

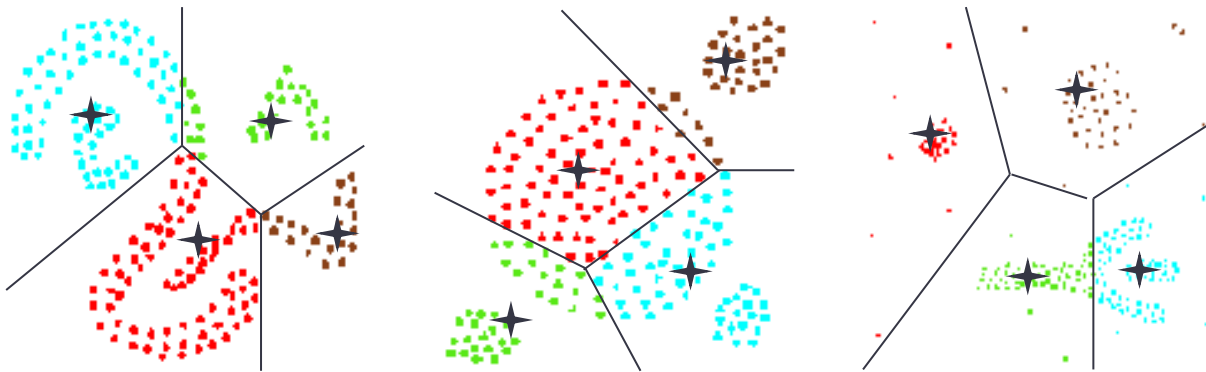
DBSCAN: Density-Based Spatial Clustering of Applications with Noise

Presented by Wondong Lee

Written by M.Ester, H.P.Kriegel, J.Sander and Xu.

Why Density-Based Clustering?

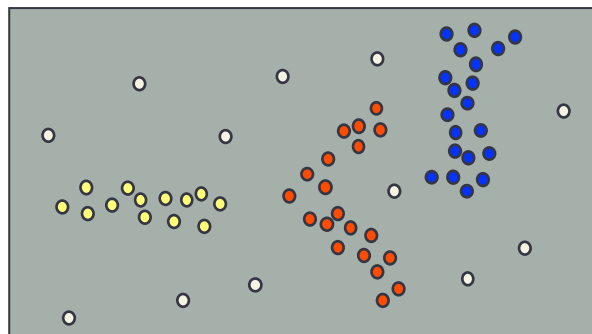
- Results of k-means algorithm for $k = 4$



➔ The result is not satisfiable!!

DBSCAN

- Relies on a density-based notion of cluster
- Discovers clusters of arbitrary shape in spatial databases with noise
- Basic Idea
 - Group together points in high-density
 - Mark as outliers → points that lie alone in low-density regions

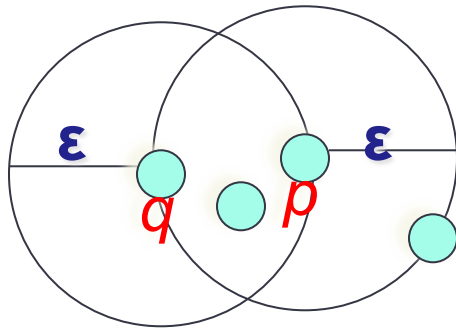


DBSCAN

- Local point density at a point p defined by two parameters
 - (1) $\epsilon \rightarrow$ radius for the neighborhood of point p :
 - ϵ -Neighborhood: all points within a radius of ϵ from the point p
$$N_{\epsilon}(p) := \{q \text{ in data set } D \mid \text{dist}(p, q) \leq \epsilon\}$$
 - (2) **MinPts** \rightarrow minimum number of points in the given neighborhood $N(p)$

High Density?

