# Course. Introduction to Machine Learning
# Work 1. Clustering Exercise
## Session 1
## Course 2020-2021

Dr. Maria Salamó Llorente

Dept. Mathematics and Informatics,

Faculty of Mathematics and Informatics,

University of Barcelona

# Contents

1. Introduction  (session 1)
2. Preprocess the data (session 1)
3. DBSCAN with sklearn (session 1)
4. K-Means (your own code) (session 2)
5. Bisecting K-Means  (your own code) (session 2)
6. K-medians, K-means++ or k-harmonic means (your own code) (session 2)
7. Fuzzy clustering (your own code) (session 3)
8. Validation techniques  (using sklearn validation metrics) (session 3)

# Introduction

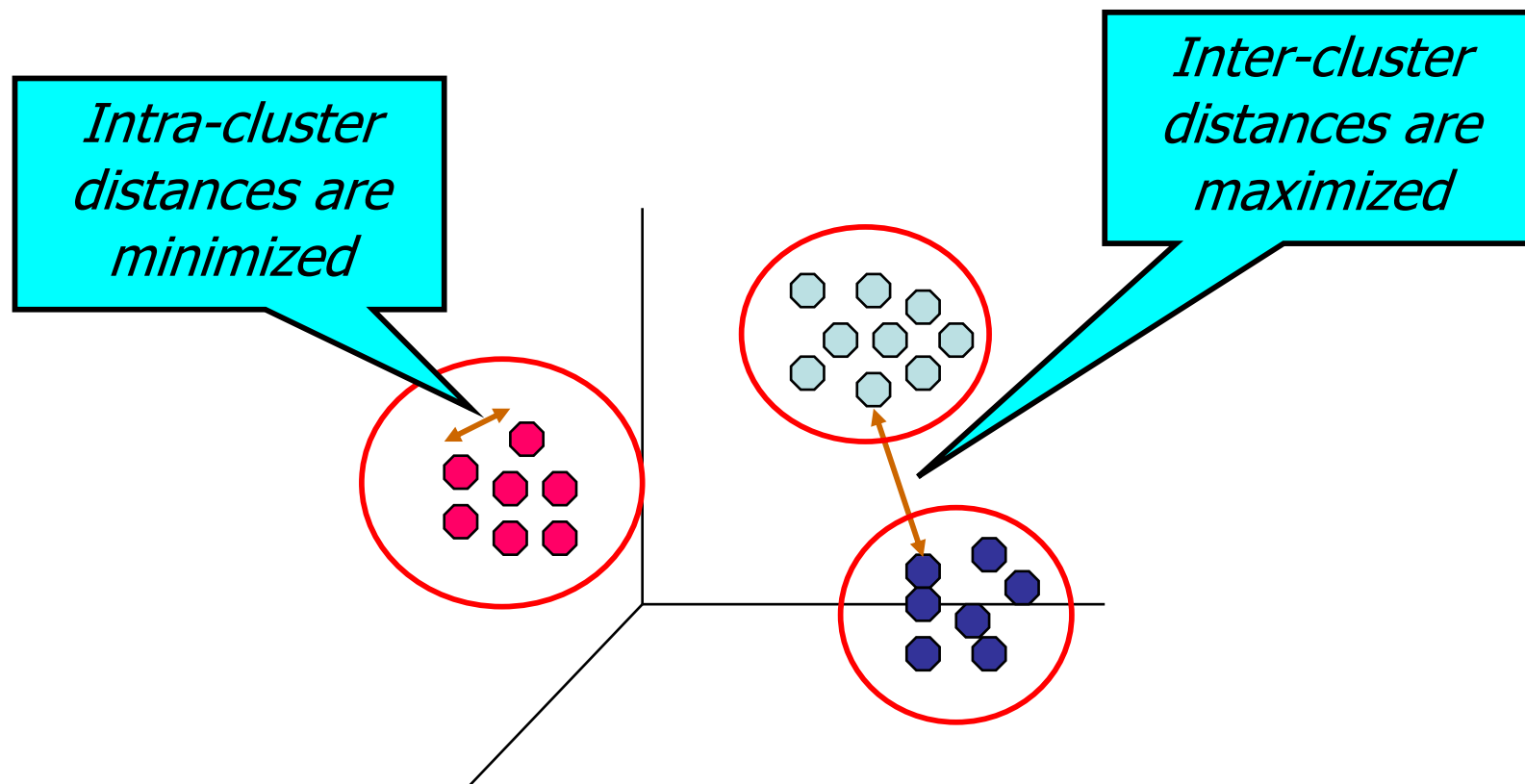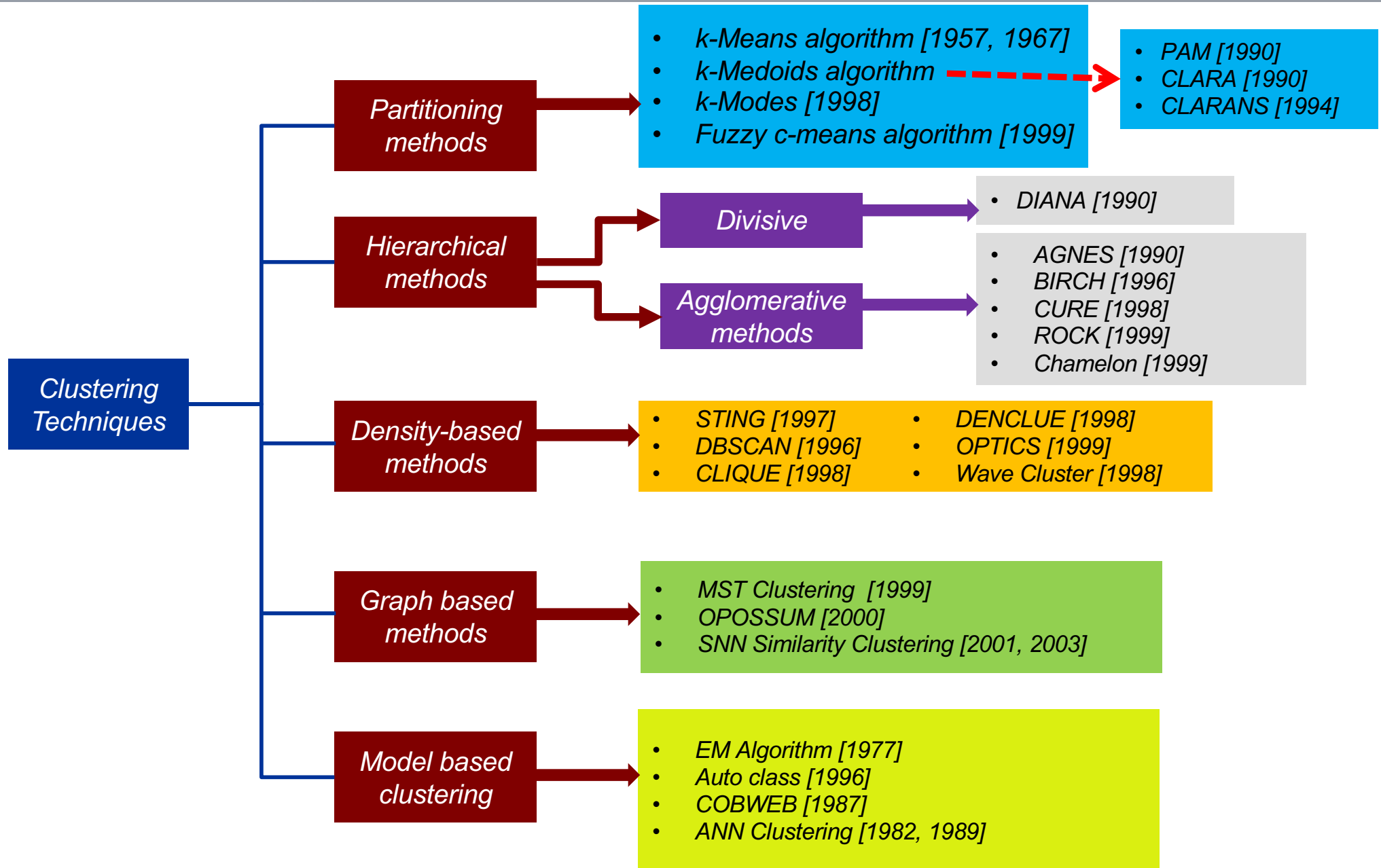- In general a **grouping** of objects such that the objects in a group (cluster) are similar (or related) to one another and different from (or unrelated to) the objects in other groups



