



# KSCHOOL

## Modelling Methodologies

Data Science

Irene Torres Valle



# ¿Who am I?



## Index

- 1 Importance of a data scientist
- 2 Types of data analysis
- 3 Some questionable conclusions
- 4 A mathematical model
- 5 Machine Learning
- 6 Types of learnings
- 7 Model metrics

## Index

8

Concepts related with ML

9

ML Canvas

10

Actual examples of ML

## Every day

2,5 trillions of bytes  
**data**



## Every minute

72 hours  
**video**

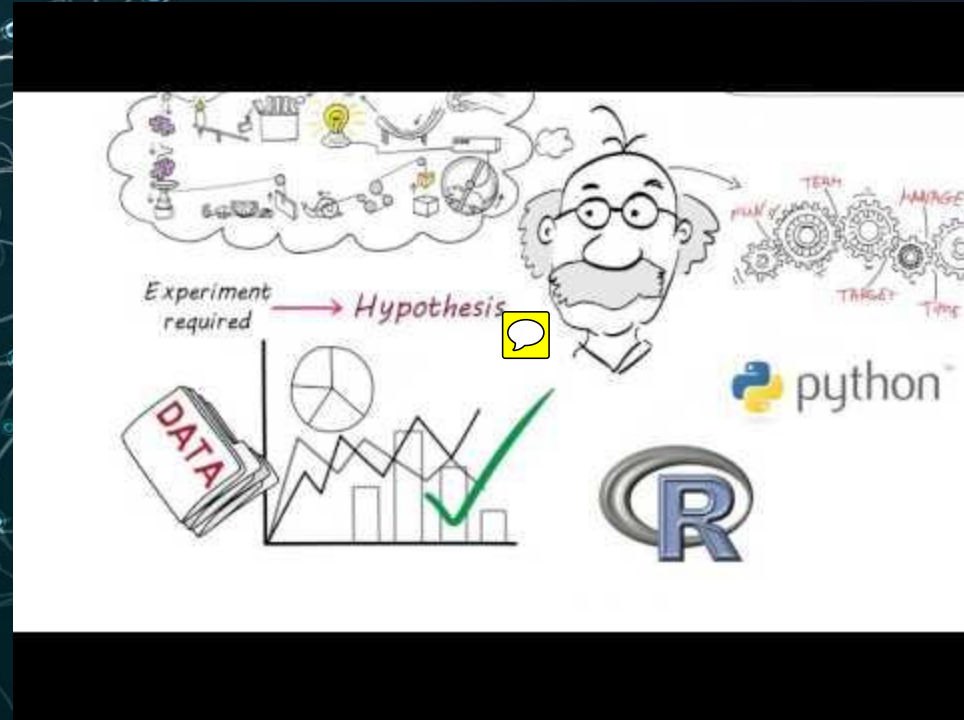
277.000  
**tweets**

200 millions  
**emails**

*“Data Scientist:  
The Sexiest Job of the  
21st Century.”*

*Harvard Business Review*

# Modelling Methodologies



# Types of data analysis



Descriptive

¿What did happen?

Diagnosis

¿Why does it happen?

Predictive

¿What will happen?

