

# Density-based Approaches

## ■ Why Density-Based Clustering methods?

- Discover clusters of arbitrary shape.
- Clusters – Dense regions of objects separated by regions of low density
- DBSCAN – the first density based clustering
- OPTICS – density based cluster-ordering
- DENCLUE – a general density-based description of cluster and clustering

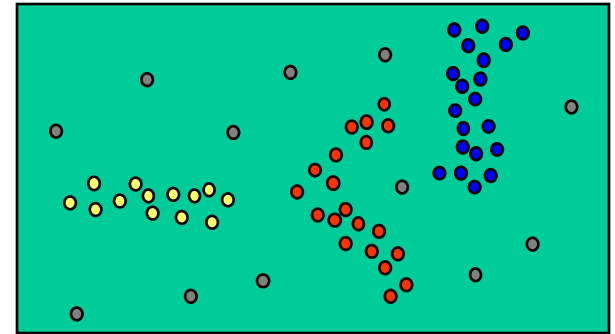
# DBSCAN: Density Based Spatial Clustering of Applications with Noise

- Proposed by Ester, Kriegel, Sander, and Xu (KDD96)
- Relies on a density-based notion of cluster: A cluster is defined as a maximal set of density-connected points.
- Discovers clusters of arbitrary shape in spatial databases with noise

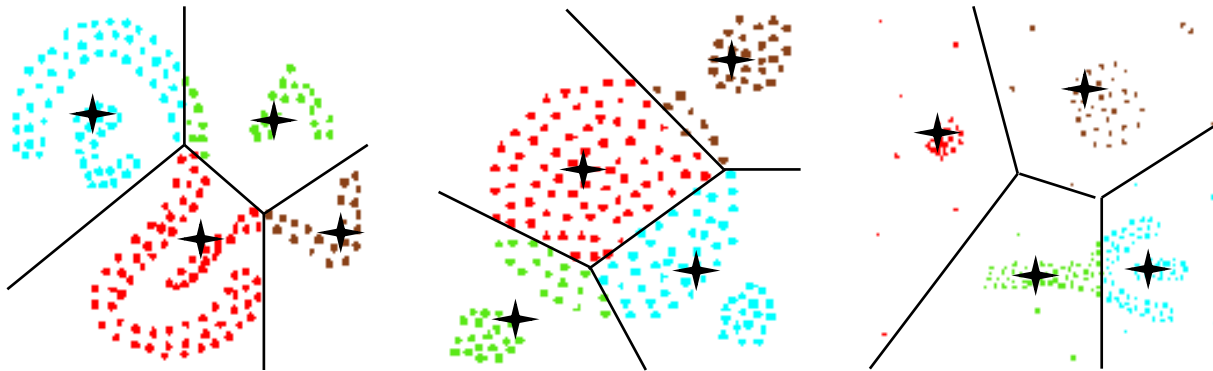
# Density-Based Clustering

## ✧ *Basic Idea:*

Clusters are dense regions in the data space, separated by regions of lower object density



## ■ Why Density-Based Clustering?



Results of a  $k$ -medoid algorithm for  $k=4$

Different density-based approaches exist (see Textbook & Papers)  
Here we discuss the ideas underlying the DBSCAN algorithm

# Density Based Clustering: Basic Concept

## ■ Intuition for the formalization of the basic idea

- For any point in a cluster, the local point density around that point has to exceed some threshold
- The set of points from one cluster is spatially connected

## ■ Local point density at a point $p$ defined by two parameters

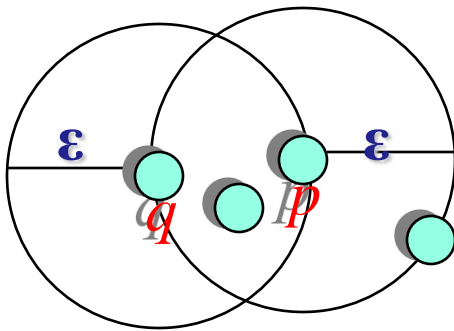
- $\varepsilon$  – radius for the neighborhood of point  $p$ :  
$$N_{\varepsilon}(p) := \{q \text{ in data set } D \mid \text{dist}(p, q) \leq \varepsilon\}$$
- $MinPts$  – minimum number of points in the given neighbourhood  $N(p)$

# $\epsilon$ -Neighborhood

- $\epsilon$ -Neighborhood – Objects within a radius of  $\epsilon$  from an object.

$$N_{\epsilon}(p) : \{q \mid d(p, q) \leq \epsilon\}$$

- “High density” -  $\epsilon$ -Neighborhood of an object contains at least *MinPts* of objects.



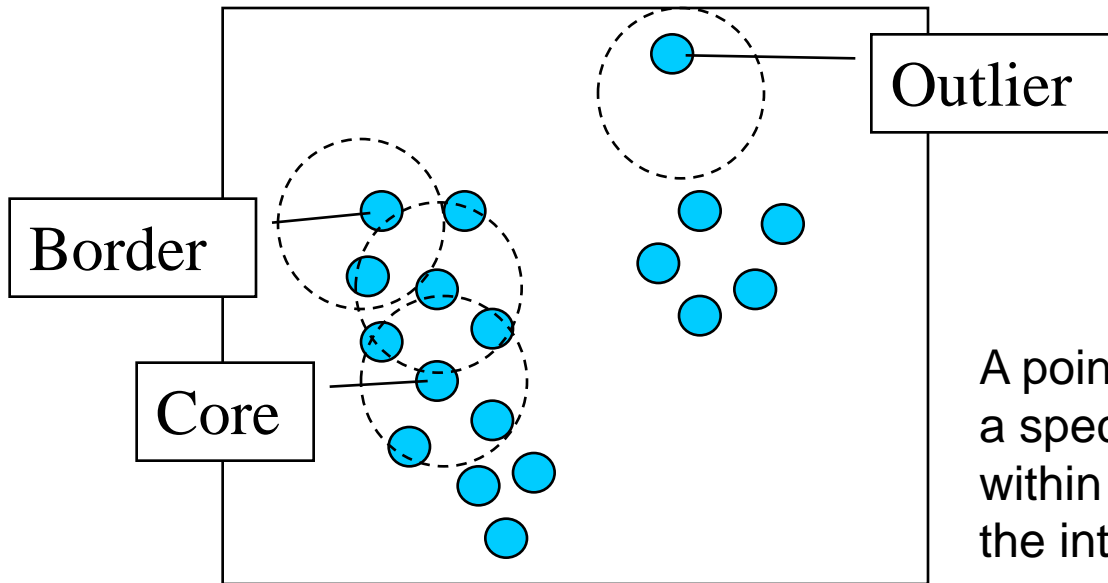
$\epsilon$ -Neighborhood of  $p$

$\epsilon$ -Neighborhood of  $q$

*Density of  $p$  is “high” (MinPts = 4)*

*Density of  $q$  is “low” (MinPts = 4)*

# Core, Border & Outlier



Given  $\epsilon$  and *MinPts*, categorize the objects into three exclusive groups.

A point is a **core point** if it has more than a specified number of points (MinPts) within  $\epsilon$ . These are points that are at the interior of a cluster.

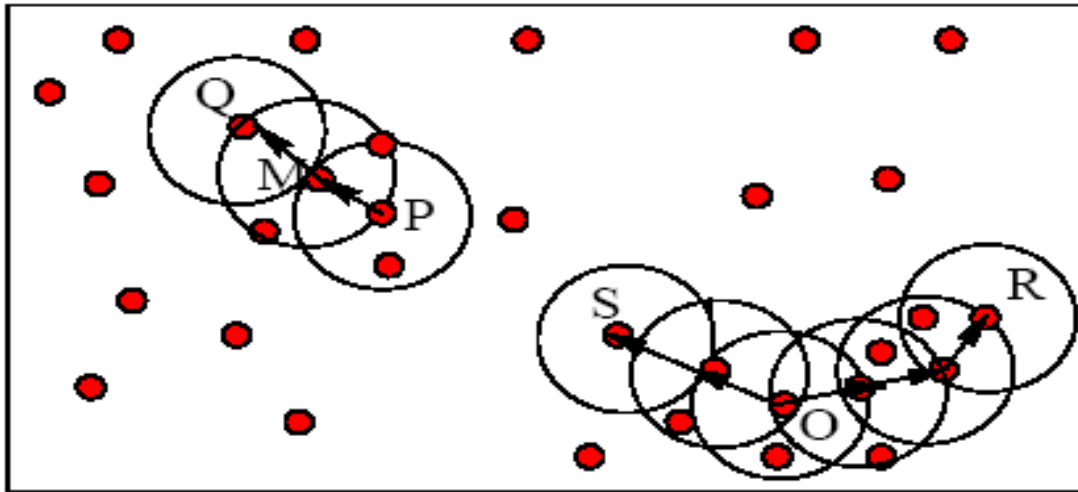
A **border point** has fewer than MinPts within  $\epsilon$ , but is in the neighborhood of a core point.

A **noise point** is any point that is not a core point nor a border point.

$\epsilon = 1\text{unit}$ ,  $\text{MinPts} = 5$

# Example

■ M, P, O, and R are core objects since each is in an Eps neighborhood containing at least 3 points



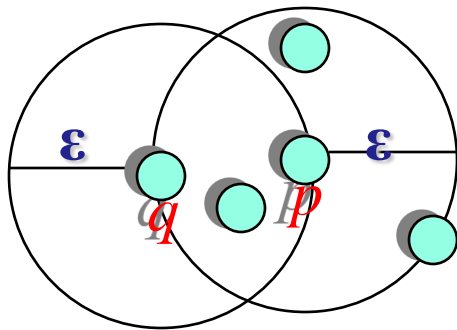
Minpts = 3

Eps=radius  
of the  
circles

# Density-Reachability

## ■ Directly density-reachable

□ An object  $q$  is directly density-reachable from object  $p$  if  $p$  is a core object and  $q$  is in  $p$ 's  $\epsilon$ -neighborhood.



MinPts = 4

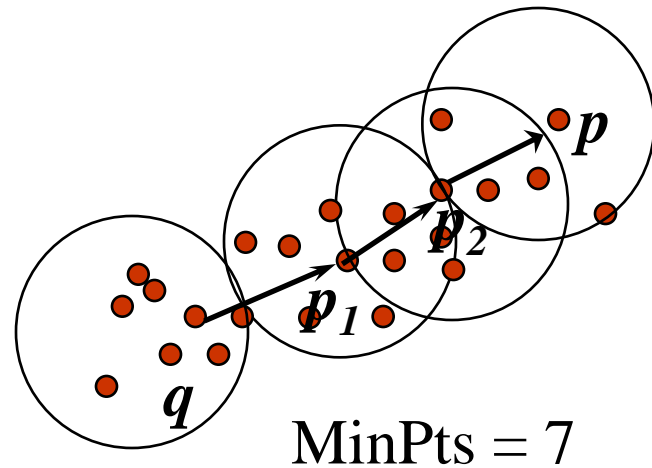
- $q$  is directly density-reachable from  $p$
- $p$  is not directly density-reachable from  $q$ ?
- Density-reachability is asymmetric.