- ϵ -Neighborhood of an point contains at least *MinPts*



ϵ -Neighborhood of p

ϵ -Neighborhood of q

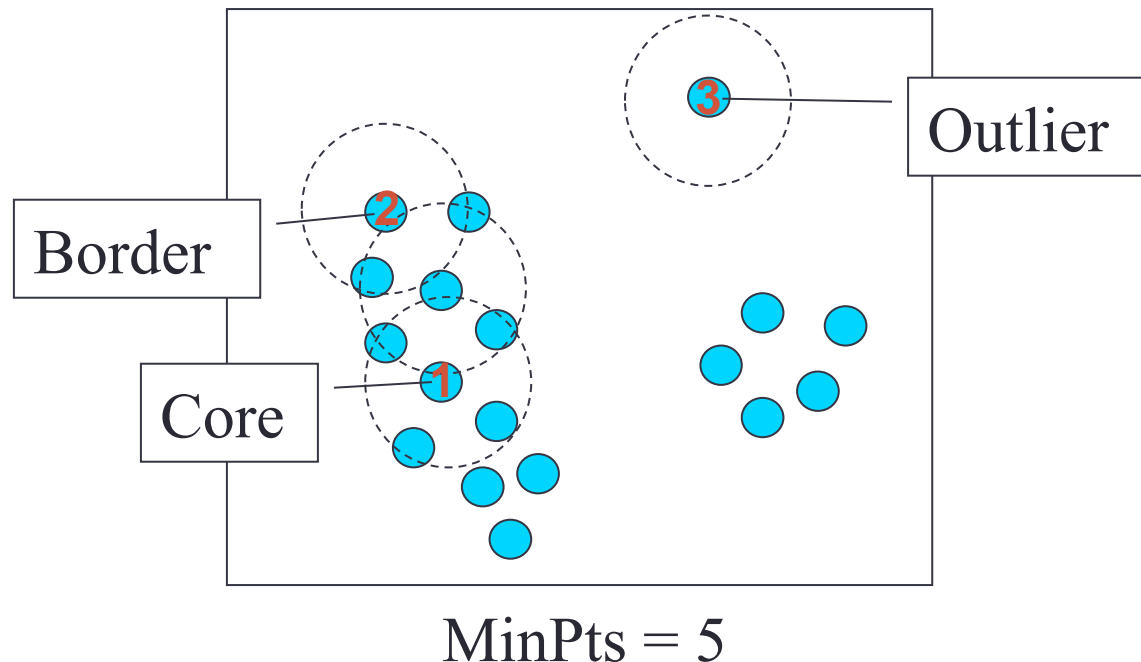
Q. When $\text{MinPts} = 4$?

Density of p is “high”

Density of q is “low”

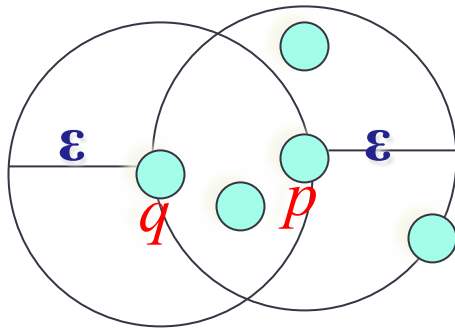
Core, Border & Outlier

- Three category for each point
 - **Core point**: if its density is **high**
 - **Border point**: density is low (**but in the neighborhood of a core point**)
 - **Noise point**: any point that is not a core point nor a border point



Density-Reachability

- Directly density-reachable
 - A point q is **directly density-reachable** from a point p :
 - If p is a **core point** and q is in p 's ϵ -neighborhood



Minpts = 4

Q. p is directly density-reachable from q ?

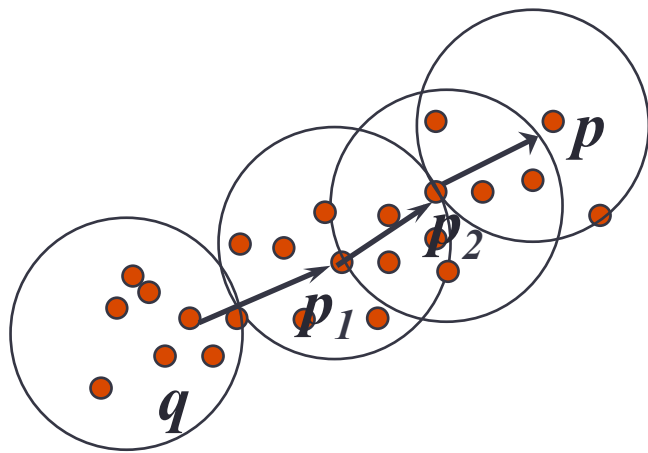
No, why?

