

# 7.02 DBSCAN Clustering

# DBSCAN Clustering

- Density-based spatial clustering of applications with noise (DBSCAN) clustering method
- It is an unsupervised method that:
  - clusters core samples (dense areas of data set)
  - denotes non-core samples (sparse portions of data set)
- Commonly used for outlier detection
- Outliers should make up  $\leq 5\%$  of the total observations
  - Enabled by adjusting model parameters accordingly

# DBSCAN Clustering

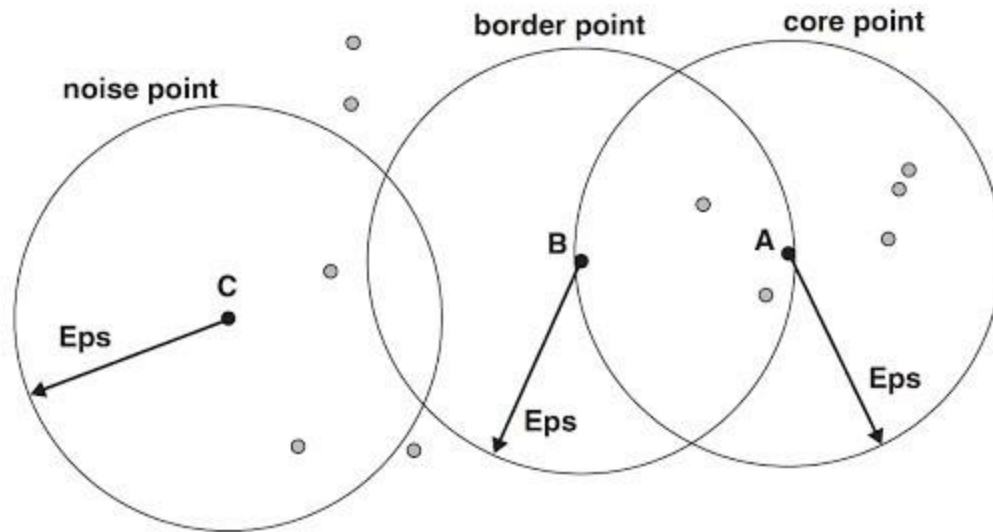
- Key Model Parameters
  - EPS
    - maximum distance between 2 samples for them to be clustered in the same neighbourhood
    - Start at a value of 0.1
  - min\_samples
    - minimum number of samples in a neighbourhood for a data point to be classified as a core point
    - Start with a very low sample size

# DBSCAN Clustering – Data Classification

- Based on these two parameters i.e., epsilon and min\_samples, we are first going to classify every point in our dataset into three categories:
  1. Core points
  2. Boundary points
  3. Noise points

# DBSCAN Clustering – Core points

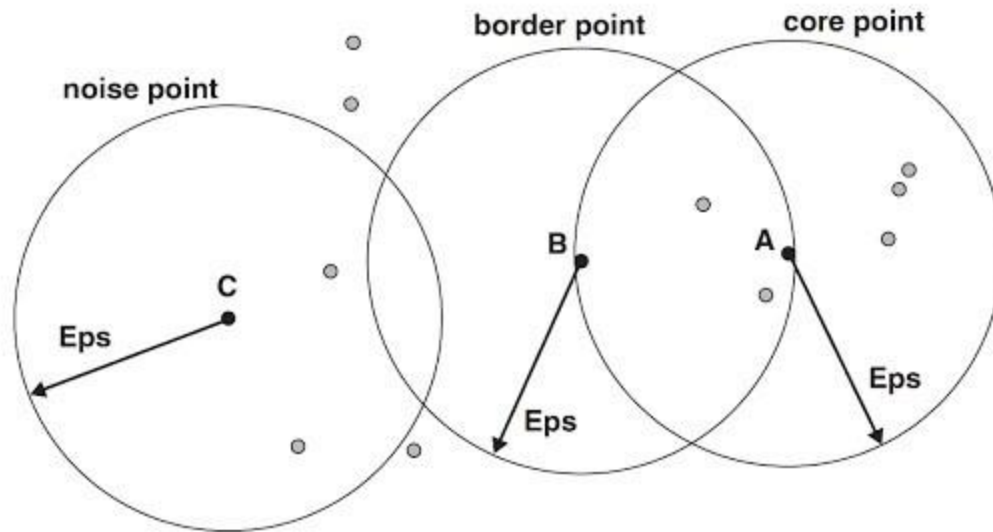
- A point is a core point when:
  - it has at least `min_samples` data points (including itself) within its  $\epsilon$ -neighborhood



Example: If we set `min_samples` = 5, then Point A satisfies this condition.

