

UNIVERSITÀ DEGLI STUDI DI BERGAMO

Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione



#### DATA SCIENCE AND AUTOMATION

Master degree in MECHATRONICS AND SMART TECHNOLOGY ENGINEERING

# Lecture 01: Introduction to data science

speaker Davide Previtali

PLACE University of Bergamo

#### Who am I



#### **Davide Previtali**

- **Currently**: Fixed-term Assistant Professor (RTD-A) @ Control systems and Automation Laboratory (University of Bergamo)
- **Studies**: Ph.D. in Engineering and Applied Sciences @ University of Bergamo. MSc in Computer Science Engineering @ University of Bergamo
- **Research topics**: lithium-ion batteries, control systems, black-box and preference-based optimization, machine learning
- <u>davide.previtali@unibg.it</u>

#### Control systems and Automation Laboratory (CAL)

- @ University of Bergamo
- Members: 7 professors, 5 Ph.D. students
- **Research topics**: control systems, advanced control, optimization, fault diagnosis, system identification, machine learning
- <u>http://cal.unibg.it/</u>
- <u>https://www.linkedin.com/school/cal-unibg/</u>





### Outline

- 1. Course introduction
- 2. Data science and the data-driven company
- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks
- 7. The data mining life cycle (CRISP-DM)



#### Outline

#### 1. Course introduction

2. Data science and the data-driven company

- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks

7. The data mining life cycle (CRISP-DM)



### **Course prerequisites**

It is **strongly suggested** to have a good knowledge of the following topics:

- Linear algebra
   Statistics
- Calculus 1 and Calculus 2

Please fill out the following questionnaire to assess your knowledge on the prerequisites:



https://forms.gle/92jNEKcmzWhMRoPq7



# **Course prerequisites**

#### How to «refresh» the prerequisites?

- Linear algebra
  - UniBg course, <u>course by Gilbert Strang @MIT</u> on YouTube, Addendum provided among the materials for this course
- Calculus 1 and Calculus 2

⇒ UniBg course

- Statistics
  - Brief review at the beginning of this course, UniBg course



A | Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione

### **Evaluation**

- Written exam, 1:30 hour
- True-false and multiple-choice questions
- Open questions

#### Up to 20 points



- Data science **project** with **discussion**
- You will receive more information on the project during the course

#### Up to 10 points



### **Educational objectives**

At the end of the course, you will be able to:

- Formulate a business problem as a data science problem
- Formulate and solve regression and classification problems
- Formulate and solve image analysis and object recognition problems
- Apply clustering and dimensionality reduction techniques
- Evaluate the goodness of a model estimated from data
- Visualize and present the results of a data science project



#### **Provided materials**

- Lessons' slides
- MATLAB code (MATLAB) (MATLAB) Focus: learn the methods, not the programming language

All the course materials are available at the following **Microsoft Team** 



[24/25] 39170-ENG - DATA SCIENCE AND AUTOMATION Generale Microsoft Teams



#### **Provided materials**

- Lessons' slides
- MATLAB code (MATLAB) (MATLAB

#### **Structure of the course**

The course will be divided in **theory** and **practice** sessions

- ✓ Theory lessons will be (mostly) of 2 hours on Mondays
- $\checkmark\,$  Practice lessons will be (mostly) of 3 hours on Fridays



#### Suggested books

 Foster Provost, Tom Fawcett.
 Data Science for Business: What you need to know about data mining and data-analytic thinking, O'Reilly Media, Inc. (2013)



T. Hastie, R. Tibshirani, J.
 Friedman. The elements of statistical learning: data mining, inference, and prediction, 2° Edition, Springer (2009)



 G. James, D. Witten, T. Hastie, R. Tibshirani. An Introduction to Statistical Learning, 2° Edition, Springer (2021)



 Andrew Gelman, Jennifer Hill.
 Data Analysis Using Regression and Multilevel/Hierarchical Models, Cambridge University Press (2006)



