

Title: Introduction

Course: Data Mining

Instructor: Claudio Sartori

Master: Data Science and Business Analytics

Master: Artificial Intelligence and Innovation

Master: Finance and Financial Technologies

Academic Year: 2023/2024

1

General information

2

2

Data in organisations

6

3

Data Mining

17

# Context

- Machine Learning
- Business Intelligence and Data Warehouse (DSBA only)
- Big Data For Industry
- Big Data Lab (DSBA only)
- Text Mining (DSBA only)
- Natural Language Processing and Applications (DTI and FFT only)

# What's in this course

- The CRISP-DM methodology for Data Mining processes (Sartori)
- Data Lake (Sartori)
- MLOps (Sartori)
- Data Collection (Sartori)
- Hands-on data mining and machine learning (Francia and Gallinucci)
- Hands-on Spark and OLAP (Francia and Gallinucci)

# Insight

Education is not the piling on of learning, information, data, facts, abilities or skills – that's training or instruction – but is rather making visible what is hidden as a seed

Thomas More<sup>1</sup>

---

<sup>1</sup> Cited by Charu C. Aggarwal in his book “Data Mining – the Textbook”

1	General information	2
2	Data in organisations	6
3	Data Mining	17

# Data, Data Mining and Machine Learning

- Data **exists** independently from Data Mining and Machine Learning
  - but you **need** Data Mining and Machine Learning techniques to derive interesting and **actionable** insights
- Data Mining and Machine Learning were created long before the dramatic increase of the amount of data available
  - the increase of the amount of data **strengthen DM and ML relevance and economic impact**

# Big Data

A new player with Data Mining and Machine Learning

- **Big Data** exists independently from Data Mining and Machine Learning
  - but you **need** Data Mining and Machine Learning techniques to **effectively analyse and use Big Data**
- Data Mining and Machine Learning were created long before the existence of Big Data
  - but using them on Big Data greatly **increase DM and ML relevance and economic impact**



Data → Information → Knowledge ⇒ better, data driven, decisions

*Data:* a collection of raw value elements

*Information:* the result of collecting and organising data

- ⇒ relationships between data items
- ⇒ context
- ⇒ meaning

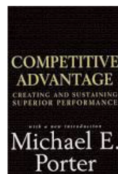
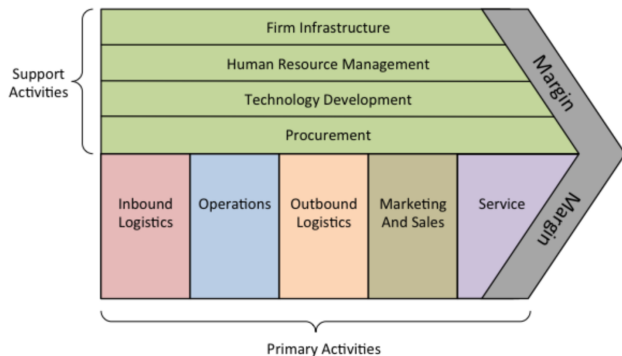
*Knowledge:* understanding information based on recognising patterns

# Increasing insights



# Where does *data* come from? 1/2

A *business process* is a set of activities that, once completed, will achieve an *organisational goal* (e.g. deliver your product to your customer)



[https://en.wikipedia.org/wiki/Value\\_chain](https://en.wikipedia.org/wiki/Value_chain)

# Where does *data* come from? 2/2

- When an event in the real world *changes the state* of the enterprise, one of the events below happens
  - a *transaction* is executed to reflect the corresponding change in the *database*
  - a signal is collected from the infrastructure and stored somewhere
- A *transaction* is a business event that generates or modifies data stored in an information system (database)
- A *signal* is the reading of a measure produced by a sensor
- Data may also be provided by *external subjects*

