

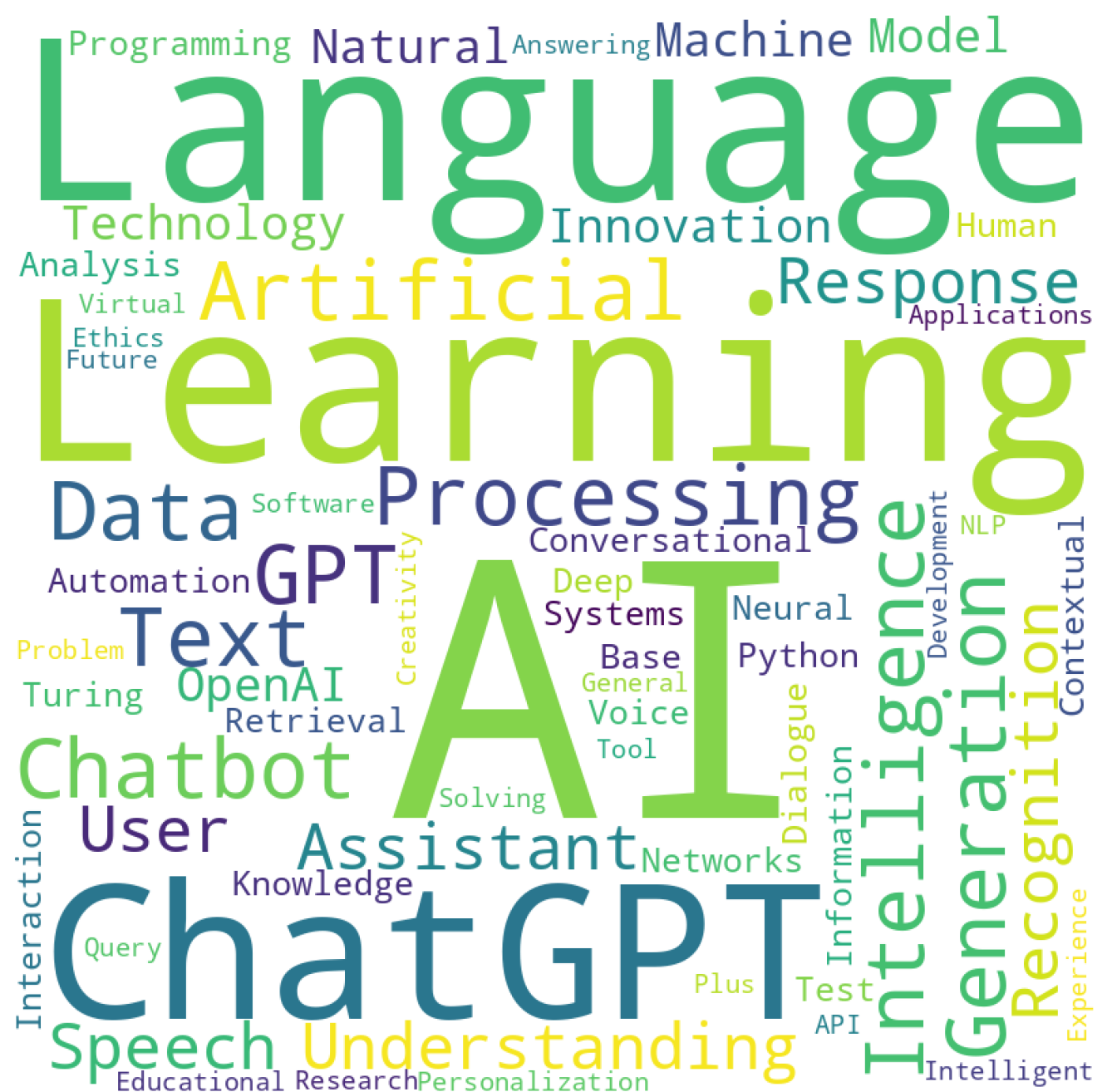
Introduction to Data Science and data Mining

STATISTICA NUMERICA - A.Y. 2024/2025

1 – Data science

A VERY SHORT TERMINOLOGY

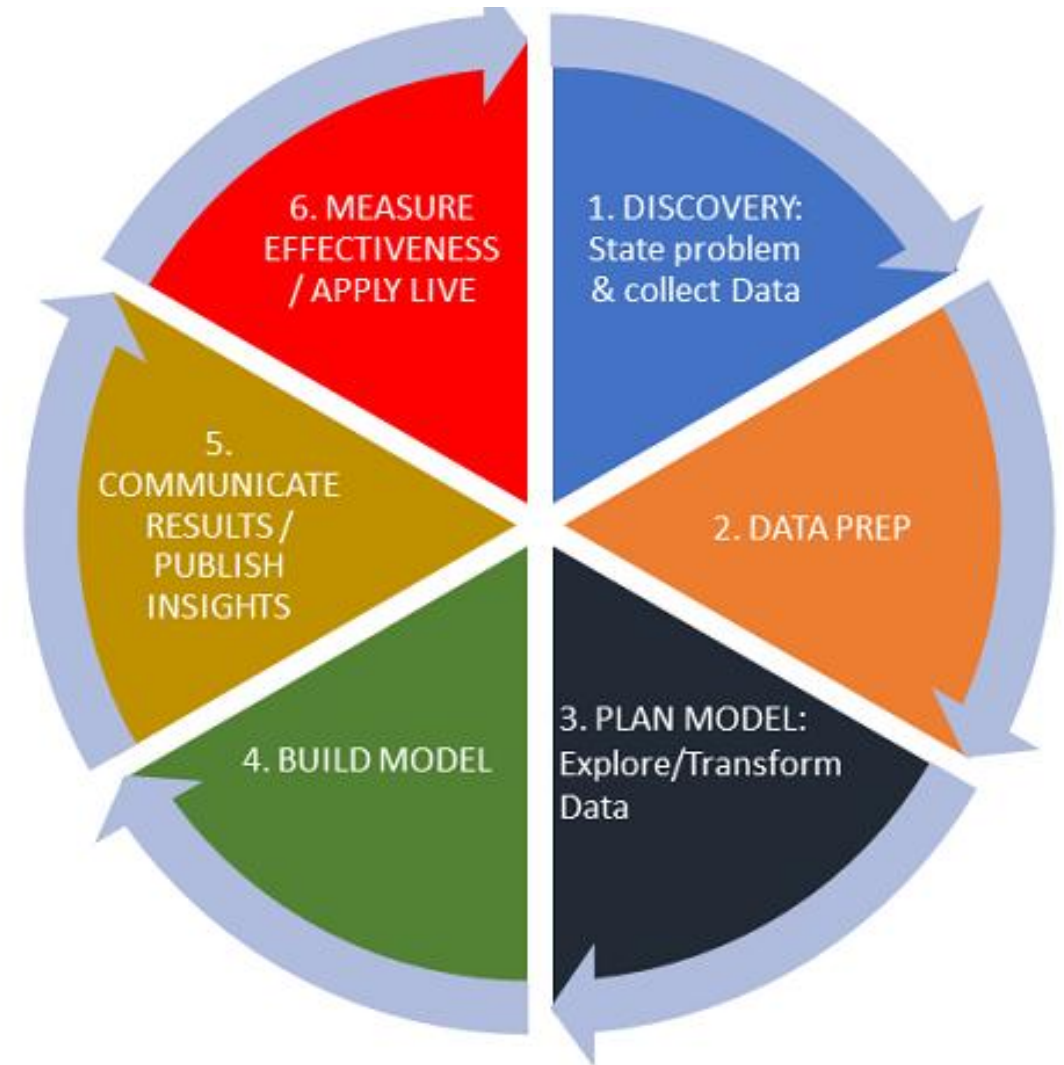
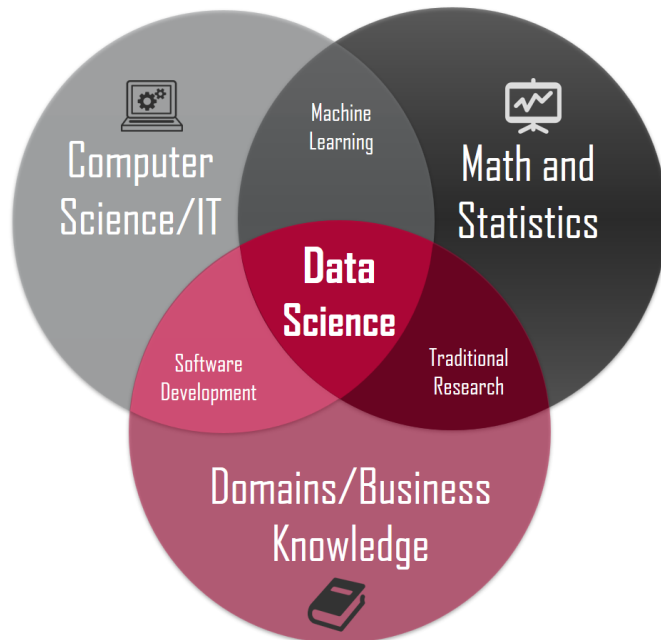
What do you
know about
Data Science?