Data Analysis Using Regression and Multilevel/Hierarchical Models ANDREW GELMAN JENNIFE HILL



#### **Suggested books**

 Cole Nussbaumer Knaflic.
 Storytelling with data: a data visualization guide for business professionals, Wiley (2015)



 I. Goodfellow, Y. Bengio, A. Courville. Deep Learning, The MIT Press (2016)



 Christopher Bishop. Pattern recognition and machine learning, Springer (2006)





#### **Interactions and feedback**

- During the course I will give you **activities** to do and **tests** to answer
- They are optional (<u>although strongly encouraged</u>) as they help you assess your level of understanding before the exam
- In addition, they will give a bonus of (at most) +3 points to the final grade

We will use the **assignments** of **Microsoft Teams** 



### **Syllabus**

- 1. Introduction to data science
- 2. Exploratory data analysis
- 3. Recap of statistics
- 4. Maximum likelihood estimation
- 5. Linear regression
- 6. Logistic regression
- 7. Bias-variance trade-off
- 8. Overfitting and regularization
- 9. Validation and cross-validation



**10.Decision trees** 

**11. Neural networks** 

**12.**Convolutional neural networks

**13.Clustering methods** 

14.Principal component analysis

### Outline

#### 1. Course introduction

#### 2. Data science and the data-driven company

- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks
- 7. The data mining life cycle (CRISP-DM)



#### «Data is the new oil»





























 UNIVERSITÀ
 Dipartimento

 DEGLI STUDI
 di Ingegneria Gestionale,

 DI BERGAMO
 dell'Informazione e della Produzione





# Data is the new oil and data science is «sexy»

The data scientist role has been deemed the **sexiest job** of the 21st century [7]

- Virtually every aspect of business is now open to **data collection** (operations, manufacturing, supply-chain management, customer behaviour, marketing campaigns)
- Collected information need to be **analyzed properly** in order to get **actionable results**
- A huge amount of data requires **specific infrastructures** to be handled
- A huge amount of data requires **computational power** to be analyzed
- We can let computers perform decisions given **past data**
- Rising of **specific job** titles



### Job positions that involve data

	Data analyst	Data scientist	Data engineer	Machine learning engineer
• • •	Data retrieval (database queries) Spot trends and patterns in the data Visualize the data and produce reports to present information to third parties	<ul> <li>Use different machine learning techniques to derive insights from data to guide business decisions</li> <li>Make predictions on products, assets and consumer behavior based on past data</li> </ul>	<ul> <li>Design and maintain data management systems</li> <li>Data collection and management</li> <li>Make data accessible to the other members of the data science team</li> </ul>	<ul> <li>Design and implementation of machine learning methods</li> <li>Extend existing machine learning frameworks and libraries</li> <li></li> </ul>
		•	•	

And many more...

#### Often, career opportunities require a **good mix** of all the aforementioned skills



### What is data science?

**Data science** is a set of fundamental principles, processes and techniques that guide the extraction of knowledge from data with the goal of **improving decision-making** 

It is an interdisciplinary academic field that is based on:

- Mathematics
- Statistics
- Machine learning and artificial intelligence
- Specialized programming

**Data mining** is the extraction of knowledge from data, via technologies that incorporate data science principles



### The data-driven company

**Data-Driven Decision-making (DDD)** refers to the practice of basing decisions on the analysis of data, rather than purely on intuition [1, 2]

- Some decisions can be made automatically (finance, recommendations)
- **Data engineering and processing** support many dataoriented business tasks but do not necessarily involve extracting knowledge or data-driven decision making
- Data, and the capability to extract useful knowledge from data, should be regarded as **key strategic asset** 
  - ✓ Need to invest to acquire the right data (even lose money)
  - ✓ Understand data science even if you will not do it





#### Anti-hippo culture



Let data drive decisions, not the <u>Highest Paid Person's Opinion</u>.



