# NOVA
# IMS
**Information Management School**

# Data Mining

## Density-based Clustering (dbscan)

**16/11/2021**

**NOVA-IMS**

Fernando Lucas Bação

bacao@isegi.unl.pt

http://www.isegi.unl.pt/fbacao

**Instituto Superior de Estatística e Gestão de Informação**
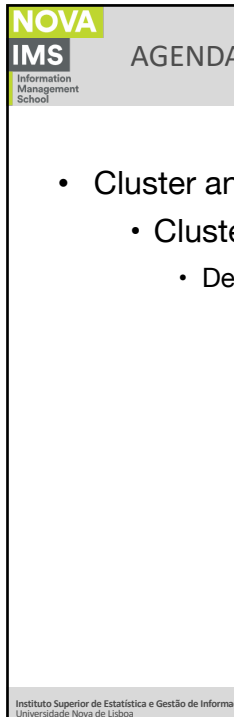Universidade Nova de Lisboa

1

---

# NOVA
# IMS
Information Management School

## AGENDA

- Cluster analysis
  - Clustering techniques
    - Density-based clustering (DBscan)

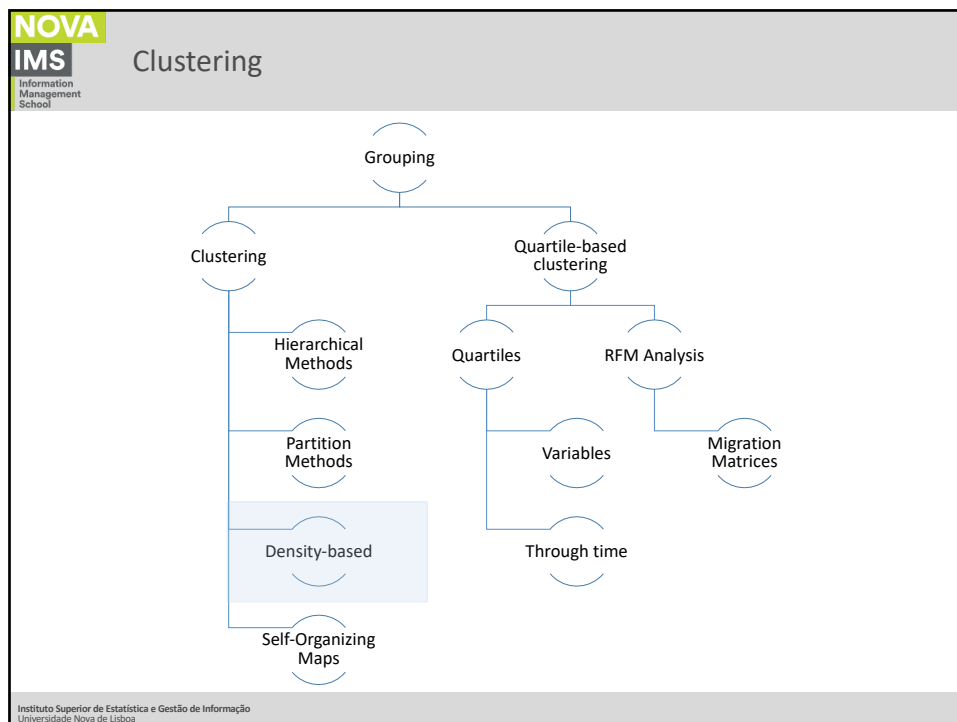Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