Data Science

DS IS A MULTIDISCIPLINARY SUBJECT

- It starts from user-specified objectives
- It exploits algorithms to extract patterns and models, with a mathematic approach



Data Science life cycle

Why/What Data Mining

We have a lot of **data** (a collection of raw value elements)

It is quite easy to extract **information** (the result of collecting and organizing data) such as

- relationships between data items
- context and meaning

We aim to something higher than information:
we need to shift from data to knowledge

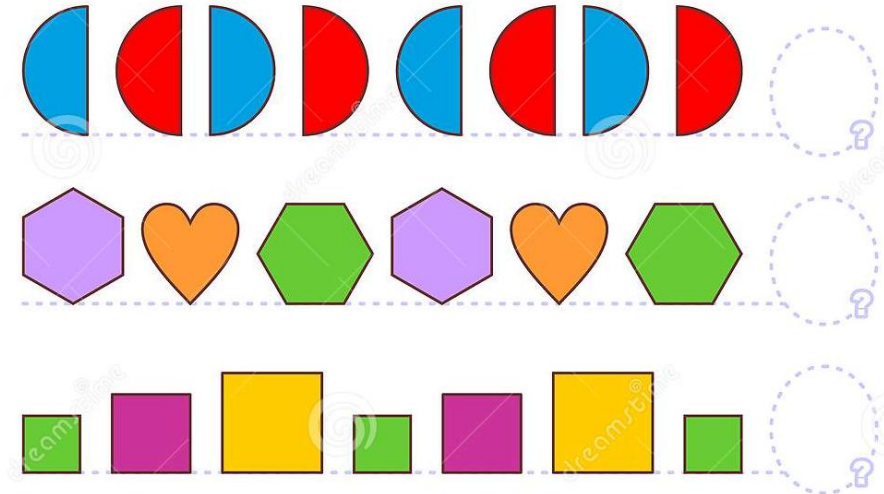
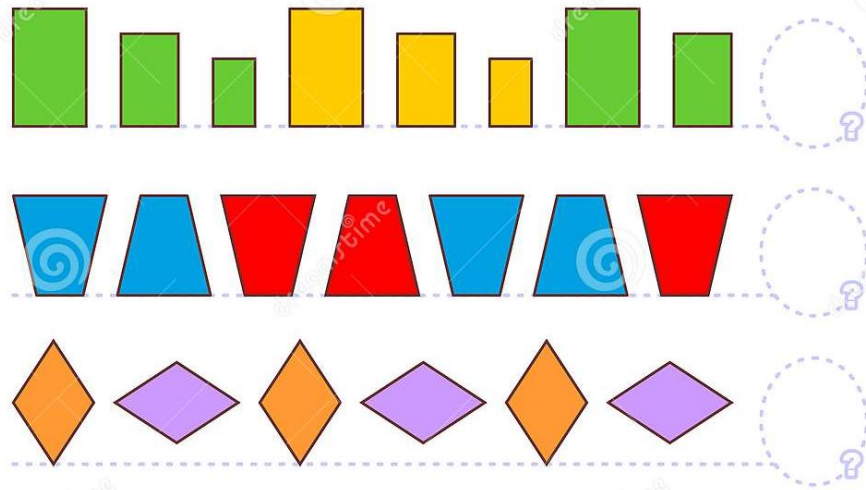
Knowledge: understanding information based on recognizing patterns

Curiosity: the DM name is wrong, because we are not mining for *data*, which are already there and available, but for *patterns*!

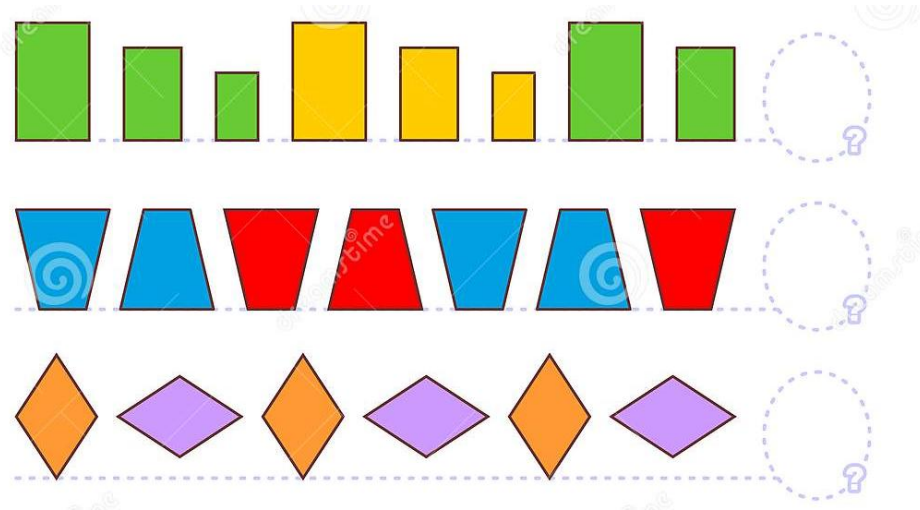


Toy examples of «mining for patterns»

WHAT COMES NEXT?



Why are we interested in patterns?
Because they allow for (future) *predictions*!



VS

Example2 - Excel

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|----|---------------|---------|---------|--------|---------|---------|--------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| | Site 1 | Site 2 | Site 3 | Site 4 | Site 5 | Site 6 | Site 7 | Site 8 | Site 9 | Site 10 | Site 11 | Site 12 | Site 13 | Site 14 | Site 15 | Site 16 | Site 17 | Site 18 | Site 19 | Site 20 | Site 21 | Site 22 | Site 23 | Site 24 | Site 25 | Site 26 |
| 1 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | January '19 | 1234 | 98 | 982 | 1234 | 98 | 982 | 2098 | 92837 | 1098 | 2098 | 92837 | 1098 | 2098 | 987 | 2098 | 1098 | 987 | 2098 | 20987 | | | | | | |
| 3 | February '19 | 23098 | 2098 | 2989 | 23098 | 2098 | 2989 | 98 | 883 | 1098 | 98 | 883 | 1098 | 1098 | 9870983 | 2098 | 1098 | 9870983 | 2098 | 987 | | | | | | |
| 4 | March '19 | 2309 | 922039 | 1234 | 2309 | 922039 | 1234 | 2098 | 982 | 2098 | 2098 | 982 | 2098 | 2098 | 92837 | 1098 | 2098 | 92837 | 1098 | 92837 | 1098 | 9870983 | | | | |
| 5 | April '19 | 2098 | 982 | 23098 | 2098 | 982 | 23098 | 922039 | 2098 | 98 | 922039 | 2098 | 98 | 98 | 883 | 1098 | 98 | 883 | 1098 | 92837 | | | | | | |
| 6 | May '19 | 2098 | 20987 | 2309 | 2098 | 20987 | 2309 | 982 | 1234 | 2098 | 982 | 1234 | 2098 | 2098 | 982 | 2098 | 982 | 2098 | 982 | 2098 | 982 | 2098 | | | | |
| 7 | June '19 | 1098 | 987 | 2098 | 1098 | 987 | 2098 | 20987 | 23098 | 98 | 20987 | 23098 | 98 | 922039 | 2989 | 98 | 922039 | 2989 | 98 | 982 | | | | | | |
| 8 | July '19 | 1098 | 9870983 | 2098 | 1098 | 9870983 | 2098 | 987 | 2309 | 2098 | 987 | 2309 | 2098 | 982 | 1234 | 2098 | 982 | 1234 | 2098 | 2989 | | | | | | |
| 9 | August '19 | 2098 | 92837 | 1098 | 2098 | 92837 | 1098 | 9870983 | 2098 | 922039 | 9870983 | 2098 | 922039 | 20987 | 23098 | 98 | 20987 | 23098 | 98 | 1234 | | | | | | |
| 10 | September '19 | 98 | 883 | 1098 | 98 | 883 | 1098 | 92837 | 2098 | 982 | 92837 | 2098 | 982 | 987 | 2309 | 2098 | 987 | 2309 | 2098 | 23098 | | | | | | |
| 11 | October '19 | 2098 | 982 | 2098 | 2098 | 982 | 2098 | 883 | 1098 | 982 | 883 | 1098 | 982 | 9870983 | 2098 | 922039 | 9870983 | 2098 | 922039 | 2309 | | | | | | |
| 12 | November '19 | 922039 | 2989 | 98 | 922039 | 2989 | 98 | 982 | 1098 | 2989 | 982 | 1098 | 2989 | 92837 | 2098 | 982 | 92837 | 2098 | 982 | 2098 | | | | | | |
| 13 | December '19 | 982 | 1234 | 2098 | 982 | 1234 | 2098 | 2989 | 2098 | 1234 | 2989 | 2098 | 1234 | 883 | 1098 | 982 | 883 | 1098 | 982 | 2098 | | | | | | |
| 14 | January '20 | 20987 | 23098 | 98 | 20987 | 23098 | 98 | 1234 | 98 | 23098 | 1234 | 98 | 23098 | 982 | 1098 | 2989 | 982 | 1098 | 2989 | 1098 | | | | | | |
| 15 | February '20 | 987 | 2309 | 2098 | 987 | 2309 | 2098 | 23098 | 2098 | 2309 | 23098 | 2098 | 2309 | 2989 | 2098 | 1234 | 2989 | 2098 | 1234 | 1098 | | | | | | |
| 16 | March '20 | 9870983 | 2098 | 922039 | 9870983 | 2098 | 922039 | 2309 | 98 | 2098 | 2309 | 98 | 2098 | 1234 | 98 | 23098 | 1234 | 98 | 23098 | 2098 | | | | | | |
| 17 | April '20 | 92837 | 2098 | 982 | 92837 | 2098 | 982 | 2098 | 2098 | 2098 | 2098 | 2098 | 2098 | 23098 | 2098 | 2309 | 23098 | 2098 | 2309 | 98 | | | | | | |
| 18 | May '20 | 883 | 1098 | 982 | 883 | 1098 | 982 | 2098 | 922039 | 1098 | 2098 | 922039 | 1098 | 2309 | 98 | 2098 | 2309 | 98 | 2098 | 2098 | | | | | | |
| 19 | June '20 | 982 | 1098 | 2989 | 982 | 1098 | 2989 | 1098 | 982 | 1098 | 1098 | 982 | 1098 | 2098 | 2098 | 2098 | 2098 | 2098 | 2098 | 922039 | | | | | | |
| 20 | July '20 | 2989 | 2098 | 1234 | 2989 | 2098 | 1234 | 1098 | 20987 | 2098 | 1098 | 20987 | 2098 | 2098 | 922039 | 1098 | 2098 | 922039 | 1098 | 2098 | | | | | | |
| 21 | August '20 | 1234 | 98 | 23098 | 1234 | 98 | 23098 | 2098 | 987 | 98 | 2098 | 987 | 98 | 1098 | 982 | 1098 | 1098 | 982 | 1098 | 98 | | | | | | |
| 22 | September '20 | 23098 | 2098 | 2309 | 23098 | 2098 | 2309 | 98 | 9870983 | 2098 | 98 | 9870983 | 2098 | 1098 | 20987 | 2098 | 1098 | 20987 | 2098 | 2098 | | | | | | |
| 23 | October '20 | 2309 | 98 | 2098 | 2309 | 98 | 2098 | 2098 | 92837 | 98 | 2098 | 92837 | 98 | 2098 | 987 | 98 | 2098 | 987 | 98 | 922039 | | | | | | |
| 24 | November '20 | 2098 | 2098 | 2098 | 2098 | 2098 | 2098 | 922039 | 883 | 2098 | 922039 | 883 | 2098 | 98 | 9870983 | 2098 | 98 | 9870983 | 2098 | 982 | | | | | | |
| 25 | December '20 | 2098 | 922039 | 1098 | 2098 | 922039 | 1098 | 2098 | 92837 | 1098 | 2098 | 92837 | 1098 | 2098 | 92837 | 98 | 2098 | 92837 | 98 | 2098 | | | | | | |
| 26 | January '21 | 1098 | 982 | 1098 | 1098 | 982 | 1098 | 98 | 883 | 1098 | 98 | 883 | 1098 | 98 | 883 | 1098 | 922039 | 883 | 2098 | 922039 | | | | | | |
| 27 | February '21 | 1098 | 20987 | 2098 | 1098 | 20987 | 2098 | 2098 | 982 | 2098 | 2098 | 982 | 2098 | 982 | 982 | 922039 | 982 | 982 | 922039 | 9870983 | | | | | | |
| 28 | March '21 | 2098 | 987 | 98 | 2098 | 987 | 98 | 922039 | 2989 | 98 | 922039 | 2989 | 98 | 20987 | 2989 | 982 | 20987 | 2989 | 982 | 92837 | | | | | | |
| 29 | April '21 | 98 | 9870983 | 2098 | 98 | 9870983 | 2098 | 982 | 1234 | 2098 | 982 | 1234 | 2098 | 987 | 1234 | 982 | 987 | 1234 | 982 | 883 | | | | | | |
| 30 | May '21 | 2098 | 92837 | 98 | 2098 | 92837 | 98 | 20987 | 23098 | 98 | 20987 | 23098 | 98 | 9870983 | 23098 | 2989 | 9870983 | 23098 | 2989 | 982 | | | | | | |