di Ingegneria Gestionale,
 dell'Informazione e della Produzione

Hippos are among the most dangerous animals in Africa. Conference rooms too. —Jonathan Rosenberg

#### The road to becoming data-driven





#### Why become data-driven?

#### Data-driven companies are

**5% more productive** [2]



1\$ invested in analytics pays back 13 \$ [3]



A | Dipartimento
 DI | di Ingegneria Gestionale,
 O | dell'Informazione e della Produzione

### Why become data-driven?

Business value created by Artificial Intelligence by 2030 [4] \$13 Trillions

Retail	\$0,8T
Travels	\$480B
Logistics	\$475B
Automotive & assembly	\$405B
Materials	\$300B
Advanced electronics & semiconductors	\$291B
Healthcare systems & services	\$267B
High tech	\$267B
Telecom	\$174B
Oil & gas	\$173B
Agriculture	\$164B

It is **difficult** to find an industrial sector **that will not benefit** from artificial intelligence in

the near future



#### Outline

- 1. Course introduction
- 2. Data science and the data-driven company

#### 3. Data and its types

- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks
- 7. The data mining life cycle (CRISP-DM)



### What are data?

We refer to **data** as any piece of information that has been collected and stored in a computer

Examples:

...

- Sensor measurements
- Customer information
- Transaction history
- Social media posts





# Types of data: structured vs unstructured

#### **Structured data**

Data that are organized following a predefined scheme and stored in tabular formats (excel sheets, SQL databases...)

<b>House area</b> [feet <sup>2</sup> ]	# bedrooms	Price [k\$]
523	1	115
645	1	150
708	2	210
	:	:

#### **Unstructured data**

Data that can have an internal structure but do not follow a predefined data model or scheme Audio files

**Text files** 

Video files

Image files





## Types of data: quantitative vs qualitative



can be ordered. Other examples: low/high income, age ranges...

Runner name	Sex	Placement	<b>Time</b> [seconds]
Orlando Dillon	М	First	14.75
Izabella Kent	F	Second	15.01
Sophia Sanders	F	Third	15.33
:	:	:	:
<u></u>			

#### Qualitative (or categorical) data

assume non-numerical values, typically belonging to pre-defined categories

Quantitative (or continuous) data

assume numerical values



**Nominal qualitative data** 

cannot be ordered
### Data are dirty

Common data problems:

• Missing values

.

...

- Unlikely values (outliers)
- Inconsistent formats

House area [feet <sup>2</sup> ]	# bedrooms	Completion date	Price [k\$]
523	1	23/06/1998	115
645	1	01/07/2000	0.001
708	unknown	19/01/1980	210
1034	3	31-Jan-2001	unknown
unknown	4	17/12/2005	355
2545	unknown	14/02/1999	440
÷	:	:	:

Typically, data must be cleaned before usage (data cleaning)



### Outline

- 1. Course introduction
- 2. Data science and the data-driven company
- 3. Data and its types

### 4. What we are going to do with data (supervised and unsupervised learning)

- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks
- 7. The data mining life cycle (CRISP-DM)



### What are we going to do with data?

In this course, we will use data for:

**LECTURE 02** 

• Descriptive analysis and visualization



• Supervised learning (in particular, regression and classification)



• **Unsupervised learning** (in particular, clustering and dimensionality reduction)



**LECTURE 13,14** 

### Supervised vs unsupervised learning

Many data science tasks can be tackled either by supervised or unsupervised learning methods

 Supervised learning: predict the values of one or more dependent variables (output(s)) based on the values of one or more independent variables (input(s))



Typically, we will focus on supervised learning problems with **only one** output

Unsupervised learning: there are <u>no</u> outputs! The goal may be to discover groups of similar entities within the data or to project the data from a high-dimensional space (#inputs > 3) down to two or three dimensions for the purpose of visualization



- **Regression\***: predict the values assumed by the **continuous output(s)** from the **input(s)** 
  - **Example**: > Predict the **prices** of houses based on their **area**