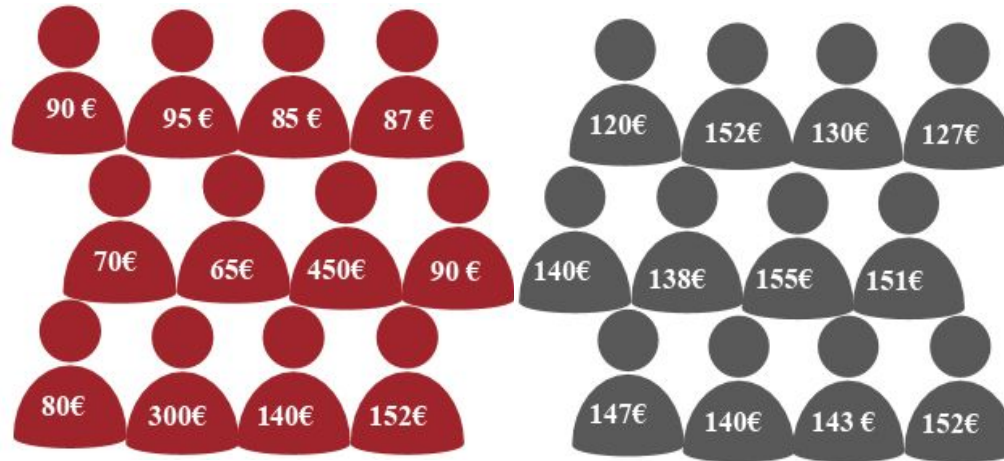
Prescriptive

¿What do I need to do?

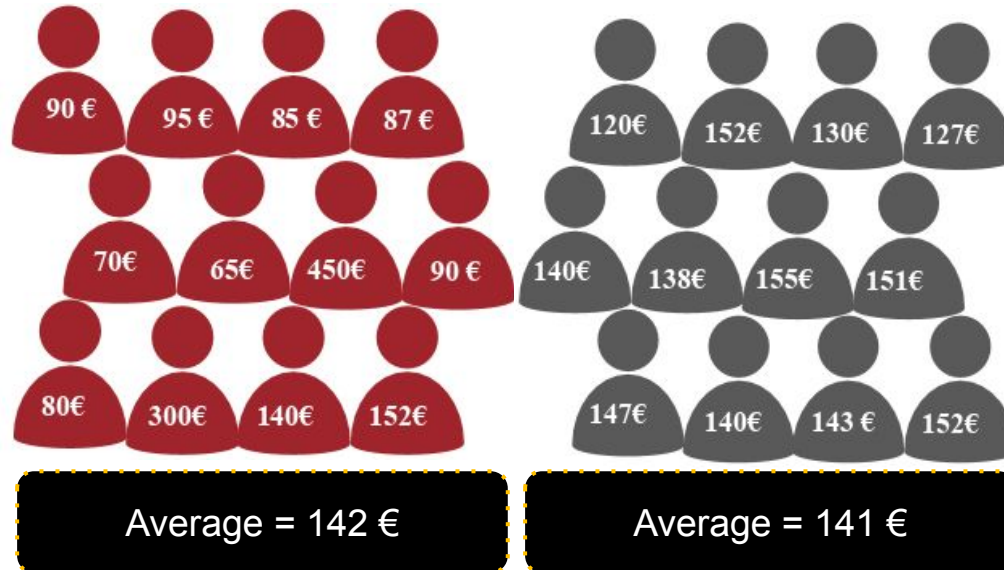
# Some questionable conclusions

Be careful with the  
arithmetical mean in  
Descriptive Analysis

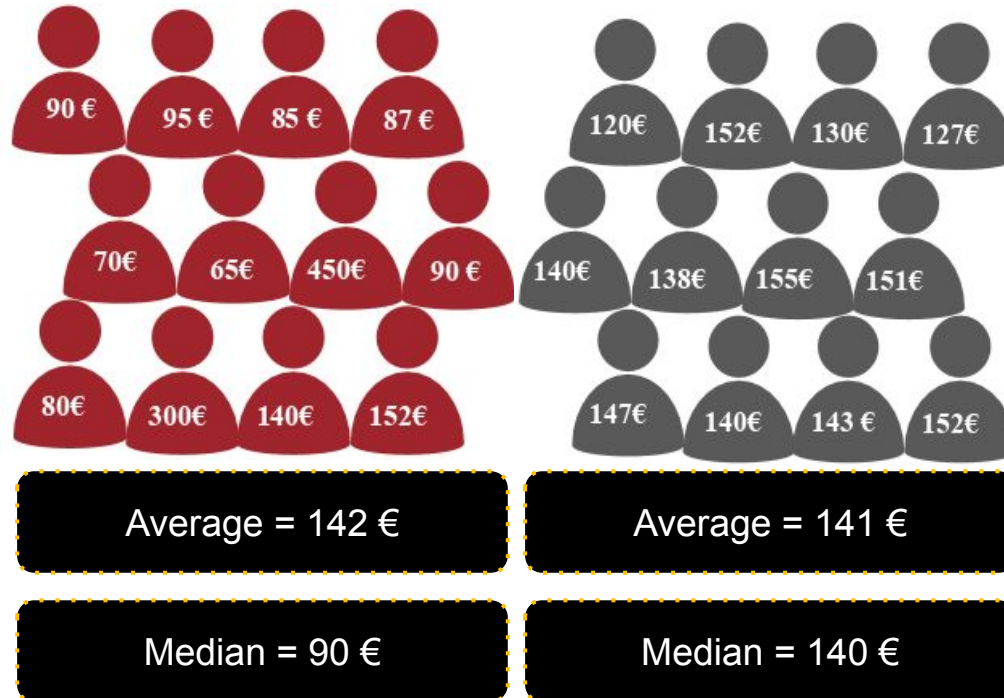
¿Which population is better?