How can we mine for patterns on huge data sets?
With algorithms running on computers!

Data structures

There are three common types of data structures, for data analysis:

1. **Unstructured data** is information that either does not have a predefined data model or is not organized in a pre-defined manner.
2. **Semi-structured data** is a form of structured data that does not conform with the formal structure of data models associated with relational databases or other forms of data tables, but nonetheless contain tags or other markers to separate semantic elements and enforce hierarchies of records and fields within the data.
Examples include JSON and XML.
3. **Structured data** is data that adheres to a pre-defined data model and is therefore straightforward to analyze. Structured data conforms to a tabular format with relationship between the different rows and columns.
Examples: Excel files or SQL databases. Each of these have structured rows and columns that can be sorted.

Structured data

A **comma-separated values (CSV)** file is a delimited text file that uses a comma to separate values. Each line of the file is a data record.

- Each record consists of one or more fields, separated by commas.
- A CSV file typically stores tabular data (numbers and text) in plain text, in which case each line will have the same number of fields.
- CSV is a common data exchange format that is widely supported by consumer, business, and scientific applications.



A **database** is an organized collection of data and it can handle very complicated queries.

- Relational databases became dominant in the 1980s. These model data as rows and columns in a series of tables, and the vast majority use SQL for writing and querying data.
- In the 2000s, non-relational databases became popular, referred to as NoSQL because they use different query languages.



Characteristics of data sets

1. Dimensionality (number of columns)
 - the difference between having a small or a large (hundreds, thousands, . . .) of attribute is also qualitative
 - see the curse of dimensionality, later
2. Sparsity (when there are many missing values)
 - Nulls in disguise (beware!): a widespread bad habit is to store zero or some special value when a piece information is not available
3. Resolution has a great influence on the results
 - the analysis of too detailed data can be affected by noise
 - the analysis of too general data can hide interesting patterns

A good understanding of the data set is required!

Data types

| <i>Data Type</i> | | <i>Description</i> | <i>Examples</i> | <i>Descriptive statistics allowed</i> |
|------------------|----------|---|--|---|
| Categorical | Nominal | The values are a set of labels, the available information allows to distinguish a label from another Operators: = and \neq | zip code, eye color, sex, ... | mode, entropy, contingency, correlation, χ^2 test |
| | Ordinal | The values provide enough information for a total ordering Operators: $<>\leq\geq$ | hardness of minerals, non-numerical quality evaluations (bad, fair, good, excellent) | median, percentiles, rank correlations |
| Numerical | Interval | The difference is meaningful Operators: $+-$ | Calendar dates, temperatures in centigrades and Fahrenheit | average, standard deviation, Pearson's correlation, F and t tests |
| | Ratio | Have a univocal definition of 0 Allow all the mathematic operations on numbers | Kelvin temperatures, masses, length, counts | geometric mean, harmonic mean, percentage variation |

The “description” and “descriptive statistics” columns are *incremental*, i.e. the properties described in a row are added to the properties described in the rows above

Data preprocessing

Data preprocessing is a data mining technique that transforms raw data into an understandable format.

It is an important step in the DM process, above all in case of:

- Presence of noise and outliers
- Missing values
- Presence of duplicates and data inconsistencies
- Curse of dimensionality (when dimensionality is very high)

It is typically a long step!

“garbage in, garbage out” (GIGO)

flawed or nonsense input data produces nonsense output

2 - Algorithms for DS

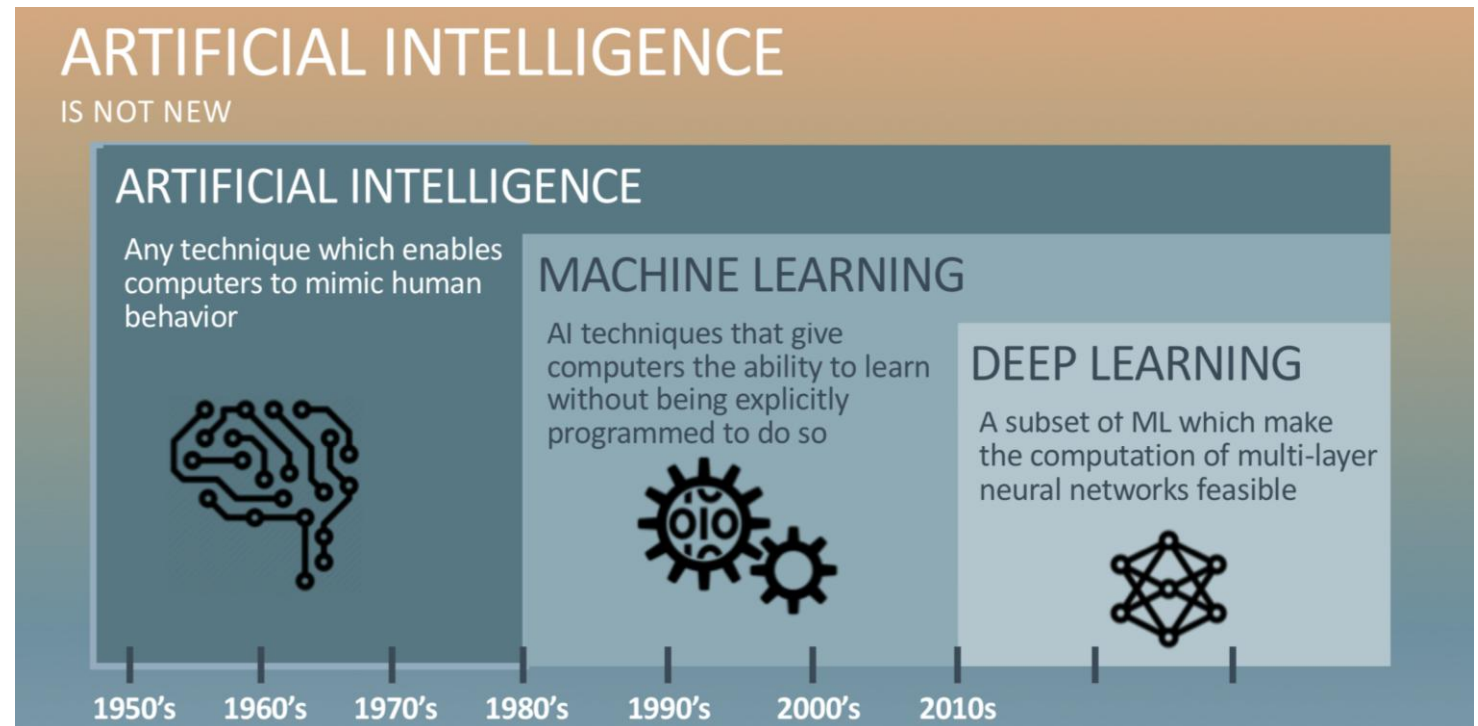
HOW TO LEARN FROM DATA

Brief history of Artificial Intelligence

In the beginning was the (pure) *statistics*. Since the 18th century, great developments in:

- Descriptive statistics
- Inferential statistics
- Statistical models

Since late '50s of XX century:





Example of AI

Soybean diseases [Michalski and Chilausky, 1980]

Soybean (Large) Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Michalski's famous soybean disease database



| | | | | | |
|----------------------------|----------------|-----------------------|-----|---------------------|------------|
| Data Set Characteristics: | Multivariate | Number of Instances: | 307 | Area: | Life |
| Attribute Characteristics: | Categorical | Number of Attributes: | 35 | Date Donated | 1988-07-11 |
| Associated Tasks: | Classification | Missing Values? | Yes | Number of Web Hits: | 120988 |

diagnosis of soybean diseases (19 different diseases)



AI approach:

starting from the expert knowledge, computer scientists tried to translate (formalised) the classification rules into code.

