

Introduction to Cluster Analysis



Fabio Stella

Department of Informatics, Systems and Communication

University of Milan-Bicocca

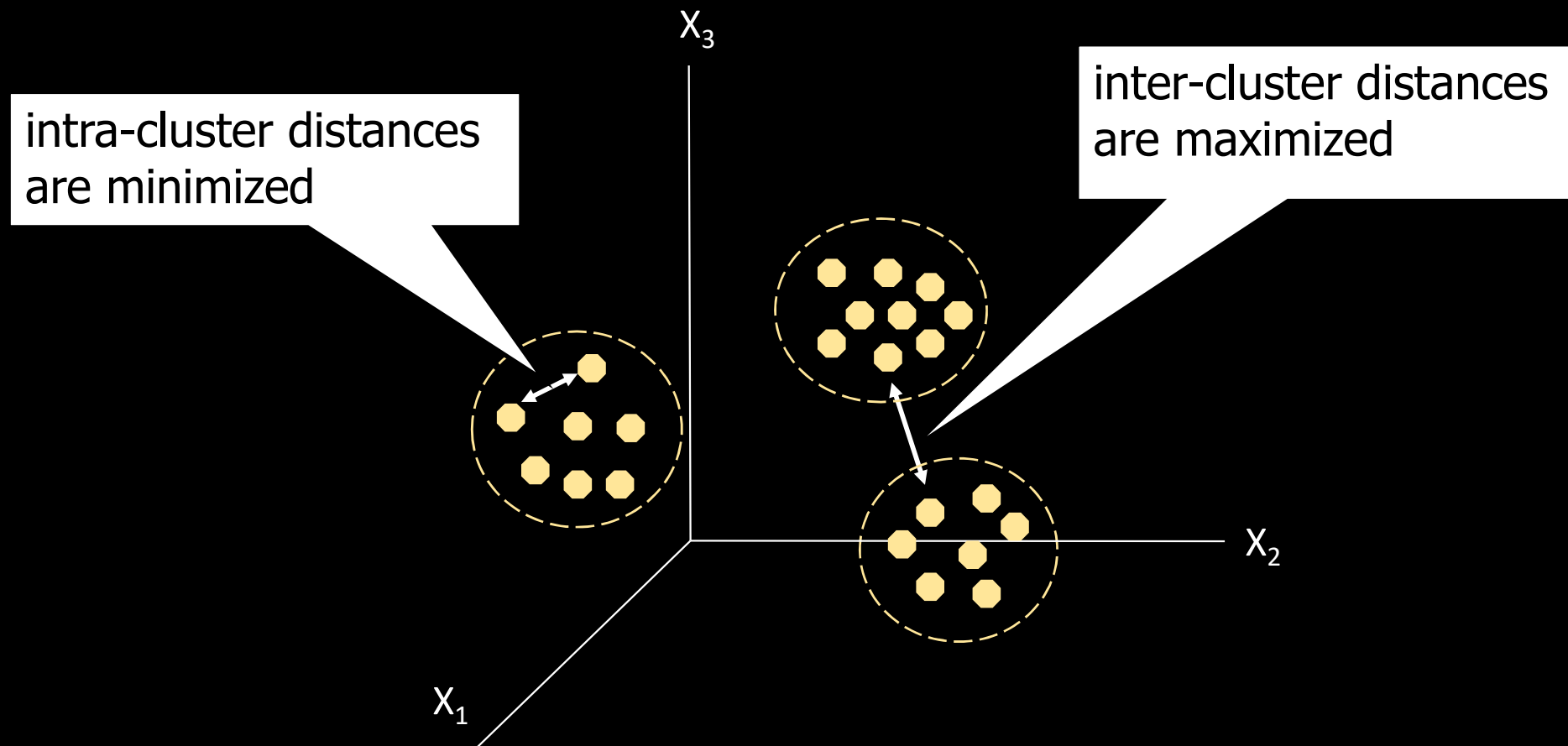
fabio.stella@unimib.it

OUTLOOK

- **CLUSTER ANALYSIS**
- **UNDERSTANDING AND SUMMARIZING**
- **TYPES OF CLUSTERING**
- **TYPES OF CLUSTERS**
- **COMPONENTS OF CLUSTER ANALYSIS**

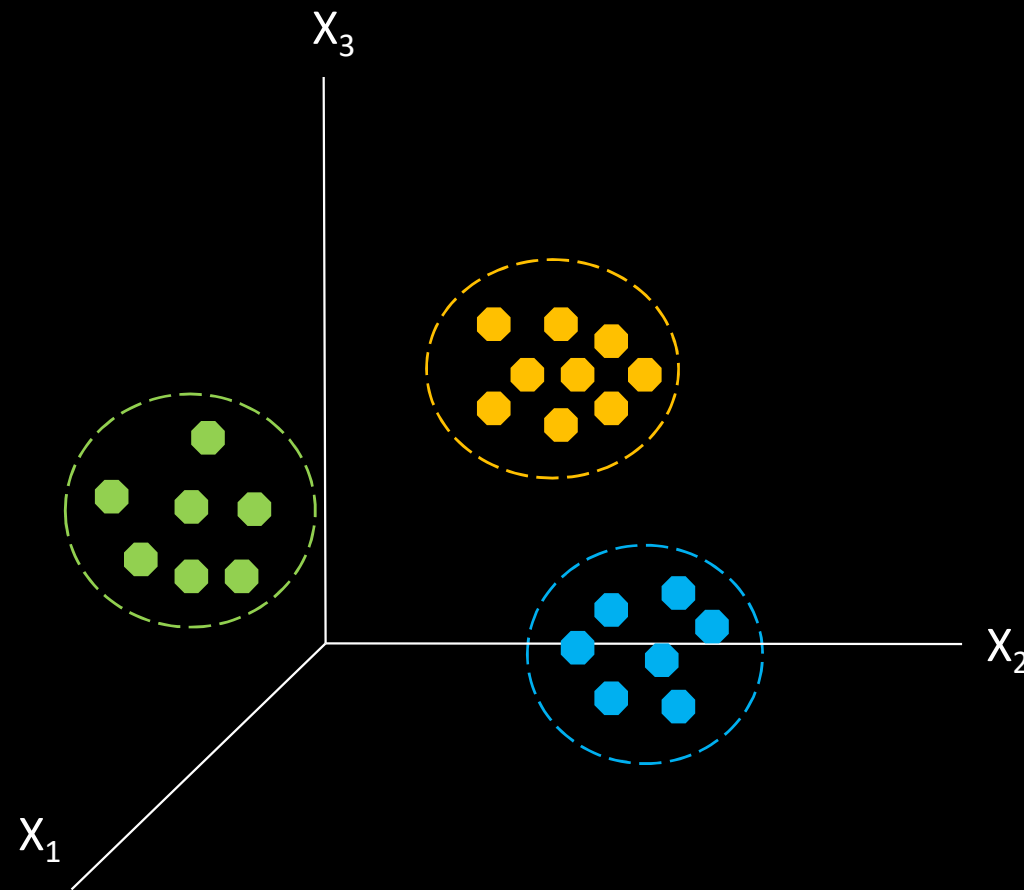
WHAT IS CLUSTER ANALYSIS?

Given a set of objects, place them in groups such that the objects in a group are similar (or related) to one another and different from (or unrelated to) the objects in other groups.



WHAT IS CLUSTER ANALYSIS?

Given a set of objects, place them in groups such that the objects in a group are similar (or related) to one another and different from (or unrelated to) the objects in other groups.