# Density-reachability

## ■ Density-Reachable (directly and indirectly):

- ❑ A point  $p$  is directly density-reachable from  $p_2$ ;
- ❑  $p_2$  is directly density-reachable from  $p_1$ ;
- ❑  $p_1$  is directly density-reachable from  $q$ ;
- ❑  $p \leftarrow p_2 \leftarrow p_1 \leftarrow q$  form a chain.



■  $p$  is (indirectly) density-reachable from  $q$

■  $q$  is not density-reachable from  $p$ ?

# Density-Connectivity

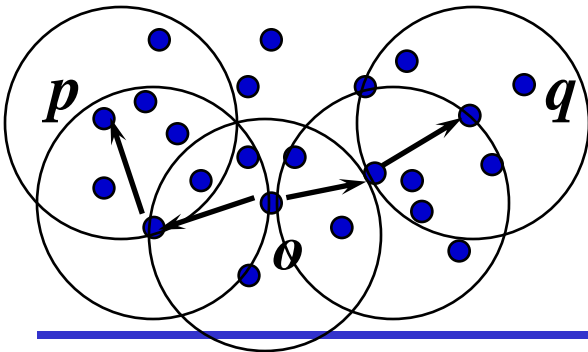
- Density-reachable is not symmetric

- not good enough to describe clusters

- Density-Connected

- A pair of points  $p$  and  $q$  are density-connected if they are commonly density-reachable from a point  $o$ .

- Density-connectivity is symmetric

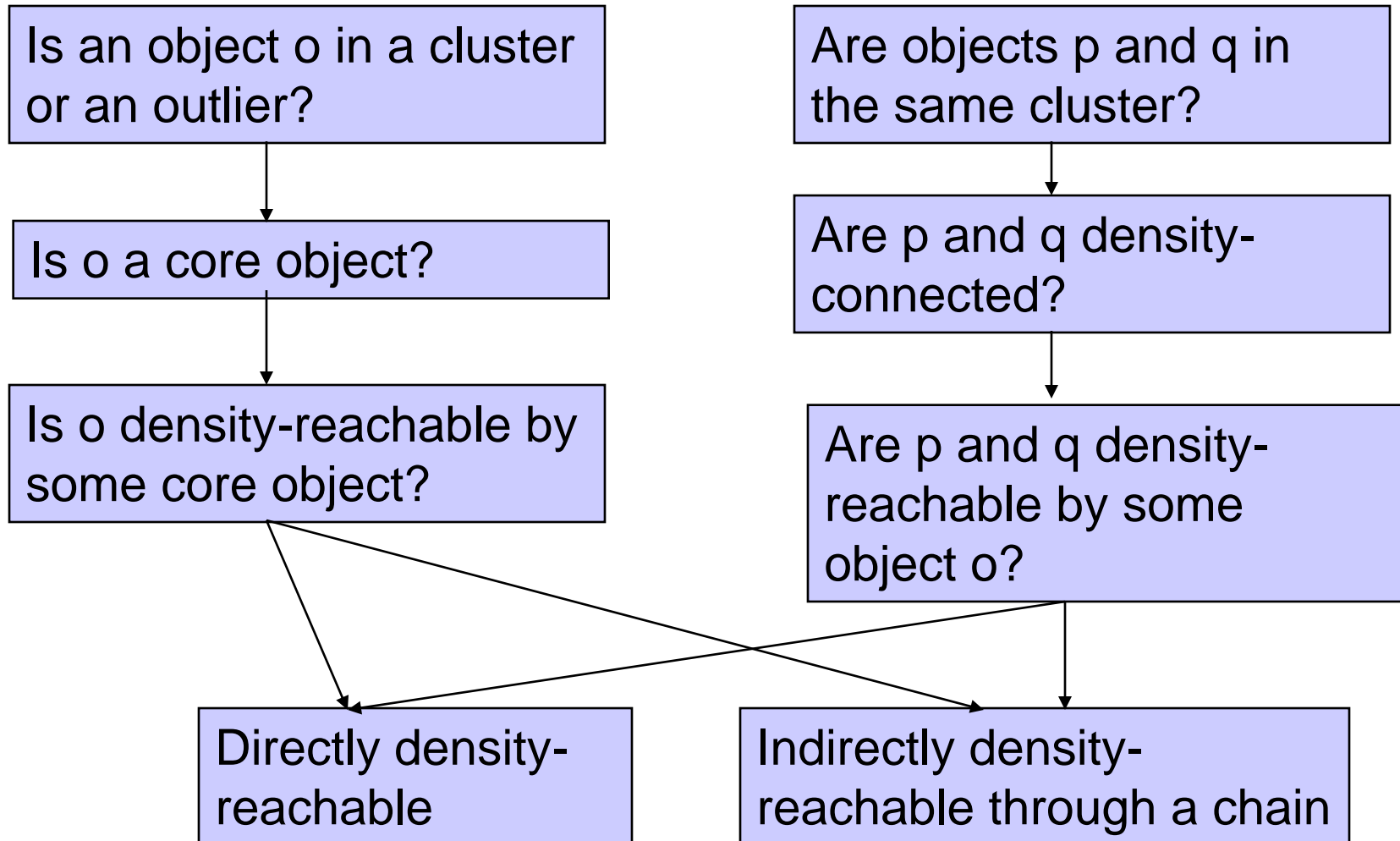


# Formal Description of Cluster

- Given a data set  $D$ , parameter  $\epsilon$  and threshold MinPts.
- A cluster  $C$  is a subset of objects satisfying two criteria:
  - **Connected:**  $\forall p, q \in C$ :  $p$  and  $q$  are density-connected.
  - **Maximal:**  $\forall p, q$ : if  $p \in C$  and  $q$  is density-reachable from  $p$ , then  $q \in C$ . (avoid redundancy)

↓  
 $P$  is a core object.

# Review of Concepts



# DBSCAN Algorithm

Input: The data set  $D$

Parameter:  $\varepsilon$ , MinPts

For each object  $p$  in  $D$

    if  $p$  is a core object and not processed then

$C = \text{retrieve all objects density-reachable from } p$

        mark all objects in  $C$  as processed

        report  $C$  as a cluster

    else mark  $p$  as outlier

    end if

End For

DBScan Algorithm

# DBSCAN: The Algorithm

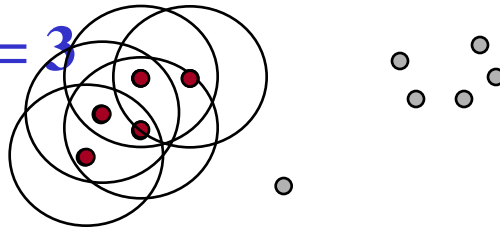
- ❑ Arbitrary select a point  $p$
- ❑ Retrieve all points density-reachable from  $p$  wrt  $Eps$  and  $MinPts$ .
- ❑ If  $p$  is a core point, a cluster is formed.
- ❑ If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database.
- ❑ Continue the process until all of the points have been processed.

# DBSCAN Algorithm: Example

## ■ Parameter

○  $\varepsilon = 2 \text{ cm}$

○  $MinPts = 3$



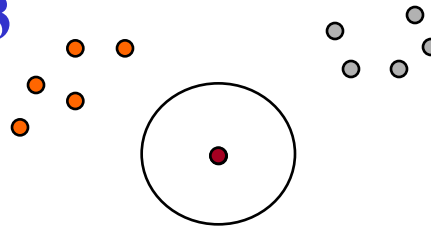
```
for each  $o \in D$  do
  if  $o$  is not yet classified then
    if  $o$  is a core-object then
      collect all objects density-reachable from  $o$ 
      and assign them to a new cluster.
    else
      assign  $o$  to NOISE
```

# DBSCAN Algorithm: Example

## ■ Parameter

○  $\varepsilon = 2 \text{ cm}$

○  $MinPts = 3$



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

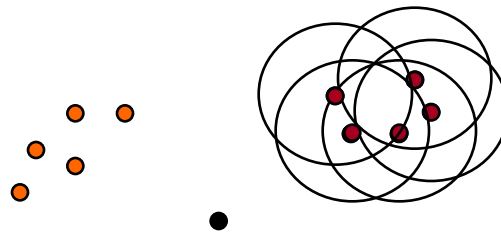


# DBSCAN Algorithm: Example

## ■ Parameter

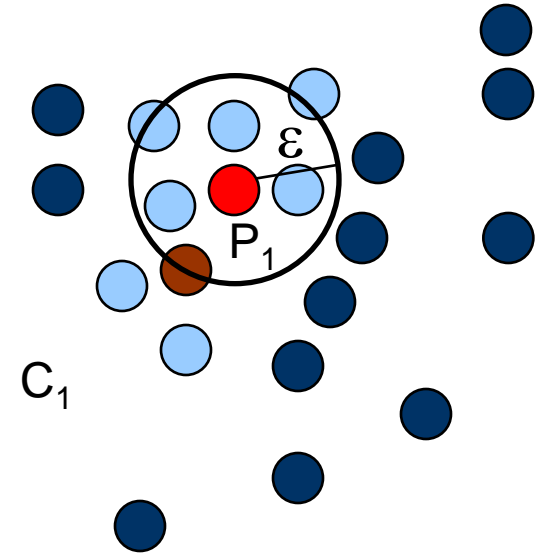
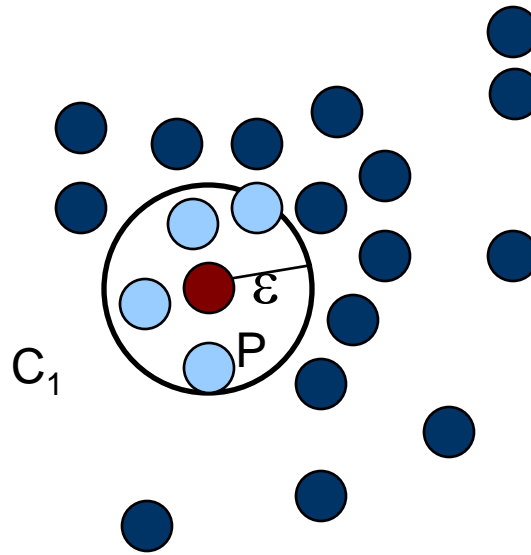
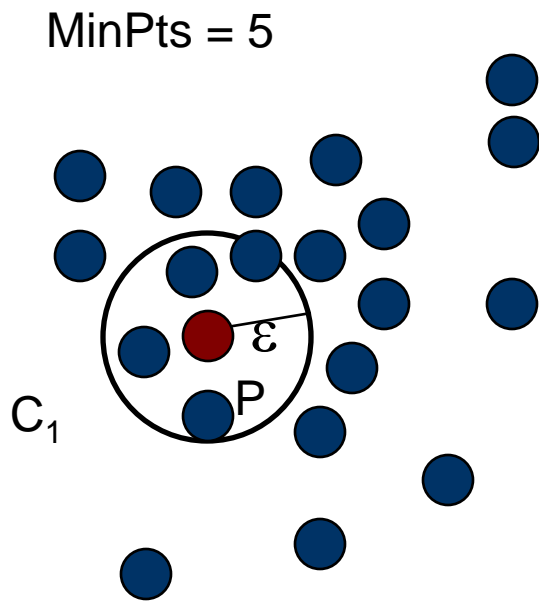
○  $\varepsilon = 2 \text{ cm}$

