



## Introduction to Data Analytics

Roberta Turra, Cineca









### Data analytics

# The **process** of extracting useful insights from raw data using algorithms that **discover** hidden patterns







## Why is it challenging









## Data typologies

- structured data
  - data matrix
  - transactional data
- 🕈 graph
  - \* web and social networks
  - molecular structures
- 🕈 ordinal data
- 🕈 spatial data
- **T** time series
- sequences
  - genetic sequences
- unstructured data
  - textual documents
  - images
  - audio and videos (multimodal)









### Data as an infrastructure

Data has become the key infrastructure for 21st century knowledge economies. Data are not the "new oil", they are rather an infrastructure and capital good that can be used across society for a theoretically unlimited range of productive purposes, without being depleted.







## CRISP-DM reference model Cross Industry Standard Process for Data Mining









## Another way of describing the process (BDVA)







## Tasks and techniques





#### Unsupervised learning

training samples have no class information guess classes or clusters in the data we are given inputs but no outputs (unlabeled data) we learn the "latent" labels

#### **Supervised learning**

use training samples with known classes to classify new data we are given examples of inputs and associated outputs we learn the relationship between them







## Different approaches to the predictive task

#### Predictive

- \* classification (the learned attribute is categorical ,"nominal")
  - Naive Bayes
  - Decision Trees
  - **Neural Networks**
  - ₹ KNN
  - Rocchio R
  - Support Vectors Machine
  - ÷ ...
- \* regression (the learned attribute is numeric)

infer how to map input to output

Statisticians: model the process that gave rise to data ML: make an accurate prediction, given the data

