

Density-based clustering

Used in many applications

- Mainly for its efficiency, resistance to noise and ability to deal with arbitrary shaped clusters

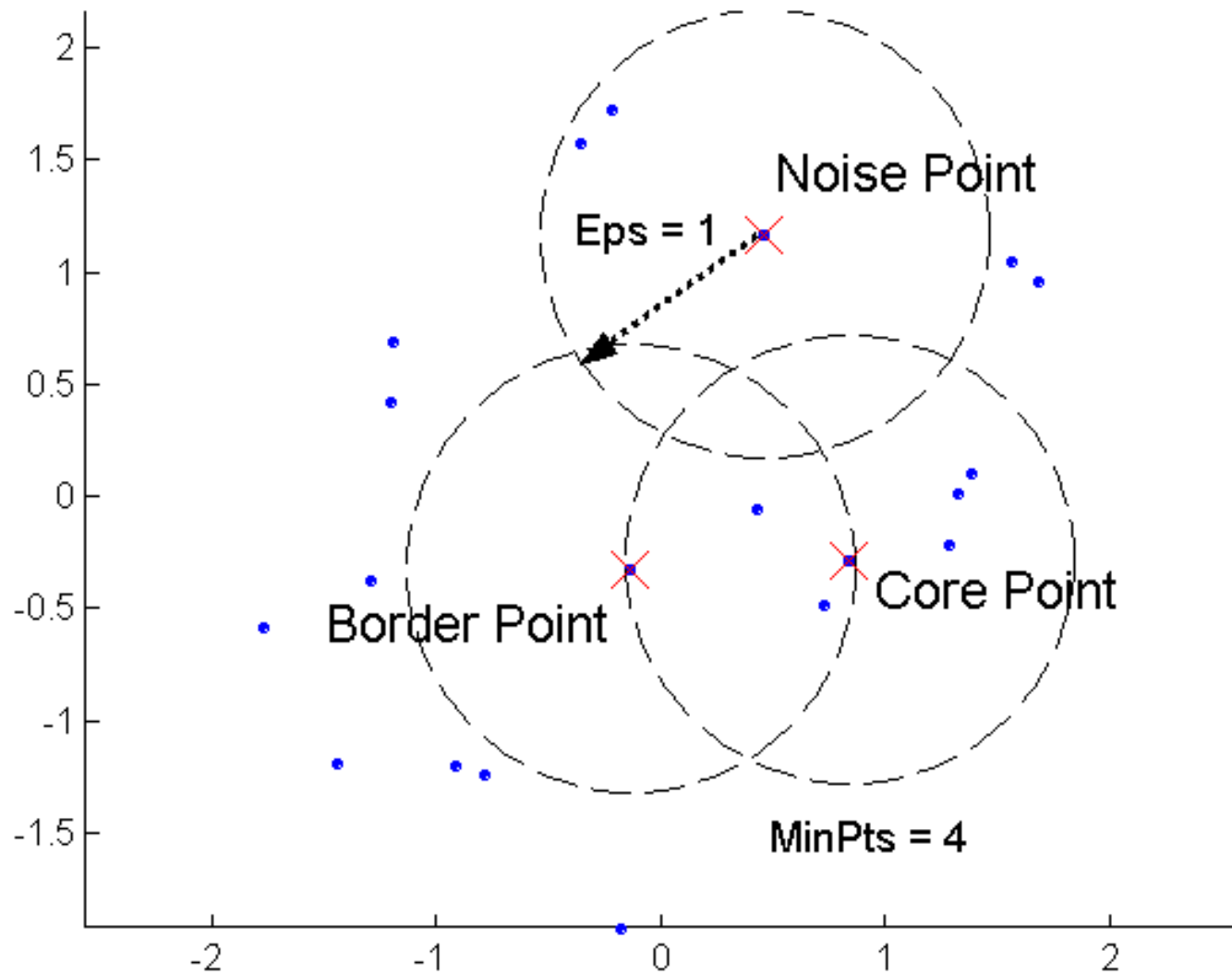
Main idea: divide noise from objects to clusters

- Objects to cluster = dense points
- Noise = low-density points

DBSCAN

- DBSCAN is a density-based algorithm.
 - Density = number of points within a specified radius (Eps)
 - A point is a **core point** if it has more than a specified number of points (MinPts) within Eps
 - ◆ These are points that are at the interior of a cluster
 - A **border point** has fewer than MinPts within Eps, but is in the neighborhood of a core point
 - A **noise point** is any point that is not a core point or a border point.