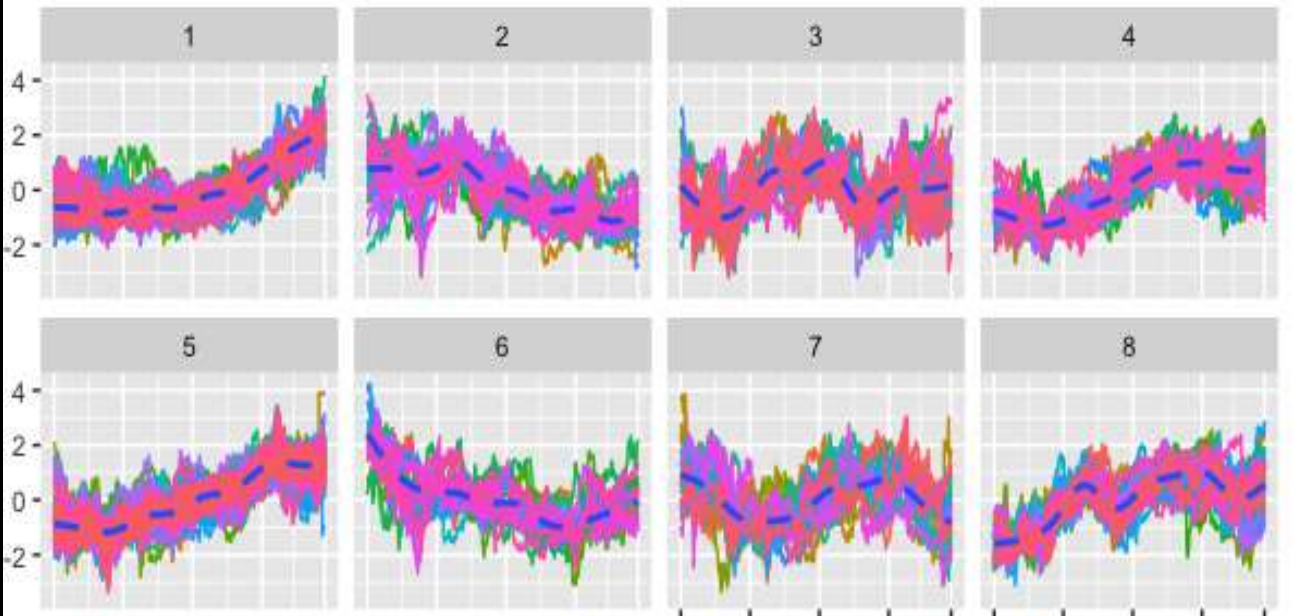
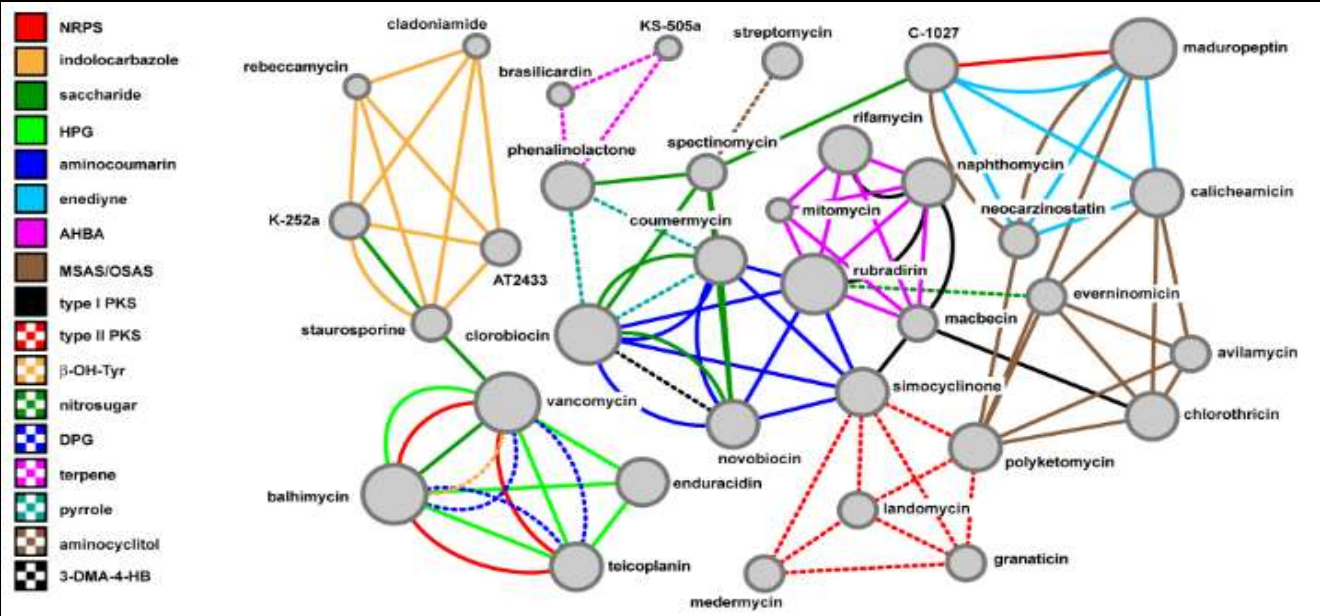


WHICH APPLICATIONS OF CLUSTER ANALYSIS?

■ UNDERSTANDING

Group related documents for browsing, group genes and proteins that have similar functionality, or group stocks with similar price fluctuations.

