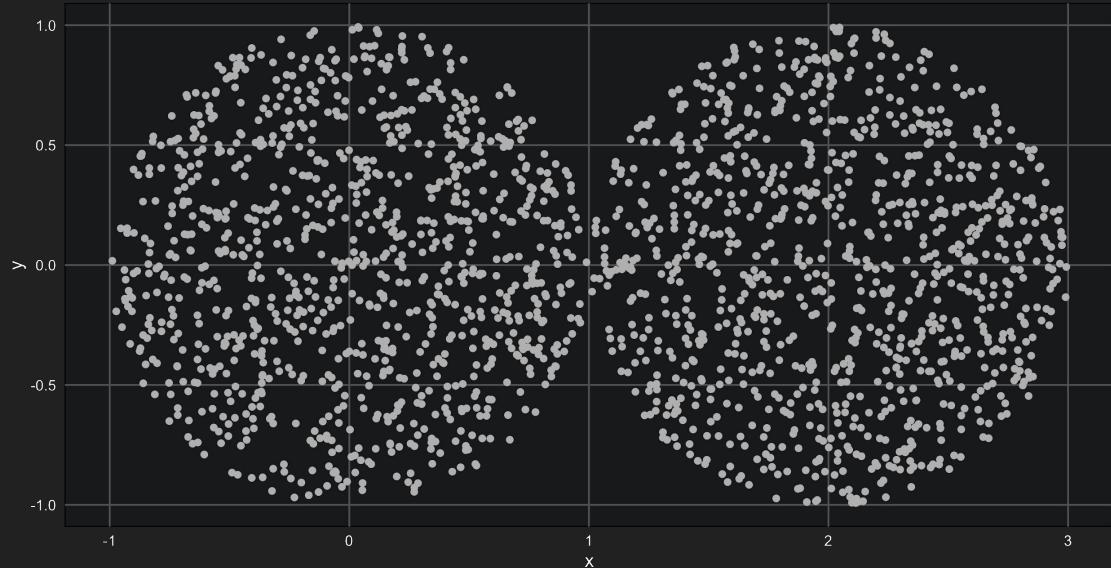


database 2

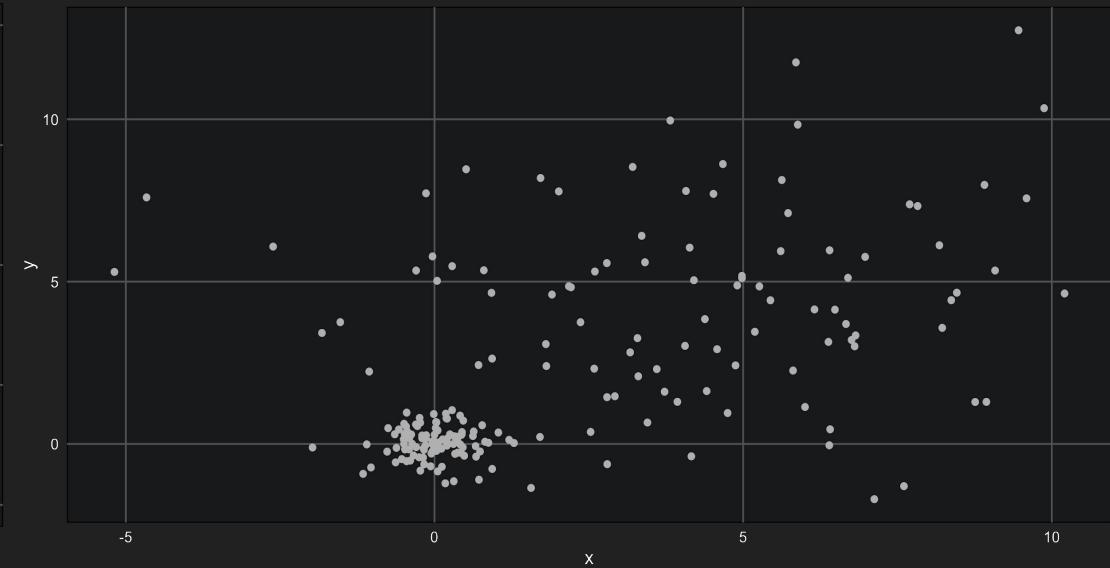