DBSCAN: Core, Border, and Noise Points



DBSCAN Algorithm

- ▮ Eliminate noise points
- ▮ Perform clustering on the remaining points

Algorithm 8.4 DBSCAN algorithm.

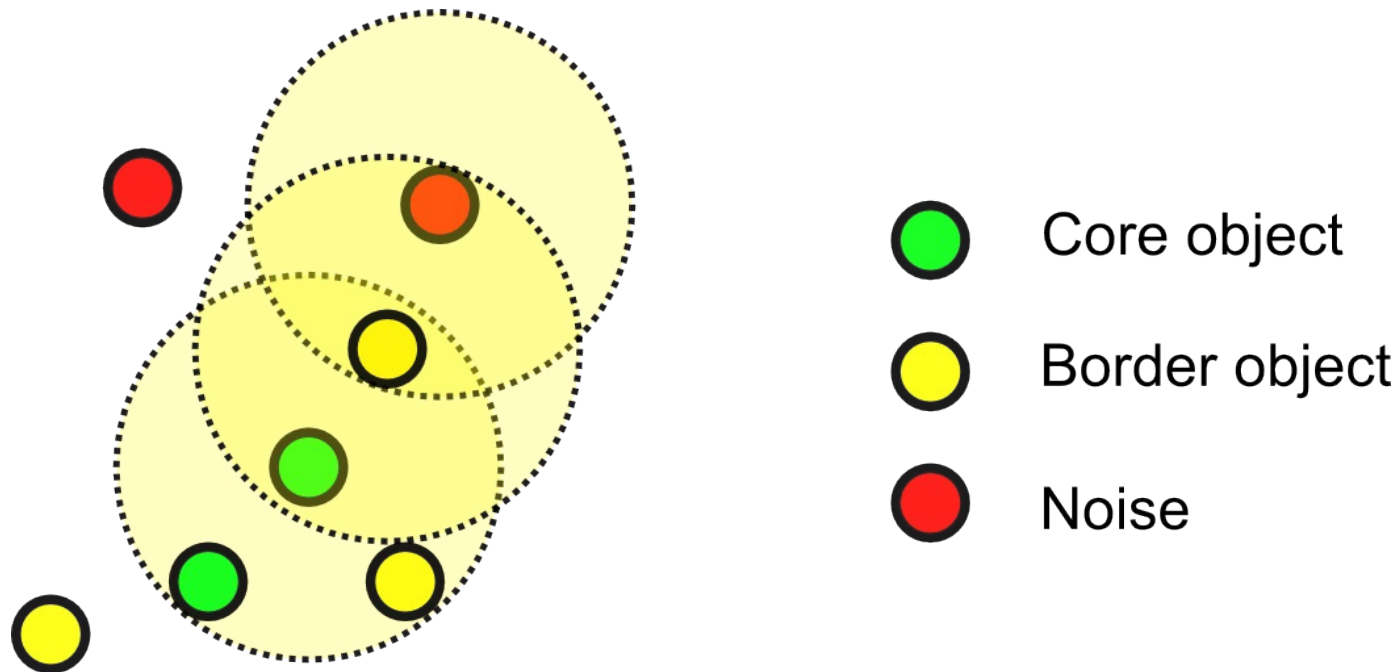
- 1: Label all points as core, border, or noise points.
- 2: Eliminate noise points.
- 3: Put an edge between all core points that are within Eps of each other.
- 4: Make each group of connected core points into a separate cluster.
- 5: Assign each border point to one of the clusters of its associated core points.

Border points can be neighbors of several core points/clusters
→ arbitrarily choose one!