○  $MinPts = 3$



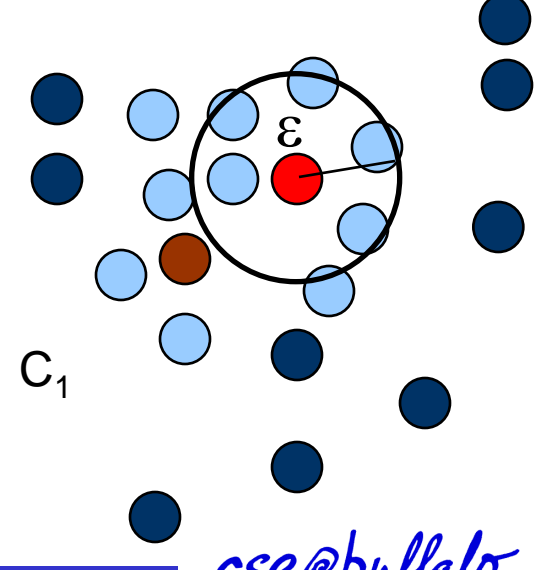
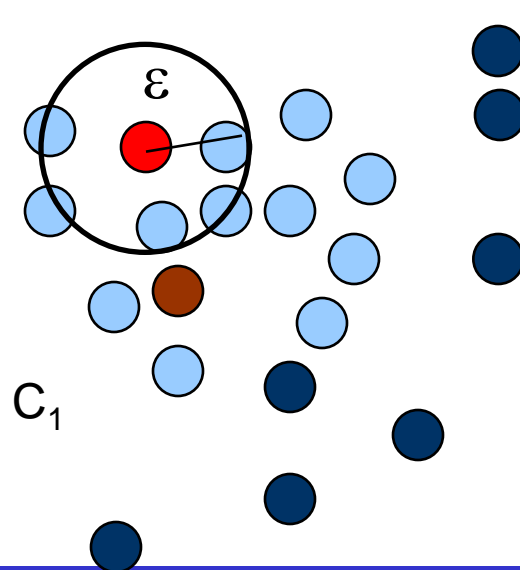
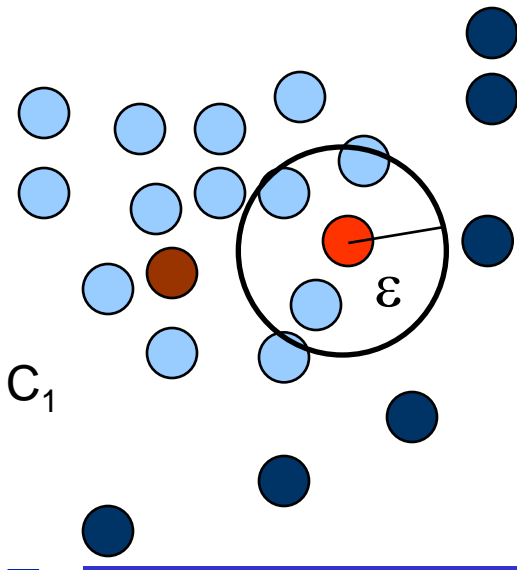
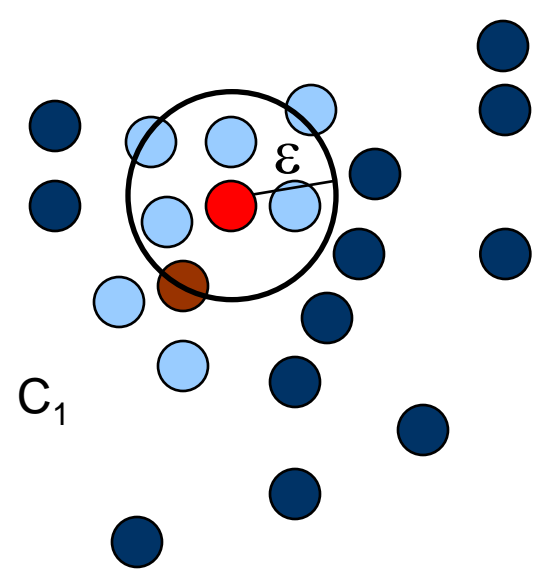
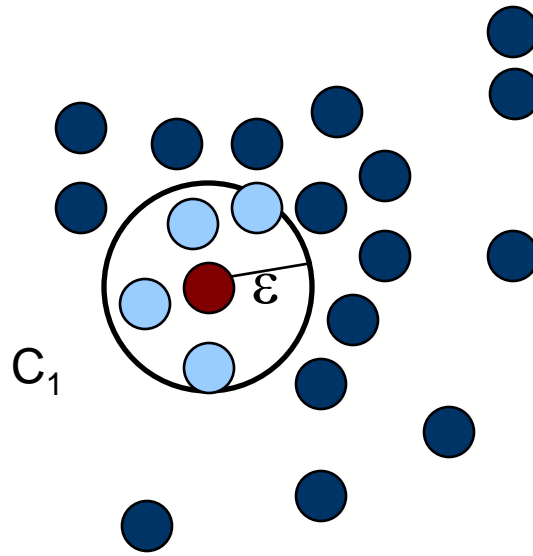
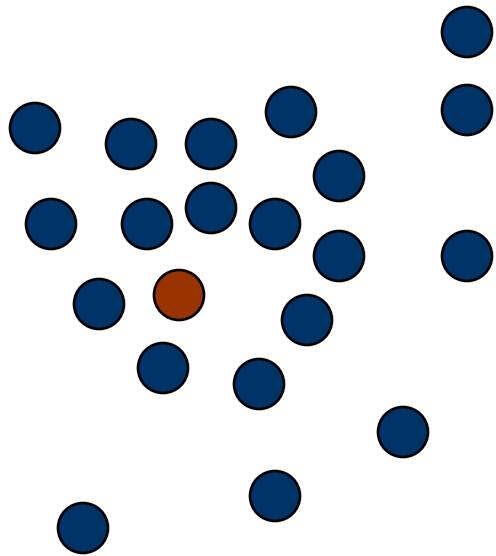
```
for each  $o \in D$  do
  if  $o$  is not yet classified then
    if  $o$  is a core-object then
      collect all objects density-reachable from  $o$ 
      and assign them to a new cluster.
    else
      assign  $o$  to NOISE
```

MinPts = 5



1. Check the  $\epsilon$ -neighborhood of  $p$ ;
2. If  $p$  has less than MinPts neighbors then mark  $p$  as outlier and continue with the next object
3. Otherwise mark  $p$  as processed and put all the neighbors in cluster  $C$

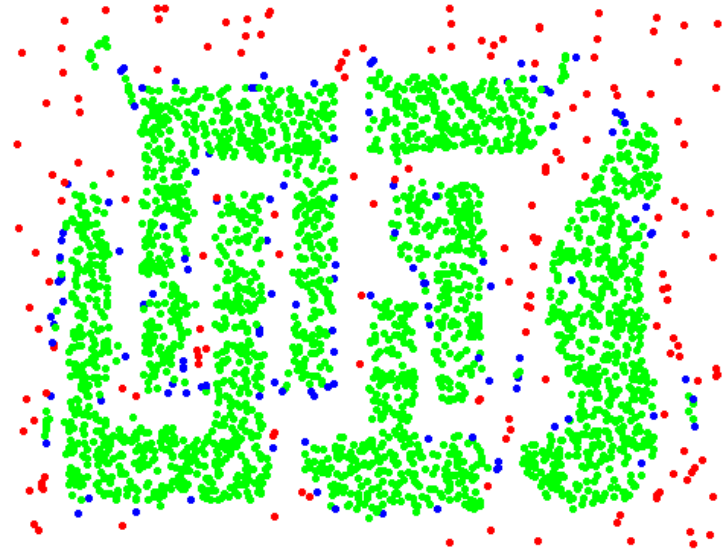
1. Check the unprocessed objects in  $C$
2. If no core object, return  $C$
3. Otherwise, randomly pick up one core object  $p_1$ , mark  $p_1$  as processed, and put all unprocessed neighbors of  $p_1$  in cluster  $C$



# Example



Original Points



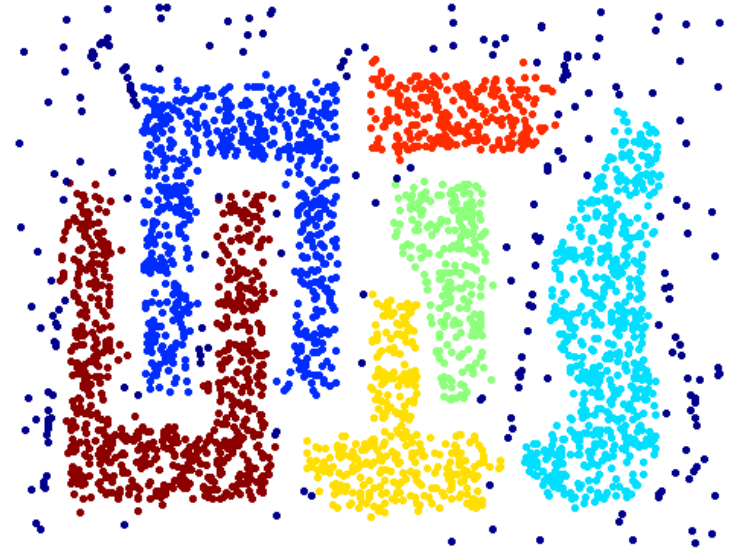
Point types: **core**,  
**border** and **outliers**