2

3



4

# Density-Based Clustering (DBSCAN)

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

5

---

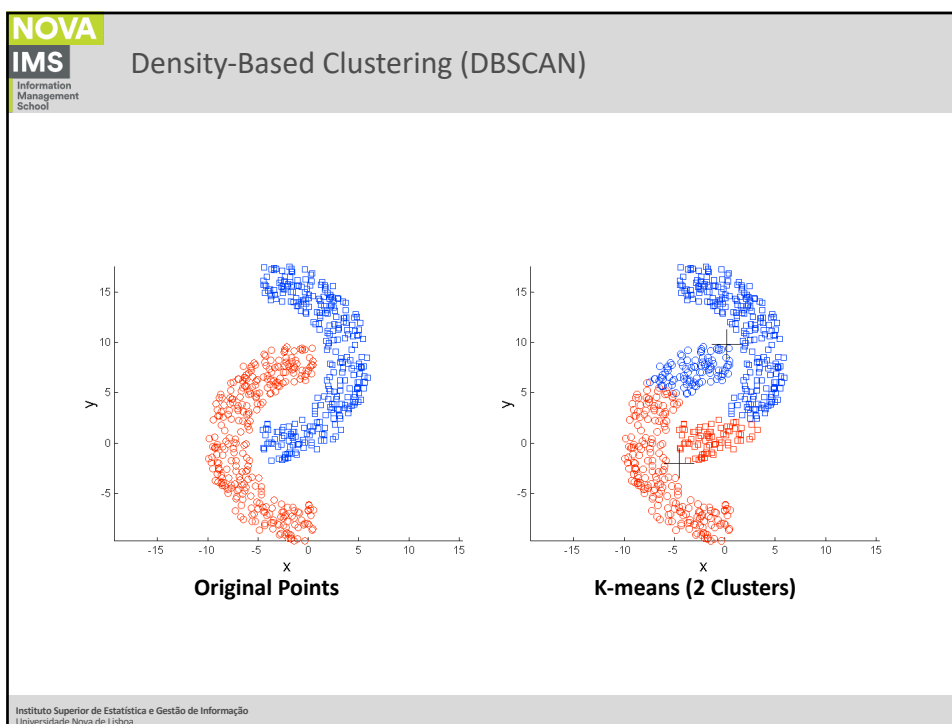## Density-Based Clustering (DBSCAN)

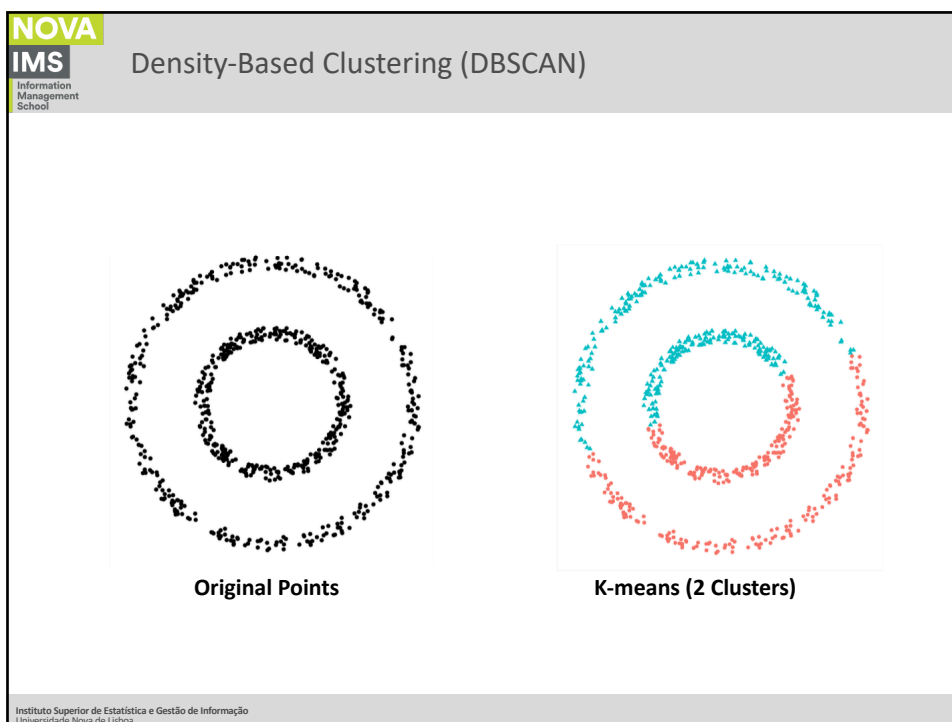- **The idea**

  - Previous clustering methods have some limitations. They are based on a particular set of assumptions that if not true, the process yields suboptimal results.

  - That is they would likely inaccurately identify non-convex regions, and where noise or outliers are included in the clusters.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

6

## Density-Based Clustering (DBSCAN)

**Original Points**

**K-means (2 Clusters)**

7

## Density-Based Clustering (DBSCAN)

**Original Points**

**K-means (2 Clusters)**

8

NOVA
IMS
Information
Management
School

Density-Based Clustering (DBSCAN)

- **The idea**

  - Density-based clustering algorithms try to **find clusters in data without assuming a particular shape**. Thus solving correctly cases like the ring example.

  - To find clusters of **arbitrary shape**, these methods model clusters as **dense regions in data space**, separated by sparse regions.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

---

NOVA
IMS
Information
Management
School

Density-Based Clustering (DBSCAN)

- **The idea**

  - We want to be able to cluster data like this:



Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

10

NOVA
IMS
Information
Management
School

Density-Based Clustering (DBSCAN)

- **Characteristics**

  - What do talk about when we talk about density?