## ¿Which population is better?



## ¿Which population is better?



## The importance of the statistical distribution of a variable

customer A

Spending A = 1.000 €  
Spending B = 1.000 €  
Spending C = 1.000 €

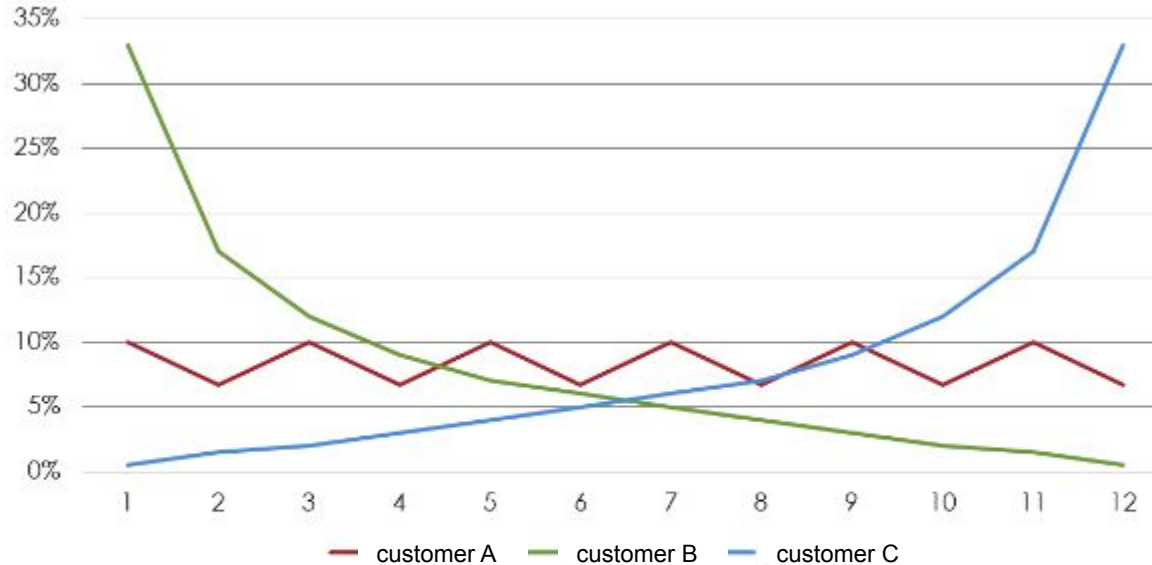
customer B

customer C

Frequency A = 12  
Frequency B = 12  
Frequency C = 12



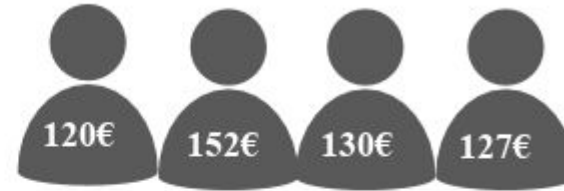
Distribution of three customers' spending along 12 months



Spending A = 1.000 €  
Spending B = 1.000 €  
Spending C = 1.000 €

Frequency A = 12  
Frequency B = 12  
Frequency C = 12

Be careful with  
comparisons



Are trully comparable?