# DBSCAN Clustering – Boundary points

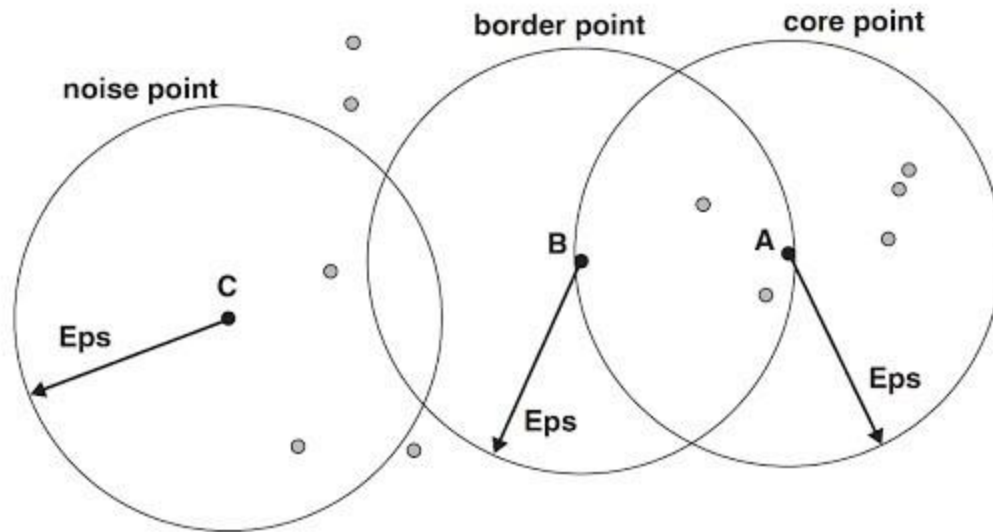
- A point is a boundary point when:
  - The number of neighbours is less than `min_samples`
  - The point is in neighbourhood of a core point



Example: Point B has less than `min_samples=5` in its neighbourhood and it is in the neighbourhood of another core point.

# DBSCAN Clustering – Noise points

- A point is a noise point when:
  - it is neither a core point nor a boundary point



Example: Point C is neither a core nor a boundary point

# Why do we need DBSCAN when we already have K-means clustering?

- **Automatic Determination of the Number of Clusters:**
  - K-means clustering requires the number of clusters to be specified in advance. This can be a limitation when the number of clusters is unknown or when the data does not naturally align with a predetermined number of clusters.
  - DBSCAN, on the other hand, does not require the number of clusters to be predefined. It automatically detects the number of clusters based on the data density and connectivity.
- **Robustness to Noise and Outliers:**
  - K-means clustering tries to assign all data points including noise points to clusters, even those that don't belong to any meaningful cluster.
  - DBSCAN, on the other hand, can effectively handle noisy data by designating them as noise points or outliers.
  - It does not force every point to belong to a cluster, allowing for a more robust identification of meaningful clusters in the presence of noise.



# Algorithmic steps for DBSCAN clustering

- The algorithm proceeds by arbitrarily picking up a point in the dataset (until all points have been visited)
- If there are at least 'min\_sample' points within a radius of ' $\epsilon$ ' to the point then we consider all these points to be part of the same cluster
- The clusters are then expanded by recursively repeating the neighborhood calculation for each neighboring point

# Why should we use DBSCAN?

- The DBSCAN algorithm should be used to find associations and structures in data that are hard to find manually but that can be relevant and useful to find patterns and predict trends.
- Clustering methods are usually used in biology, medicine, social sciences, archaeology, marketing, characters recognition, management systems, among others.

# Why should we use DBSCAN?

- Example: Suppose we have an online store, and we want to improve our sales by recommending relevant products to our customers.
  - We don't know exactly what our customers are looking for but based on a data set we can predict and recommend a relevant product to a specific customer.
  - We can apply DBSCAN to our e-commerce database and find clusters based on the products that the users have bought.
  - Using this clusters we can find similarities between customers, for example,
    - Customer A has bought 1 pen, 1 book and 1 scissors
    - Customer B has bought 1 book and 1 scissors
    - Then we can recommend 1 pen to the customer B