  - What is our physical intuition behind density?

    - We say a point **p** is in a dense region if there are **many points** in the neighborhood of **p**. In this sense, the density around a point can be measured by the number of points surrounding it.

    - To measure the density around a point **p** we use the tradition topological definition of neighborhood. The $\varepsilon$-*neighborhood* of a point **p** is the space within a radius $\varepsilon > 0$ centered in **p**.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

11

NOVA
IMS
Information
Management
School

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

  - DBSCAN takes two input parameters:

    - $\varepsilon$ - the radius defining the neighborhood

    - *MinPts* - the minimum of points in $\varepsilon$ -neighborhood

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

12

## Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

    - DBSCAN importante concepts:

        - If $z$ is a point that have at least *MinPts* in its $\varepsilon$ -neighborhood is called **core point**.

        - $x$ is **border point**, if the number of its neighbors is less than *MinPts*, but it belongs to the $\varepsilon$ -neighborhood of some core point $z$

        - If a point is neither a **core** nor a **border** point, then it is called a **noise point** or an **outlier**.

13

## Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density



- Assuming MinPts = 6.
    - **x** is a **core point** because $\varepsilon$ -neighborhood (x) = 6,
    - **y** is a **border point** because $\varepsilon$ –neighborhood (y) < *MinPts*, but it belongs to the $\varepsilon$ -neighborhood of the core point x.
    - **z** is a **noise point**.

14

**NOVA**
**IMS**
Information
Management
School

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

    - We define 3 terms, required for understanding the DBSCAN algorithm:

        - **Direct density reachable**: A point "A" is directly density reachable from another point "B" if:

            - i) "A" is in the $\varepsilon$ –neighborhood of "B" and

            - ii) "B" is a core point.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

15

**NOVA**
**IMS**
Information
Management
School

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

    - We define 3 terms, required for understanding the DBSCAN algorithm:

        - **Density reachable**: A point "A" is density reachable from "B" if there are a set of core points leading from "B" to "A.

        - **Density connected**: Two points "A" and "B" are density connected if there are a core point "C", such that both "A" and "B" are density reachable from "C".

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

16

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

17



Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density
  - A **density-based cluster** is defined as a **group of density connected points**.

18

## Density-Based Clustering (DBSCAN)

- **Algorithm**

    1. DBSCAN begins with an arbitrary data point that has not been visited. The $\varepsilon$ – neighborhood of this point is extracted (all points which are within the $\varepsilon$ distance are neighborhood points).

    2. If there are a sufficient number of points (according to *MinPts*) within this neighborhood then the clustering process starts and the current data point becomes the first point in the new cluster. Otherwise, the point will be labeled as noise. In both cases that point is marked as "visited".

    3. For this first point in the new cluster, the points within its $\varepsilon$ distance neighborhood also become part of the same cluster. This procedure of making all points in the $\varepsilon$ –neighborhood belong to the same cluster is then repeated for all of the new points that have been just added to the cluster group.

    4. This process of steps 2 and 3 is repeated until all points in the cluster are determined i.e all points within the $\varepsilon$ –neighborhood of the cluster have been visited and labeled.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

19

## Density-Based Clustering (DBSCAN)



epsilon = 1.00
minPoints = 4

Restart    Pause

https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

20

## Density-Based Clustering (DBSCAN)

- **Advantages**
  - It does not require a pre-set number of clusters at all.
  - It identifies outliers as noises (not affected by),
  - it can finds arbitrarily sized and arbitrarily shaped clusters quite well.

- **Disadvantages**
  - It doesn't perform well when the clusters are of varying density
  - Setting of the distance threshold $\varepsilon$ and *MinPts* for identifying the neighborhood points will vary from cluster to cluster when the density varies
  - Doesn't work well in high-dimensional spaces

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

21

# Questions?

22

22

## Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

    1. For each point $x_i$, compute the distance between $x_i$ and the other points. Finds all neighbor points within $\varepsilon$ –neighborhood of the starting point $x_i$. Each point, with a neighbor count greater than or equal to *MinPts*, is marked as **core point** or visited.

    2. For each **core point**, if it's not already assigned to a cluster, create a new cluster. Find recursively all its density connected points and assign them to the same cluster as the **core point**.

    3. Iterate through the remaining unvisited points in the dataset.

    Those points that do not belong to any cluster are treated as outliers or noise.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

23