# DATA MINING

# 2016/2017

Will become *Big Data Computing* in 2017/2018

DATA

NEEDS

TOOLS

Data Scientist

# DATA

1. Technological progress
   - Growth of storage capacity
   - Growth of communication bandwidth
   - Growth of computing power
   - Reduction of ICT costs
2. Digital universe:
   - Integration of digital technologies in *all sorts* of human activities
   - Scientific research: biology (e.g., genomics); physics (e.g., particle physics, LHC); astronomy (e.g., wide-field survey telescopes); climate monitoring; etc.
   - Exponential growth of data (doubles every 2 years)
3. Data can be either unstructured (e.g., textual data) or structured (e.g., database records, networks)

# DATA: digital universe



Source: *The Digital Universe of Opportunities*, by ICD (2014)

1 ZettaByte (ZB) = $10^{21}$B

# DATA: the four V's



Source: IBM Big Data & Analytics Hub

# NEEDS

- Either predetermined or induced by avilability of data/tools
- Examples:
  - Retailing: product improvement; recommendation systems
  - Banking/Finance: fraud detection; risk prediction; financial forecast
  - Telecommunications: user profiling; quality-of-service improvement
  - Science: validation of methods/results
  - Medicine: diagnosis/therapy improvement
  - Social impact: analysis of social networks; smart cities
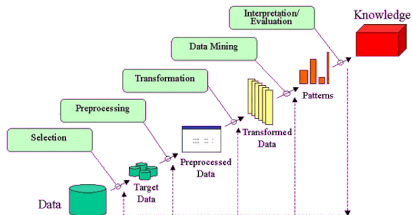
# TOOLS

- Statistics
- Machine Learning
- Algorithmics: tradeoffs between accuracy and space-time efficiency; novel computing/programming frameworks
- High-performance computing
- Natural Language Processing
- Visualization

# What is Data Mining?

Given a potentially large dataset, discover (through automatic procedues) patterns, models, properties, which are

- *Valid*
- *Useful*
- *Unexpected* (i.e., not explicitly contained in the data) and *previously unknown*
- *Understandable* (to humans)

Process of Knowledge Discovery in Databases

# Presentation of the Course

## What will we learn?

1. A sample of key primitives for data analysis
   - Rigorous setting
   - Algorithmic solutions with focus on large inputs
2. Novel computing/programming frameworks for (large) data analysis: theory and practice

## More specific contents

1. Association Analysis
2. Computational Frameworks: MapReduce/Spark, Streaming
3. Clustering
4. Graph Analytics
5. Mining primitives for data streams
6. Similarity Search

# Administative issues

## Evaluation

- Written test: 75%
- Project: 25%

## Online tools

- **Moodle:** registration (by march 12, password: DM1617), forum, results of written tests
- **Uniweb:** official exam lists, final grades
- **Course web site** (http://www.dei.unipd.it/~capri/DM): complete info, slides, links, solutions