*Intra-cluster distances are minimized*

*Inter-cluster distances are maximized*

# Taxonomy of Clustering Algorithms

UNIVERSITAT DE BARCELONA

**Clustering Techniques**

**Partitioning methods**
- k-Means algorithm [1957, 1967]
- k-Medoids algorithm
- k-Modes [1998]
- Fuzzy c-means algorithm [1999]

- PAM [1990]
- CLARA [1990]
- CLARANS [1994]

**Hierarchical methods**

*Divisive*
- DIANA [1990]

*Agglomerative methods*
- AGNES [1990]
- BIRCH [1996]
- CURE [1998]
- ROCK [1999]
- Chamelon [1999]

**Density-based methods**
- STING [1997]
- DBSCAN [1996]
- CLIQUE [1998]
- DENCLUE [1998]
- OPTICS [1999]
- Wave Cluster [1998]

**Graph based methods**
- MST Clustering [1999]
- OPOSSUM [2000]
- SNN Similarity Clustering [2001, 2003]

**Model based clustering**
- EM Algorithm [1977]
- Auto class [1996]
- COBWEB [1987]
- ANN Clustering [1982, 1989]

- You need to implement the code using Python 3.6 and Pycharm IDE

- Packages allowed in this exercise:
  - arff_loader
  - numpy
  - pandas
  - scipy
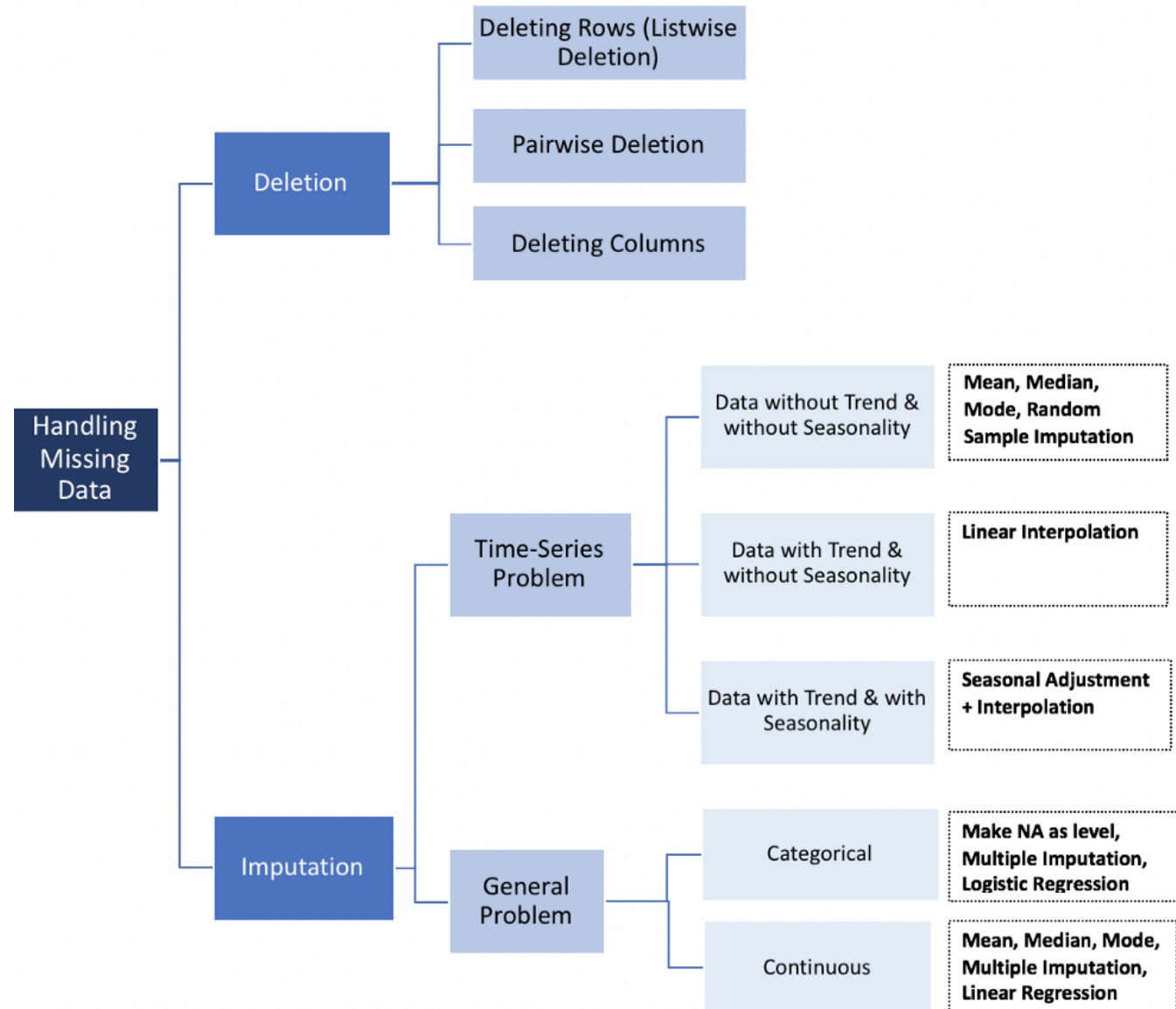  - sklearn (only for some parts)
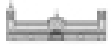  - matplotlib
  - seaborn