$\epsilon = 10$ , MinPts = 4

# When DBSCAN Works Well



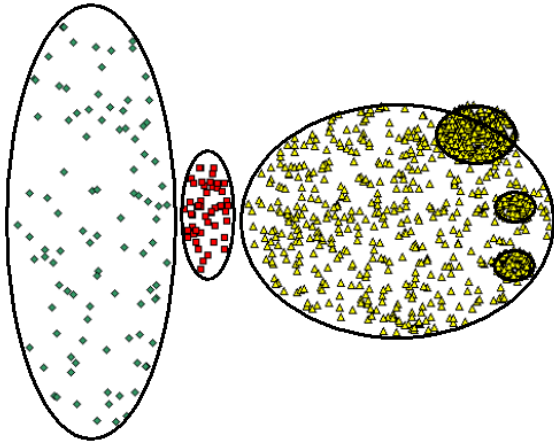
Original Points



Clusters

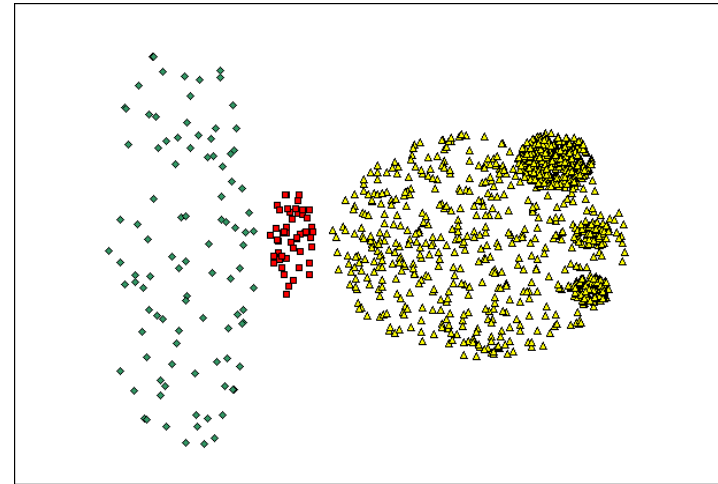
- Resistant to Noise
- Can handle clusters of different shapes and sizes

# When DBSCAN Does NOT Work Well

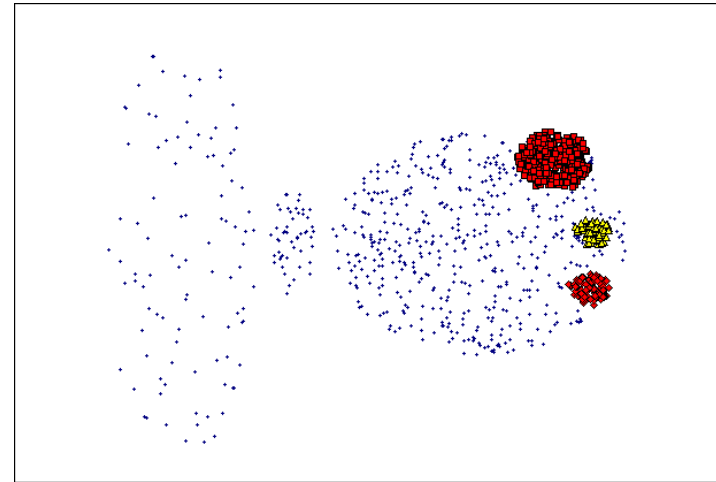


Original Points

- Cannot handle Varying densities
- sensitive to parameters



(MinPts=4, Eps=9.92).



(MinPts=4, Eps=9.75)

# DBSCAN: Sensitive to Parameters

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

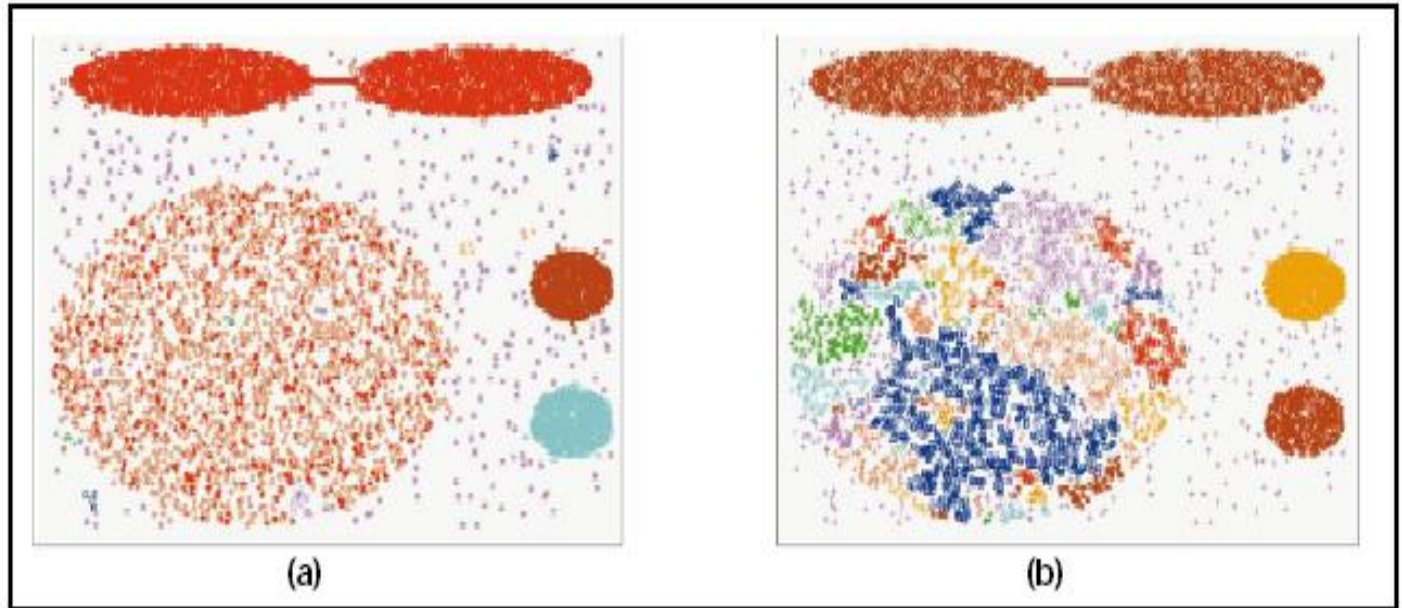
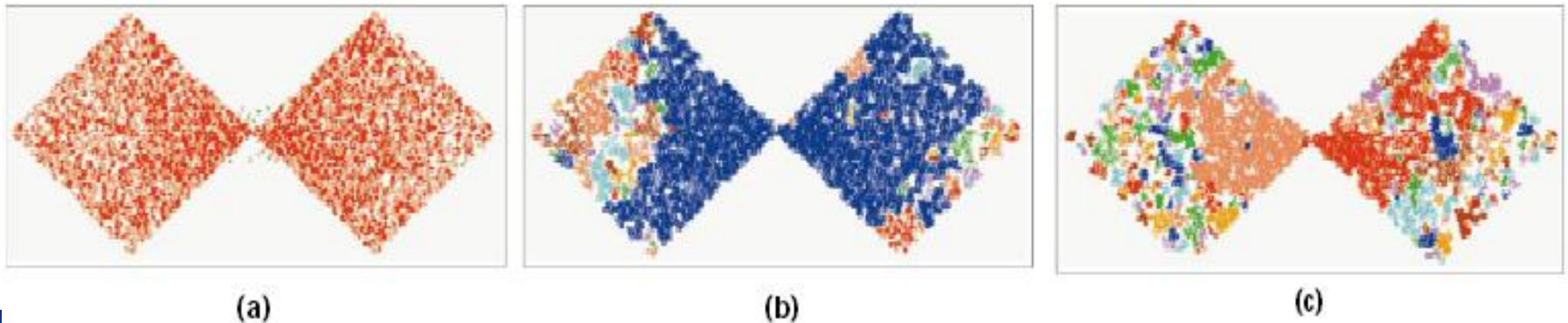


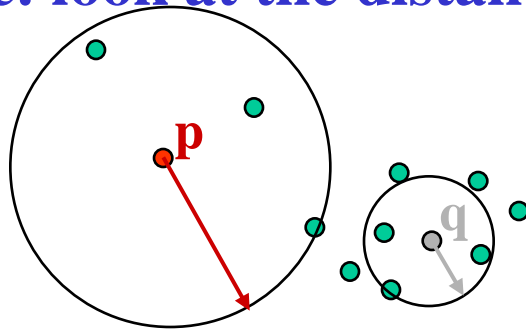
Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.







# Determining the Parameters $\varepsilon$ and *MinPts*

- **Cluster:** Point density higher than specified by  $\varepsilon$  and *MinPts*
- **Idea:** use the point density of the least dense cluster in the data set as parameters – but how to determine this?
- **Heuristic:** look at the distances to the  $k$ -nearest neighbors



$3\text{-distance}(p) :$  

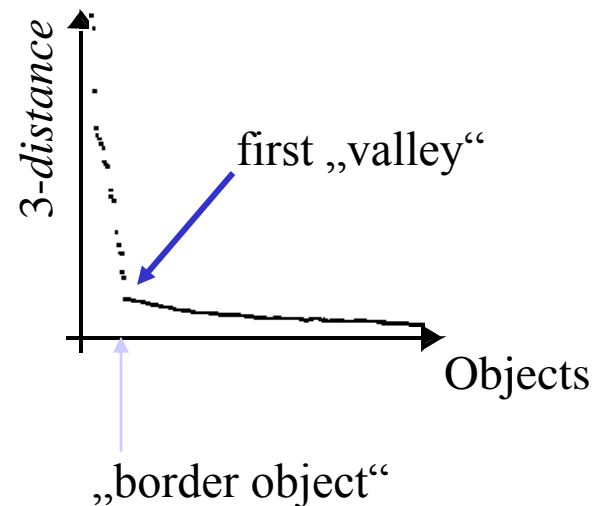
$3\text{-distance}(q) :$  

- **Function  $k\text{-distance}(p)$ :** distance from  $p$  to the its  $k$ -nearest neighbor
- **$k\text{-distance plot}$ :**  $k$ -distances of all objects, sorted in decreasing order