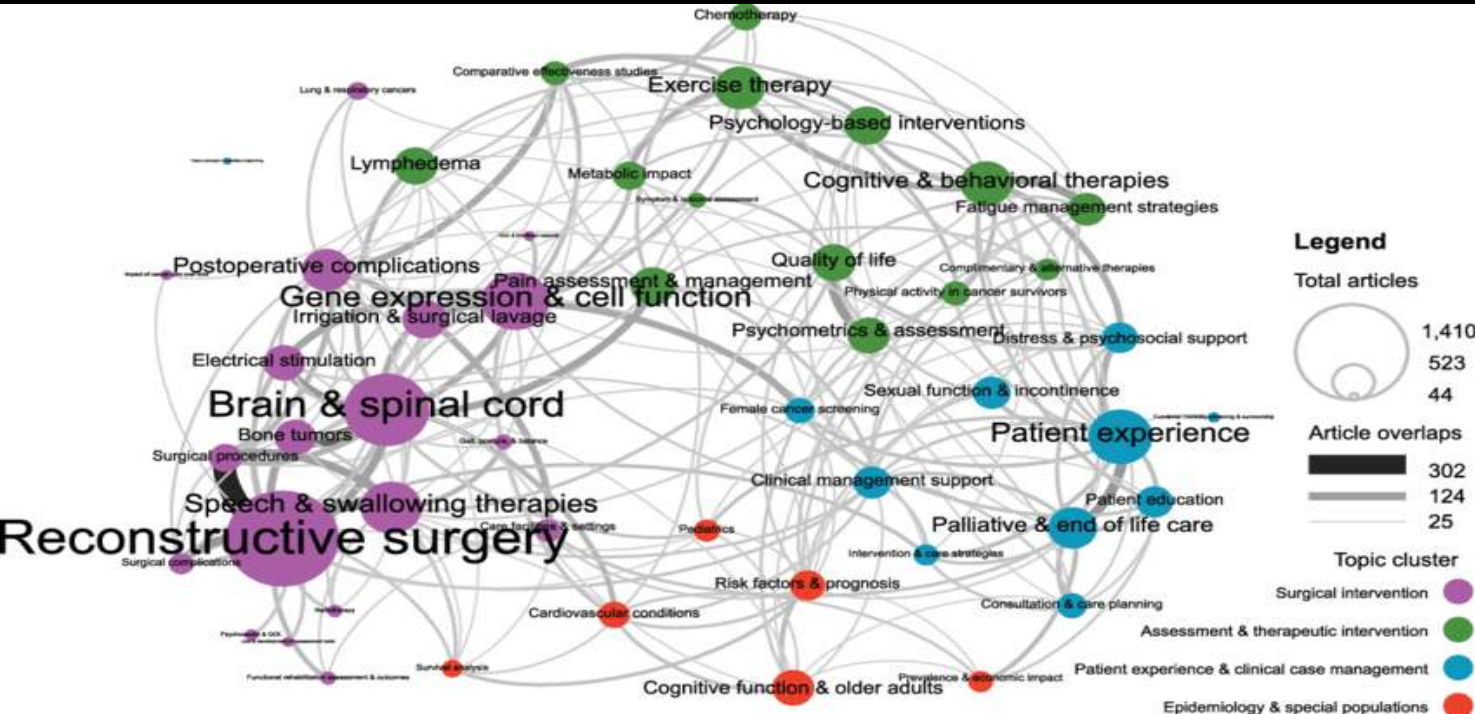
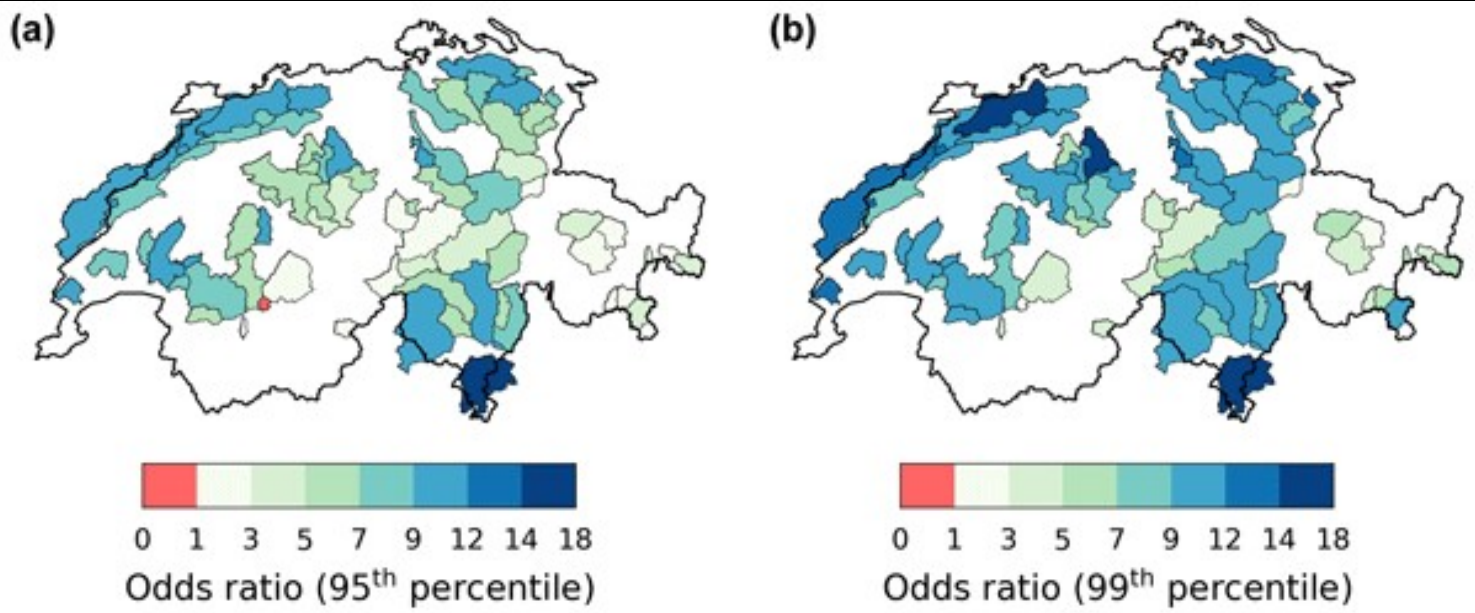
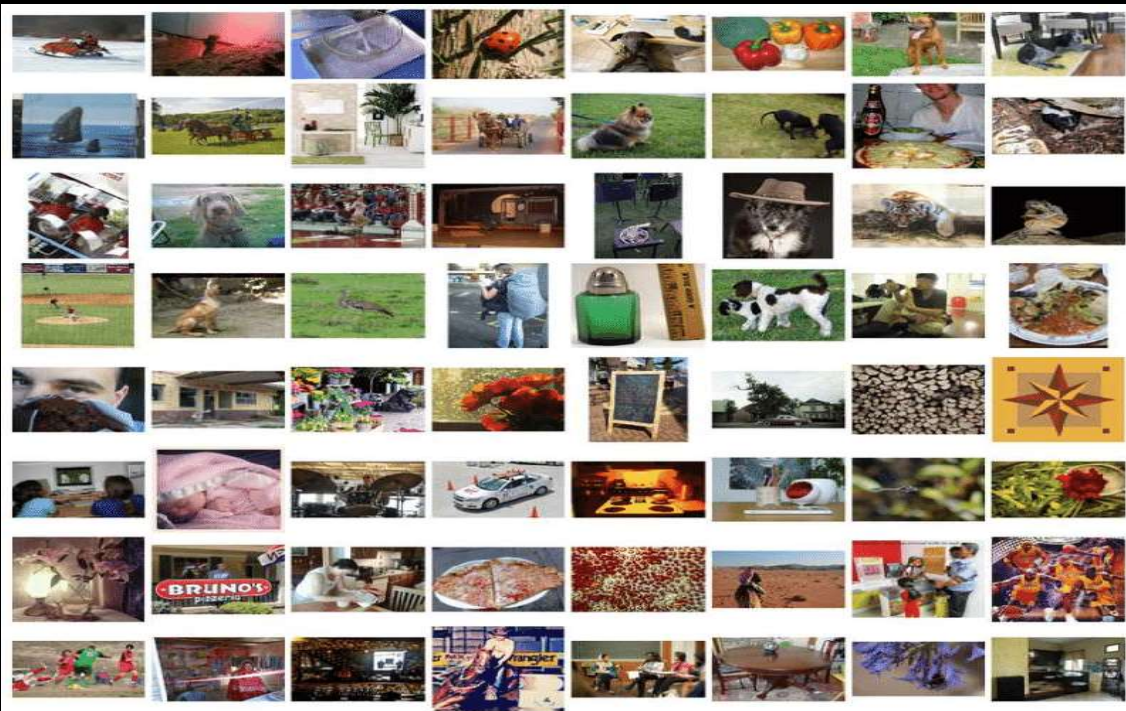
1. [Apache Solr - The Apache Software Foundation](#)
Solr is highly reliable, scalable and fault tolerant, providing distributed indexing, replication and load-balanced querying, automated failover and recovery, ...
<http://lucene.apache.org/solr/> [Bing, Google, Yahoo]
2. [Apache Solr - Features](#)
Solr is a standalone enterprise search server with a REST-like API. You put documents in it (called "indexing") via JSON, XML, CSV or binary over HTTP.
<http://lucene.apache.org/solr/features.html> [Bing, Google]
3. [FrontPage - Solr Wiki](#)
Nov 12, 2014 ... Official documentation for the latest release of Solr can be found on the Solr website. Of particular note is the Solr Reference Guide which is ...
<https://wiki.apache.org/solr/> [Google]
4. [Solr Quick Start - Apache Lucene - The Apache Software Foundation](#)
This document covers getting Solr up and running, ingesting a variety of data sources into multiple collections, and getting a feel for the Solr administrative and ...
<http://lucene.apache.org/solr/quickstart.html> [Google]
5. [Apache Solr - Wikipedia, the free encyclopedia](#)
Solr (pronounced "solar") is an open source enterprise search platform, written in Java, from the Apache Lucene project. Its major features include full-text ...
https://en.wikipedia.org/wiki/Apache_Solr [Google]
6. [Basic Solr Concepts - SolrTutorial.com](#)
Basic Solr Concepts. In this document, we'll cover the basics of what you need to know about Solr in order to use it. Indexing. Solr is able to achieve fast search ...



