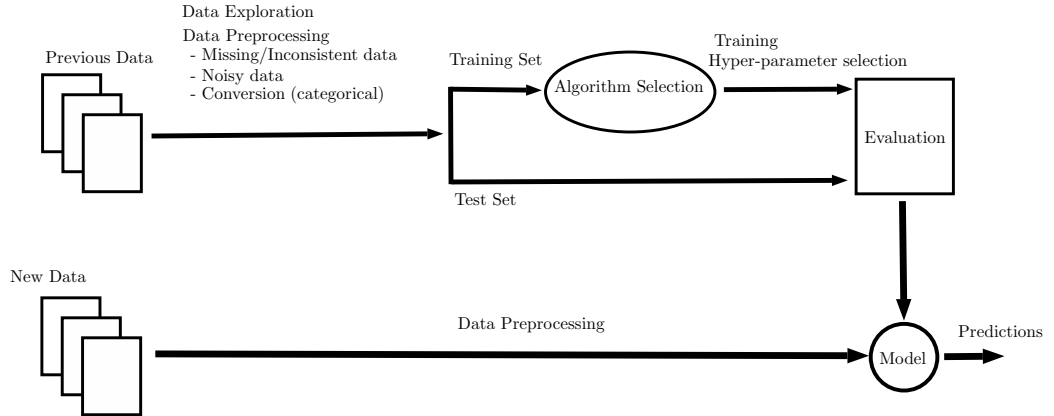POLIMI GRADUATE SCHOOL OF MANAGEMENT

# DATA PREPARATION

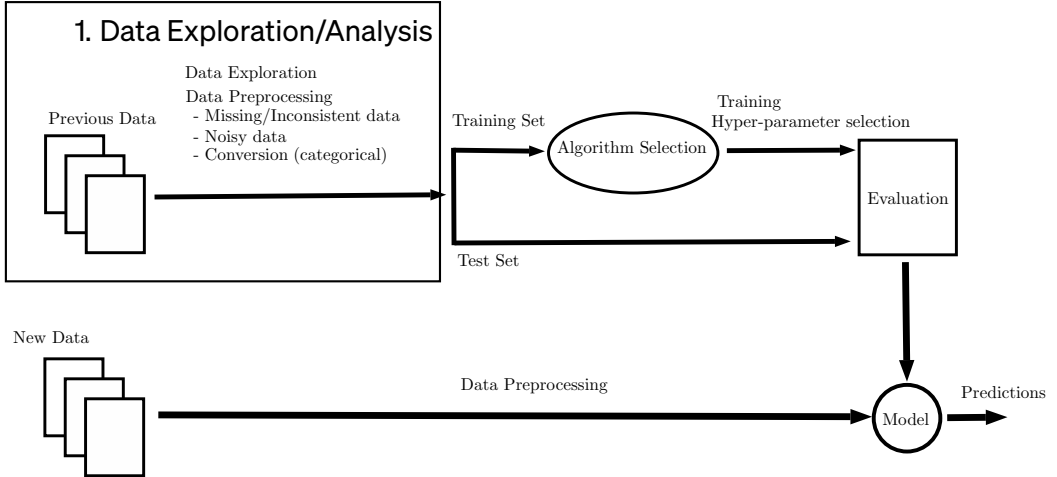PERCORSO EXECUTIVE DATA SCIENCE AND BUSINESS ANALYTICS