Examples:

If leaf condition = normal and stem condition = abnormal
 and stem cankers = below soil line and canker lesion color = brown
then: diagnosis is rhizoctonia root rot

If leaf malformation = absent and stem condition = abnormal
 and stem cankers = below soil line and canker lesion color = brown
then: diagnosis is rhizoctonia root rot

Results:

the diagnosis accuracy obtained by the rules alone, without expert assistance, is 72%

Problems:



1. the elicitation of rules from expert is difficult and time consuming (35 attributes = huge number of “if” cases!);

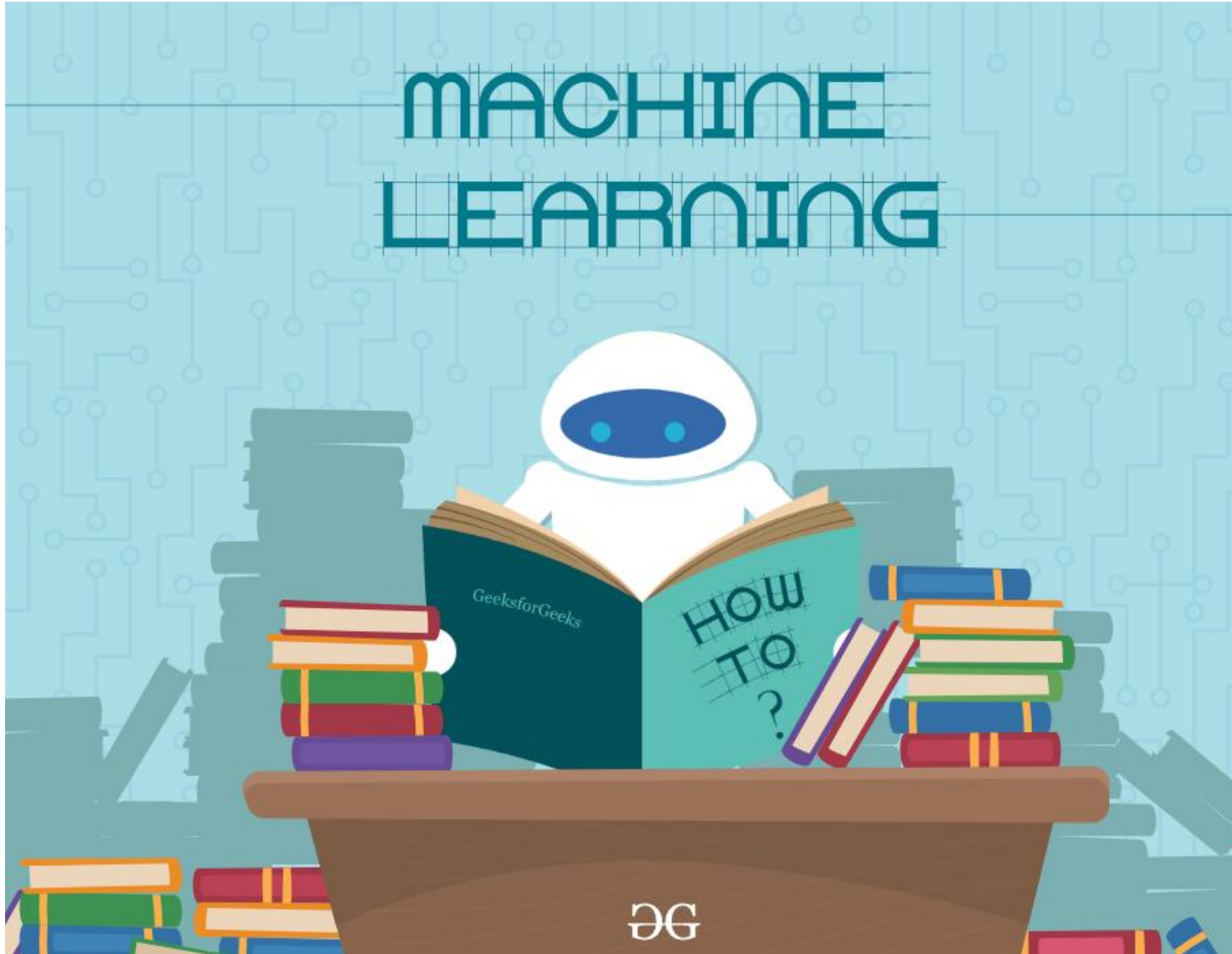
If leaf condition = normal and stem condition = abnormal
 and stem cankers = below soil line and canker lesion color = brown
then: diagnosis is rhizoctonia root rot

If leaf malformation = absent and stem condition = abnormal
 and stem cankers = below soil line and canker lesion color = brown
then: diagnosis is rhizoctonia root rot

2. the rules are not independent, and they should be carefully checked;
3. the achieved accuracy (72%) is not satisfactory.

Conclusion:

the rules are not able to capture all the expert knowledge *by being told*



Alternative approach



The same data set has been processed by a machine learning algorithm, to generate its classification rules.

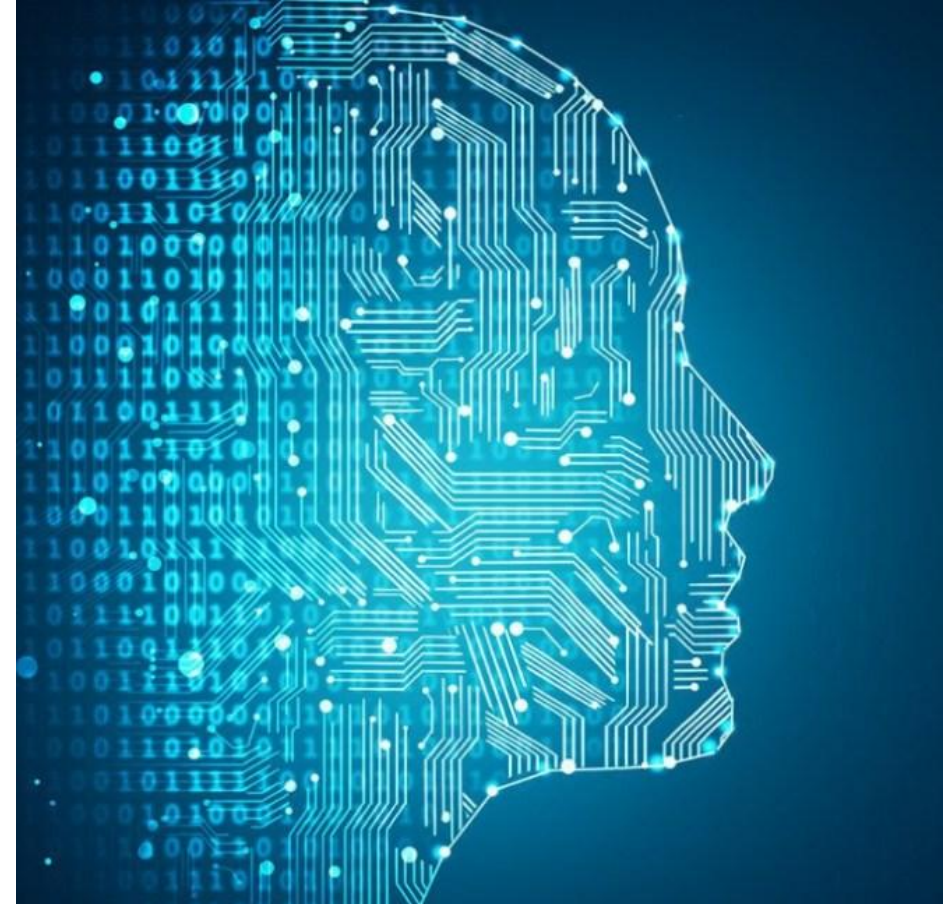
The new set of rules got an accuracy of 97.5%, comparable with that of a junior expert.

Machine Learning

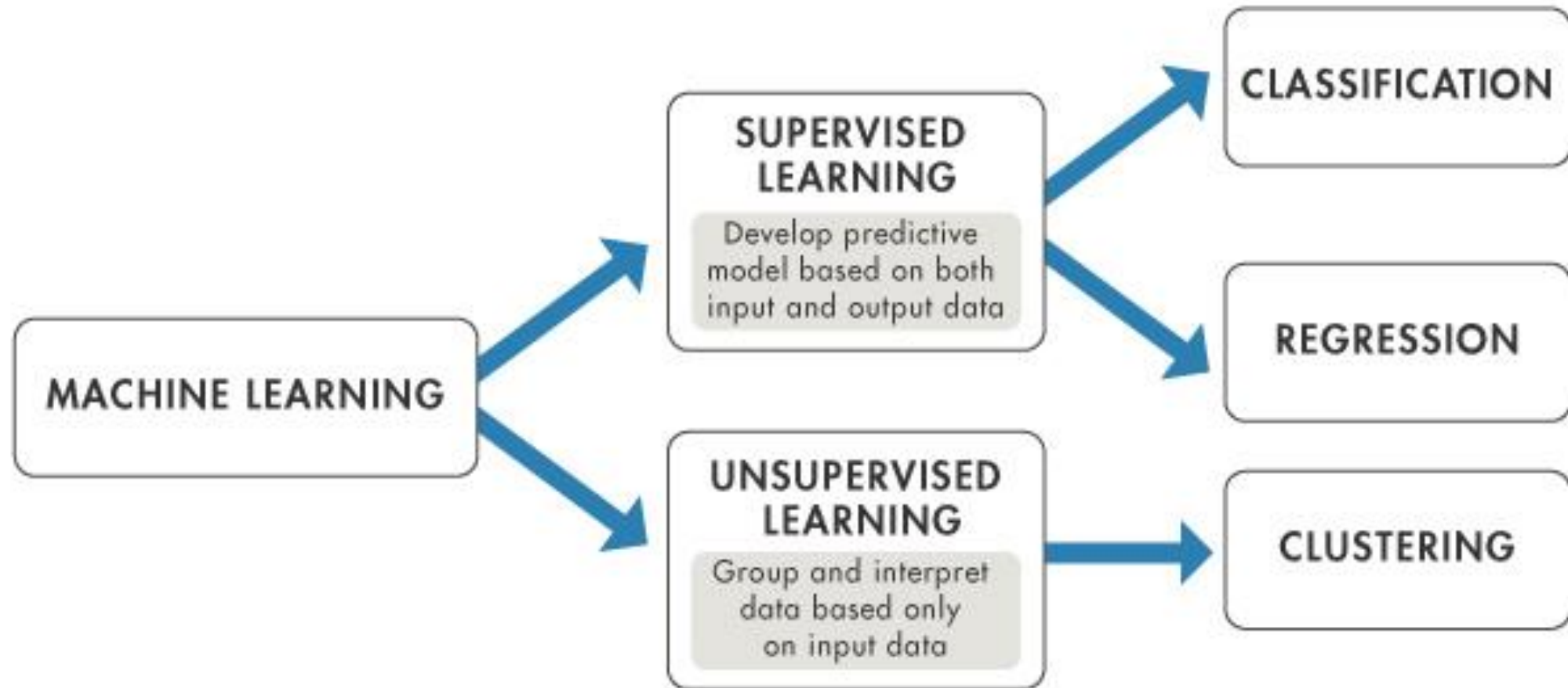
Machine Learning is a field of study that gives computers the ability to learn without being explicitly programmed (*learning from examples*)