Q. Density-reachability is **asymmetric**

Density-Reachability

- Density-reachable

- A point p is **density-reachable** from a point q
 - If there is a chain of points p_1, \dots, p_n , with $p_1 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i



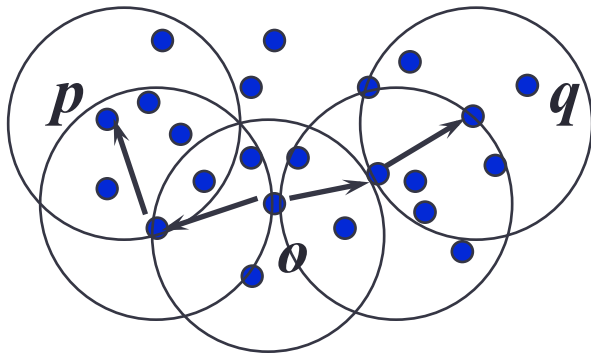
- p_1 is directly density-reachable from q
- p_2 is directly density-reachable from p_1
- p is directly density-reachable from p_2
- There is a chain from q to p ($q \rightarrow p_1 \rightarrow p_2 \rightarrow p$)

Q. q is density-reachable from p ?

No, why?

Density-Connectivity

- Density-connected
 - A pair of points p and q are density-connected
 - If they are commonly **density-reachable** from a point o



MinPts = 7

Q. o is density-reachable from p ?

Yes, why?

Q. Density-connectivity is **symmetric**

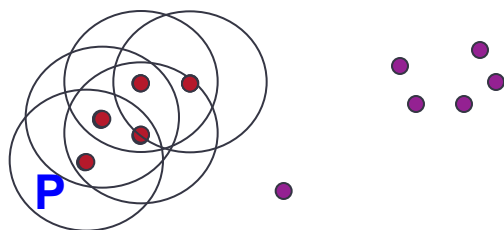
Formal Description of Cluster

- Given a data set D , parameter ε and *MinPts*,
- A cluster C is a subset of D satisfying two criteria:
 - **Maximality**
 - $\forall p, q$ if $p \in C$ and if q is **density-reachable** from p , then also $q \in C$
 - **Connectivity**
 - $\forall p, q \in C$, p and q are **density-connected**
- **Note:** cluster contains *core points* as well as *border points*

)

Parameter

- $\varepsilon = 2, MinPts = 3$



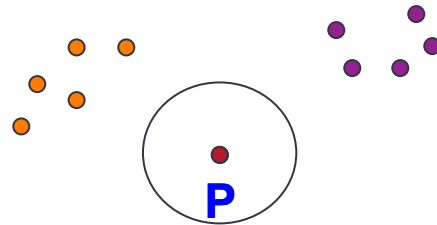
```

if  $p$  is not yet classified then
  if  $p$  is a core-point then
    collect all points density-reachable from  $p$ 
    and assign them to a new cluster.
  else
    else assign  $p$  to NOISE
    assign  $p$  to NOISE
  
```

example)

Parameter

- $\varepsilon = 2$, $MinPts = 3$

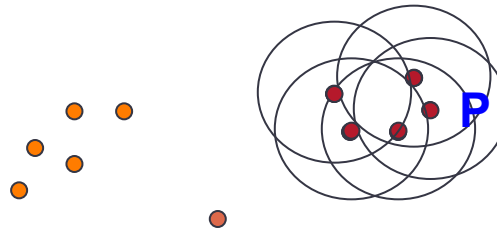


```
for  $p \in D$  do
  if  $p$  is not yet classified then
    if  $p$  is a core-point then
      collect all points density-reachable from  $p$ 
      and assign them to a new cluster.
    else
      assign  $p$  to NOISE
```

example)

Parameter

- $\varepsilon = 2$, $MinPts = 3$

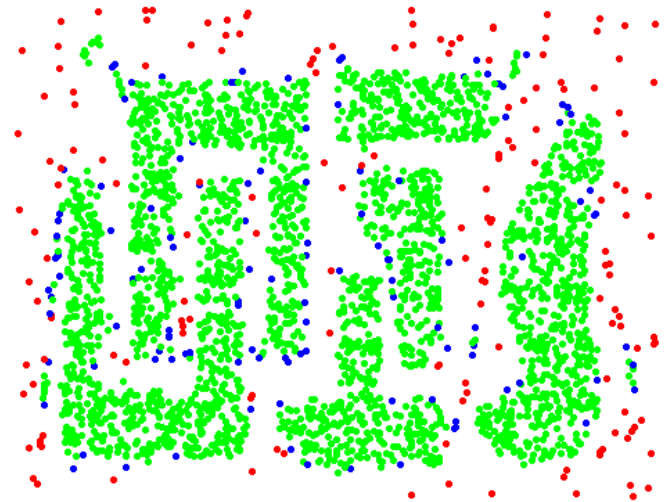


```
for  $p \in D$  do
  if  $p$  is not yet classified then
    if  $p$  is a core-point then
      collect all points density-reachable from  $p$ 
      and assign them to a new cluster.
    else
      assign  $p$  to NOISE
```