# Preprocess the data

- You need to read the .arff file
  - You can implement your own code or use scipy.io.arff.loadarff
- Data needs pre-processing
  - Features may contain **different ranges**
    - Normalize or Standarize the machine learning data
  - Features may have **different types**
    - Categorical, Numerical, and mix-type data
  - Data may contain **missing values**
    - Use the median (for example)

- **To deal with different ranges**
  - Normalize or scale features

- Alternatives
  - **Standardisation**: Standardisation replaces the values by their Z scores. `sklearn.preprocessing.scale`
  - **Mean normalisation**: This distribution will have values between **-1 and 1** with **μ=0**.
    `sklearn.preprocessing.StandardScaler`
  - **Min-Max scaling**: This scaling brings the value between 0 and 1. `sklearn.preprocessing.MinMaxScaler`
  - **Unit vector**: Scaling is done considering the whole feature vector to be of unit length.
    `sklearn.preprocessing.Normalizer`

- **To deal with different types**
- Alternatives
  - **Label encoding**: convert to a number
    `sklearn.preprocessing.LabelEncoder`
  - **One hot encoding**: where a categorical variable is converted into a binary vector, each possible value of the categorical variable becomes the variable itself with default value of zero and the variable which was the value of the categorical variable will have the value 1.
    `sklearn.preprocessing.OneHotEncoder`

- **To deal with missing values**

# DBSCAN
# Density-Based Clustering

Using sklearn

# Some Links

- http://www2.cs.uh.edu/~ceick/7363/Papers/dbscan.pdf

- https://youtu.be/sKRUfsc8zp4

- https://youtu.be/6jl9KkmgDlw

- https://scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html#sklearn.cluster.DBSCAN

- Clustering based on density (local cluster criterion), such as **density-connected points** or based on an explicitly constructed density function
- Major features:
  - Discover clusters of arbitrary shape
  - Handle noise
  - One scan
  - Need density parameters
- Several interesting studies:
  - <u>DBSCAN</u>: Ester, et al. (KDD'96)
  - <u>DENCLUE</u>: Hinneburg & D. Keim  (KDD'98/2006)
  - <u>OPTICS</u>: Ankerst, et al (SIGMOD'99).
  - <u>CLIQUE</u>: Agrawal, et al. (SIGMOD'98)

- DBSCAN is a Density-Based Clustering algorithm

- Reminder: In density based clustering we partition points into dense regions separated by not-so-dense regions.

- Important Questions:
  – How do we measure density?
  – What is a dense region?