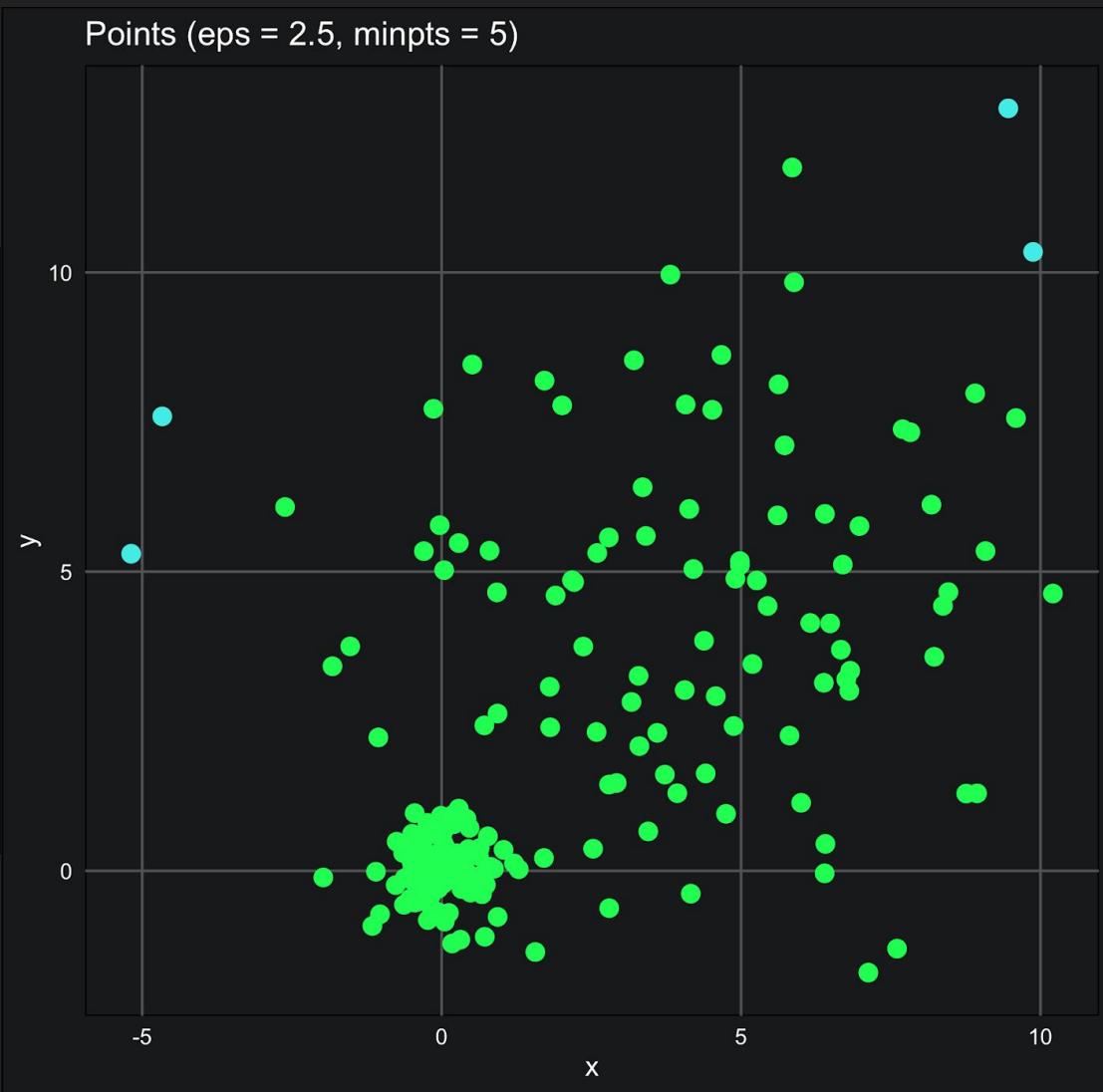
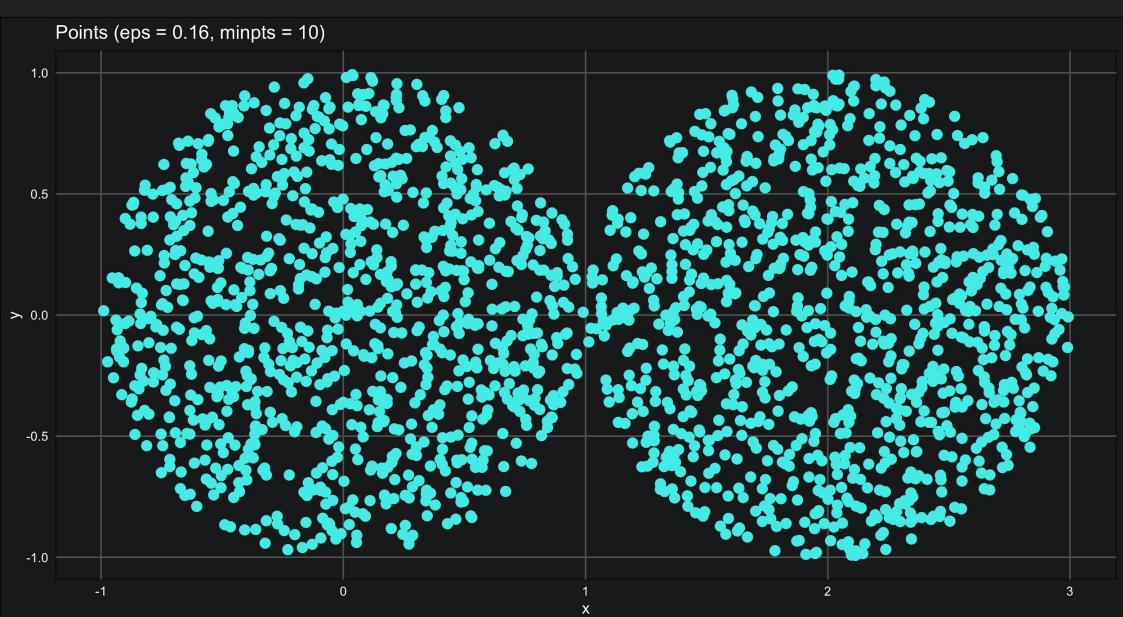
database 3