In particular:

- ML was born in 80s
- Since the early '90s we talk about “Data Mining processes with Machine Learning”

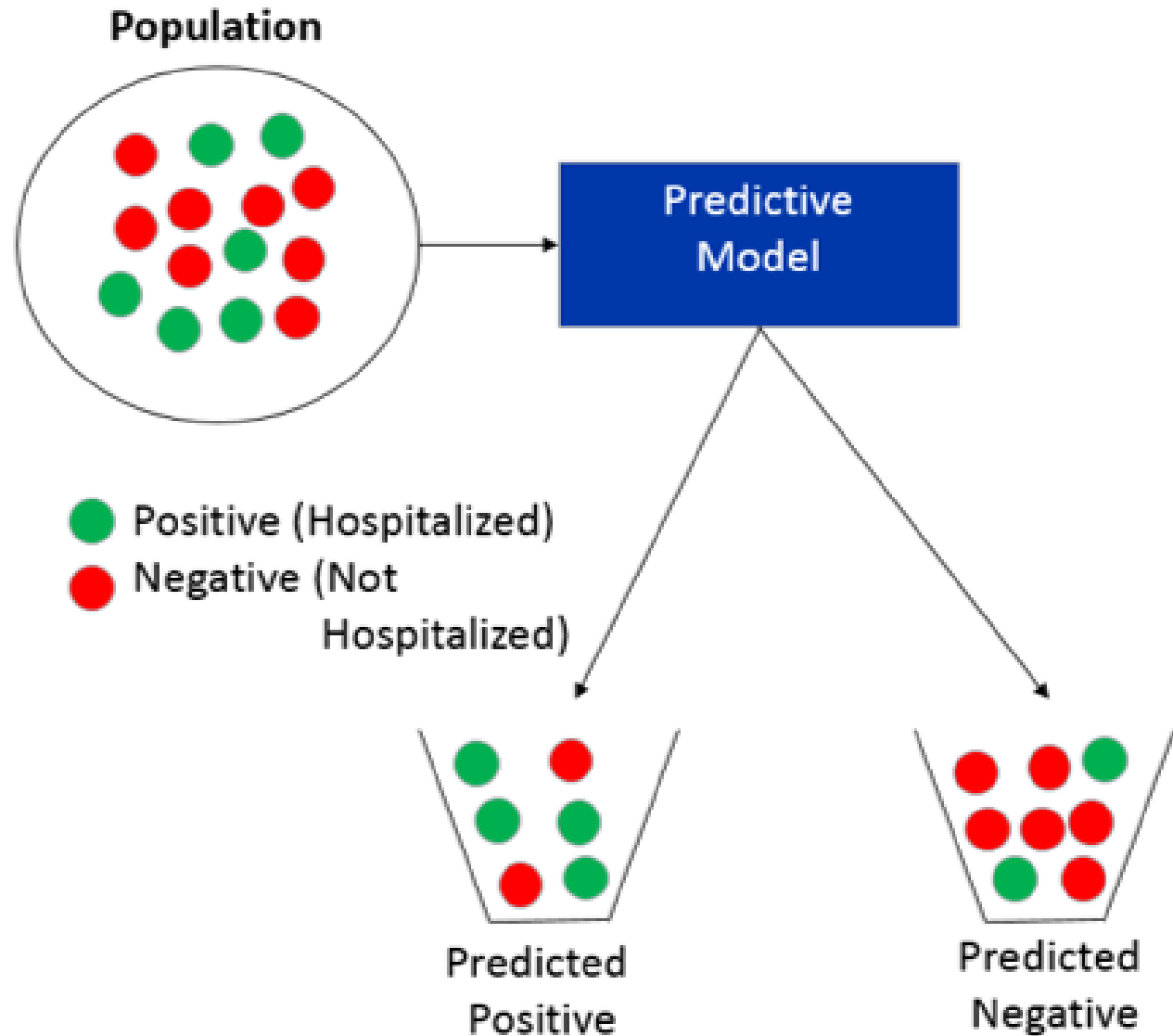


Algorithms for approximations (i.e. for learning from available data)



Classification

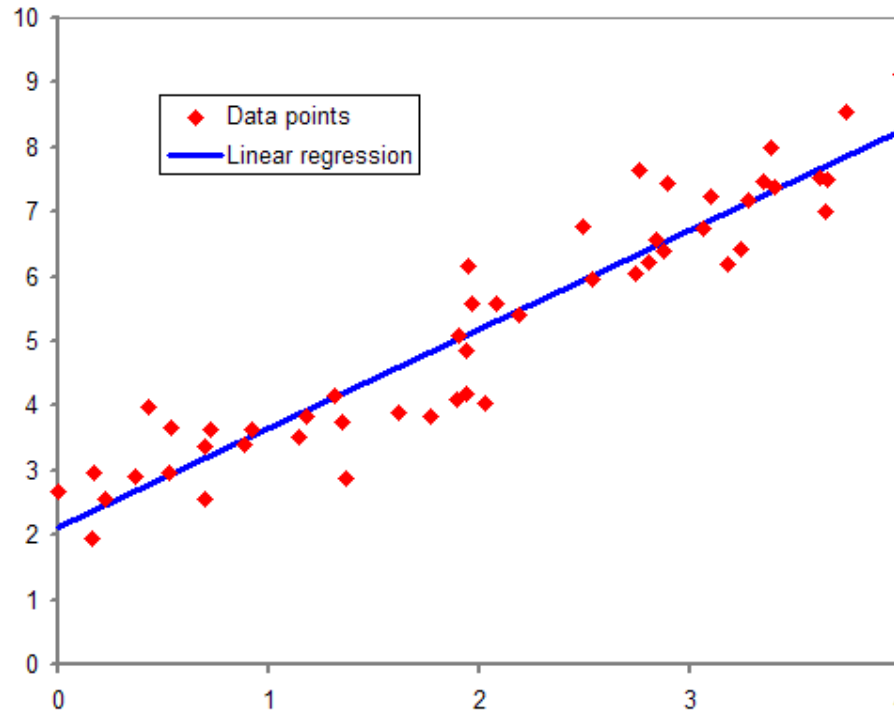
Classification is a supervised task that aims at predicting one class over two or more possible classes for a particular item.



Regression

regression, value estimation

Regression analysis is a set of statistical processes for estimating the relationships between a dependent variable (often called the 'outcome variable') and one or more independent variables (often called 'predictors', 'covariates', or 'features').



x = independent variable
 y = dependent variable,
 $y \approx f(x)$

Given all the data points (x_i, y_i) , we look for f

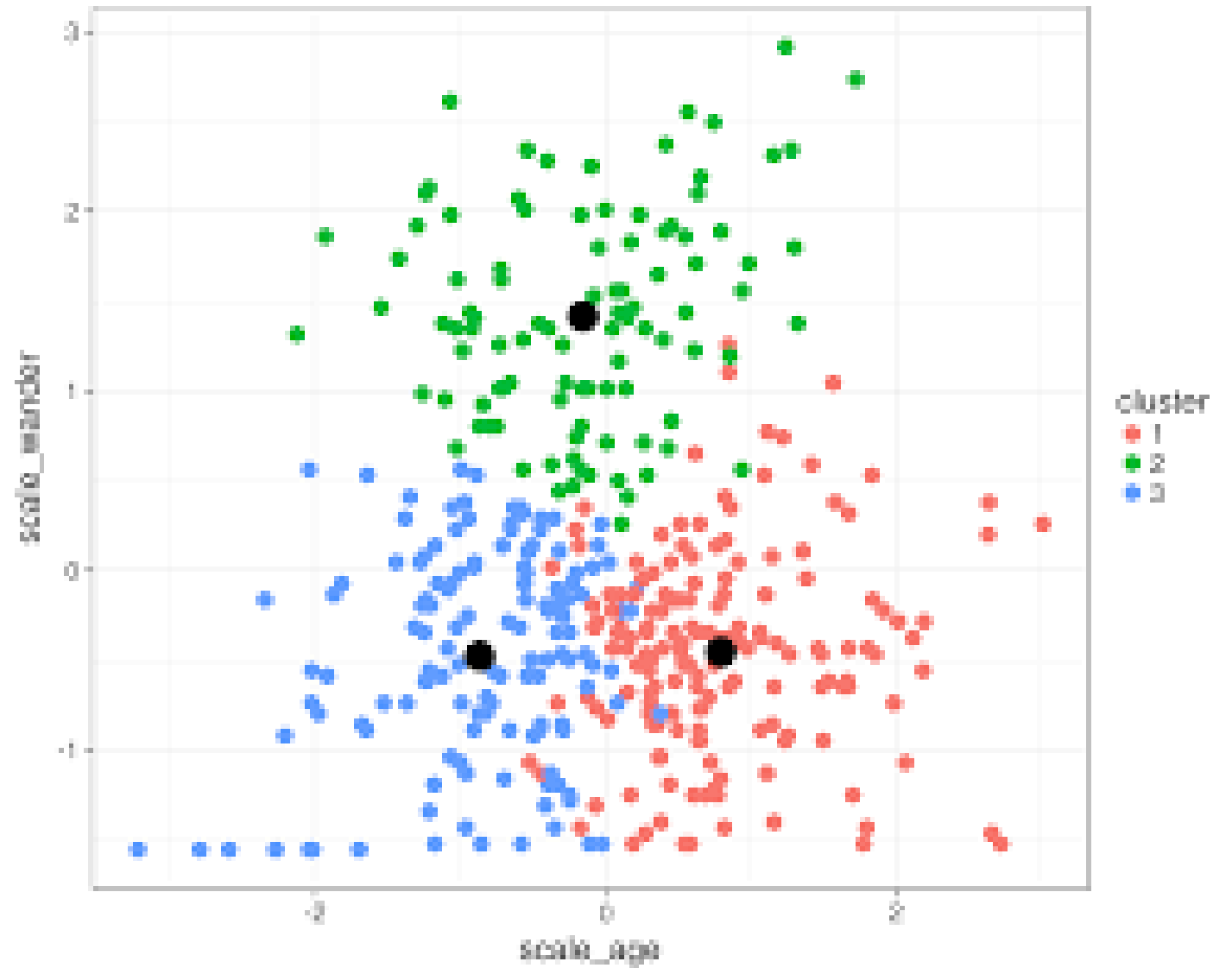
Remarks

- given a set of **numeric** attribute values for an individual, it estimates the value of another numeric attribute
- it is related to *classification*, but the methods are completely different
- is primarily used
 - for **prediction** and forecasting
 - to infer causal relationships between two variables