# Determining the Parameters $\varepsilon$ and *MinPts*

## ■ Example *k*-distance plot

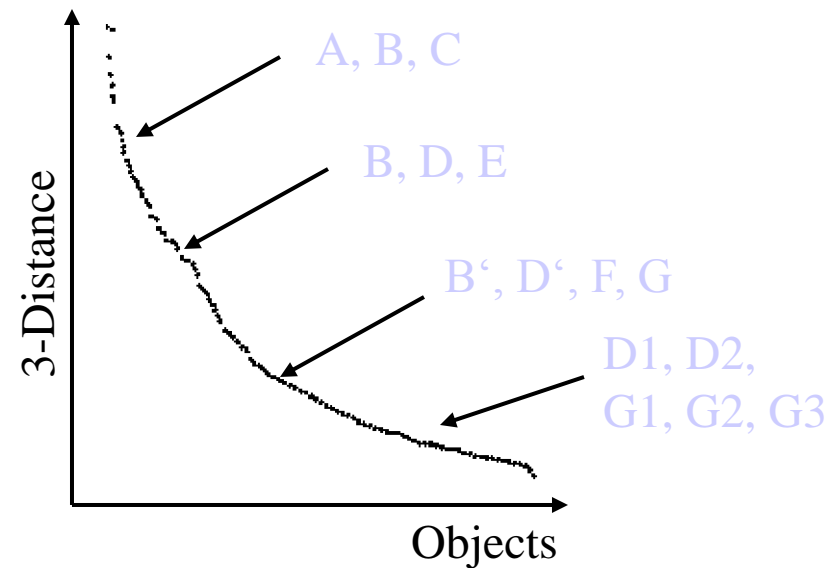
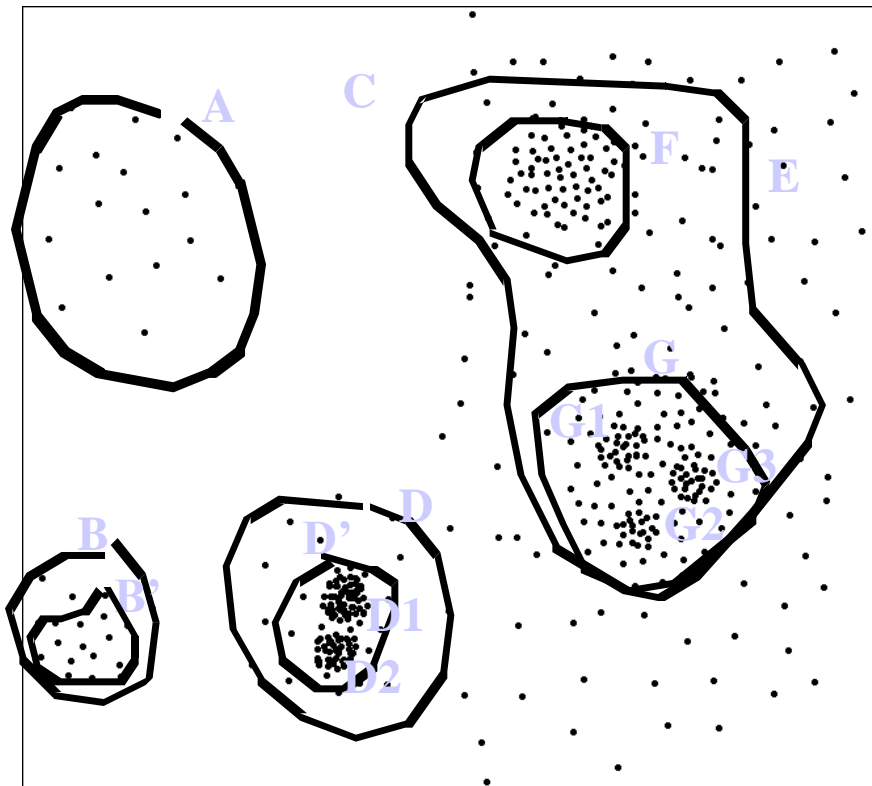


## ■ Heuristic method:

- ❑ Fix a value for *MinPts* (default:  $2 \times d - 1$ )
- ❑ User selects “border object” *o* from the *MinPts*-distance plot;  $\varepsilon$  is set to *MinPts*-distance(*o*)

# Determining the Parameters $\varepsilon$ and *MinPts*

## ■ Problematic example



# Density Based Clustering: Discussion

## ■ Advantages

- ❑ Clusters can have arbitrary shape and size
- ❑ Number of clusters is determined automatically
- ❑ Can separate clusters from surrounding noise
- ❑ Can be supported by spatial index structures

## ■ Disadvantages

- ❑ Input parameters may be difficult to determine
- ❑ In some situations very sensitive to input parameter setting

# OPTICS: Ordering Points To Identify the Clustering Structure

## ■ DBSCAN

- ❑ Input parameter – hard to determine.
- ❑ Algorithm very sensitive to input parameters.

## ■ OPTICS – Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)

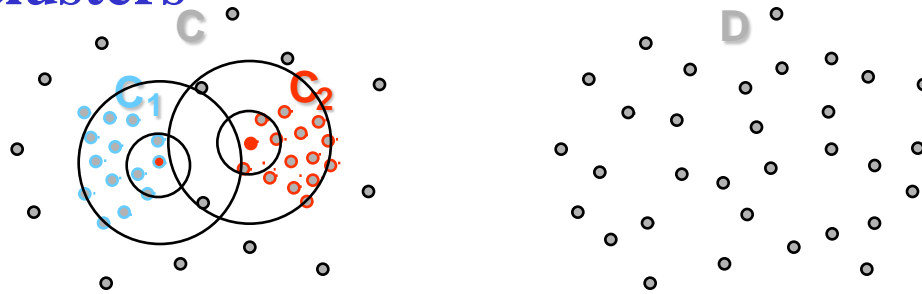
- ❑ Based on DBSCAN.
- ❑ Does not produce clusters explicitly.
- ❑ Rather generate an ordering of data objects representing density-based clustering structure.

# OPTICS con't

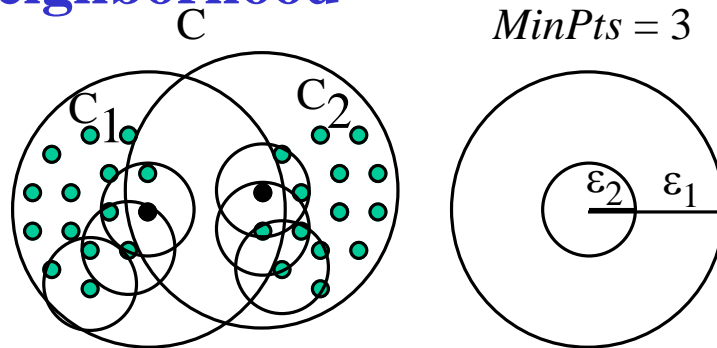
- Produces a special order of the database wrt its density-based clustering structure
- This cluster-ordering contains info equiv to the density-based clusterings corresponding to a broad range of parameter settings
- Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
- Can be represented graphically or using visualization techniques

# Density-Based Hierarchical Clustering

- *Observation:* Dense clusters are completely contained by less dense clusters



- *Idea:* Process objects in the “right” order and keep track of point density in their neighborhood



# Core- and Reachability Distance

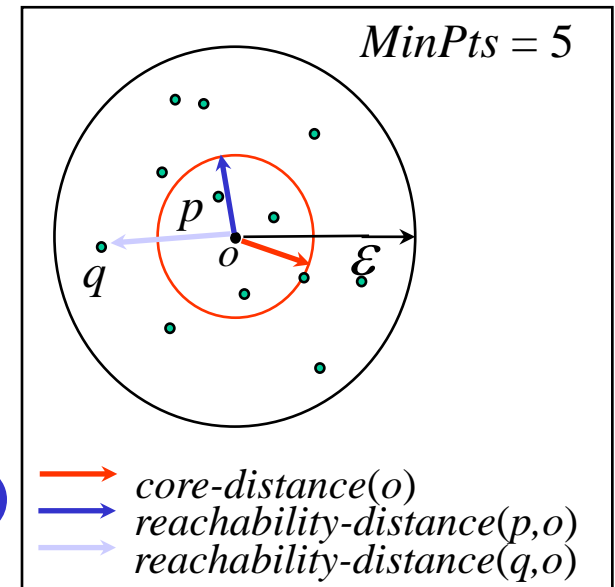
■ Parameters: “generating” distance  $\varepsilon$ , fixed value  $MinPts$

■  $core-distance_{\varepsilon, MinPts}(o)$

“smallest distance such that  $o$  is a core object”  
(if that distance is  $\leq \varepsilon$ ; “?” otherwise)

■  $reachability-distance_{\varepsilon, MinPts}(p, o)$

“smallest distance such that  $p$  is  
directly density-reachable from  $o$ ”  
(if that distance is  $\leq \varepsilon$ ; “?” otherwise)



# OPTICS: Extension of DBSCAN

- Order points by shortest *reachability distance* to guarantee that clusters w.r.t. higher density are finished first. (for a constant MinPts, higher density requires lower  $\epsilon$ )



# The Algorithm OPTICS

## ■ Basic data structure: `controlList`

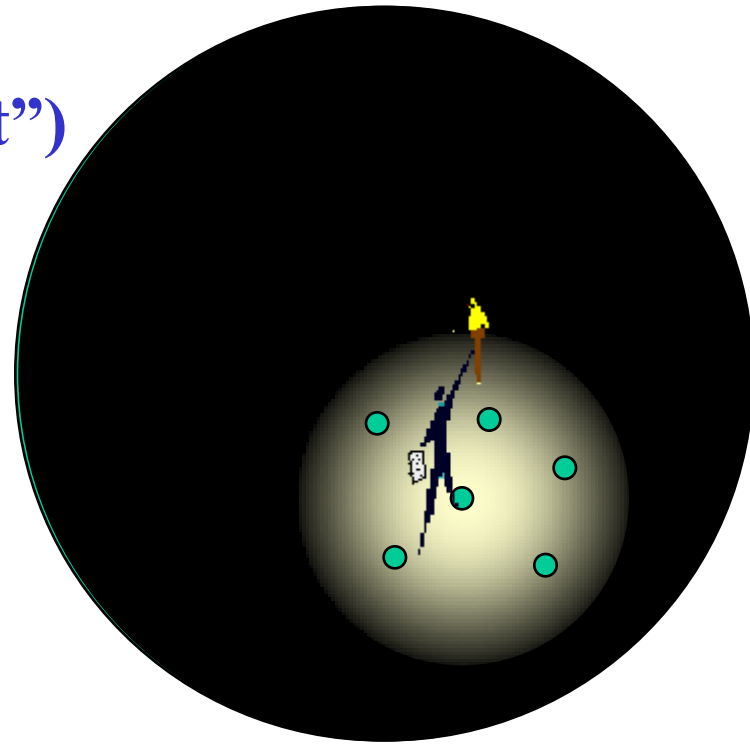
- Memorize shortest reachability distances seen so far  
    (“distance of a jump to that point”)

## ■ Visit each point

- Make always a shortest jump

## ■ Output:

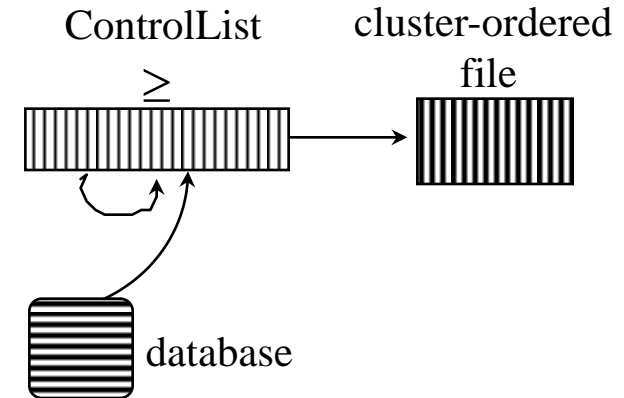
- order of points
- core-distance of points
- reachability-distance of points



# The Algorithm OPTICS

✴ *ControlList* ordered by reachability-distance.

