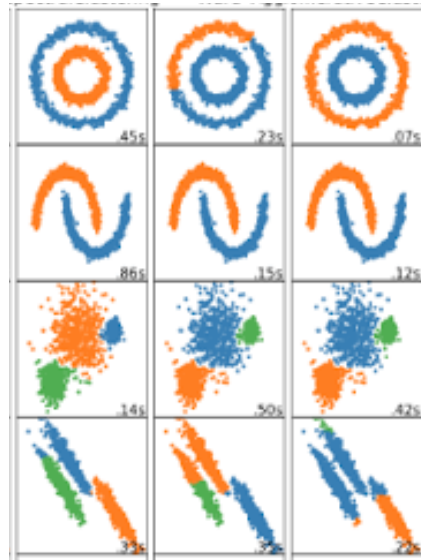# Density-Based Clustering Methods

- Partitioning methods (K-means, PAM clustering) work for finding spherical-shaped clusters or convex clusters.

- In other words, they are suitable only for compact and well-separated clusters.

- For data set containing nonconvex clusters k-means algorithm has difficulties for identifying these clusters with arbitrary shapes.

# Density-Based Clustering Methods

- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
  - Discover clusters of arbitrary shape
  - Handle noise
  - One scan
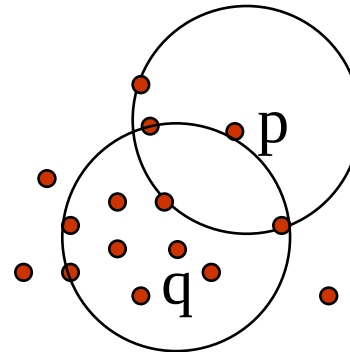  - Need density parameters as termination condition

The key idea is that for each point of a cluster, the neighborhood of a given radius has to contain at least a minimum number of points.

# Density-Based Clustering: Basic Concepts

- Two parameters*:*

  - *Eps*: Maximum radius of the neighbourhood

  - *MinPts*: Minimum number of points in an Eps-neighbourhood of that point

- $N_{Eps}(p)$: {q belongs to D | dist(p,q) ≤ Eps}

- Directly density-reachable: A point *p* is directly density-reachable from a point *q* w.r.t. *Eps*, *MinPts* if

  - *p* belongs to $N_{Eps}(q)$

  - core point condition:
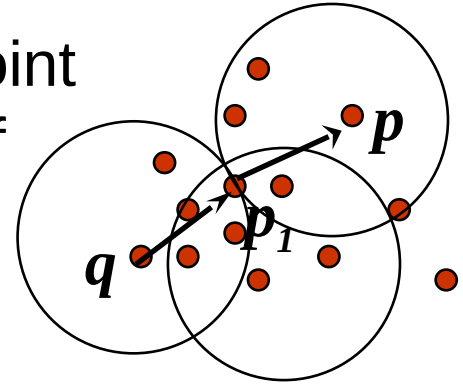
    $$|N_{Eps}(q)| ≥ MinPts$$

MinPts = 5

Eps = 1 cm

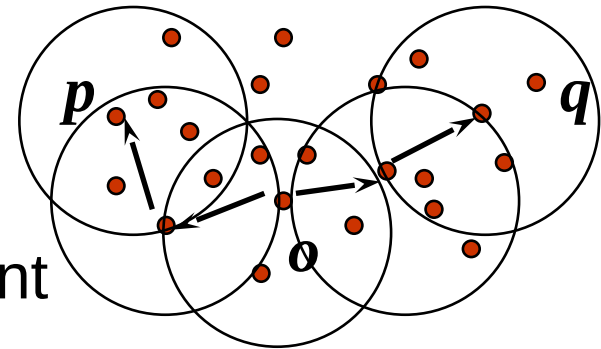# Density-Reachable and Density-Connected

- Density-reachable:

  - A point *p* is density-reachable from a point *q* w.r.t. *Eps*, *MinPts* if there is a chain of points $p_1$, …, $p_n$, $p_1 = q$, $p_n = p$ such that $p_{i+1}$ is directly density-reachable from $p_i$