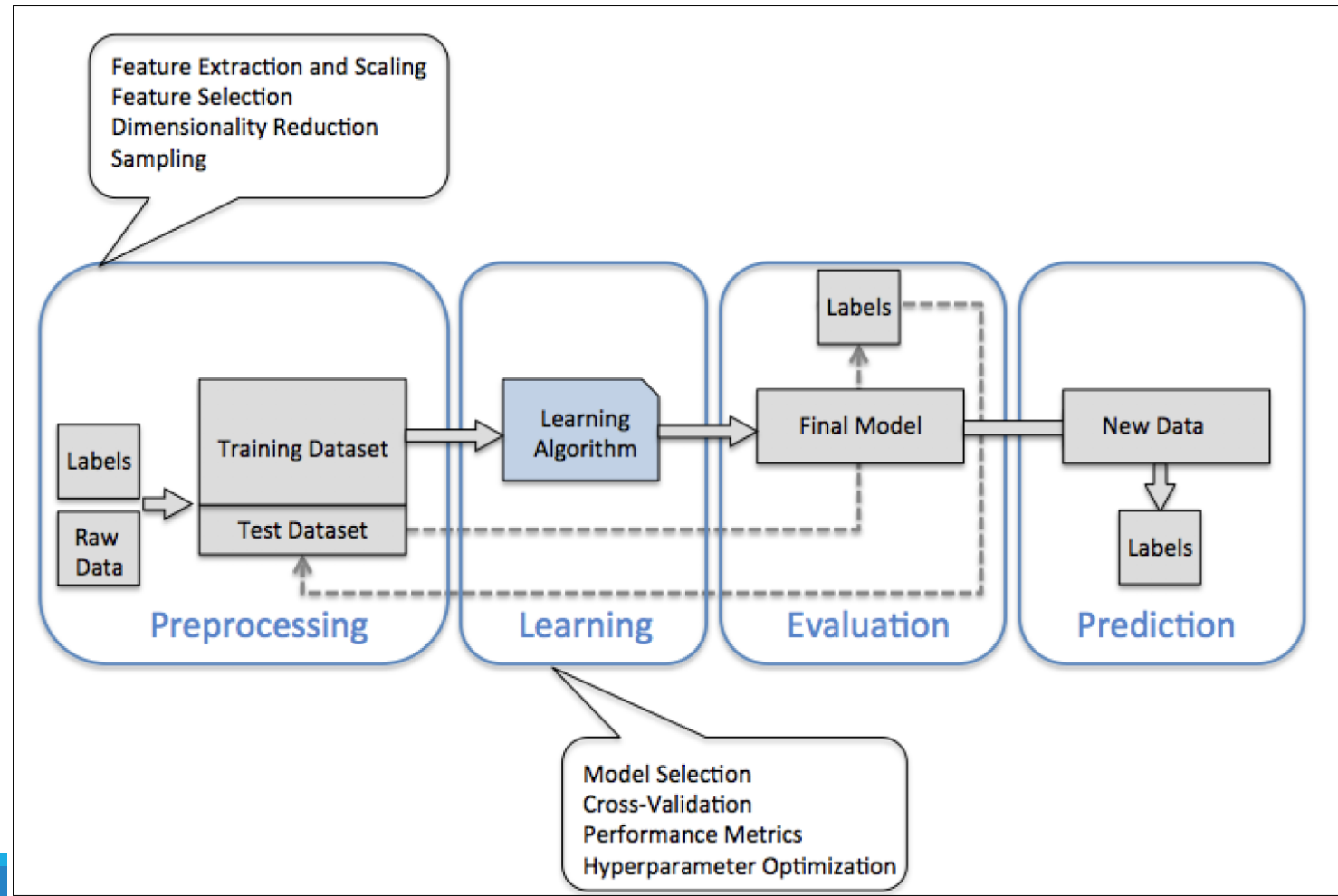
Clustering

clustering , cluster analysis

Clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups.



Steps for ML

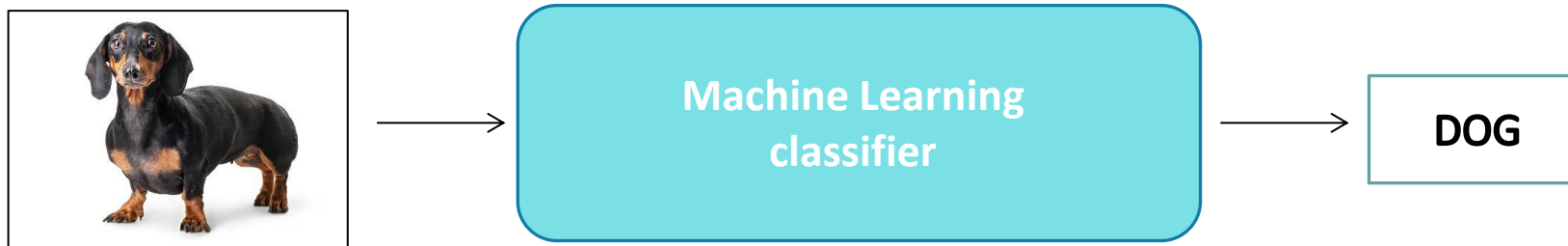


ML definition: what does it mean?

“Machine Learning is that branch of computer science that gives computers the *ability to learn on their own* without being explicitly programmed to do so”

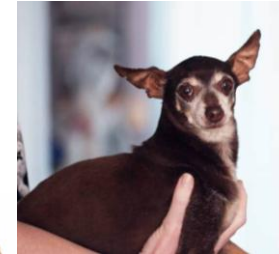
Let's try to understand it with an example of image recognition.
In our setting:

- we have many images (containing cats or dogs) already labelled;
- we want to automatically classify new images as “CAT” or “DOG”, through a ML program.



What we have:

1. Dog images



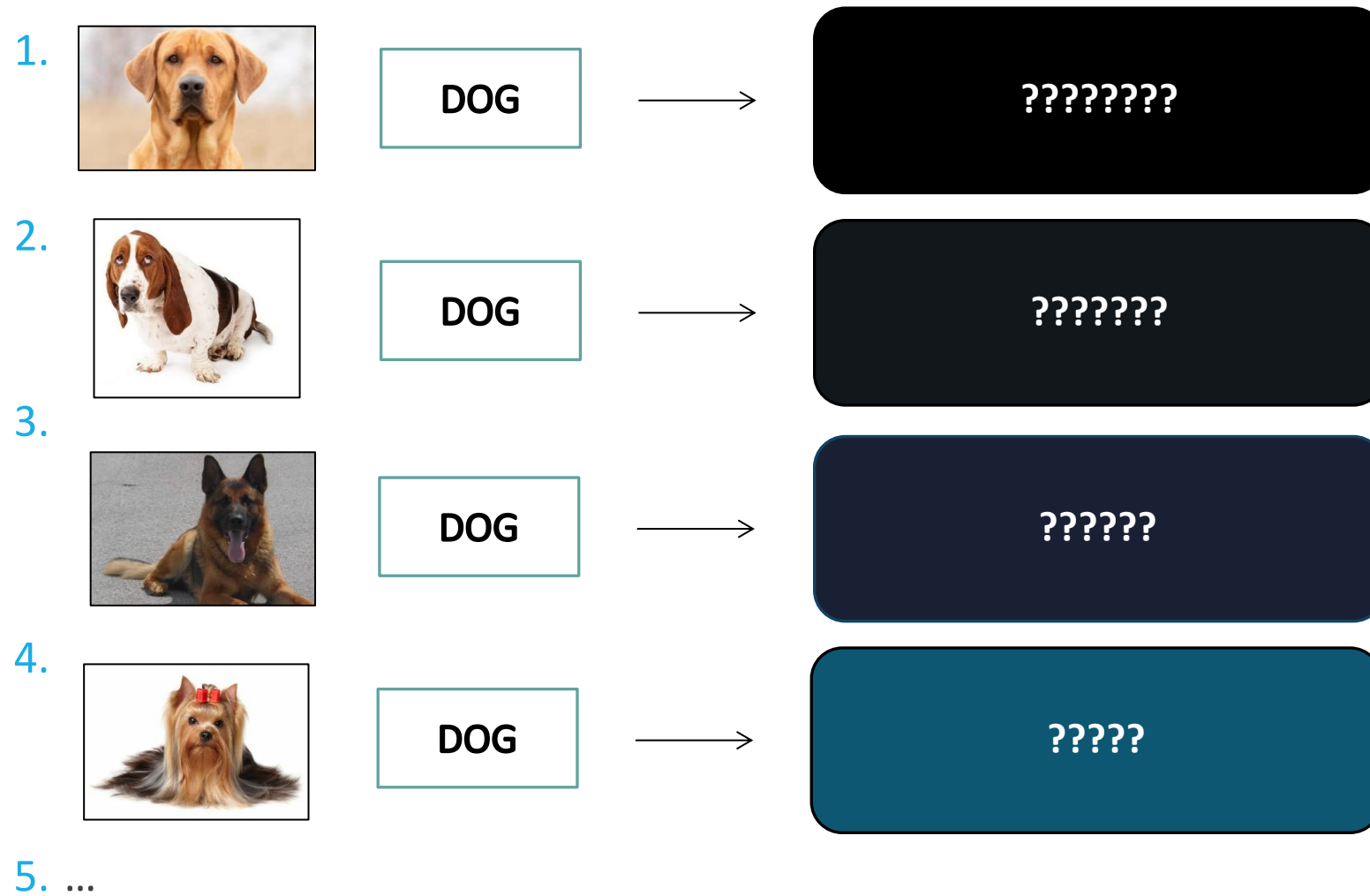
2. Cat images



3. ML theory

????????