Example



Original Points

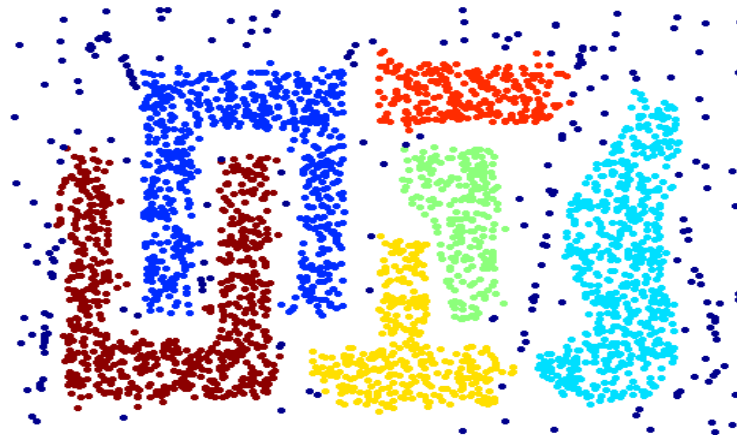


Point types: **core**,
border and **outliers**

$\epsilon = 10$, MinPts = 4

When DBSCAN Works Well

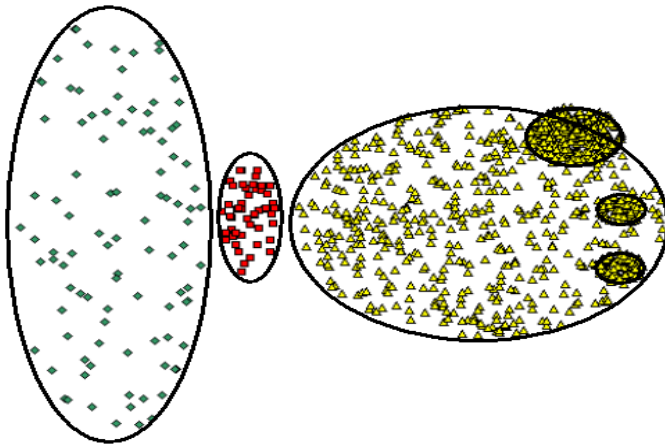
- Resistant to Noise
- Can handle clusters of different shapes and sizes