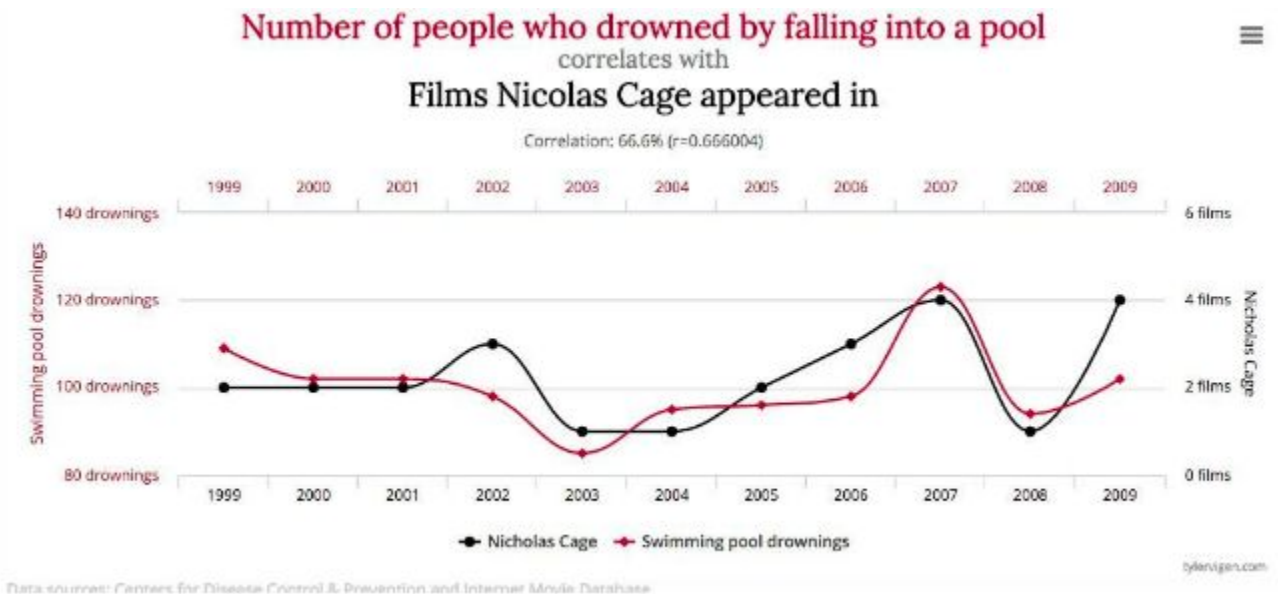
- Let's imagine that one of them are women and the other one are men.
- Let's imagine that one of them are customers and the other one are employees.

It is necessary the same starting point.  
In another case, we will have the obvious, instead of learnings.



Relationship between  
variables does not  
imply causal relation

When Nicolas Cage acts, the number of people who drown by falling into a pool increase.



# A mathematical model

A **mathematical model** is a description of a system using mathematical concepts and language. A model may help to explain a system and to study the effects of different components, and to make predictions about behaviour.

## Example of a model:

Objective

How could I measure the maximum heart rate?