- House area **# bedrooms** Price [k\$] [feet<sup>2</sup>] 523 115 1 645 150 708 210 2  $\varphi \in \mathbb{R}$  $\in \mathbb{R}$  $\boldsymbol{\varphi} \in \mathbb{R}^{2 \times 1}$ \*: covered in this course : supervised : unsupervised
- Predict the prices of houses based on their area and number of bedrooms

- Classification\*: predict the values assumed by the <u>categorical</u> output(s) from the input(s)
  - **Example**: > Develop an application that recognizes cats in **images**



Classification\*: predict the values assumed by the <u>categorical</u> output(s) from the input(s)

**Example**: > Distinguish cats from dogs based on their **height** and **weight** 



43 /91

- **Causal modeling**: identify which **inputs** (**causes**) actually influence the **outputs** (**effects**) and, possibly, to what extent
  - **Example**: > Did a particular marketing campaign influence the consumers to purchase our product?

Causal modeling typically involves substantial investments in data, such as randomized controlled experiments (**A/B tests**) and sophisticated methods for drawing causal observation data (**"counterfactual" analysis**)

What would be the difference in sales if we used an advertisement instead of another?

**Technical note**: regression and classification are based on correlation, causal modeling is based on causality

: supervised

Causal modeling: identify which inputs (causes) actually influence the outputs (effects)

and, possibly, to what extent



Ice Cream Sales vs. Shark Attacks

#### **Correlation does not imply causation!**

If we take a look at the data representing monthly ice cream sales and monthly shark attacks around the United States each year, we can see that the two variables are highly correlated

 Does this mean that consuming ice cream causes shark attacks? No! The more likely explanation is that more people consume ice cream and get in the ocean when it's warmer outside, explaining the high correlation

\*: covered in this course

: supervised

: unsupervised 45 /91

- **Clustering\***: organize the data into different groups based on their similarity
  - Example: ➤ Understand which types of customers are similar to each other by grouping individuals according to several characteristics → personalized marketing campaigns



- **Co-occurrence grouping**: find associations between different entities (characterized by a set of **features**) based on transactions involving them
  - **Example**: > What items are commonly purchased together? (market basket analysis)



Clustering looks at the similarity between entities based on their features, co-occurrence grouping considers the similarity of entities based on their appearing together in transactions (e.g., "a keyboard is not similar to a mouse, although they are typically bought together")

\*: covered in this course

: supervised

- **Profiling**: find the typical behavior of an individual, group or population
  - **Example**: > What is the typical credit card usage of a customer segment?
    - Profile the typical wait time of customers who call into a call center



48 /91

- **Link prediction**: predict connections between entities in a network, usually by suggesting that a link should exist, and possibly also estimating the strength of the link
  - **Example**: > Friend recommendations in social networks



\*: covered in this course

: supervised

 Dimensionality reduction\*: take a large dataset (many inputs and, possibly, many outputs) and replace it with a smaller dataset, retaining as much information as possible



- **Similarity matching**: find similar entities based on data known about them
  - **Example**: > Recommendation systems



#### Inputs:

- Song titles
- Song genres
- Audio signals
- :
- User ratings
- •

# **Output**: none (in this example)

Clustering is used for exploratory data analysis ("can we partition the data into different groups of similar entities?"), similarity matching has the specific goal of finding similar entities

\*: covered in this course

: supervised

: unsupervised 51 /91

### Data science tasks vs algorithms

### Data science task

(the problem that we are trying to solve, what we are trying to do) Regression, classification, ...



### Algorithm (or method)

(how we solve it, a sequence of operations to follow) Neural networks, KNN, K-means

clustering, ...