  - *Except last all other must be a core point*

- Density-connected

  - A point *p* is density-connected to a point

  *q* w.r.t. *Eps*, *MinPts* if there is a point *o* such that both, *p* and *q* are density-reachable from *o* w.r.t. *Eps* and *MinPts*
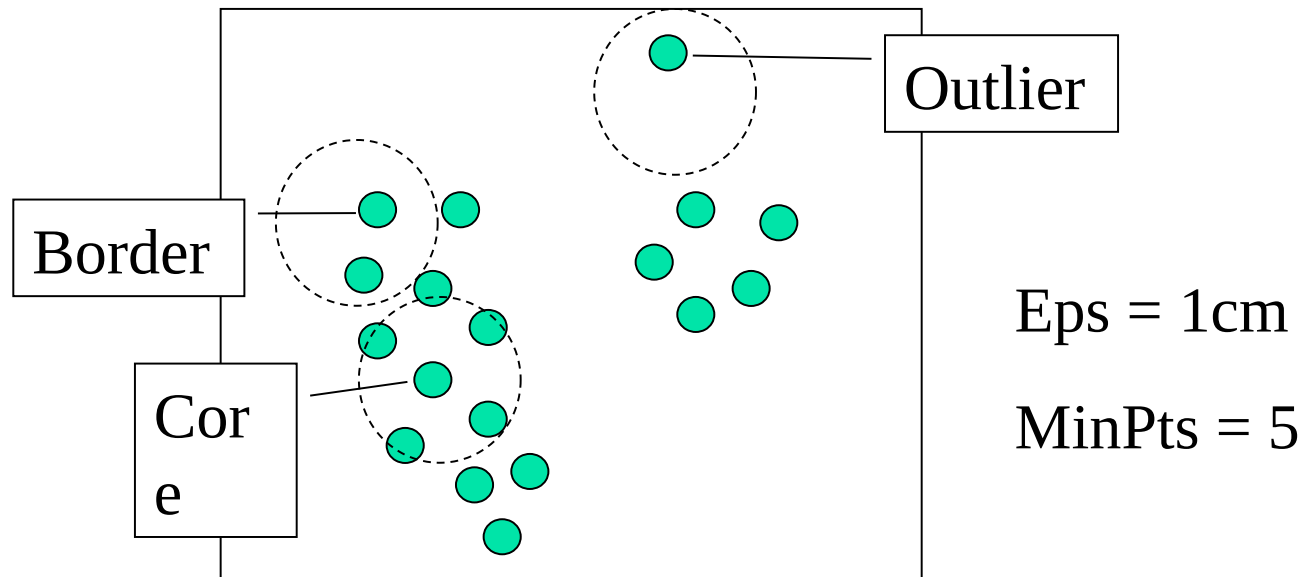
# DBSCAN: Density-Based Spatial Clustering of Applications with Noise

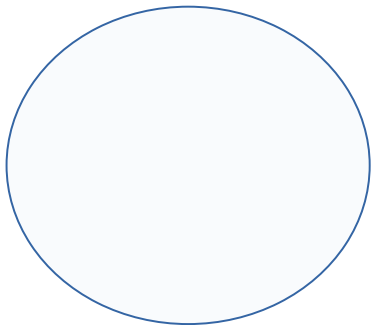Core Point: A point is a core point if it has more or equal to MinPts points within eps.

Border Point: A point which has fewer than MinPts within eps but it is in the neighborhood of a core point.

Noise or outlier: A point which is not a core point or border point.

Outlier

Border

Cor e

Eps = 1cm

MinPts = 5

# DBSCAN: The Algorithm

- Compute Neighbors of each points and identify core points

- Joint neighboring core points into clusters

- For each non-core points do:

  – Add to a neighboring core points if possible //border

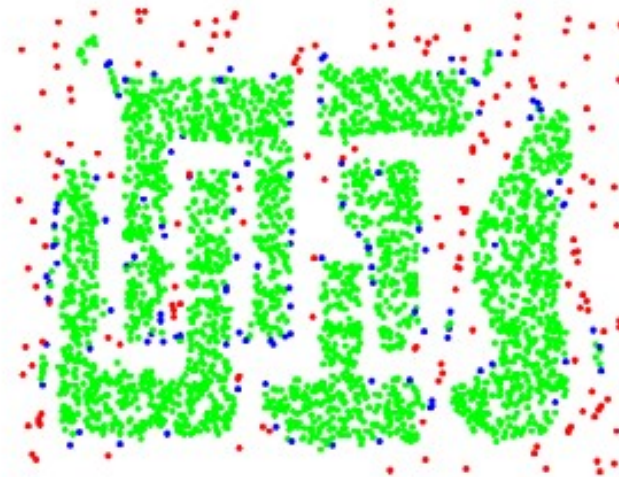  – Otherwise add it to noise //noise point