# Structured vs unstructured decisions

<i>Structured</i>		<i>Unstructured</i>	
<i>Description</i>	<i>Example</i>	<i>Description</i>	<i>Example</i>
Made under an established situation	Hiring a new employee	Made under an emergent situation	Fire breakout
Programmed	Start the monthly payment of salaries	Unplanned	Opportunity for financial investment
Fully understood	When a bank customer makes huge fund movements ask him the reason	Unclear or uncertain	Necessary to acquire information to understand which operation is to be performed
Routine task	Hiring new personnel in a given sector	Sudden One-shot situation	Dealing with a labor strike
Specified process	Manufacturing something	General processes	Managing security for IT equipment
Well defined methodology	Possible withdraw of funds from international accounts according to currency rates	Decisions relying on knowledge and/or expertise and on analysis of information	What new market segment could be targeted

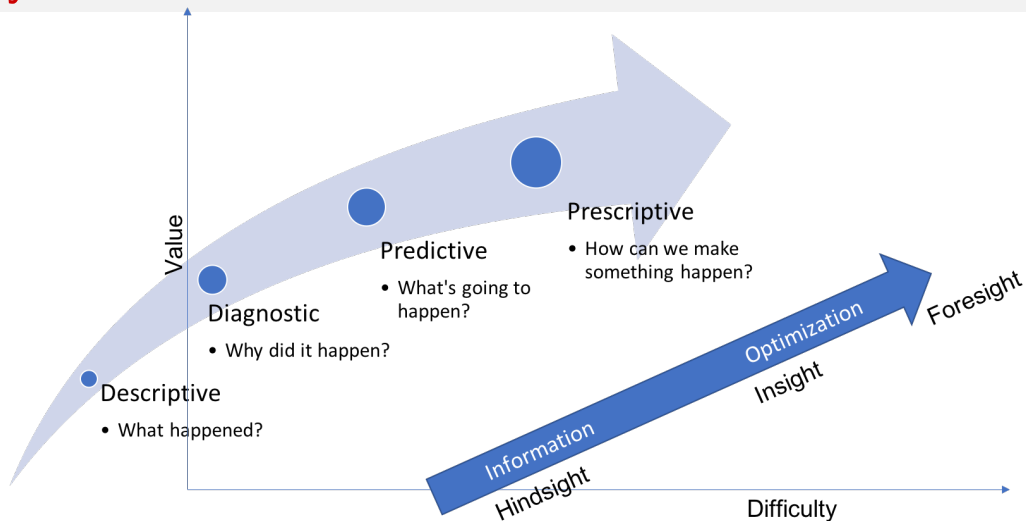
# Analytics vs Data Mining

**Analytics** – Structured decisions driven by data

**Data Mining** – Unstructured decisions driven by data

Sometimes they can provide insights in order to define a new structured decision

# Analytics



# Analytics

- descriptive
  - **aggregate** data with DB techniques, **understand data**, descriptive statistics and **unsupervised** machine learning
- diagnostic
  - descriptive + **domain knowledge**, **understand causes**
- predictive
  - calculate the most probable value of a variable in a future time, given the **history** of a set (sequence) of variables
- prescriptive
  - **suggest** actions to be taken to obtain the desired effect, **choose** among options and strategies, **optimize**



1

General information

2

2

Data in organisations

6

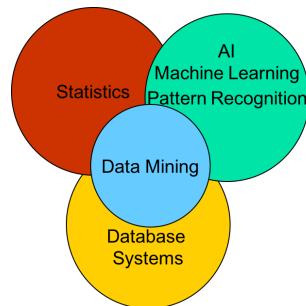
3

Data Mining

17

# Data Mining Origins

- The sizes of the circles to not reflect the relative importance/size of the topics
- Many textbooks referring either to *machine learning* or *data mining* have a significant overlap, sometimes the separation between the two topics is a little *fuzzy*



# Data Mining $\Leftrightarrow$ Machine Learning

In the following we will use the topic names as follows

- **Data mining** is the discovery process described in page 20
- **Machine learning for data mining** is the core of learning models and algorithms which allow to extract actionable patterns from data

Looking at the literature

- Machine learning includes also other concepts and methods which are not used for data mining
- Data mining books frequently include also *learning models* which are not traditionally covered in machine learning literature
  - Look [here](#) for a comprehensive list of data mining topics

# The Data Mining Process – attach labels to numbers

Internal data

Selection and pre-processing

Machine Learning

Knowledge

Interpretation and evaluation

Prepared data

Data Lake

Patterns and models

External data

Measure

Data Warehouse

Take action

Data Sources

Change