DBSCAN

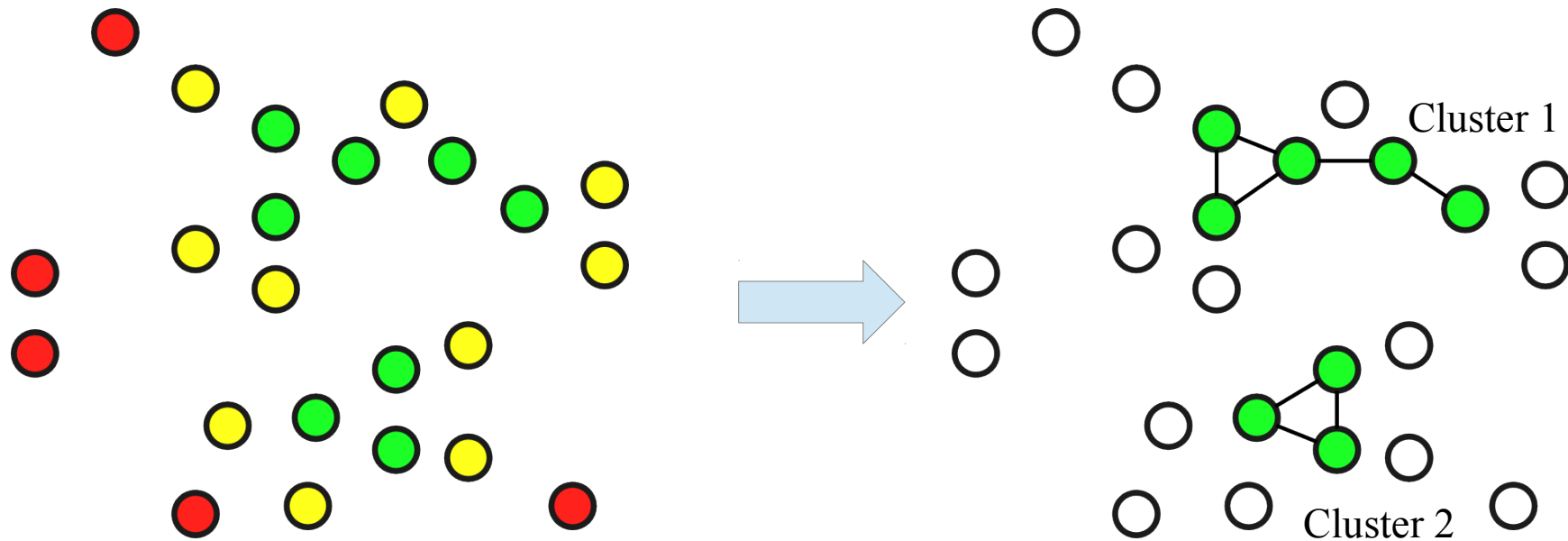
Step 1: label points as core (dense), border and noise

- Based on thresholds R (radius of neighborhood) and min_pts (min number of neighbors)