Mauricio Soto - mauricioabel.soto@polimi.it
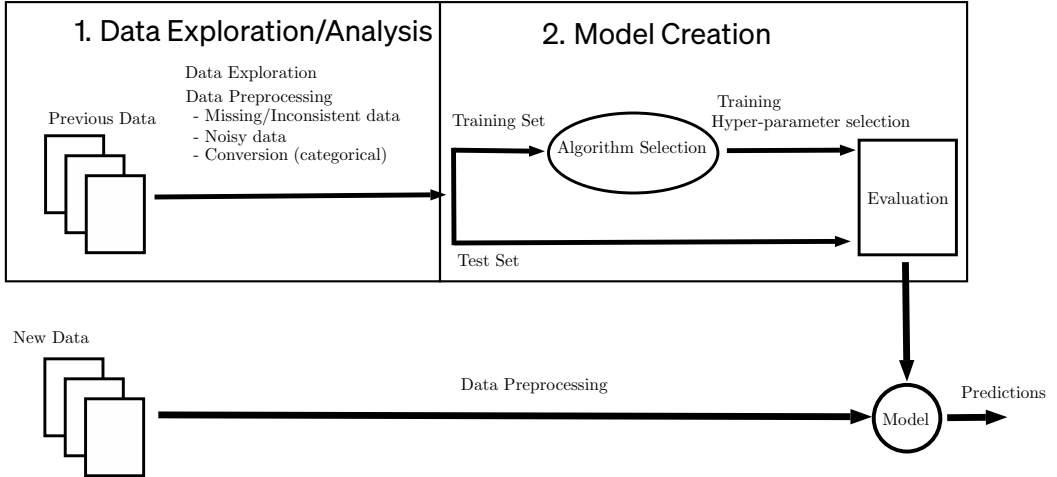
# WORKFLOW



Data Exploration
Data Preprocessing
- Missing/Inconsistent data
- Noisy data
- Conversion (categorical)

Previous Data

Training Set

Algorithm Selection

Training
Hyper-parameter selection

Evaluation

Test Set

New Data

Data Preprocessing

Model

Predictions

# WORKFLOW



POLIMI GRADUATE SCHOOL OF MANAGEMENT

# WORKFLOW



**1. Data Exploration/Analysis**

Data Exploration
Data Preprocessing
- Missing/Inconsistent data
- Noisy data
- Conversion (categorical)

Previous Data

**2. Model Creation**

Training Set

Algorithm Selection

Training
Hyper-parameter selection

Evaluation
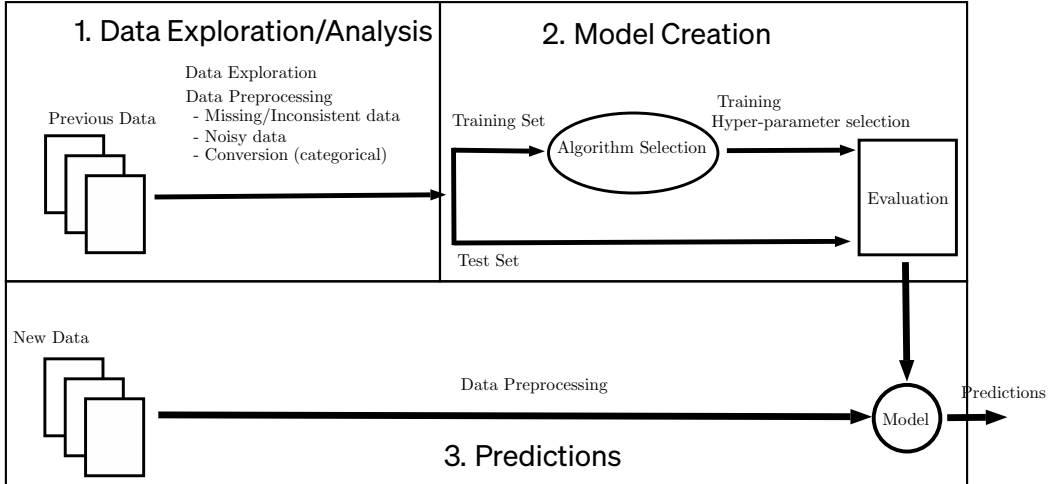
Test Set

New Data

Data Preprocessing

Model → Predictions

# WORKFLOW

# INCOMPLETE DATA

- ► Inspection
- ► Elimination
- ► Identification
- ► Replacement
  - mean value of numerical attributes
  - mean value of the target class
  - value estimated sing statistical inference

# WHAT IS AN OUTLIER AND HOW TO RECOGNIZE IT

https://pollev.com/mauriciosoto

# NOISY DATA

- ► Univariate
  - Normal-like distribution
  
  $$[\bar{\mu} - 2\bar{\sigma}, \bar{\mu} + 2\bar{\sigma}]$$
  
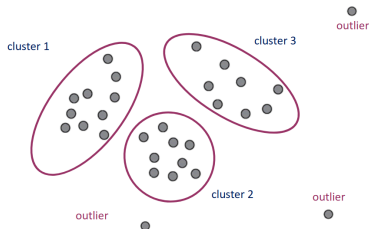    contains about 96% of the data
  - In the general case, **Tchebysheff Theorem** states taht for $\gamma > 1$
  
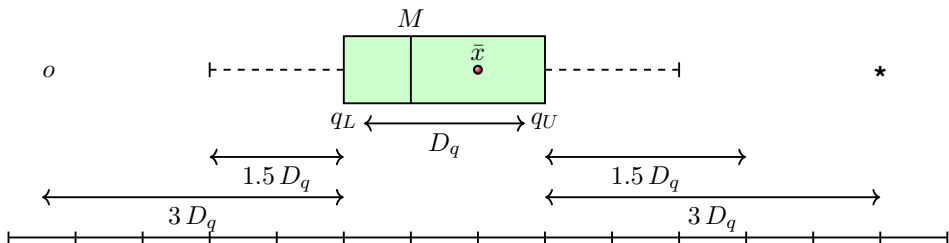  $$[\bar{\mu} - \gamma\bar{\sigma}, \bar{\mu} + \gamma\bar{\sigma}]$$
  
    contains $1 - 1/\gamma^2$ proportion of the observations
- ► Multi variate
  - Clustering techniques

# BOX-PLOT



- $D_q = q_U - q_L = q_{0.75} - q_{0.25}$
- internal lower edge= $q_L - 1.5\,D_q$
- external lower edge= $q_L - 3\,D_q$

# DATA TRANSFORMATION

- **Decimal Scaling**

$$x'_{ij} = \frac{x_{ij}}{10^k}$$

- **Min-Max** in the interval $[x'_{\min,j}, x'_{\max,j}]$

$$x'_{ij} = \frac{x_{ij} - x_{\min,j}}{x_{\max,j} - x_{\min,j}}(x'_{\max,j} - x'_{\min,j}) + x'_{\min,j}$$

- **z-index**

$$x'_{ij} = \frac{x_{ij} - \bar{\mu}_j}{\bar{\sigma}_j}$$

# DATA REDUCTION

- ▶ **Sampling**
  - Simple sampling
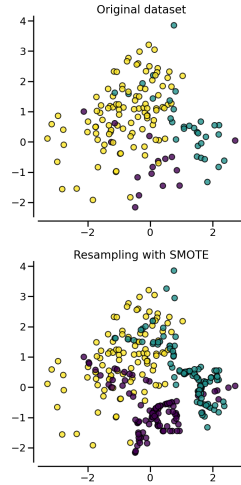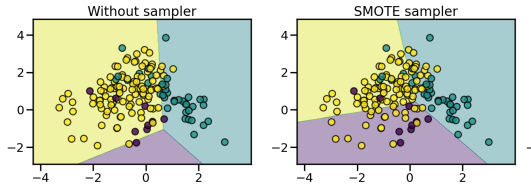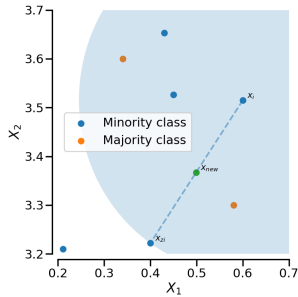  - Stratified sampling
- ▶ **Selection**
  - Filter methods
  - Wrapper methods
  - Embedded methods
- ▶ **Discretization, Aggregation**
- ▶ **Projection** (ex. PCA)

# DATA UNBALANCE - SMOTE

https://imbalanced-learn.org/

# THANK YOU

POLIMI GRADUATE SCHOOL OF MANAGEMENT