- Different data science tasks can be solved by the same algorithms *K*-means clustering can be used both for clustering and similarity matching
- Different algorithms can solve the same data science task A regression problem can be solved by the linear regression method, neural networks and *K*NN

In this course, we will study methods for solving different data science tasks



### **Syllabus**

- 1. Introduction to data science
- 2. Exploratory data analysis
- 3. Recap of statistics
- 4. Maximum likelihood estimation
- 5. Linear regression (regression)
- 6. Logistic regression (classification)
- 7. Bias-variance trade-off
- 8. Overfitting and regularization
- 9. Validation and cross-validation

: supervised

10.Decision trees (regression and classification)
11. Neural networks (regression, classification, dimensionality reduction...)
12.Convolutional neural networks (regression, classification, ...)
13.Clustering methods (clustering)
14.Principal component analysis (dimensionality reduction)

: unsupervised

### Outline

- 1. Course introduction
- 2. Data science and the data-driven company
- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)

### 5. Static and dynamical models in supervised learning

6. From business problems to data science tasks

7. The data mining life cycle (CRISP-DM)



## **Models in supervised learning**

Most supervised learning methods rely on mathematical **models** that describe the relationship between the **inputs** and the **outputs** 



Supervised learning methods estimate  ${\mathcal M}$  from data



## **Models in supervised learning**

We view both S and  $\mathcal{M}$  as mathematical functions that map **inputs** (**features**) to **outputs** (**targets**)



The goal of supervised learning methods is to learn a function  $\hat{f}(\cdot)$  that approximates  $f(\cdot)$  well **on the whole domain** of  $\varphi$ 



## Models in supervised learning

We view both S and  $\mathcal{M}$  as mathematical functions that map **inputs** (**features**) to **outputs** (**targets**)



The goal of supervised learning methods is to learn a function  $\hat{f}(\cdot)$  that approximates  $f(\cdot)$  well **on the whole domain** of  $\varphi$ 



### **Dataset notation**

Before moving on, we introduce the following notation that we will use for any dataset





V(t) $\varphi(i)$ We can view each voltage/current measurement by itself (i.e. as an observation  $(\varphi(i), y(i))$  in its own right), we do not need to consider V(t) and I(t) as signals "The time t can be omitted"

y(i)

#### **Example:** Ohm's law

### Static systems (and models)

A system whose **outputs** can be determined directly from the inputs is said to be a static system ("memoryless" system)

 $\boldsymbol{\varphi}(i)$  $f(\cdot)$ Inputs Outputs

The output I(t) at time t only depends on the input V(t) at the same time instant



## Static systems (and models)

Static systems need **not** describe **only** physics phenomena

House area [feet <sup>2</sup> ]	# bedrooms	Price [k\$]
523	1	115
645	1	150
708	2	210
:	•	:

 $f(\cdot)$ : mapping from house area and # bedrooms to price



 $f(\cdot)$ : mapping from image to label



### Learning static systems

In the **regression** setting, the simplest model that can be used to describe static systems is the **linear model** 

$$y(i) = \theta_0 + \theta_1 \varphi_1(i) + \dots + \theta_{d-1} \varphi_{d-1}(i) + \epsilon(i) = \sum_{j=0}^{d-1} \theta_j \varphi_j(i) + \epsilon(i)$$
  

$$i - \text{th observation} = \frac{\varphi(i)^{\mathsf{T}} \theta}{1 \times d \ d \times 1 \ 1 \times 1} + \epsilon(i) \cdot \varphi_0 = 1$$
  

$$(\varphi(i) = [\varphi_0 \ \varphi_1(i) \ \dots \ \varphi_{d-1}(i)]^{\mathsf{T}} \in \mathbb{R}^{d \times 1}$$
  

$$\cdot \theta = [\theta_0 \ \theta_1 \ \dots \ \theta_{d-1}]^{\mathsf{T}} \in \mathbb{R}^{d \times 1}$$
  

$$\cdot y(i) \in \mathbb{R}$$

- The vector  $\boldsymbol{\theta}$  is called **parameters vector**  $\rightarrow$  to be found by minimizing a cost function
- The vector  $\varphi(i)$  is called **features vector** for the *i*-th observation  $\rightarrow$  attributes of entities
- The quantity  $\epsilon(i)$  is the **error** due to not perfect explanation of y(i) using  $\varphi(i)$