- DBSCAN:
  – Density at point p: number of points within a circle of radius Eps
  – Dense Region: A circle of radius Eps that contains at least MinPts points

## Characterization of points

- 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 points belong in a dense region and 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

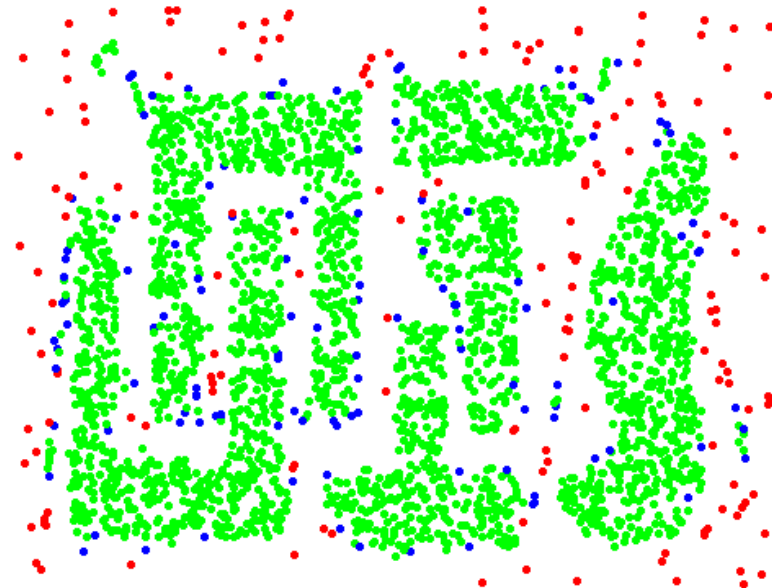- A **noise point** is any point that is not a core point or a border point

Original Points

Point types: core, border and noise

Eps = 10, MinPts = 4

- Parameters must be specified by the user
  - $\varepsilon$ = physical distance (radius),
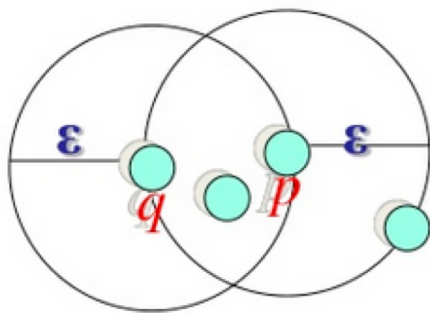  - *minPts* = desired minimum cluster size

## *minPts*

  - derived from the number of dimensions *D* in the data set, as *minPts* ≥ *D* + *1*
  - *minPts* = *1* does not make sense, as then every point on its own will already be a cluster
  - *minPts* must be chosen at least 3. Larger is better.
  - larger the dataset, the larger the value of *minPts* should be chosen

## $\varepsilon$

  - value can be chosen by using a k-distance graph
  - If $\varepsilon$ is chosen much too small, a large part of the data will not be clustered
  - If too high value, majority of objects will be in the same cluster
  - In general, small values of $\varepsilon$ are preferable

- Ɛ-Neighborhood: Objects within a radius of Ɛ from an object (epsilon-neighborhood)

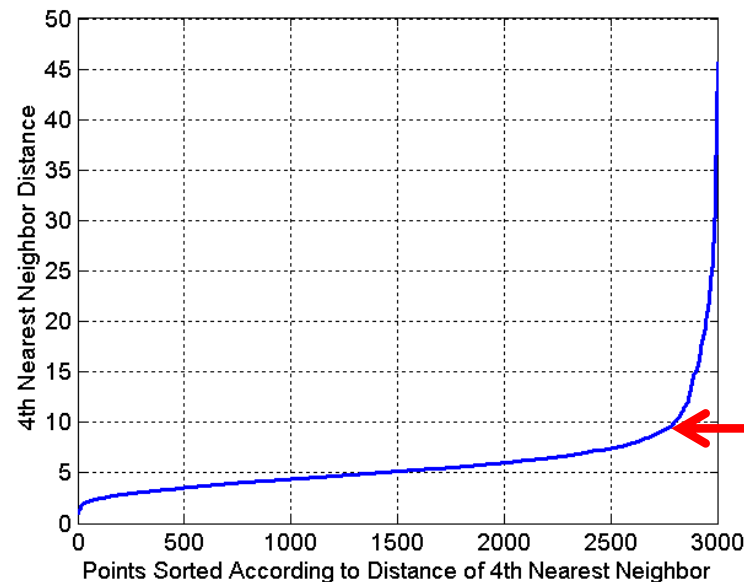- Core objects: Ɛ-Neighborhood of an object contains at least *MinPts* of objects