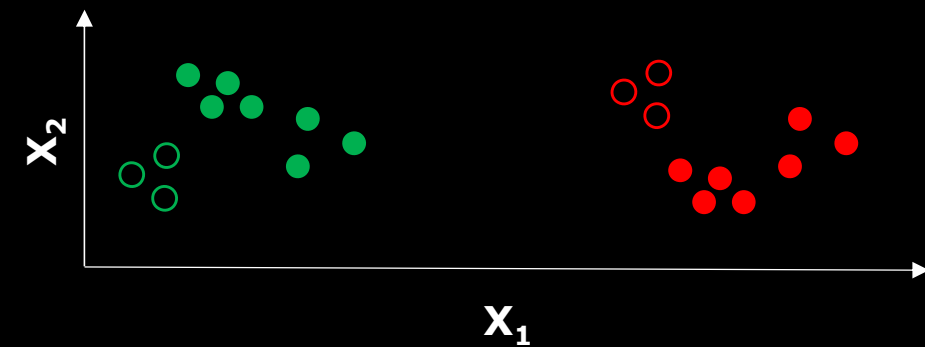
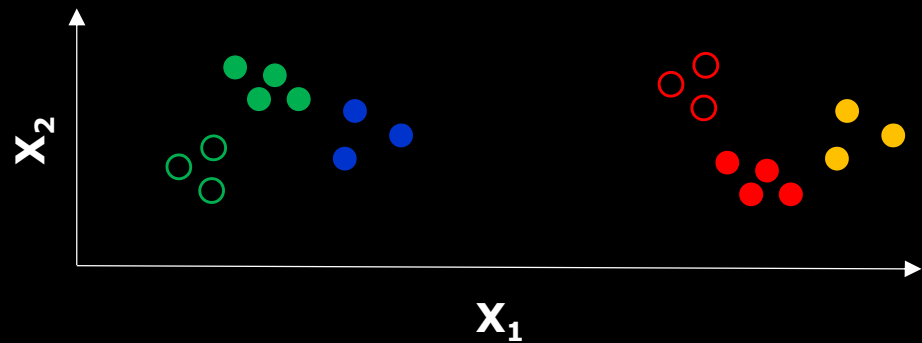
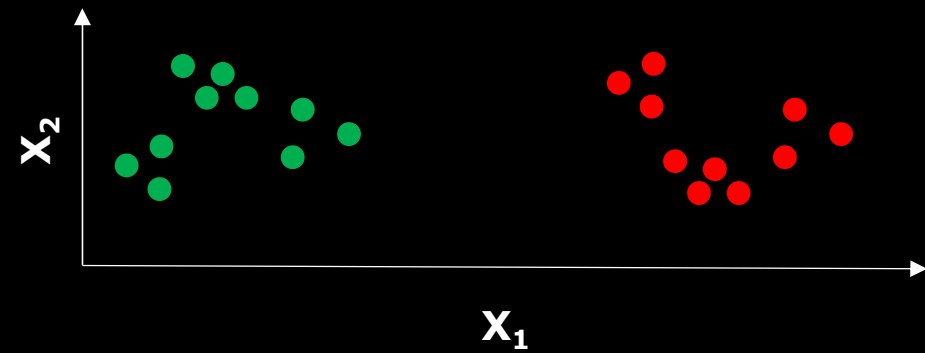
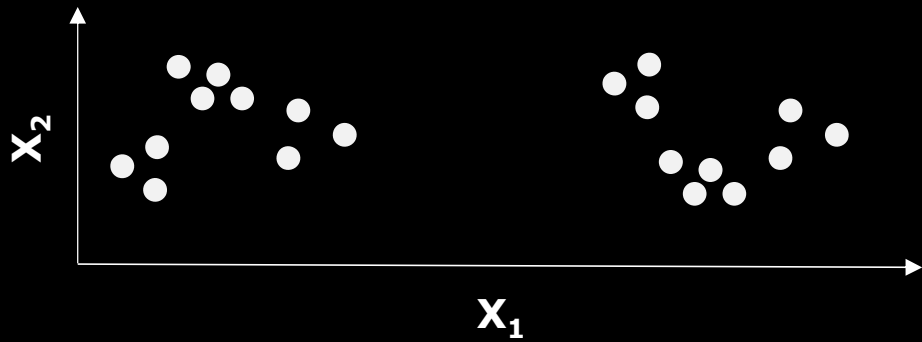
WHICH APPLICATIONS OF CLUSTER ANALYSIS?

SUMMARIZATION

Reduce the size of large data sets.



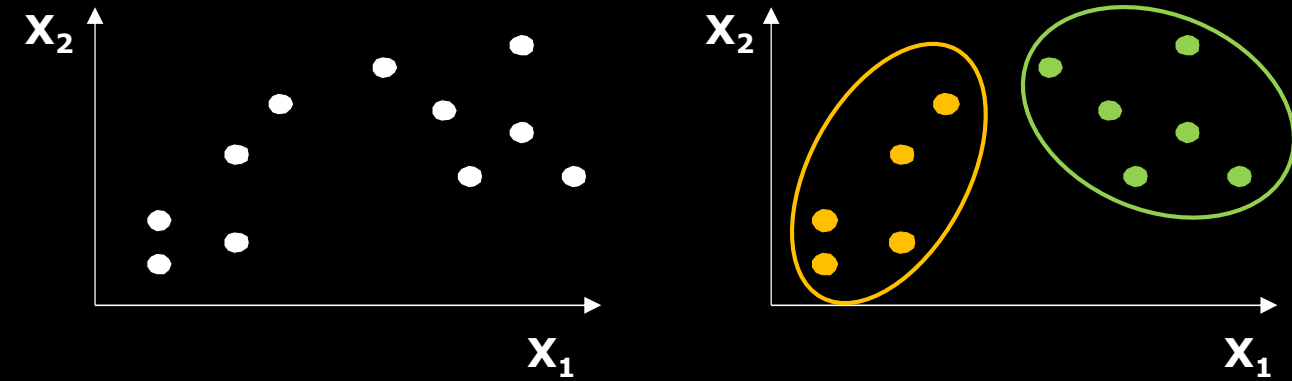
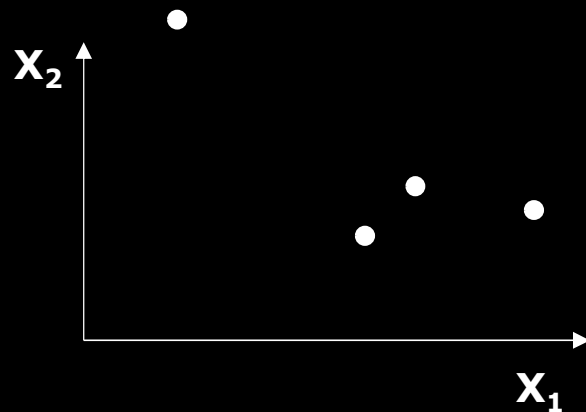
WHAT DO WE MEAN BY CLUSTER?



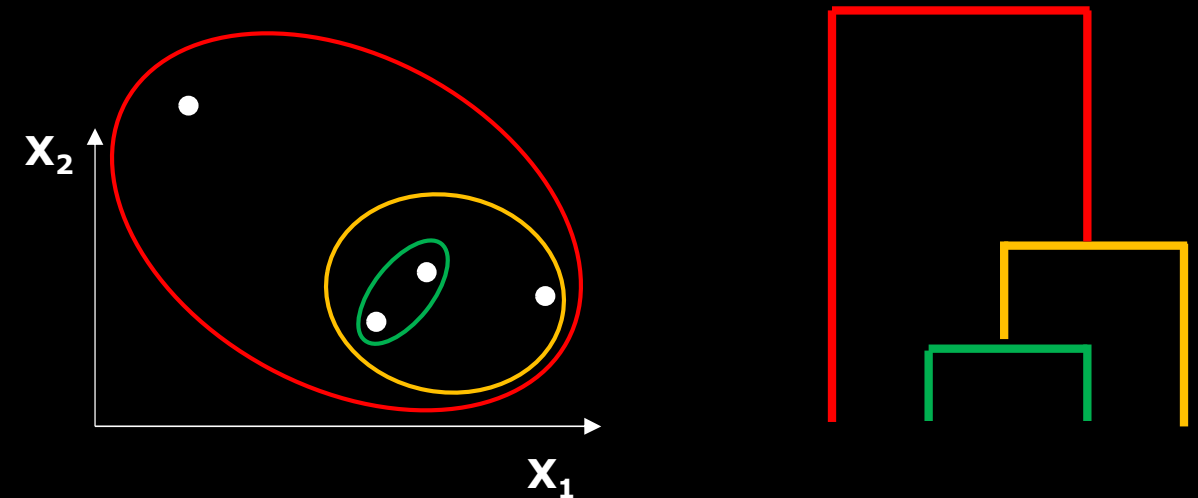
The definition of cluster is imprecise and the best definition depends on the nature of data and the desired results.

TYPES OF CLUSTERING

- A **CLUSTERING** is a set of clusters
- Important distinction between **PARTITIONAL** and **HIERARCHICAL** sets of clusters
 - **PARTITIONAL CLUSTERING**: a division of data objects into non-overlapping subsets
 - **HIERARCHICAL CLUSTERING**: a set of nested clusters organized as a hierarchical tree



PARTITIONAL CLUSTERING



HIERARCHICAL CLUSTERING

OTHER DISTINCTIONS BETWEEN SETS OF CLUSTERS

- **EXCLUSIVE** versus **NON-EXCLUSIVE**

- In non-exclusive clusterings, points may belong to multiple clusters
- Can belong to multiple classes or could be '**BORDER**' points

- **FUZZY CLUSTERING** (ONE TYPE OF NON-EXCLUSIVE)

- In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
- Weights must sum to 1
- Probabilistic clustering has similar characteristics

- **PARTIAL** versus **COMPLETE**

- In some cases, we only want to cluster some of the data
- Data can contain **OUTLIERS** or **ANOMALOUS OBSERVATIONS**

TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function

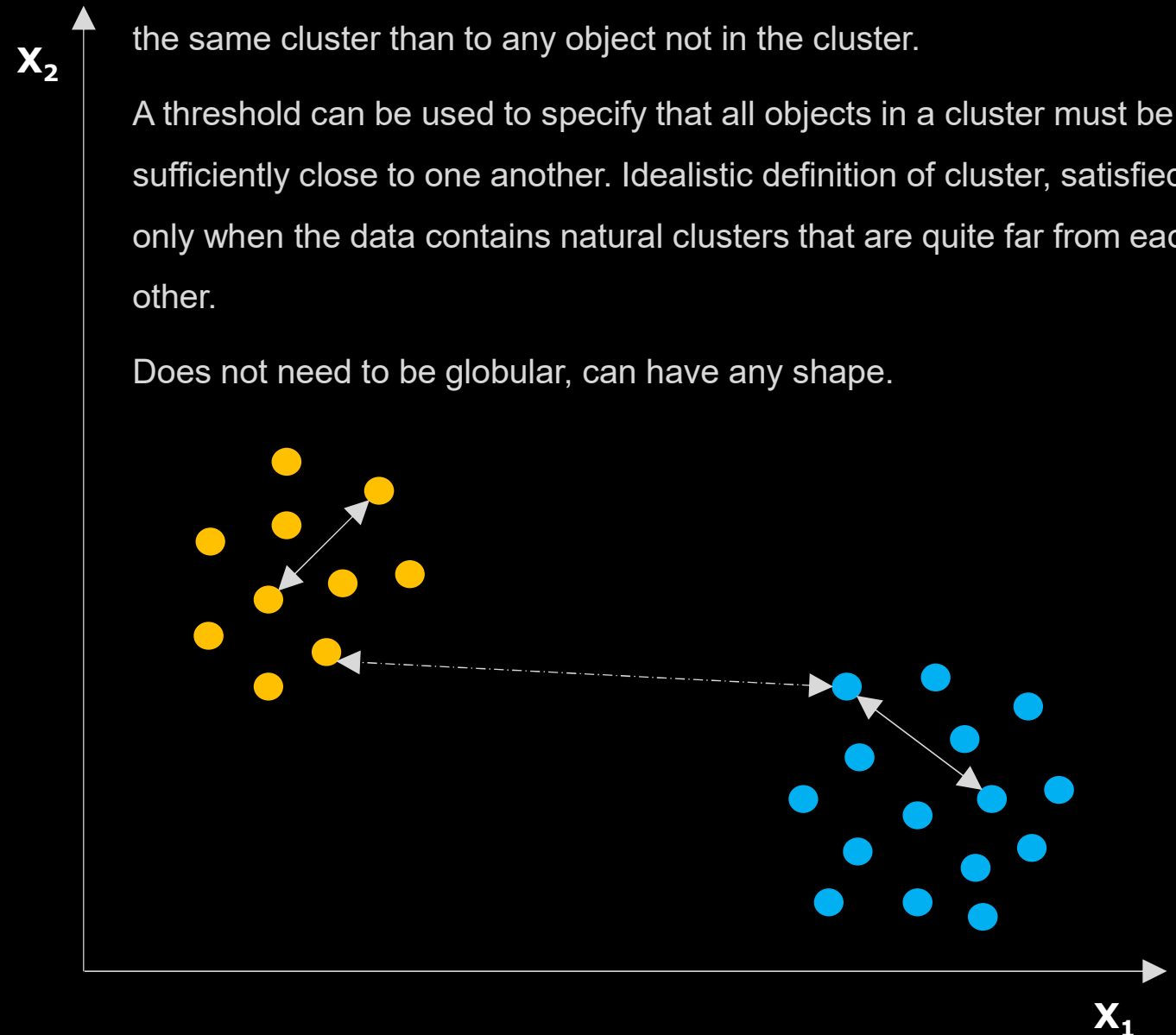
TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function

Each object in the cluster is closer (more similar) to every other object in the same cluster than to any object not in the cluster.

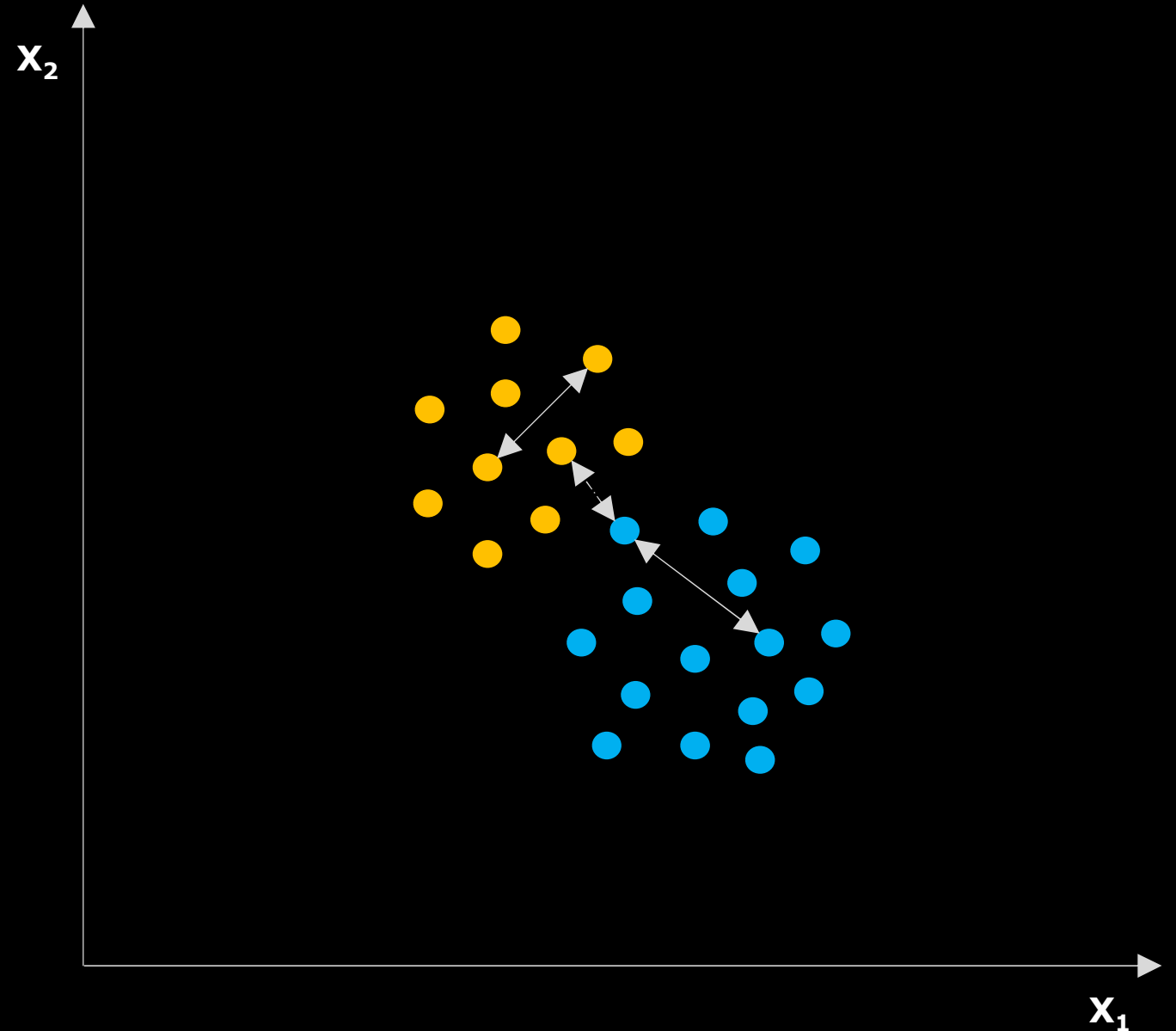
A threshold can be used to specify that all objects in a cluster must be sufficiently close to one another. Idealistic definition of cluster, satisfied only when the data contains natural clusters that are quite far from each other.

Does not need to be globular, can have any shape.