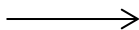
We need to train the ML classifier with the original labelled images



6.

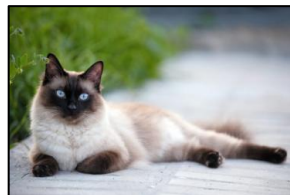


CAT

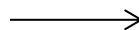


????

7.



CAT

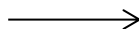


???

8.



CAT

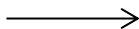


??

9.



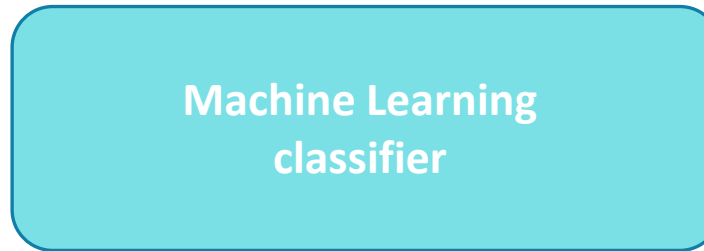
CAT



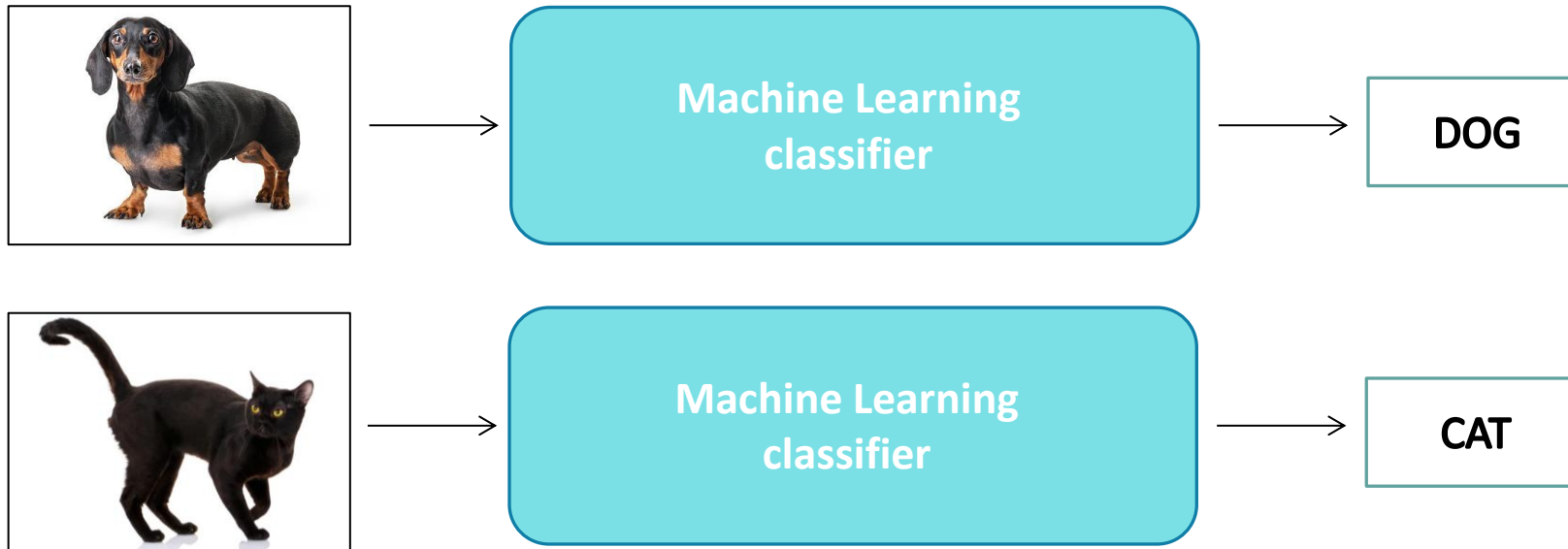
?

10. ...

After this phase, the Machine Learning classifier will be defined:

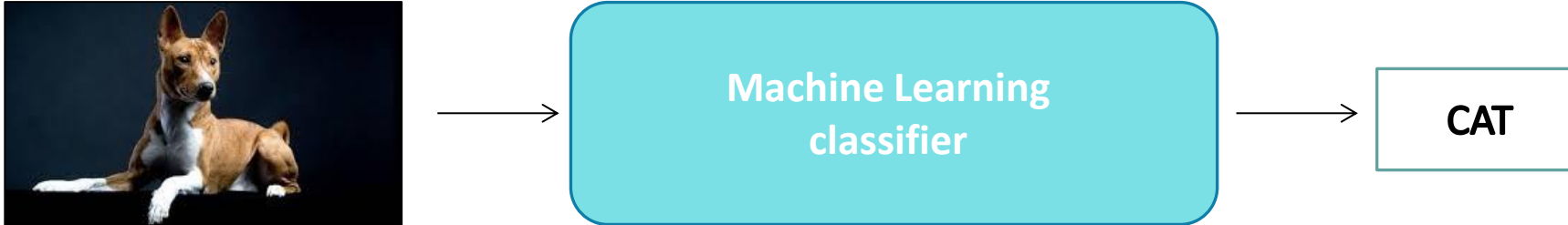


Thus, it will be usable on new, unseen images:

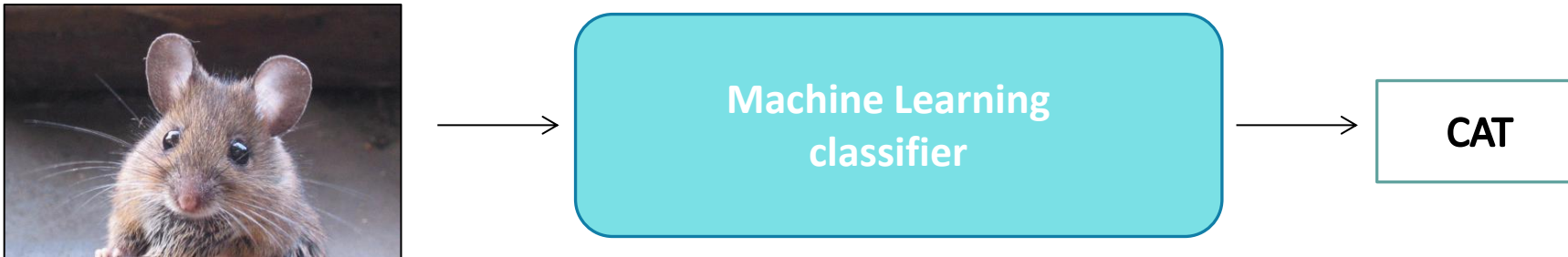


BUT pay attention, because there may be errors:

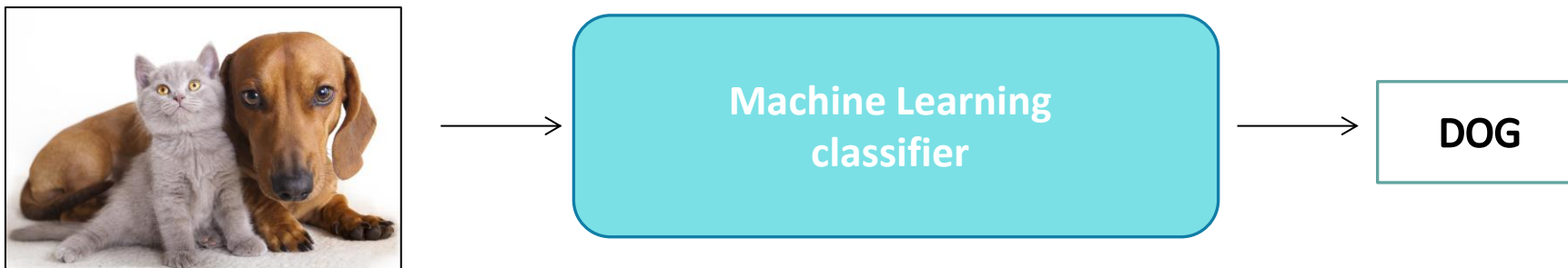
- on some unseen images



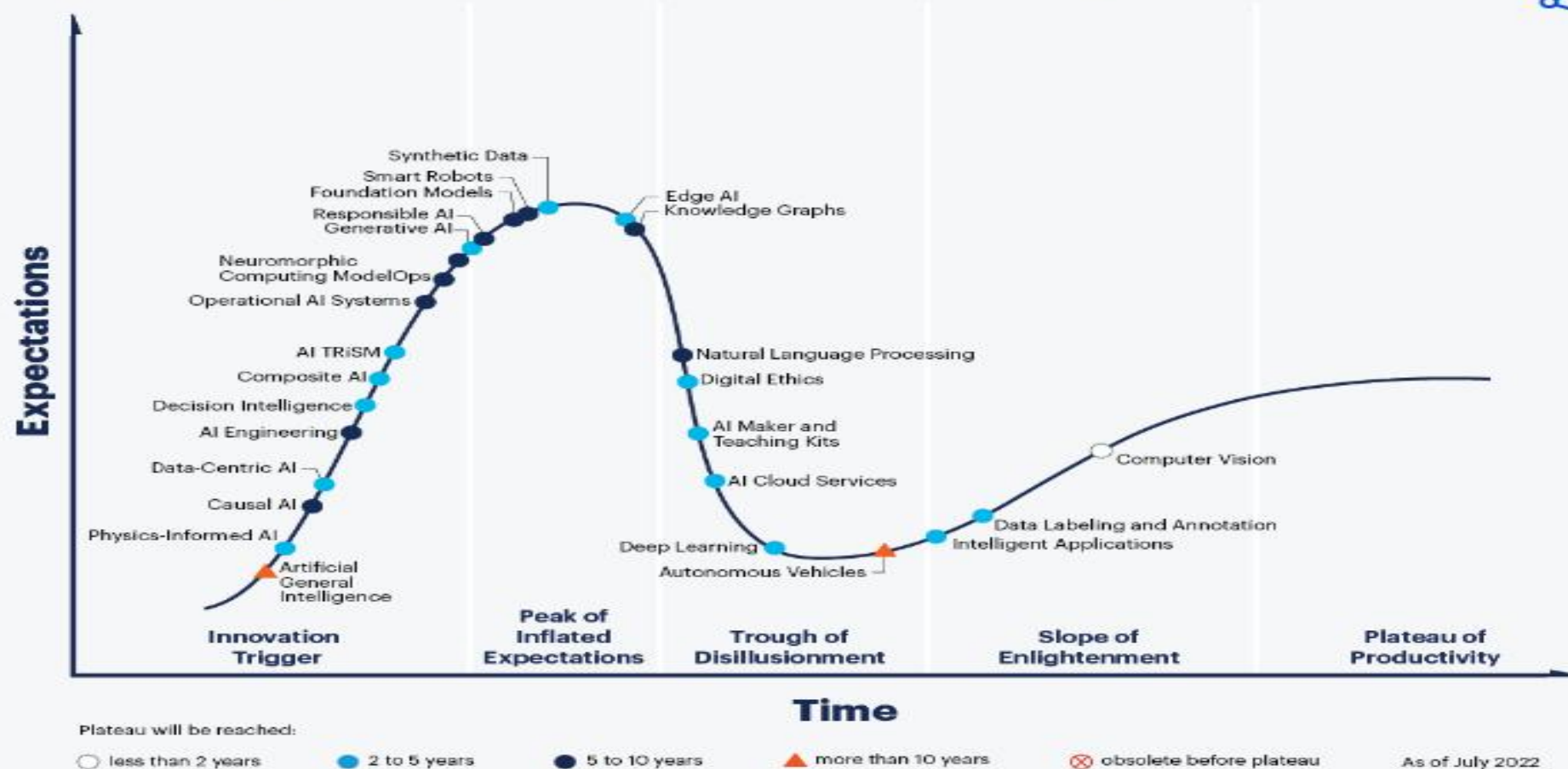
- on unseen objects (other than cats or dogs)



- on more than 1 subject



Hype Cycle for Artificial Intelligence, 2022

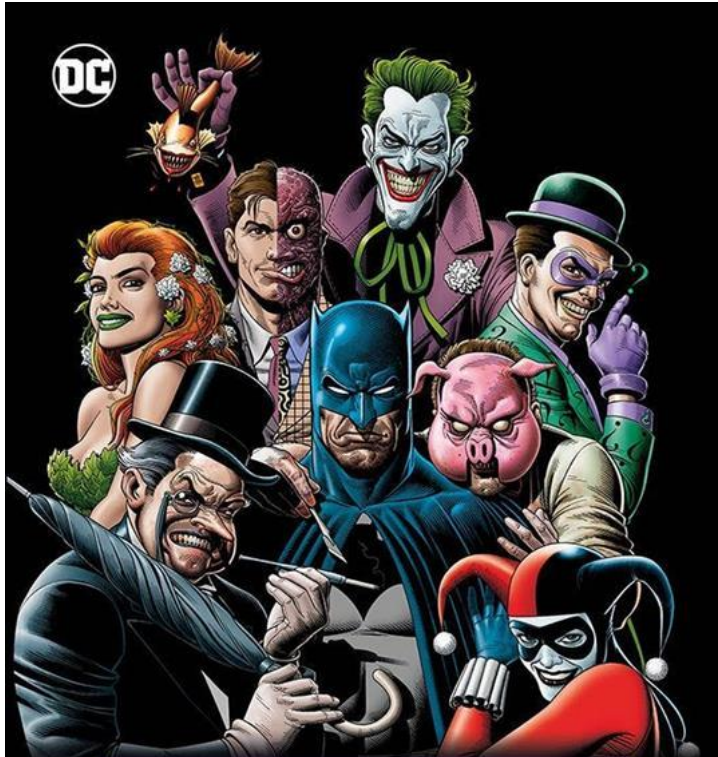


[gartner.com](https://www.gartner.com)

Sources: Gartner
© 2022 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and Hype Cycle are registered trademarks of Gartner, Inc. and its affiliates in the U.S. 1957302

Gartner

Example of decision tree

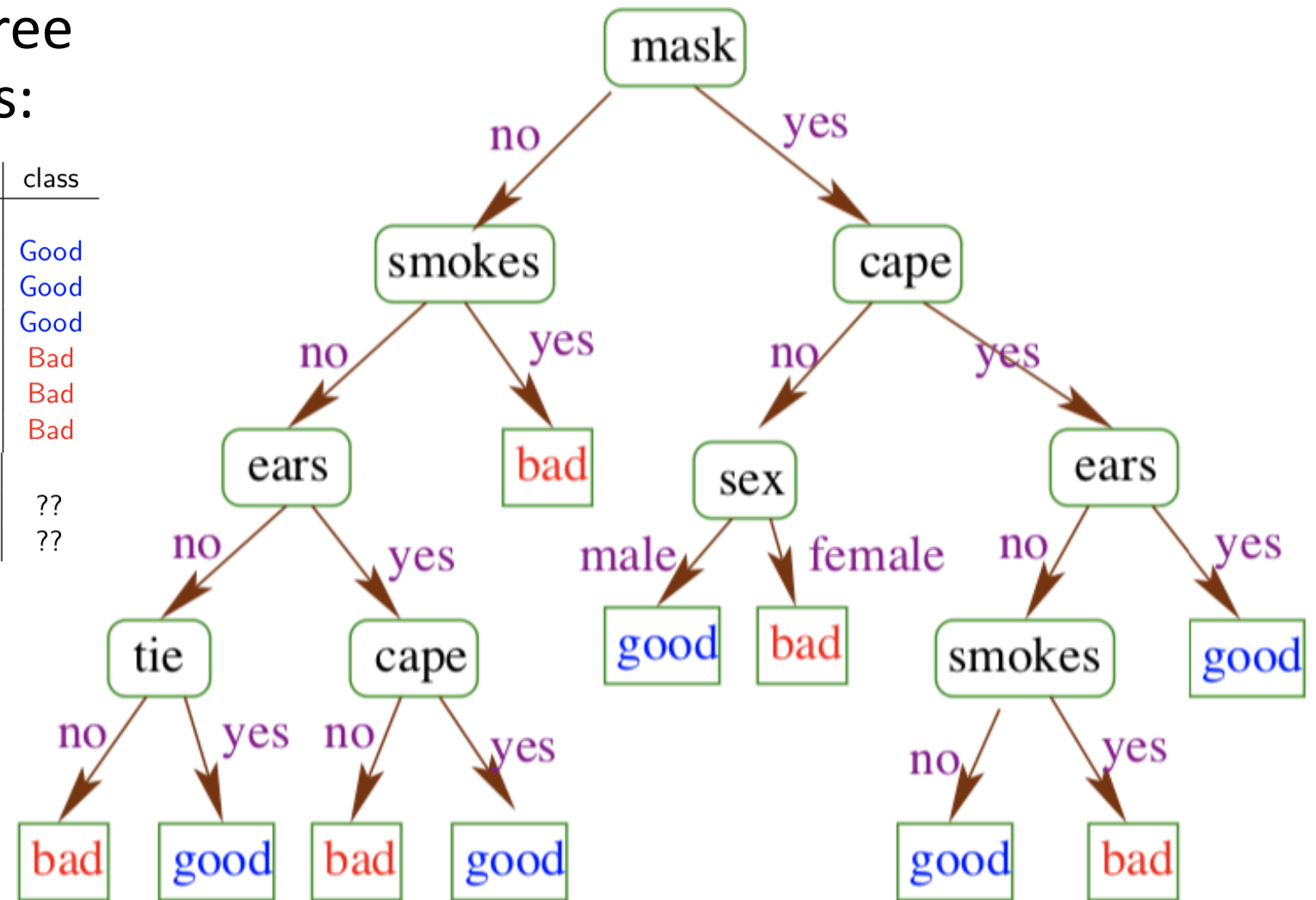


We want to identify characters as good or bad, from their appearance based on:

| | sex | mask | cape | tie | ears | smokes | class |
|----------|---------------|------|------|-----|------|--------|-------|
| | training data | | | | | | |
| batman | male | yes | yes | no | yes | no | Good |
| robin | male | yes | yes | no | no | no | Good |
| alfred | male | no | no | yes | no | no | Good |
| penguin | male | no | no | yes | no | yes | Bad |
| catwoman | female | yes | no | no | yes | no | Bad |
| joker | male | no | no | no | no | no | Bad |

After training a decision tree on the available 6 samples:

| | sex | mask | cape | tie | ears | smokes | class |
|---------------|--------|------|------|-----|------|--------|-------|
| training data | | | | | | | |
| batman | male | yes | yes | no | yes | no | Good |
| robin | male | yes | yes | no | no | no | Good |
| alfred | male | no | no | yes | no | no | Good |
| penguin | male | no | no | yes | no | yes | Bad |
| catwoman | female | yes | no | no | yes | no | Bad |
| joker | male | no | no | no | no | no | Bad |
| test data | | | | | | | |
| batgirl | female | yes | yes | no | yes | no | ?? |
| riddler | male | yes | no | no | no | no | ?? |



Testing:

- Batgirl:
good (correctly classified)
- Riddler:
good (incorrectly classified)