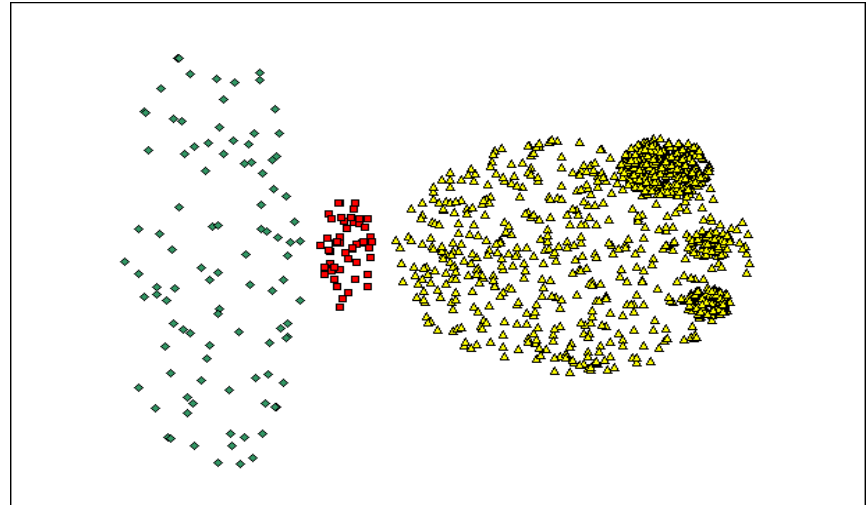
Clusters

When DBSCAN Does Not Work Well

- Cannot handle varying densities



Original Points

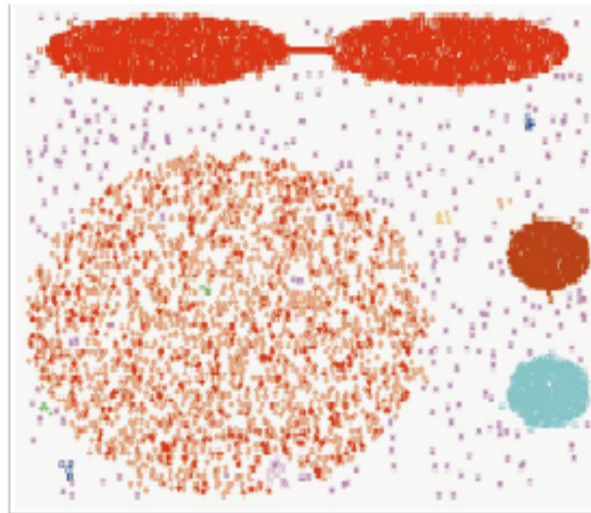


($\epsilon = 9.92$, MinPts=4)

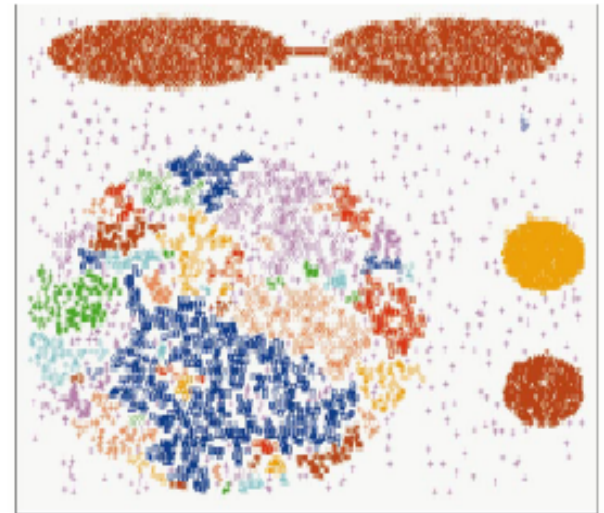
When DBSCAN Does Not Work Well

- Sensitive to parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

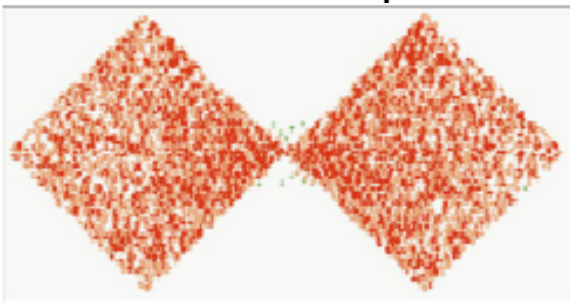


(a)

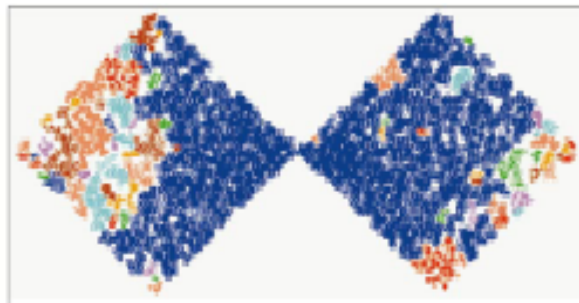


(b)

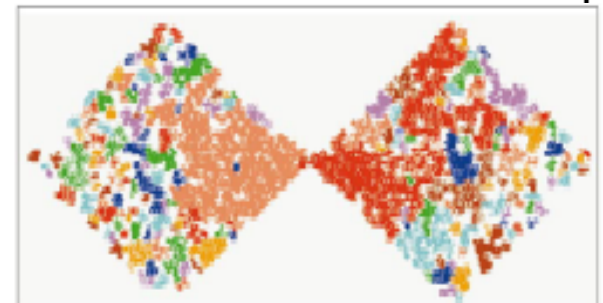
Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



(a)

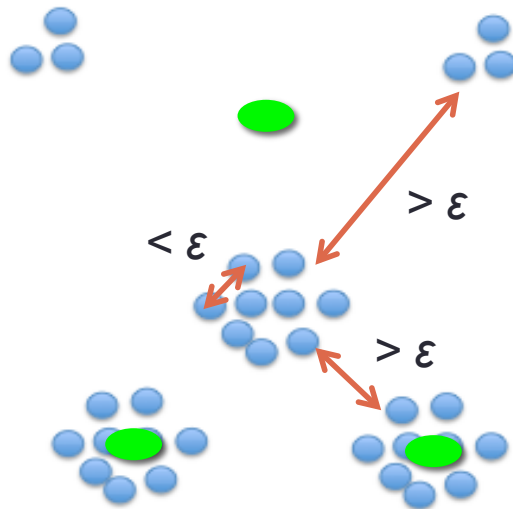


(b)