Touching/Overlapping Clusters



Different Density Clusters



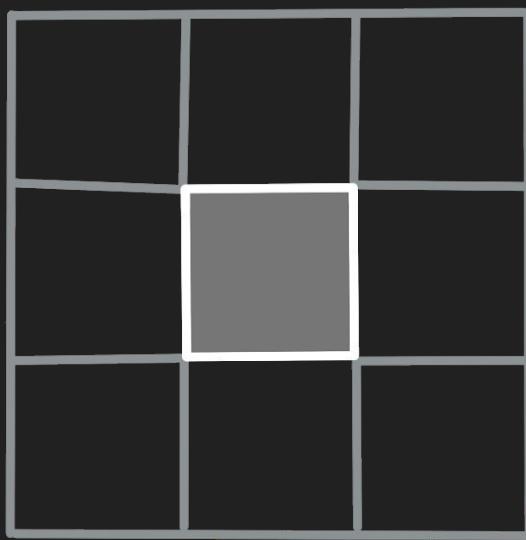


Choosing Minimum Points

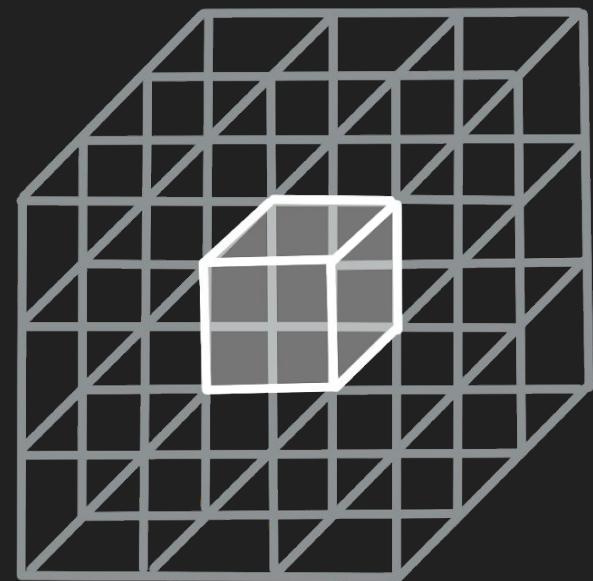
- Domain Knowledge + Distance Metrics
- More rows = larger min_pts
- More noise = larger min_pts
- More features = larger min_pts



(a)



(b)



(c)

Choosing Epsilon

- Elbow method
(k-dist)
- Domain Knowledge

Elbow Method for Choosing eps