### Learning static systems

To "learn" means to estimate the values of the parameters in  $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_{d-1} \end{bmatrix}^{\mathsf{T}}$ 

**<u>Key idea</u>**: find the values of  $\theta$  that **minimize** a "cost" (or "loss"), i.e. an "error" or "something bad"  $\rightarrow$  it is good to minimize something bad

• This is achieved through **optimization** 

A typical cost in the regression setting is the following

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{y}(i) - \boldsymbol{\varphi}(i)^{\mathsf{T}} \boldsymbol{\theta})^2 = \frac{1}{N} \sum_{i=1}^{N} \epsilon(i)^2$$

With this cost, we are **minimizing the sum of the squared errors** between the observed

outputs (i.e. those reported in our dataset) and the outputs estimated by the linear model



### Learning static systems

### **Scalar (single) parameter** $\theta$

<u>Multiple parameters  $\theta$ </u>



This rationale is followed by the **linear regression method** 

$$\hat{y}(i) = \hat{f}(\boldsymbol{\varphi}(i)) = \boldsymbol{\varphi}(i)^{\top} \hat{\boldsymbol{\theta}}$$



### **Dynamical systems (and models)**

A system whose **outputs** (at a certain time instant) cannot be determined directly from the **inputs** (at the same time instant) is said to be a **dynamical system** 



Dynamical models are mathematical models that describe the future evolution of the variables involved as a **function of their past trend** 

Dynamical systems usually involve the **time**: the **outputs** y(t) at a certain time t **depend** on the outputs at previous times

This dependency on the past endows the model with a **"memory"** (i.e. the dynamics)



### Dynamical systems (and models)

This dependency on the past endows the model with a "memory" (i.e. the dynamics)





### Learning dynamical systems

- In the control systems community, the concept of learning dynamical systems is typically referred to as **system identification**
- In any case, system identification belongs to the supervised learning framework:
  - Instead of dealing with datasets of observations, each observation representing an entity (e.g., a house), we have at our disposal datasets of signals (e.g., the electric motor's voltage and current signals). In this context, an observation is the collection of measurements of signals of interest at a certain time
  - Models used to describe dynamical systems must embed the time dependency in some way (e.g., we can employ transfer function models, state-space models, ...)



### Learning dynamical systems

- System identification methods are **<u>not covered</u>** in this course but are based on the concepts that we are going to learn in the following lectures
- If you are interested in system identification, these UniBg courses are what you are looking for:
  - [38020] Identificazione dei modelli e analisi dei dati
  - > [38095] Adaptive learning, estimation and supervision of dynamical systems
- A good book on the topic is:



Michel Verhaegen, Vincent Verdult

Filtering and system identification: a least squares approach

Cambridge University Press (2007)



### Machine Learning (ML), Artificial Intelligence (AI), Data Science and System Identification



All in all, we need a model to **better understand the phenomena** that are of our interest. <u>Models are useful for:</u>

- Decision-making: suppose that we are testing a new vaccine. We have two groups of people. We give the vaccine to the first group (test group) and a placebo to the second one (control group). Then, we measure some variables from the patients. How can we determine if the vaccine was effective or not?
- **Communication:** a model allows to communicate to third parties the main insights and results of your analysis



All in all, we need a model to **better understand the phenomena** that are of our interest. <u>Models are useful for:</u>

• **Prediction:** forecast the values that the output variables will assume based on the values assumed by the inputs variables and on which we have no data about

House area [feet <sup>2</sup> ]	# bedrooms	Price [k\$]
523	1	115
645	1	150
708	2	210
:	:	:

low much does a 600 feet<sup>2</sup> house with 2 pedrooms cost?