p6

p1

p2

p3

p4

p5

Original Points

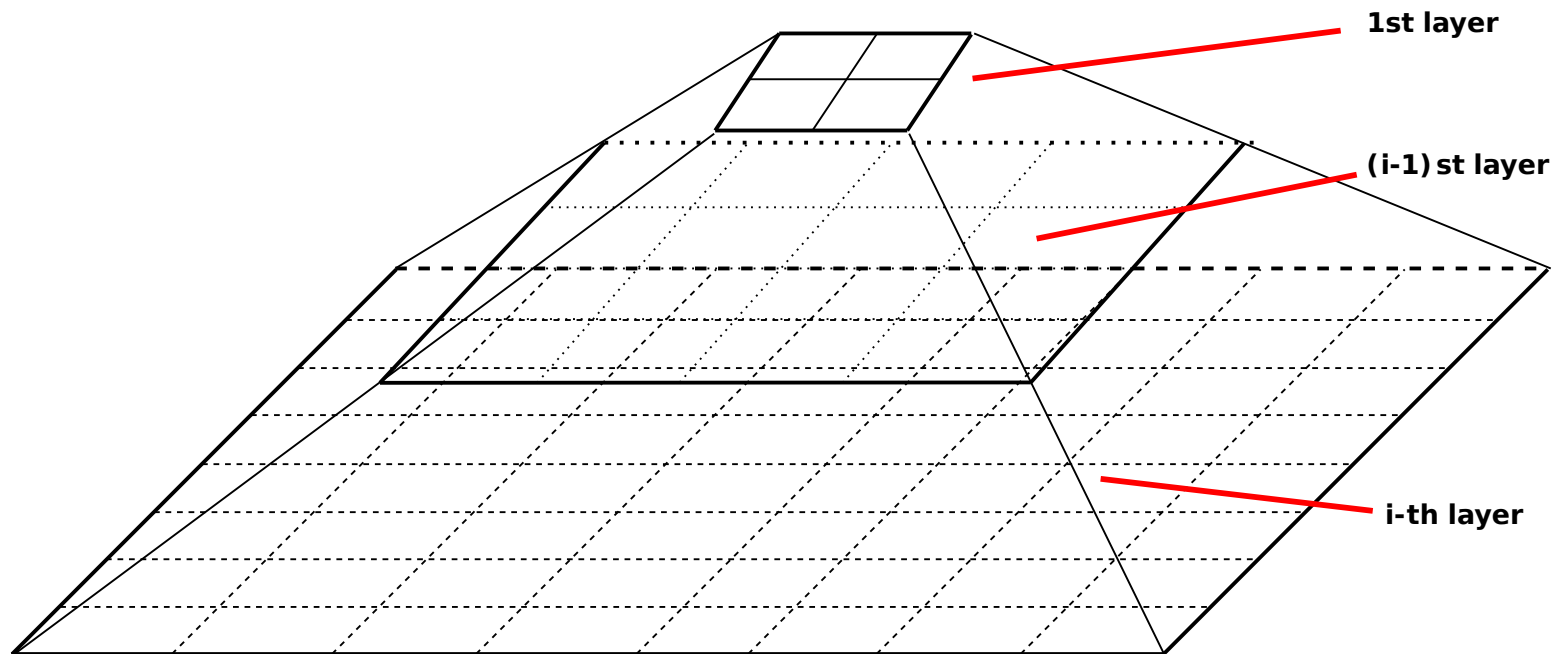Point types: core, border and noise

Eps = 10, MinPts = 4

# Grid-Based Methods

# Grid-Based Clustering Method

- Using multi-resolution grid data structure
- Several interesting methods
  - STING (a STatistical INformation Grid approach) by Wang, Yang and Muntz (1997)
  - WaveCluster by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
    - A multi-resolution clustering approach using wavelet method
  - CLIQUE: Agrawal, et al. (SIGMOD'98)
    - Both grid-based and subspace clustering

# STING: A Statistical Information Grid Approach

- Wang, Yang and Muntz (VLDB'97)
- The spatial area is divided into rectangular cells
- There are several levels of cells corresponding to different levels of resolution

1st layer

(i-1) st layer

i-th layer

# The STING Clustering Method

- Each cell at a high level is partitioned into a number of smaller cells in the next lower level

- Statistical info of each cell is calculated and stored beforehand and is used to answer queries

- Parameters of higher level cells can be easily calculated from parameters of lower level cell
  - *count*, *mean*, *s*, *min*, *max*
  - type of distribution—*normal*, *uniform*, etc.

- Use a top-down approach to answer spatial data queries

- Start from a pre-selected layer—typically with a small number of cells

- For each cell in the current level compute the confidence interval

# STING Algorithm and Its Analysis

- Remove the irrelevant cells from further consideration
- When finish examining the current layer, proceed to the next lower level
- Repeat this process until the bottom layer is reached
- Advantages:
  - Query-independent, easy to parallelize, incremental update

- Disadvantages:
  - All the cluster boundaries are either horizontal or vertical, and no diagonal boundary is detected