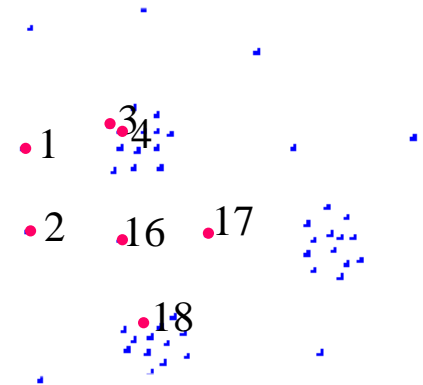
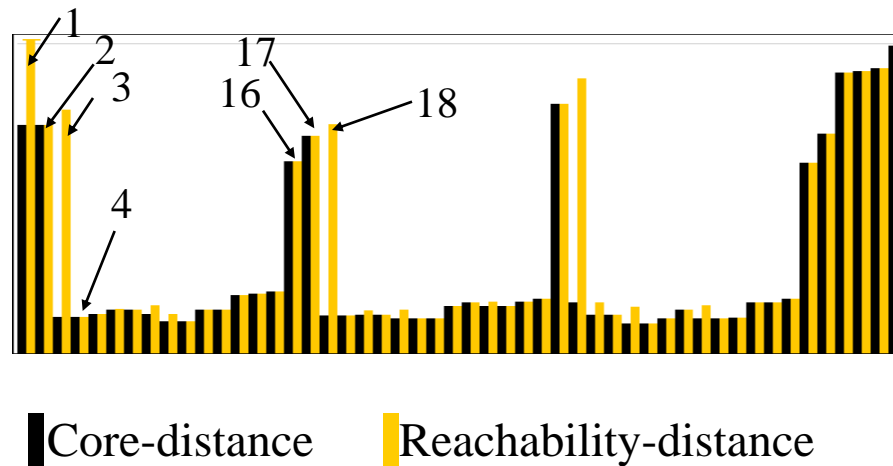
```
foreach  $o \in \text{Database}$ 
  // initially,  $o.\text{processed} = \text{false}$  for all objects  $o$ 
  if  $o.\text{processed} = \text{false}$ ;
    insert ( $o$ , “?”) into ControlList;
  while ControlList is not empty
    select first element ( $o$ ,  $r\_dist$ ) from ControlList;
    retrieve  $N_\epsilon(o)$  and determine  $c\_dist = \text{core-distance}(o)$ ;
    set  $o.\text{processed} = \text{true}$ ;
    write ( $o$ ,  $r\_dist$ ,  $c\_dist$ ) to file;
    if  $o$  is a core object at any distance  $\leq \epsilon$ 
      foreach  $p \in N_\epsilon(o)$  not yet processed;
        determine  $r\_dist_p = \text{reachability-distance}(p, o)$ ;
        if  $(p, \_) \notin \text{ControlList}$ 
          insert ( $p$ ,  $r\_dist_p$ ) in ControlList;
        else if  $(p, \text{old\_}r\_dist) \in \text{ControlList}$  and  $r\_dist_p < \text{old\_}r\_dist$ 
          update ( $p$ ,  $r\_dist_p$ ) in ControlList;
```



# OPTICS: Properties

■ “Flat” density-based clusters wrt.  $\epsilon^* \leq \epsilon$  and *MinPts* afterwards:

- ❑ Starts with an object  $o$  where  $c\text{-dist}(o) \leq \epsilon^*$  and  $r\text{-dist}(o) > \epsilon^*$
- ❑ Continues while  $r\text{-dist} \leq \epsilon^*$

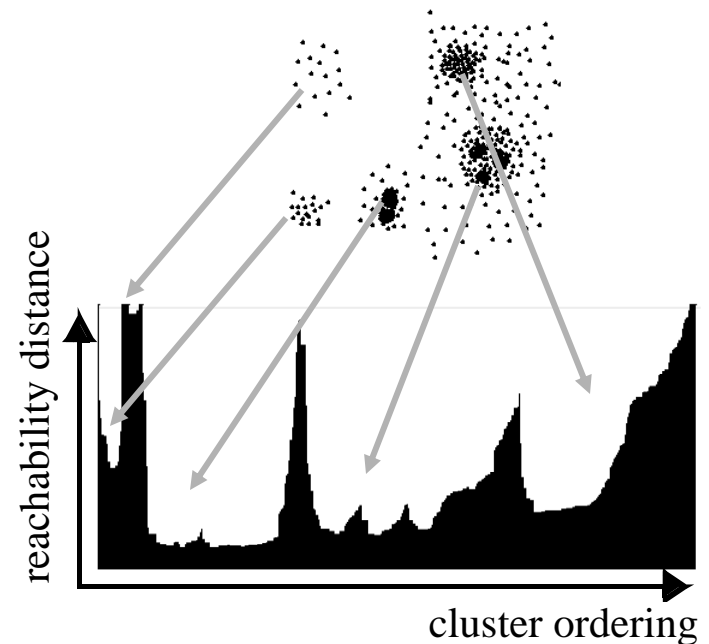
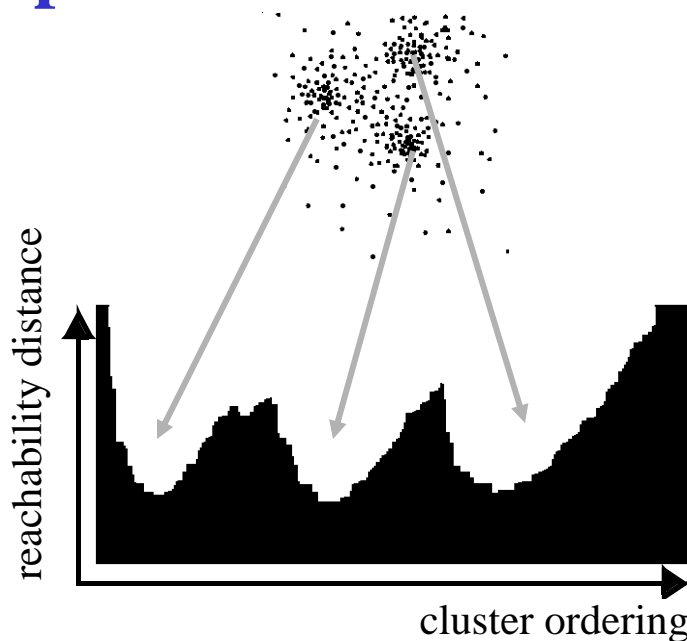


■ Performance: approx. runtime( DBSCAN( $\epsilon$ , *MinPts*) )

- ❑  $O(n * \text{runtime}(\epsilon\text{-neighborhood-query}))$ 
  - without spatial index support (worst case):  $O(n^2)$
  - e.g. tree-based spatial index support:  $O(n * \log(n))$

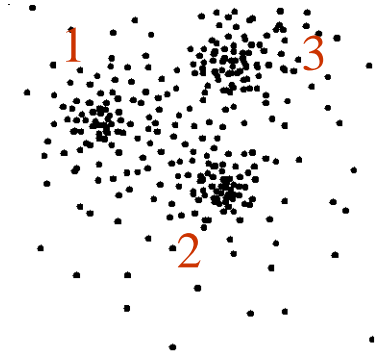
# OPTICS: The Reachability Plot

- represents the density-based clustering structure
- easy to analyze
- independent of the dimension of the data

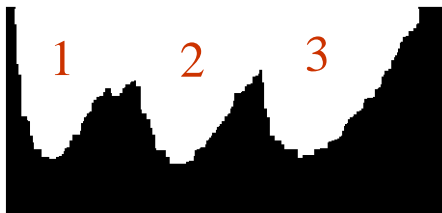


# OPTICS: Parameter Sensitivity

- Relatively insensitive to parameter settings
- Good result if parameters are just “large enough”