All in all, we need a model to **better understand the phenomena** that are of our interest. **Models are useful for**:

• Inference: understand how changes in the inputs affect the outputs

House area [feet <sup>2</sup> ]	# bedrooms	Price [k\$]
523	1	115
645	1	150
708	2	210
	÷	÷

- Does increasing house area increase the house price (and by how much)?
- Is # bedrooms actually associated with the price of a house?

**Prediction vs inference**: prediction is not necessarily concerned with the structure of the model  $\hat{f}(\cdot)$  and its complexity ( $\hat{f}(\cdot)$  can be seen as a black-box) while inference uses the model to understand the relationship between each input and each output



All in all, we need a model to **better understand the phenomena** that are of our interest. <u>Models are useful for:</u>

• **Simulation:** we can simulate, with a computer, the response (outputs) of a model due to certain inputs. By looking at the model's response, we can get a better grasp of the modeled system





Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione
### Why do we need models?

All in all, we need a model to **better understand the phenomena** that are of our interest. <u>Models are useful for:</u>

• **Control:** often, in control engineering, we need a model of a system to design a controller that limits the deviation of the controlled variables y(t) from the reference variables s(t) ( $t \in \mathbb{R}$  represents the time)





### Why do we need models?

All in all, we need a model to **better understand the phenomena** that are of our interest. <u>Models are useful for:</u>

• **Fault diagnosis:** we can check the presence of faults by comparing signals that come from the real system with those simulated by the estimated model



### Outline

- 1. Course introduction
- 2. Data science and the data-driven company
- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning

### 6. From business problems to data science tasks

### 7. The data mining life cycle (CRISP-DM)



# Business problems as data science tasks

Each data-driven project is **unique**. First and foremost, **decompose** the business problem

into data science subtasks that can be solved by **existing methods** 





# Business problems as data science tasks

- Spam e-mail detection system Classification
- Credit approval Classification
- Fraud detection Profiling
- Recognize objects in images Classification
- Find the relationship between house prices and house sizes Regression
- Predict the stock market Regression

- Market segmentation Clustering
- Market basket analysis
  Co-occurrence
  grouping
- Language models (word2vec) Similarity matching
- Social network analysis
  Link
  prediction
- Low-order data representations Dimensionality reduction
- Movies recommendation Similarity matching
- A/B testing Causal modeling



Focus on data science and machine learning projects that are **valuable** and **feasible** 

Think about automating **tasks** rather than automating **jobs** 

What are the main **drivers** of the business values?

What are the main **pain points** in your business?





#### **MANUFACTURING LINE MANAGER**

#### Data science



• Optimize production yield

#### **Machine learning**



• Automatic visual inspection



Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

#### **RECRUITING**



• Optimize recruiting process



• Automatic resume screening



#### MARKETING



#### **Machine learning**



• A/B testing websites

• Recommendation system



A Dipartimento
 D di Ingegneria Gestionale,
 D dell'Informazione e della Produzione

### Outline

- 1. Course introduction
- 2. Data science and the data-driven company
- 3. Data and its types
- 4. What we are going to do with data (supervised and unsupervised learning)
- 5. Static and dynamical models in supervised learning
- 6. From business problems to data science tasks

### 7. The data mining life cycle (CRISP-DM)



### **CRISP-DM process**

**Cross Industry Standard Process for Data Mining** (CRISP-DM)

#### **Iteration is the rule** rather than the exception:

- Business understanding
- Data understanding
- Data preparation
- Modeling
- Evaluation
- Deployment





# **CRISP-DM: Business understanding**

Cast the business problem into one or more data science problems

- Regression
- Classification
- Causal modeling
- Clustering
- Co-occurrence grouping
- Profiling
- Link prediction
- Dimensionality reduction
- Similarity matching

#### Think carefully about the **use scenario**:

- What exactly do we want to do?
- **How** exactly would we do it?
- What parts of this use scenario constitute possible data mining models?