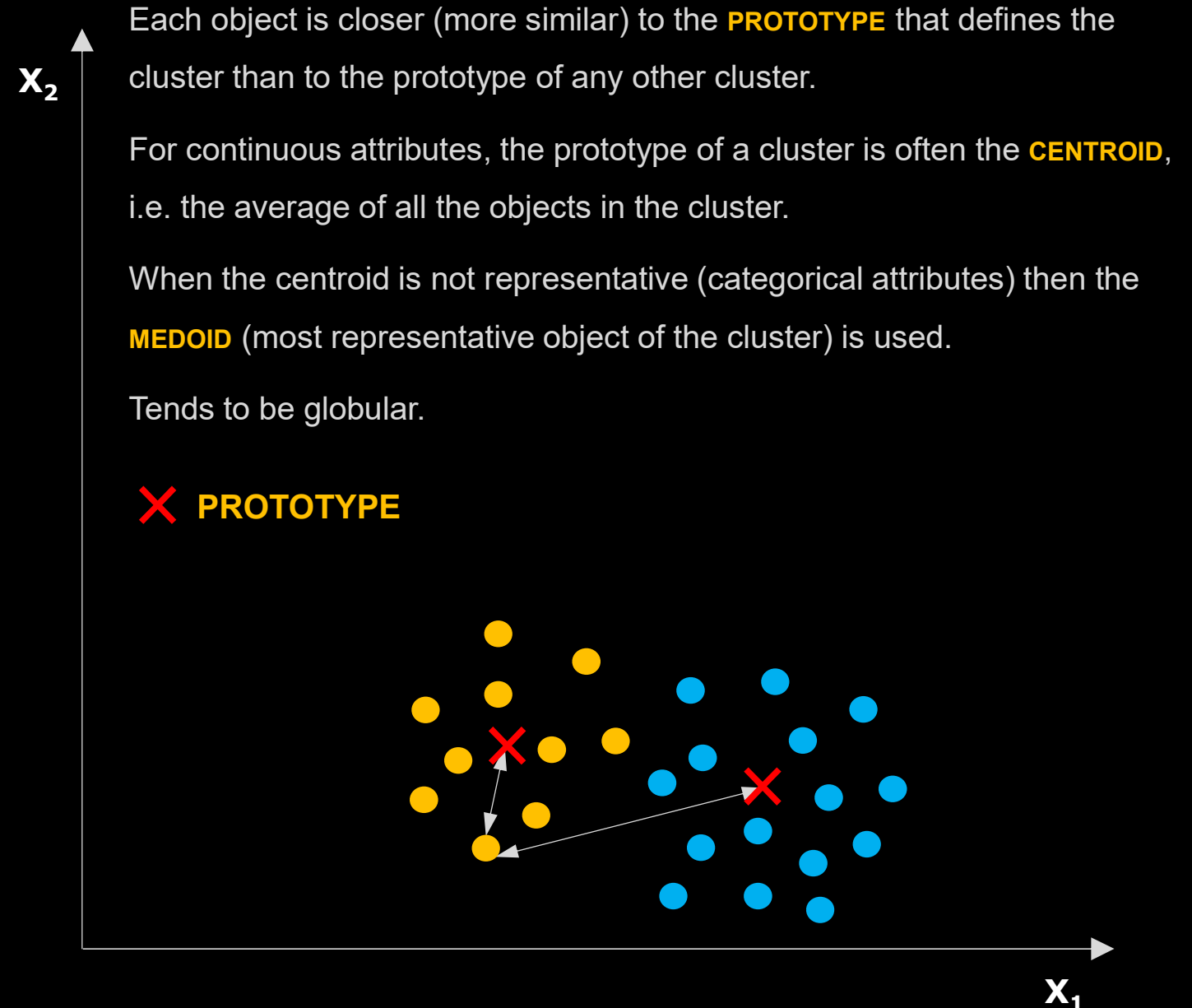
TYPES OF CLUSTERS

- Not well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function



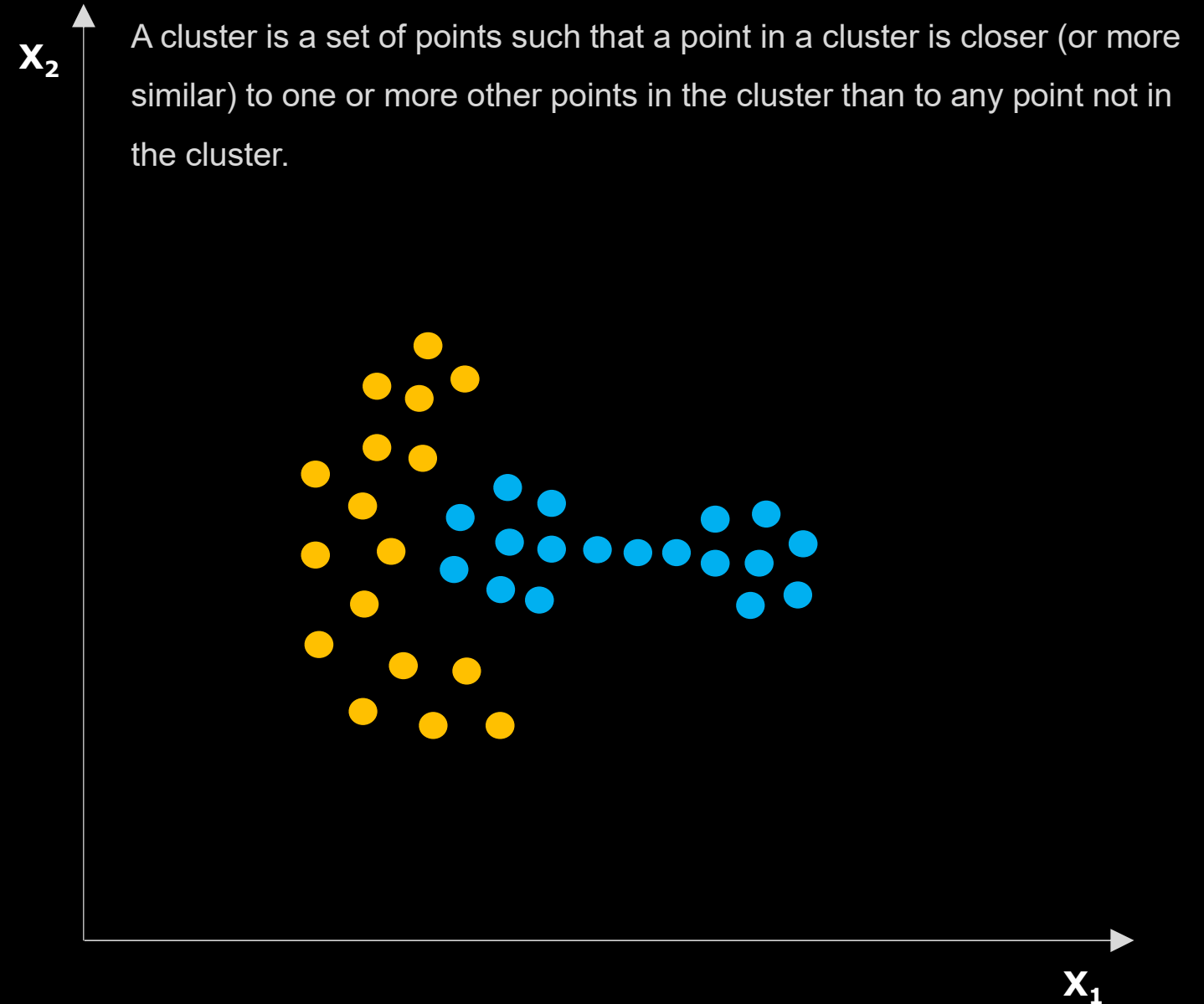
TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function



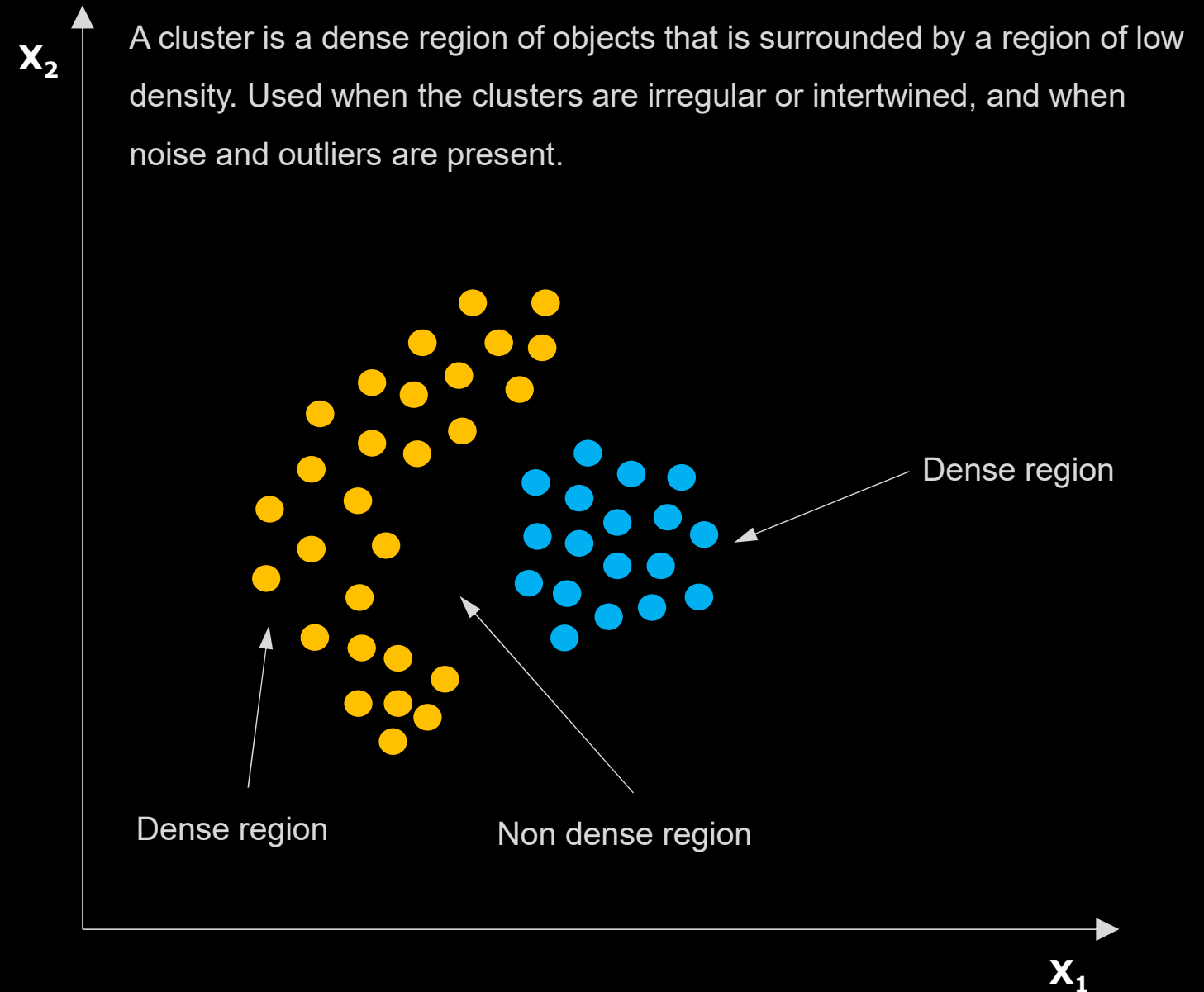
TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function



TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- **Density-based clusters**
- Described by an Objective Function



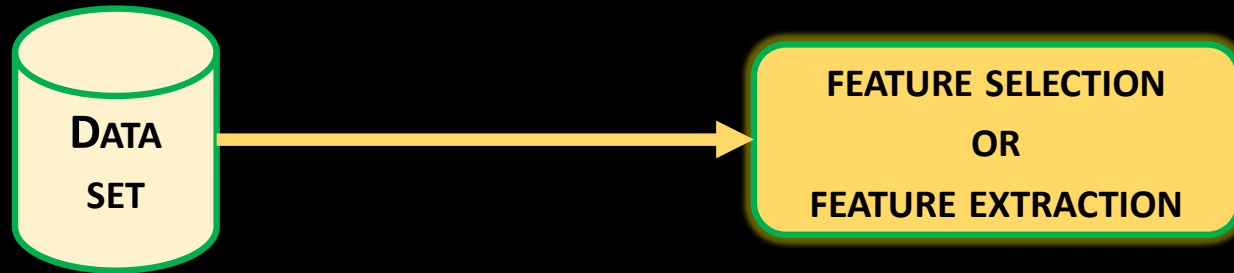
TYPES OF CLUSTERS

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- Described by an Objective Function
 - Finds clusters that minimize or maximize an objective function.
 - Enumerate all possible ways of dividing the points into clusters and evaluate the 'goodness' of each potential set of clusters by using the given objective function. (NP Hard).
 - Can have global or local objectives.
 - Hierarchical clustering algorithms typically have local objectives
 - Partitional algorithms typically have global objectives
 - A variation of the global objective function approach is to fit the data to a parameterized model.
 - Parameters for the model are determined from the data.
 - Mixture models assume that the data is a 'mixture' of a number of statistical distributions.

CHARACTERISTICS OF THE INPUT DATA ARE IMPORTANT

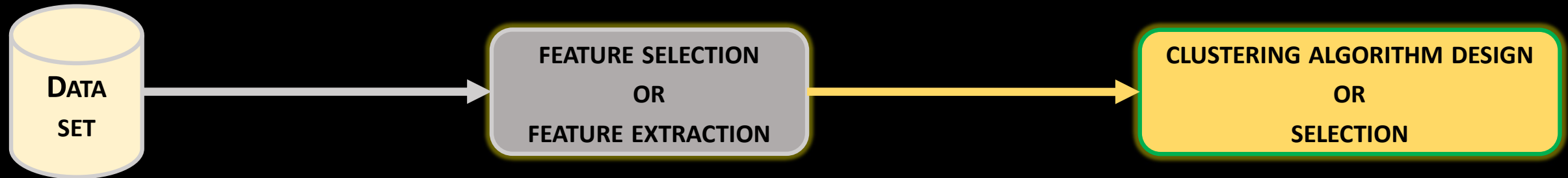
- Type of PROXIMITY or DENSITY MEASURE
 - Central to clustering
 - Depends on data and application
- Data characteristics that affect proximity and/or density are
 - DIMENSIONALITY
 - Sparseness
 - ATTRIBUTE TYPE
 - SPECIAL RELATIONSHIPS IN THE DATA
 - For example, autocorrelation
 - DISTRIBUTION OF THE DATA
- NOISE AND OUTLIERS
 - Often interfere with the operation of the clustering algorithm
- Clusters of DIFFERING SIZES, DENSITIES, and SHAPES

COMPONENTS OF CLUSTER ANALYSIS



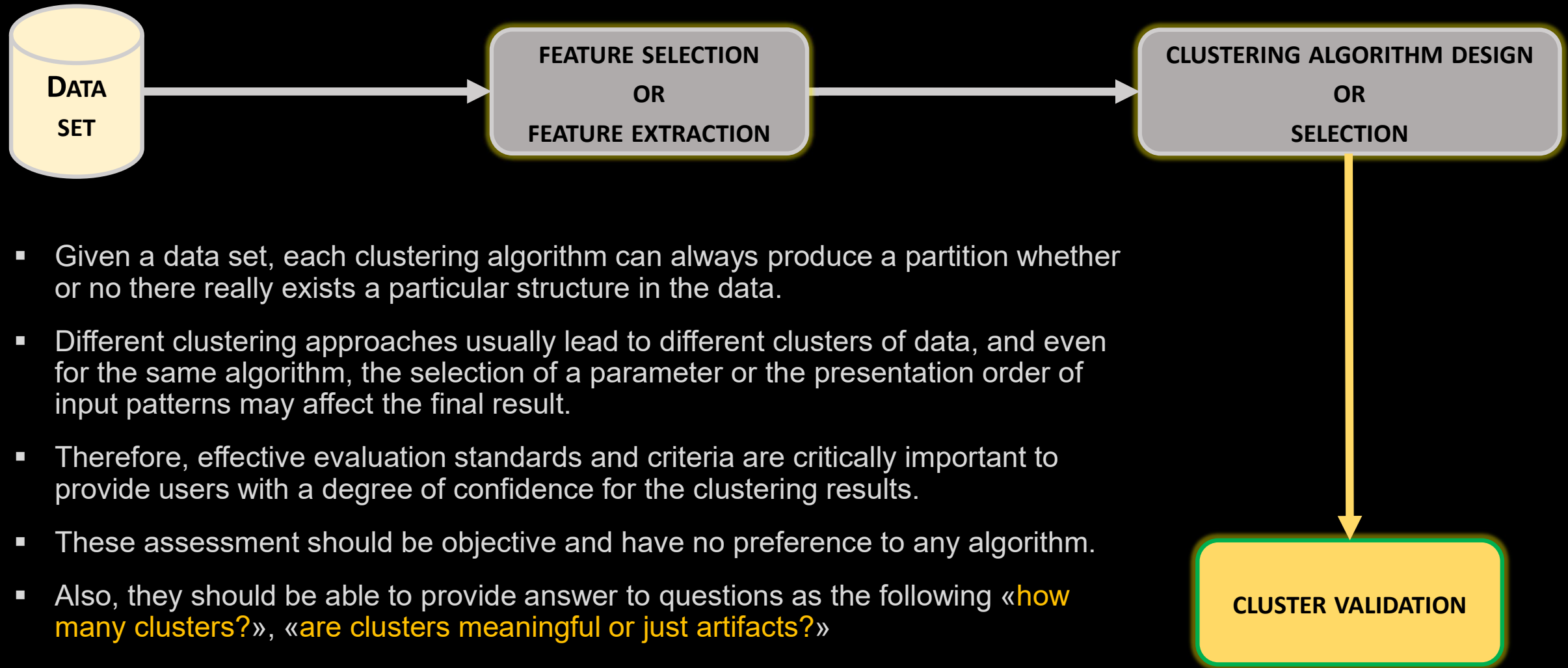
- **FEATURE SELECTION** assures the retention of the meaning of the original attributes.
- **FEATURE EXTRACTION** is capable of producing features that could be of better use in uncovering the data structure. However, it may generate features that are difficult to interpret.
- Ideal features should be of use in distinguishing patterns belonging to different clusters, immune to noise, and easy to obtain and interpret.

COMPONENTS OF CLUSTER ANALYSIS

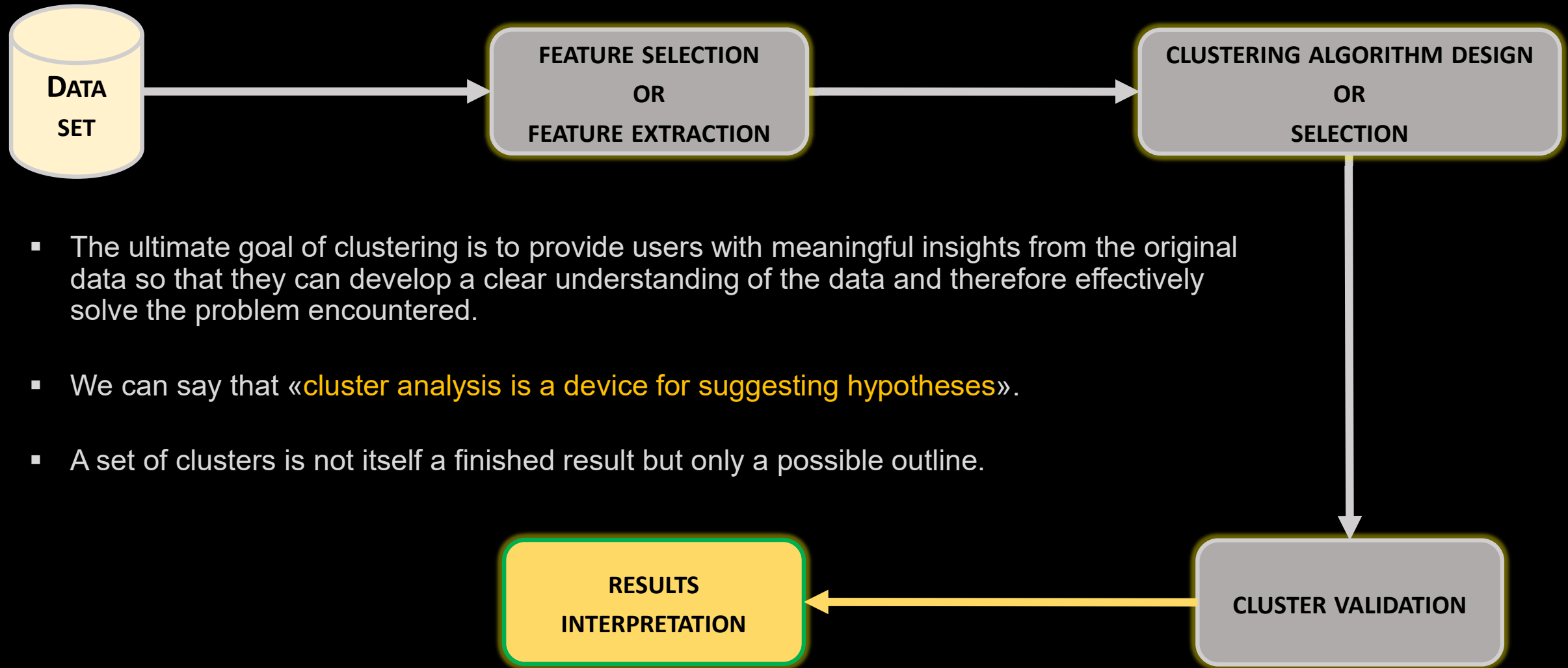


- Determining the **PROXIMITY** measure and constructing a **CRITERION FUNCTION**.
- Once a **PROXIMITY MEASURE** is determined, clustering can be formulated as an **OPTIMIZATION PROBLEM** with a specific **OBJECTIVE FUNCTION**.

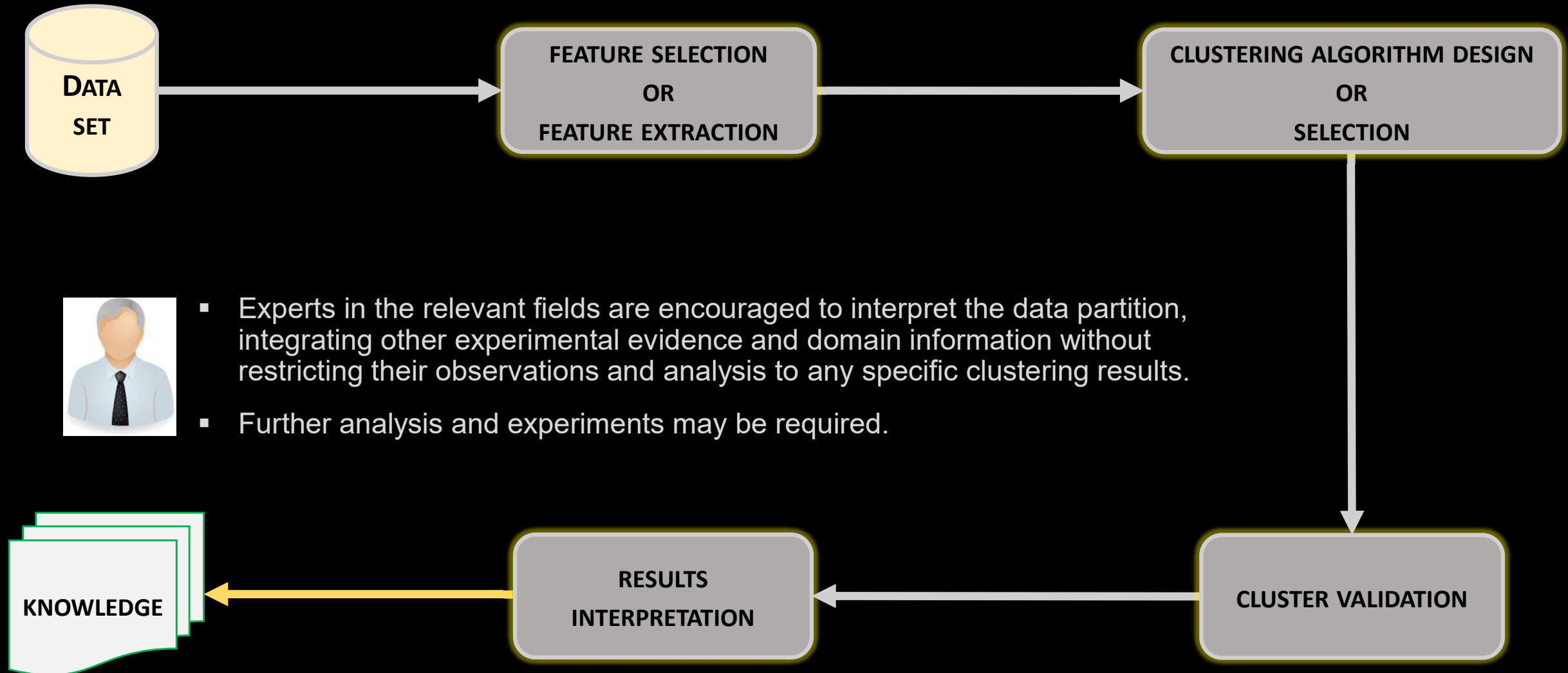
COMPONENTS OF CLUSTER ANALYSIS



COMPONENTS OF CLUSTER ANALYSIS



COMPONENTS OF CLUSTER ANALYSIS



RECAP

- **CLUSTER ANALYSIS**
- **UNDERSTANDING AND SUMMARIZING**
- **TYPES OF CLUSTERING**
- **TYPES OF CLUSTERS**
- **COMPONENTS OF CLUSTER ANALYSIS**