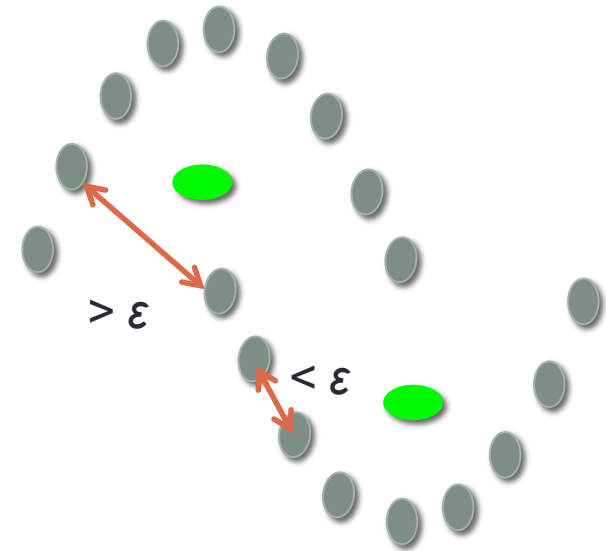
(c)

K-means VS DBSCAN



(1) When $k = 3$
MinPts = 4

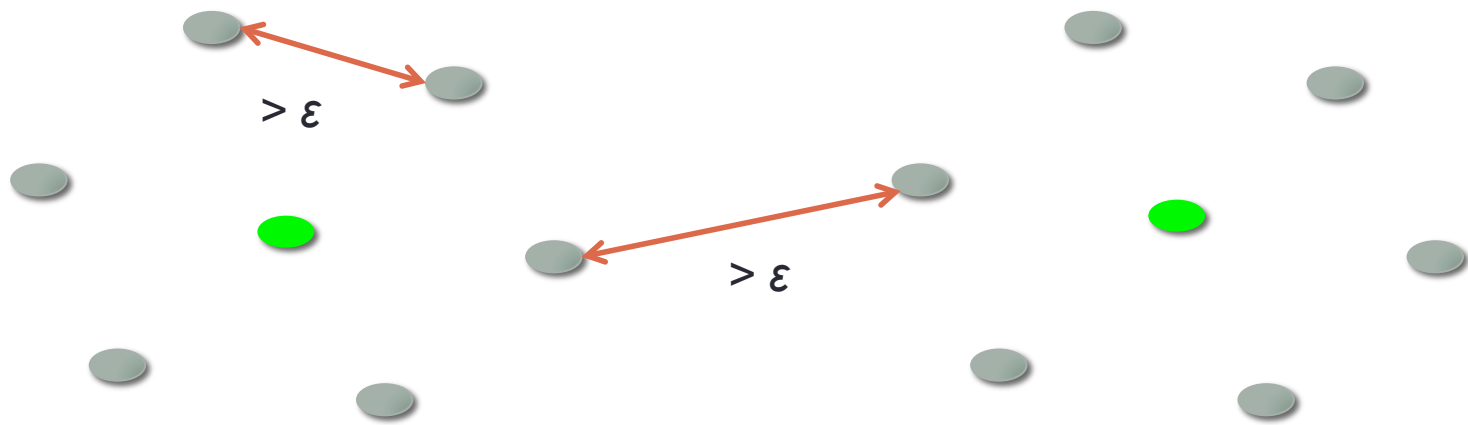
● : Initial center



(2) When $k = 2$
MinPts = 3

Winner is DBSCAN

K-means VS DBSCAN



● : Initial center

(1) When $k = 2$
MinPts = 3

Winner is K-means

Thank you for attention

Any Questions?



Reference

- Comparing Clustering Algorithm
 - <http://www.cise.ufl.edu/~jmishra/clustering/DataMiningPresentation.ppt>
- Density-Based Clustering
 - <http://www.cse.buffalo.edu/faculty/azhang/cse601/density-based.ppt>
- Ester, Martin, et al. "A density-based algorithm for discovering clusters in large spatial databases with noise." *Kdd*. Vol. 96. No. 34. 1996.