Ɛ-Neighborhood of $p$
Ɛ-Neighborhood of $q$

$p$ is a core object (MinPts = 4)

$q$ is not a core object

1. Create a graph whose nodes are the points to be clustered
2. For each core-point c create an edge from c to every point p in the **ε-neighborhood** of c
3. Set N to the nodes of the graph;
4. If N does not contain any core points terminate
5. Pick a core point c in N
6. Let X be the set of nodes that can be reached from c by going forward;
   1. create a cluster containing X∪{c}
   2. N=N/(X∪{c})
7. Continue with step 4

*Remark: points that are not assigned to any cluster are outliers;*

- 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
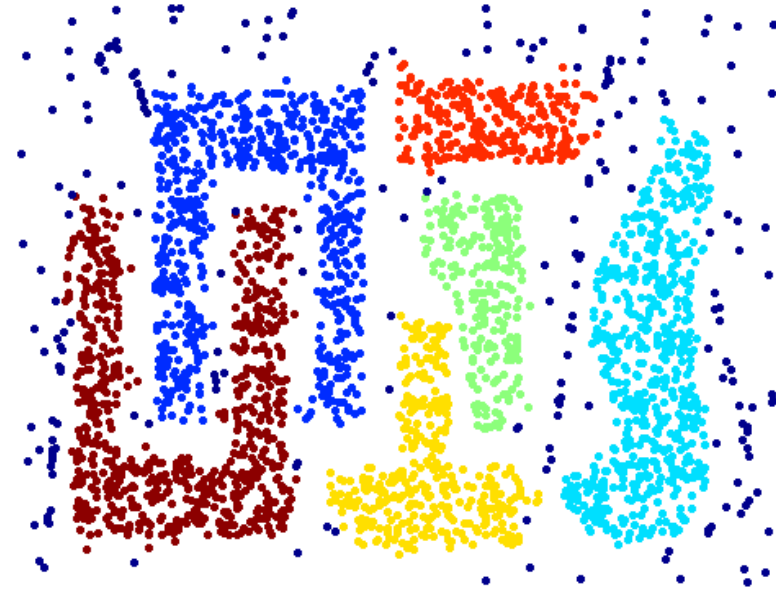- Find the distance d where there is a "knee" in the curve
  - Eps = d, MinPts = k



*Eps ~ 7-10*
*MinPts = 4*

*Original Points*

*Clusters*

- Resistant to Noise

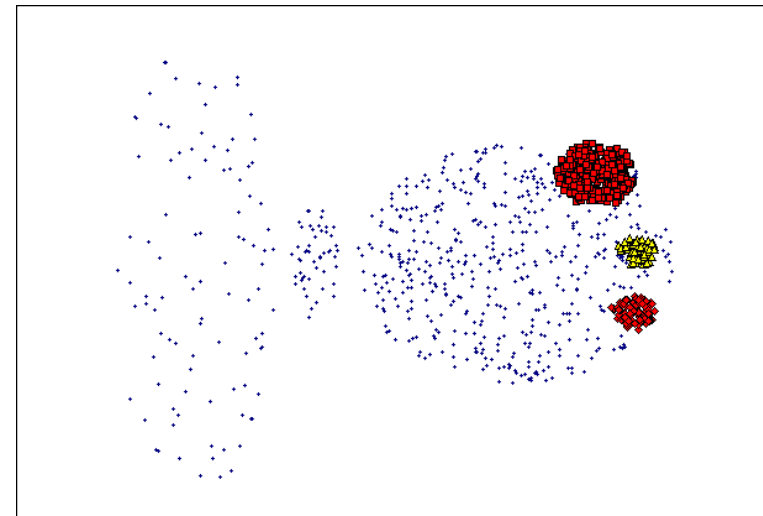- Can handle clusters of different shapes and sizes

*(MinPts=4, Eps=9.75).*

*Original Points*

- Varying densities
- High-dimensional data

*(MinPts=4, Eps=9.92)*

- Time Complexity: $O(n^2)$
  - For each point it has to be determined if it is a core point
  - Can be reduced to `O(n*log(n))` in lower dimensional spaces by using efficient data structures (where n is the number of objects to be clustered);

- Space Complexity: $O(n)$