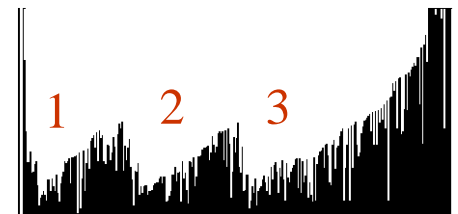
$MinPts = 10, \epsilon = 10$



$MinPts = 10, \epsilon = 5$



$MinPts = 2, \epsilon = 10$



# An Example of OPTICS

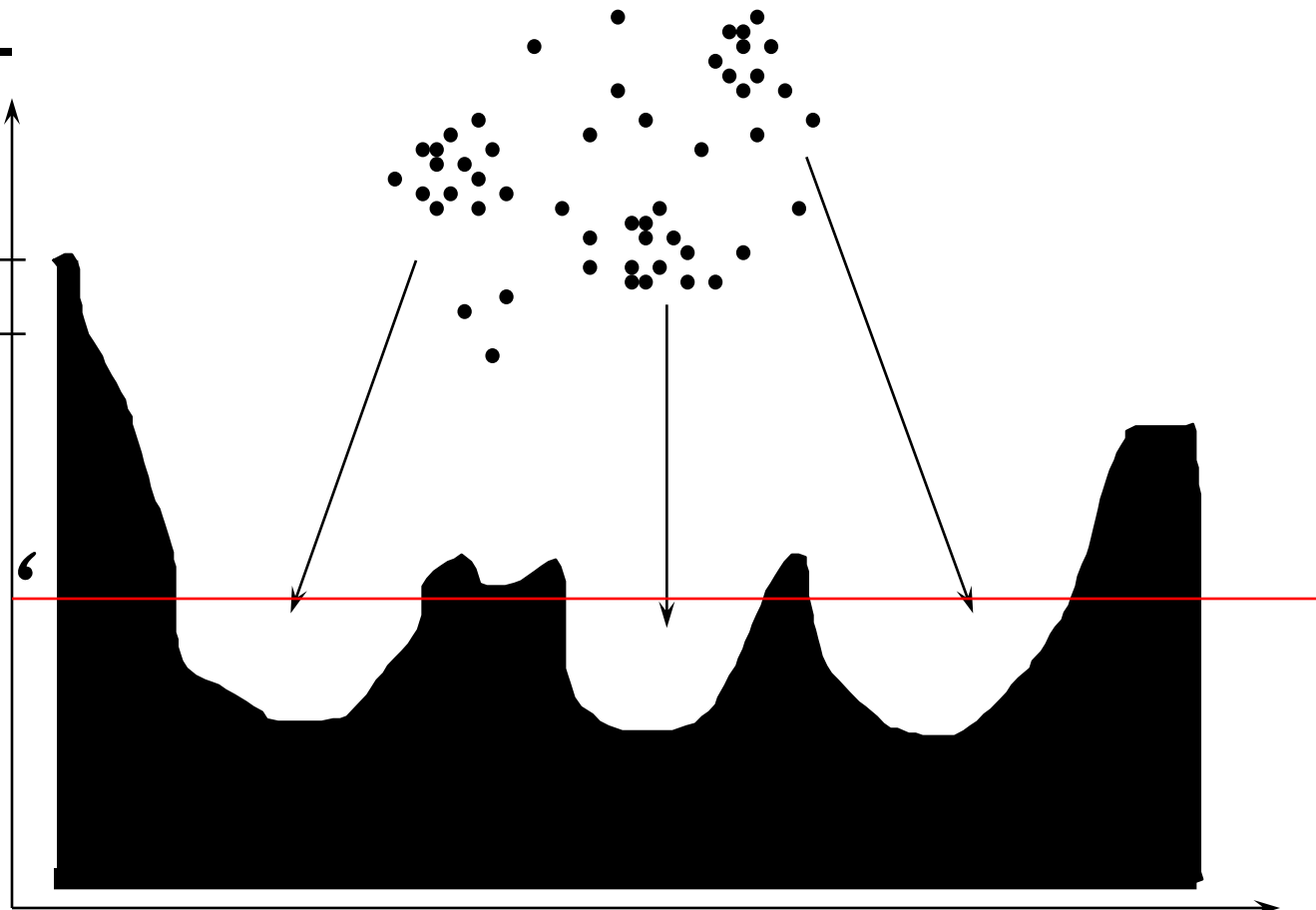
neighboring objects stay close to each other in a linear sequence.

Reachability-distance

undefined

$\epsilon$

$\epsilon'$

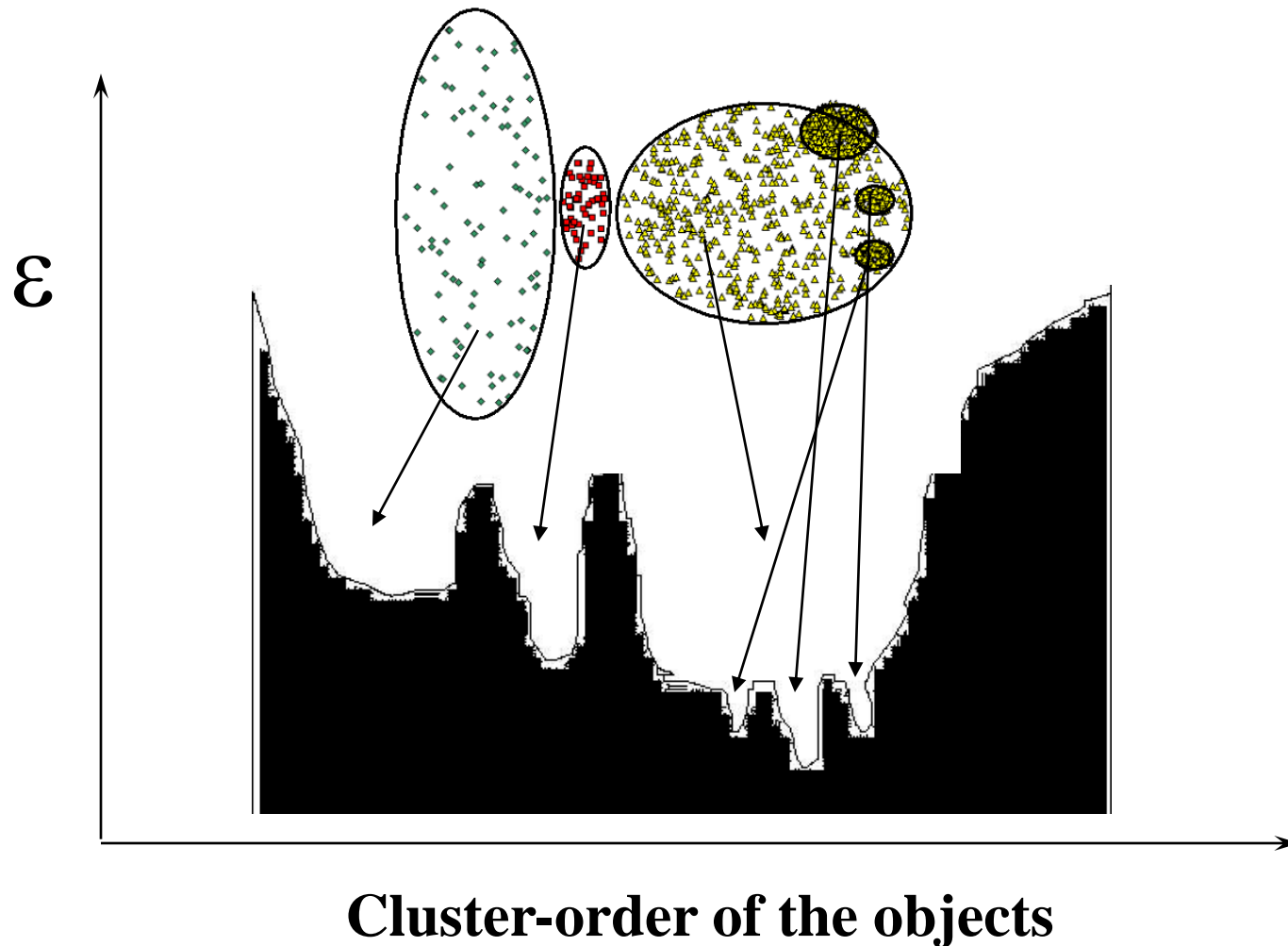


Cluster-order of the objects

# DBSCAN VS OPTICS

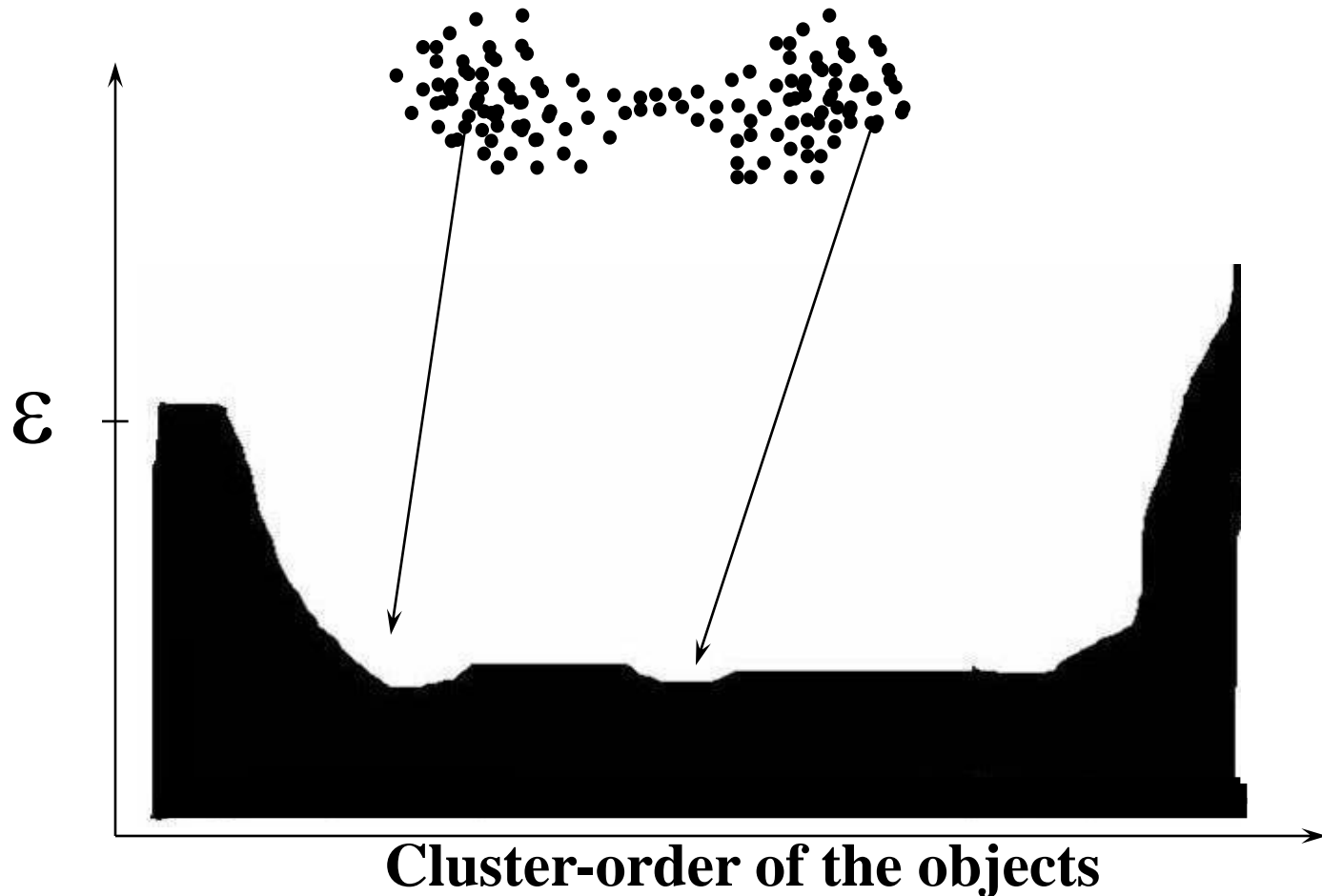
	DBSCAN	OPTICS
Density	Boolean value (high/low)	Numerical value (core distance)
Density- connected	Boolean value (yes/no)	Numerical value (reachability distance)
Searching strategy	random	greedy

# When OPTICS Works Well





# When OPTICS Does NOT Work Well



# DENCLUE: using density functions

■ DENSity-based CLUstEring by Hinneburg & Keim (KDD'98)