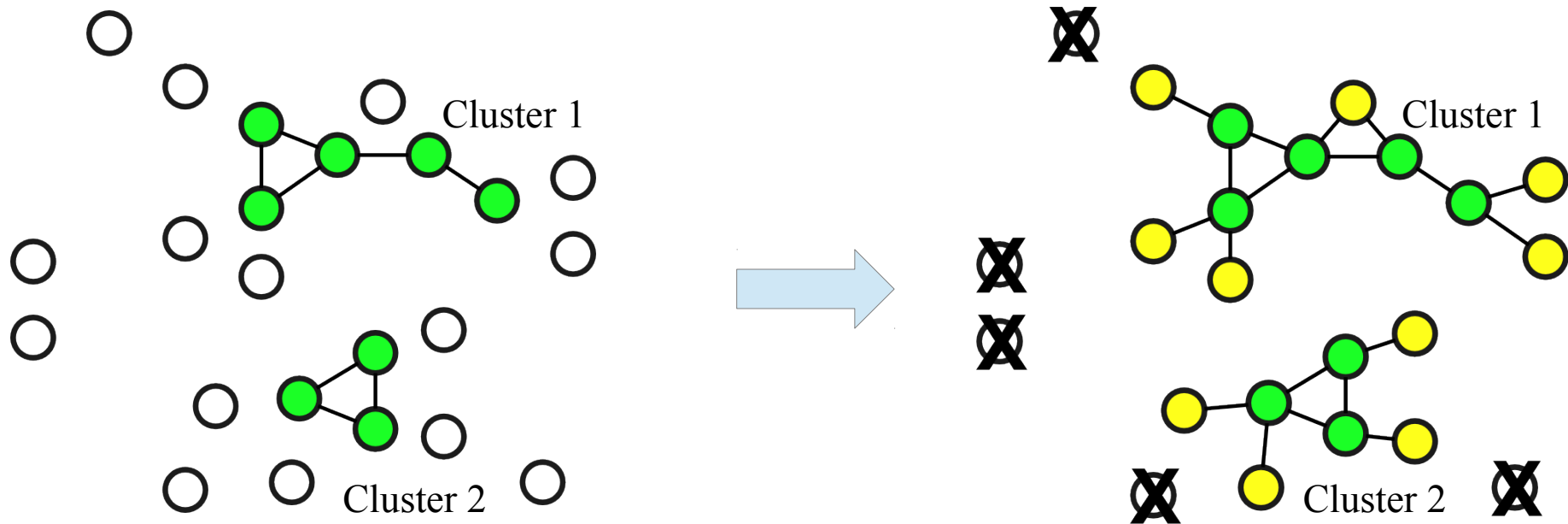
DBSCAN

Step 2: connect core objects that are neighbors, and put them in the same cluster



DBSCAN

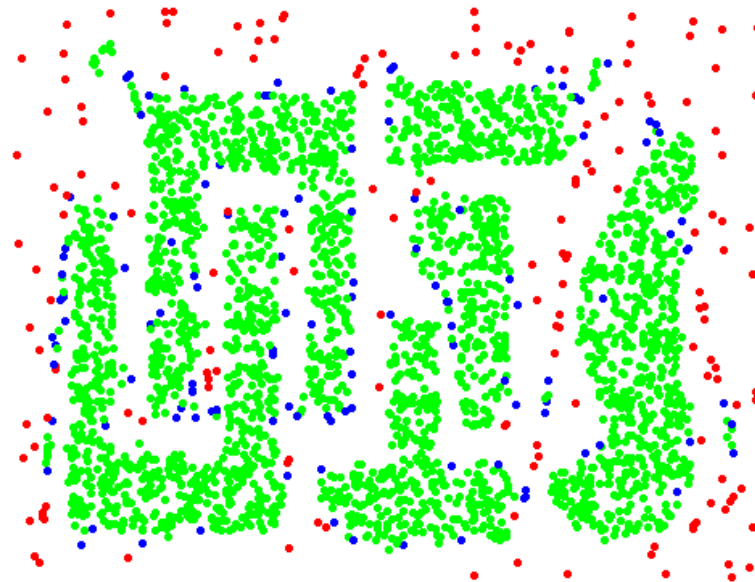
Step 3: associate border objects to (one of) their core(s), and remove noise



DBSCAN: Core, Border and Noise Points



Original Points



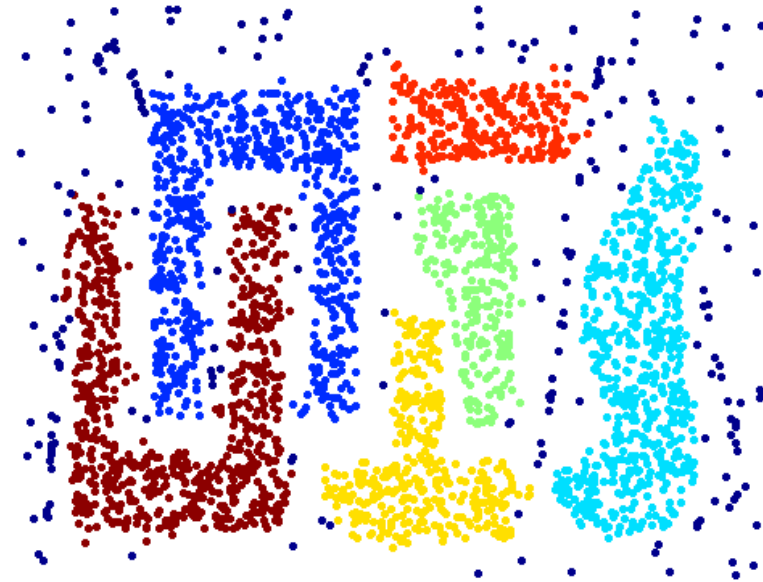
Point types: **core**,
border and **noise**

Eps = 10, MinPts = 4

When DBSCAN Works Well



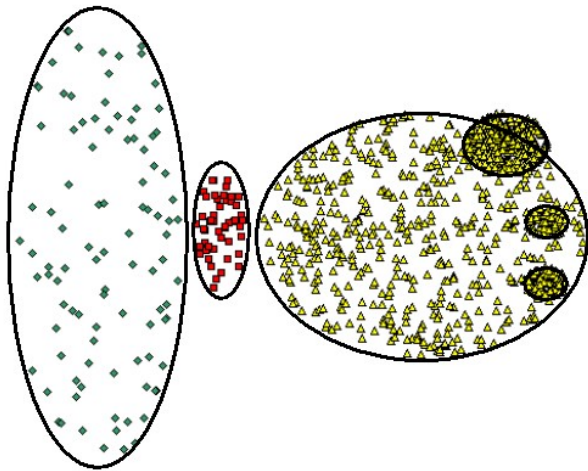
Original Points



Clusters

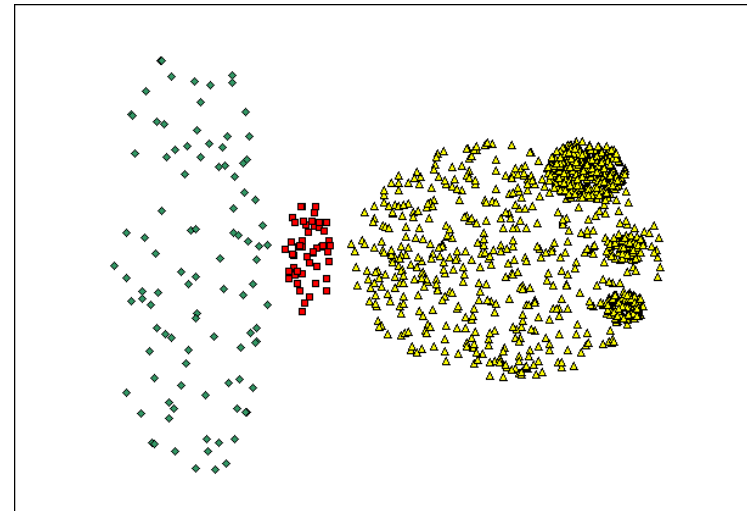
- **Resistant to Noise**
- **Can handle clusters of different shapes and sizes**

When DBSCAN Does NOT Work Well

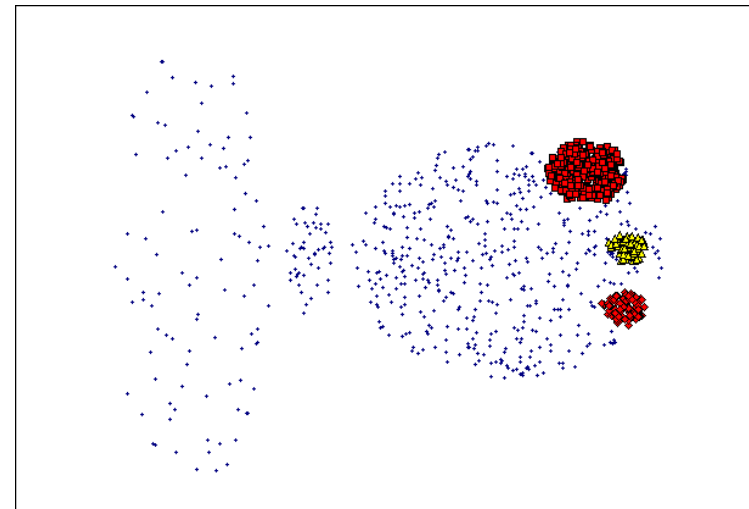


Original Points

- Varying densities
- High-dimensional data



(MinPts=4, Eps=9.92)



(MinPts=4, Eps=9.75).

DBSCAN: Determining EPS and MinPts

- Idea is that for points in a cluster, their k^{th} nearest neighbors are at roughly the same distance
- Noise points have the k^{th} nearest neighbor at farther distance
- So, plot sorted distance of every point to its k^{th} nearest neighbor