### **CRISP-DM: Data understanding**

Identify the available and needed data

**Costs/benefits** of acquiring each source of data

Are the data at our disposal **related to the business problem**?

Can we use a **proxy** for the data that we do not have?

As data understanding progresses, the **solution paths** may differ





# **CRISP-DM:** Data preparation

Clean and prepare the data for usage

Usually, data mining algorithms require **data in a specific format** which is different from the one that is readily available

• Convert string to numbers, infer missing data, import data from excel files, ...

# Data preprocessing/cleaning/labeling (most of data science project time is spent here) [5]

Pay attention to not use historical data that **will not be available** when decisions need to be made





# **CRISP-DM: Modeling**

Estimate a mathematical model to extract patterns from data

In most cases, standard algorithms can be directly applied on data

The aim is to find a model that performs well **on unseen data** 

The type of the model is chosen based on:

- What data science **task** we want to solve
- Performance measures
- Availability of **libraries** for deployment





### **CRISP-DM: Evaluation**

Assess the validity of the results

We could find patterns that exist only in the particular dataset that we have at our disposal **(overfitting)** 

Does the model **satisfy** the original business goals?

The devised solution and the model's decisions should be **comprehensible** by the stakeholders

Usually, evaluation is performed **before deploying**. In this case, build environments that **closely mimic** the real use scenario





# **CRISP-DM: Deployment**

Put the model (or the data mining steps) into production

Usually requires to **re-code** the model, to make it compatible with existing technologies



This step can require a notable **investment in time**. Usually, the data science team builds a prototype that is then passed on to the development team

For this reason, it is suggested to **include a member of the development team** in the early phases of the data science project

Deployment can involve not only the final model, but also **previous phases** (data collection, model building, evaluation)



# Workflow of a machine learning project





# Workflow of a data science project

Mix clay

<u>Optimize a</u> <u>manufacturing line</u>

- 1. Collect data
- 2. Analyze data
  - Iterate many times to get good insights
- 3. Deploy model
  - Deploy changes
  - Re analyze new data periodically

→	D->	- Fe		>	
Shape mug		Add glaze		Fire kiln	Final inspection
	<b>Clay batch</b> # 001 034			Supplier	<b>Mixing time</b> [minutes]
				Supplier 1	35
				Supplier 1	22
		109		Supplier 2	28
	Mug batch	; #	Humidity	Temperature in kiln [F]	<b>Time in kiln</b> [hours]
	001		0.002%	1410	22
	034		0.003%	1520	24
	109		0.002%	1420	22





UNIVERSITÀ DEGLI STUDI DI BERGAMO

| Dipartimento | di Ingegneria Gestionale, | dell'Informazione e della Produzion

### References

- 1. Provost, Foster, and Tom Fawcett. *"Data Science for Business: What you need to know about data mining and data-analytic thinking"*. O'Reilly Media, Inc., 2013. **Chapters 1-2**.
- 2. Brynjolfsson, E., Hitt, L. M., and Kim, H. H. *"Strength in numbers: How does data driven decision making affect firm performance?".* Tech. rep., available at SSRN: <u>http://ssrn.com/abstract=1819486</u>, 2011
- 3. Nucleus Research, 2014. <u>http://bit.ly/XQFDbv</u>.
- 4. <u>Notes from the AI frontier: Modeling the impact of AI on the world economy</u>, 2018.
- 5. Pyle, D. "Data Preparation for Data Mining". Morgan Kaufmann, 1999. Chapter 1.
- 6. G. James, D. Witten, T. Hastie, R. Tibshirani. *"An Introduction to Statistical Learning"*. 2° Edition, Springer, 2021. **Chapters 1-2**.
- 7. <u>Data scientist: The Sexiest Job the 21<sup>st</sup> Century</u>, 2012.
- 8. <u>Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020,</u> with forecasts from 2021 to 2025, 2022.
- 9. <u>Correlation does not imply causation: 5 real-world examples</u>, 2021.



Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione