

# Data Analytics

Giuseppe Magazzù

2021 - 2022

# Contents

<b>1</b>	<b>Introduzione</b>	<b>1</b>
1.1	Definizioni . . . . .	2
1.2	Pre-Processing . . . . .	3
1.3	Data Cleaning . . . . .	3

# Chapter 1

## Introduzione

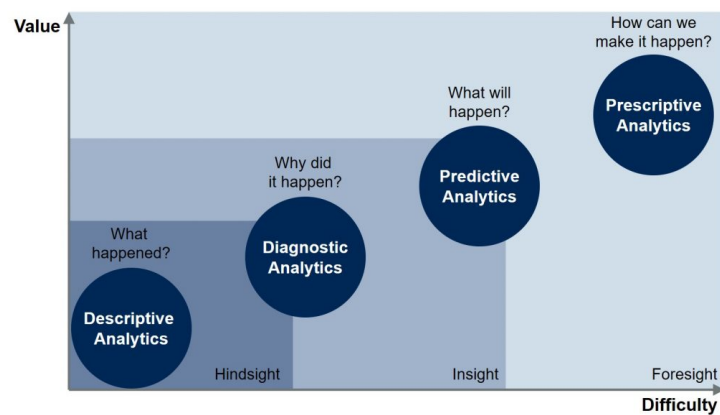


Figure 1.1: Diversi tipi di analisi per valore e difficoltà [3].

## 1.1 Definizioni

Un **istanza** (instance, item, record) é un esempio descritto da un insieme finito di attributi. Il numero di attributi può variare per alcune istanze.

Un **attributo** (attribute, field, variable) é una misura di un aspetto di un'istanza.

Tipi di attributi:

- **Quantità nominali:** i valori sono **simboli** distinti. Non hanno relazioni come ordinamento o distanza.  
(e.g. attributo: "outlook", valori: "sunny", "cloudy", and "rainy").
- **Quantità ordinali:** i valori hanno una relazione d'ordine, ma non di distanza.  
(e.g. attributo: "temperature", valori: "hot" > "mild" > "cold").
- **Quantità d'intervallo:**
- **Quantità di rapporto:**

Una **classe** (class, label) rappresenta un gruppo di istanze che condividono delle caratteristiche comuni.

## Propositionalization

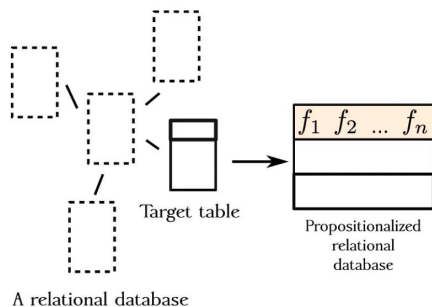


Figure 1.2: Processo di propositionalization [2]

## 1.2 Pre-Processing

I dati nel mondo reale sono **incompleti**, **rumorosi** e **inconsistenti**. Per ottenere dell'analisi di qualità é necessario effettuare prima delle operazioni sui dati.

- **Data Cleaning**: sostituire valori mancanti, smussare dati rumorosi, identificare o rimuovere outliers e risolvere inconsistenze.
- **Data Integration**: integrazione di diversi dati.
- **Data Transformation**: normalizzare o aggregare i dati
- **Data Reduction**: feature selection, feature extraction

## 1.3 Data Cleaning

### Dati Mancanti

Alcuni dati possono non essere stati calcolati o possono non essere disponibili per malfunzionamenti o per errori umani.

L'assenza di questi dati **complica l'analisi** poiché non tutti i metodi di analisi non gestiscono questo problema, inoltre comporta una **perdita di efficacia** nell'estrarre dei pattern.

Categorie di valori mancanti:

- **Missing Completely At Random (MCAR)**:
- **Missing At Random (MAR)**:
- **Not Missing At Random (NMAR)**:

Gestione dei valori mancanti:

- **Ignorare** le istanze o gli attributi con valori mancanti. Praticabile solo se ci sono pochi esempi mancanti poiché introdurrebbe un bias.
- **Convertire** i valori mancanti in un nuovo valore ("missing", "?", "NA").
- **Imputare** i valori mancanti basandosi sul resto del dataset.

### Metodi di Imputazione

- **Most Common (MC) Value**  
Assunzione: ogni attributo ha una distribuzione normale.
  - Valori **continui**: rimpiazza con la media dell'attributo nel dataset
  - Valori **discreti**: rimpiazza con il valore più frequente dell'attributo nel dataset

- **Concept Most Common (CMC) Value**

Assunzione: ogni attributo ha una distribuzione normale per tutte le istanze che appartengono alla stessa classe.

I valori mancanti vengono rimpiazzati con il valore medio/più frequente delle istanze della stessa **classe**.

- **K-Nearest Neighbors**

Le istanze vengono disposte in uno spazio metrico e i valori mancanti vengono imputati considerando le k istanze più vicine.

## Dati Rumorosi

Alcuni dati possono avere errori dovuti a **strumenti difettosi**, **errori umani** o **di calcolo**, errori durante la **trasmissione** dei dati o **limitazioni tecnologiche**.

Questi errori introducono del “rumore” all'interno dei dati che può essere rimosso usando tecniche di **data smoothing**. Queste tecniche riducono il rumore e rendono i pattern più identificabili, tuttavia si riduce la quantità di dati da analizzare e inoltre gli outliers possono alterare l'analisi.

## Binning

I dati vengono **ordinati** e **partizionati** in bin. Quindi ogni bin si può smussare con media, mediana dei valori all'interno o utilizzando gli estremi.

- **Equal-width** (distance) partitioning: viene diviso il range in N intervalli di uguale dimensione.
- **Equal-depth** (frequency) partitioning: viene diviso il range in N intervalli, ognuno dei quali contiene approssimativamente lo stesso numero di esempi.

## Dati Sbilanciati

Esistono molti problemi di classificazione in cui una classe ha una distribuzione fortemente sbilanciata, ovvero che il numero di osservazioni per una classe è molto inferiore a un'altra (e.g. fraud detection, disease diagnosis, natural disaster, etc.). Quindi risulta difficile ottenere buoni valori di accuracy su entrambe le classi.

Un possibile approccio è quello di bilanciare i dati del train set.

- **Oversampling**: aggiungere istanze alla **classe minoritaria** tramite campionamento con rimpiazzo (i.e. duplicare alcuni valori) fino a ottenere lo stesso numero di istanze per classe. Bilancia le classi, ma non fornisce nuove informazioni al modello.
- **Undersampling**: rimuovere randomicamente istanze dalla **classe maggioritaria** fino a ottenere lo stesso numero di istanze per classe.

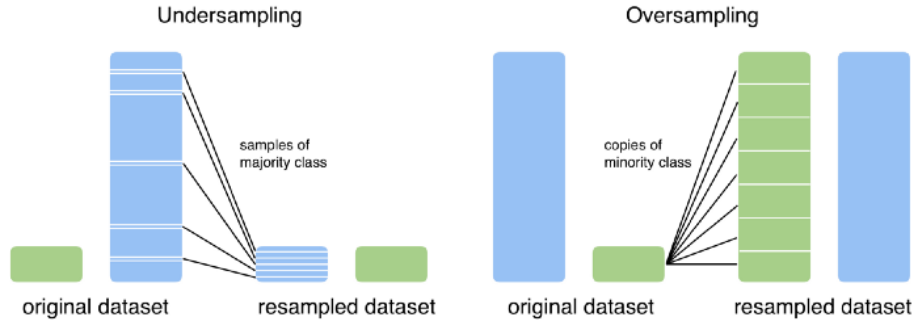


Figure 1.3: Rappresentazione del funzionamento del random resampling.

### Synthetic Minority Oversampling Technique (SMOTE)

SMOTE è una tecnica di oversampling che genera esempi sintetici della classe minoritaria a partire dai dati esistenti. Dato un esempio della classe minoritaria vengono selezionati i  $k$  esempi più vicini, viene scelto uno a caso tra questi e viene generato un numero esempio tra questi due.

### Tomek Links

Tomek Links è una tecnica di undersampling che rimuove gli esempi della classe maggioritaria che appartengono a un Tomek Link. Un **Tomek Link** é una coppia d'istanze  $(E_i, E_j)$  di classi diverse per cui non esiste nessun'altra istanza che sia più vicina a uno dei due.

La collezione di Tomek Links nel dataset definisce le frontiere delle classi.

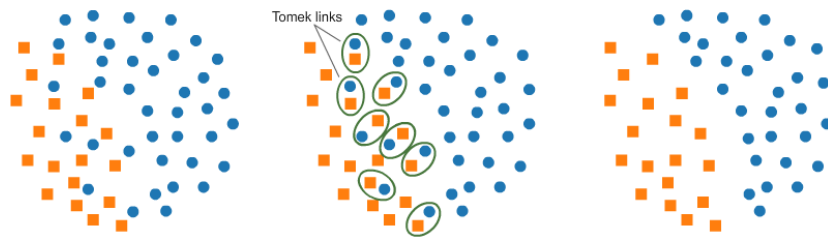


Figure 1.4: Esempio di undersampling tramite Tomek Links [1].

# Bibliography

- [1] Rahul Agarwal. *The 5 most useful Techniques to Handle Imbalanced datasets*. [Online; accessed 14/03/2022]. 2020. URL: [https://mlwhiz.com/images/imbal/1\\_hubf0730b098fff787d09b5f9aa956817e\\_24275\\_500x0\\_resize\\_box\\_2.png](https://mlwhiz.com/images/imbal/1_hubf0730b098fff787d09b5f9aa956817e_24275_500x0_resize_box_2.png).
- [2] Nada Lavrač, Blaž Škrlj, and Marko Robnik-Šikonja. "Propositionalization and embeddings: two sides of the same coin". In: *Machine Learning* 109.7 (2020), pp. 1465–1507.
- [3] Jason McNellis. *Gartner Four Analytic Types*. [Online; accessed 13/03/2022]. 2019. URL: <https://blogs.gartner.com/jason-mcnellis/files/2019/11/GartnerFourAnalyticTypesV5.jpg>.