■ Major features

- ❑ Solid mathematical foundation
- ❑ Good for data sets with large amounts of noise
- ❑ Allows a compact mathematical description of arbitrarily shaped clusters in high-dimensional data sets
- ❑ Significantly faster than existing algorithm (faster than DBSCAN by a factor of up to 45)
- ❑ But needs a large number of parameters

# Denclude: Technical Essence

- Model density by the notion of influence
- Each data object exert influence on its neighborhood.
- The influence decreases with distance
- Example:
  - Consider each object is a radio, the closer you are to the object, the louder the noise
- Key: Influence is represented by mathematical function

# Denclue: Technical Essence

■ Influence functions: (influence of  $y$  on  $x$ ,  $\sigma$  is a user given constant)

□ Square :  $f^y_{square}(x) = 0$ , if  $\text{dist}(x,y) > \sigma$ ,  
1, otherwise

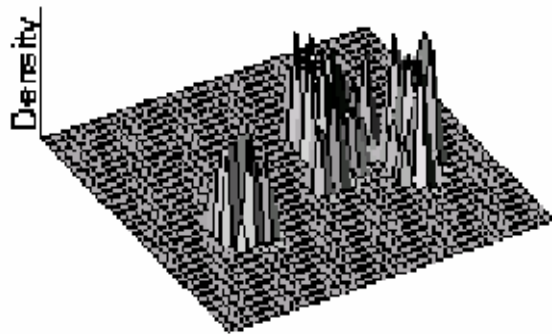
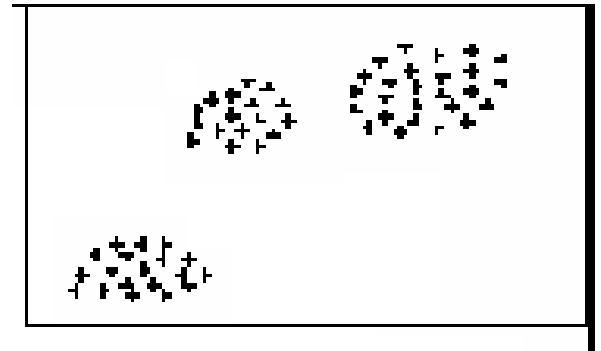
□ Gaussian:

$$f^y_{Gaussian}(x) = e^{-\frac{d(x,y)^2}{2\sigma^2}}$$

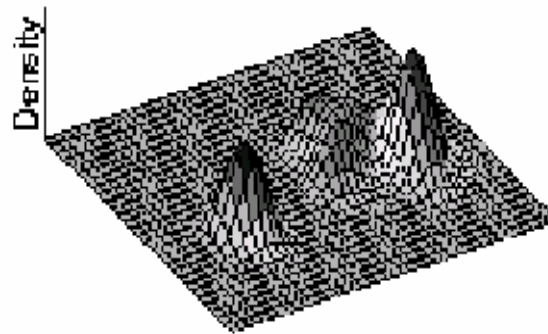
# Density Function

- Density Definition is defined as the sum of the influence functions of all data points.

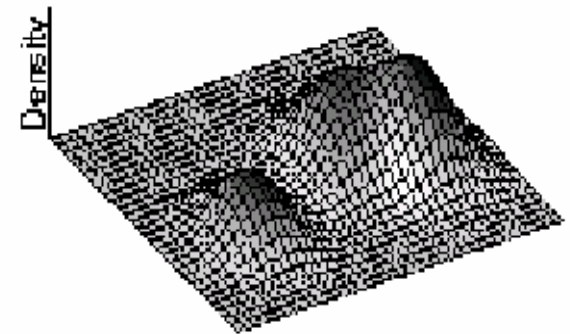
$$f_{Gaussian}^D(x) = \sum_{i=1}^N e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$



(a)  $\sigma = 0.2$



(b)  $\sigma = 0.6$



(d)  $\sigma = 1.5$

# Gradient: The steepness of a slope

## ■ Example

$$f_{\text{Gaussian}}(x, y) = e^{-\frac{d(x, y)^2}{2\sigma^2}}$$

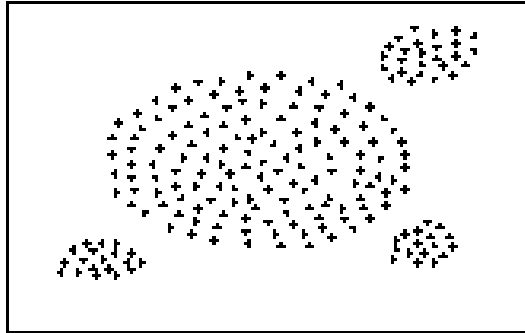
$$f_{\text{Gaussian}}^D(x) = \sum_{i=1}^N e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

$$\nabla f_{\text{Gaussian}}^D(x, x_i) = \sum_{i=1}^N (x_i - x) \cdot e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

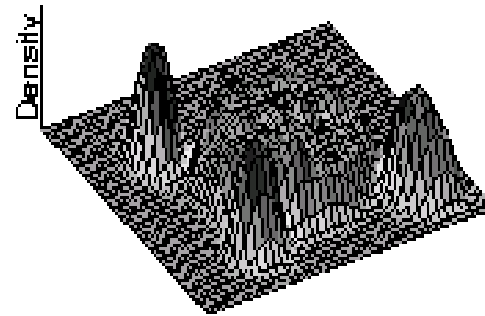
# Denclue: Technical Essence

- Clusters can be determined mathematically by identifying density attractors.
- Density attractors are local maximum of the overall density function.

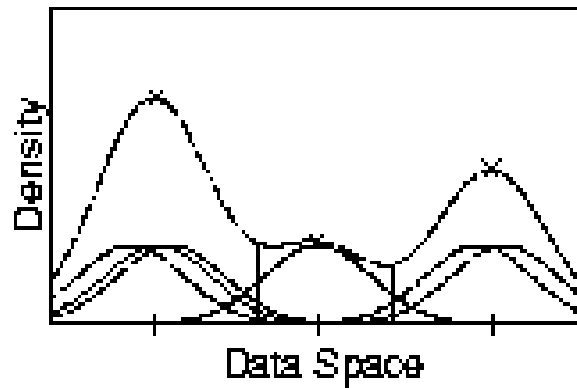
# Density Attractor



(a) Data Set



(c) Gaussian





# Cluster Definition

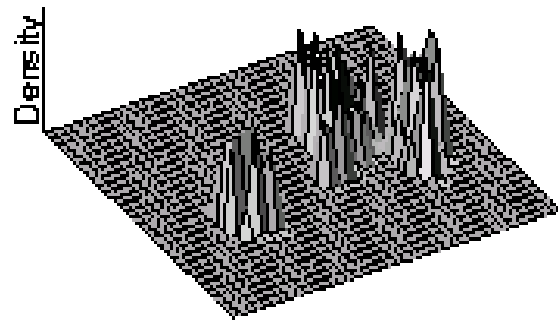
## ■ Center-defined cluster

- A subset of objects attracted by an attractor  $x$
- $\text{density}(x) \geq \xi$

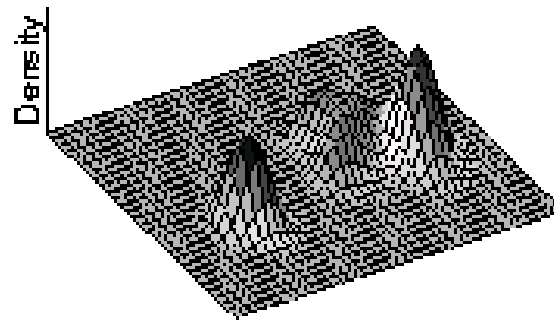
## ■ Arbitrary-shape cluster

- A group of center-defined clusters which are connected by a path  $P$
- For each object  $x$  on  $P$ ,  $\text{density}(x) \geq \xi$ .

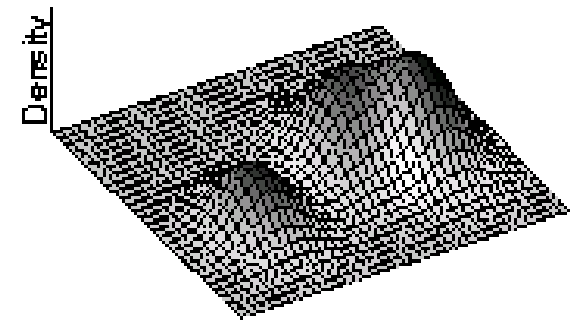
# Center-Defined and Arbitrary



(a)  $\sigma = 0.2$

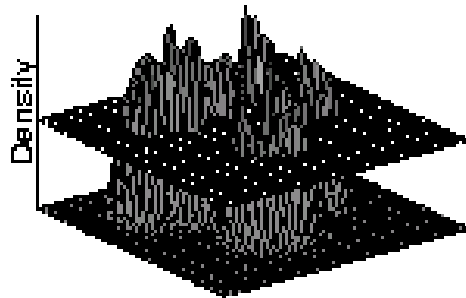


(b)  $\sigma = 0.6$



(d)  $\sigma = 1.5$

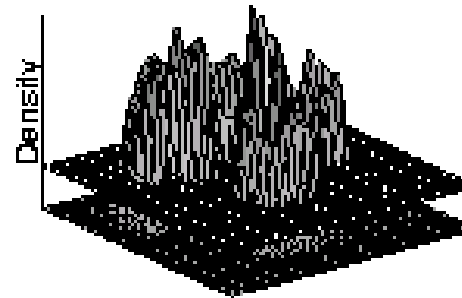
Figure 3: Example of Center-Defined Clusters for different  $\sigma$



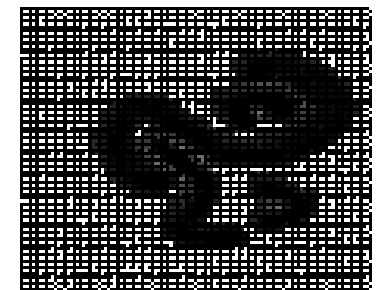
(a)  $\xi = 2$



(b)  $\xi = 2$



(c)  $\xi = 1$



(d)  $\xi = 1$

Figure 4: Example of Arbitrary-Shape Clusters for different  $\xi$

# DENCLUE: How to find the clusters

- Divide the space into grids, with size  $2\sigma$
- Consider only grids that are highly populated
- For each object, calculate its density attractor using hill climbing technique
  - Tricks can be applied to avoid calculating density attractor of all points
- Density attractors form basis of all clusters

# Features of DENCLUE

## ■ Major features

### □ Solid mathematical foundation

- Compact definition for density and cluster
- Flexible for both center-defined clusters and arbitrary-shape clusters

### □ But needs a large number of parameters

- $\sigma$ : parameter to calculate density
- $\xi$ : density threshold
- $\delta$ : parameter to calculate attractor