## Example of a model:

Data that I  
have

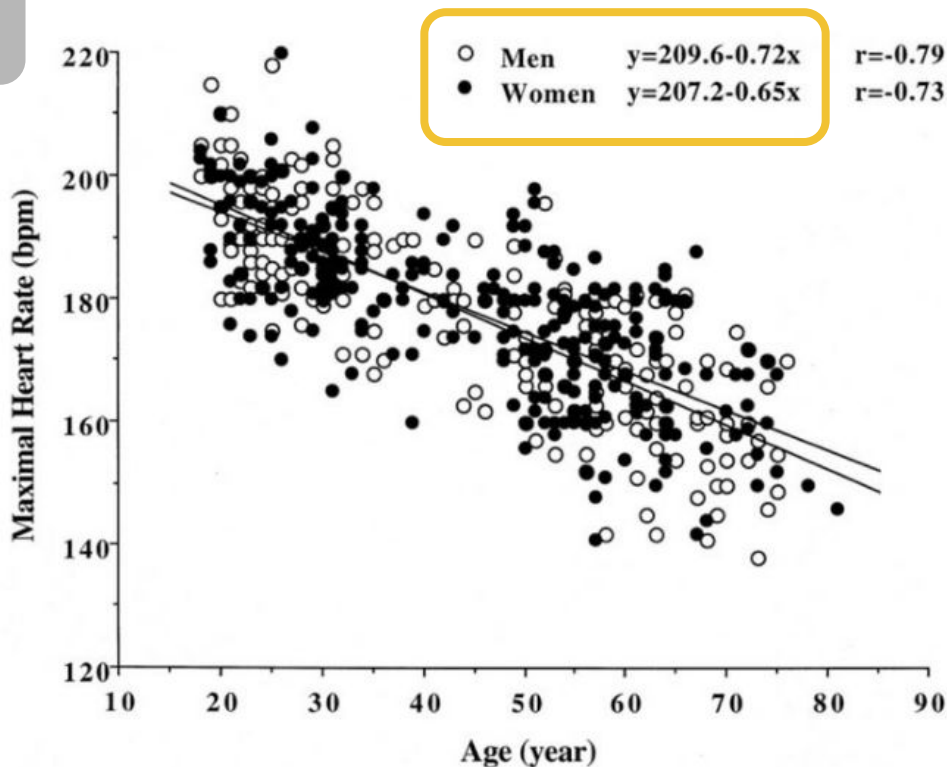
	Sexo	Edad	FCM
1	Mujer	37	185
2	Mujer	71	163
3	Hombre	25	210
4	Mujer	22	198
5	Hombre	48	179

each line  
represents  
a person

...

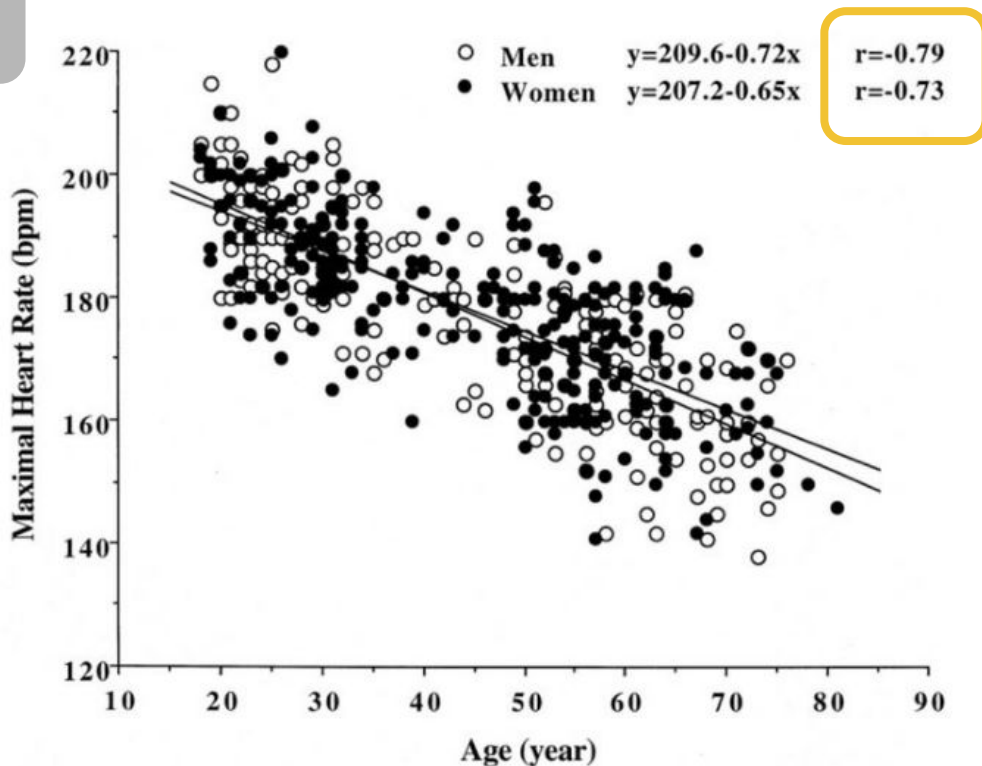
## Example of a model:

## Results



## Example of a model:

## Results

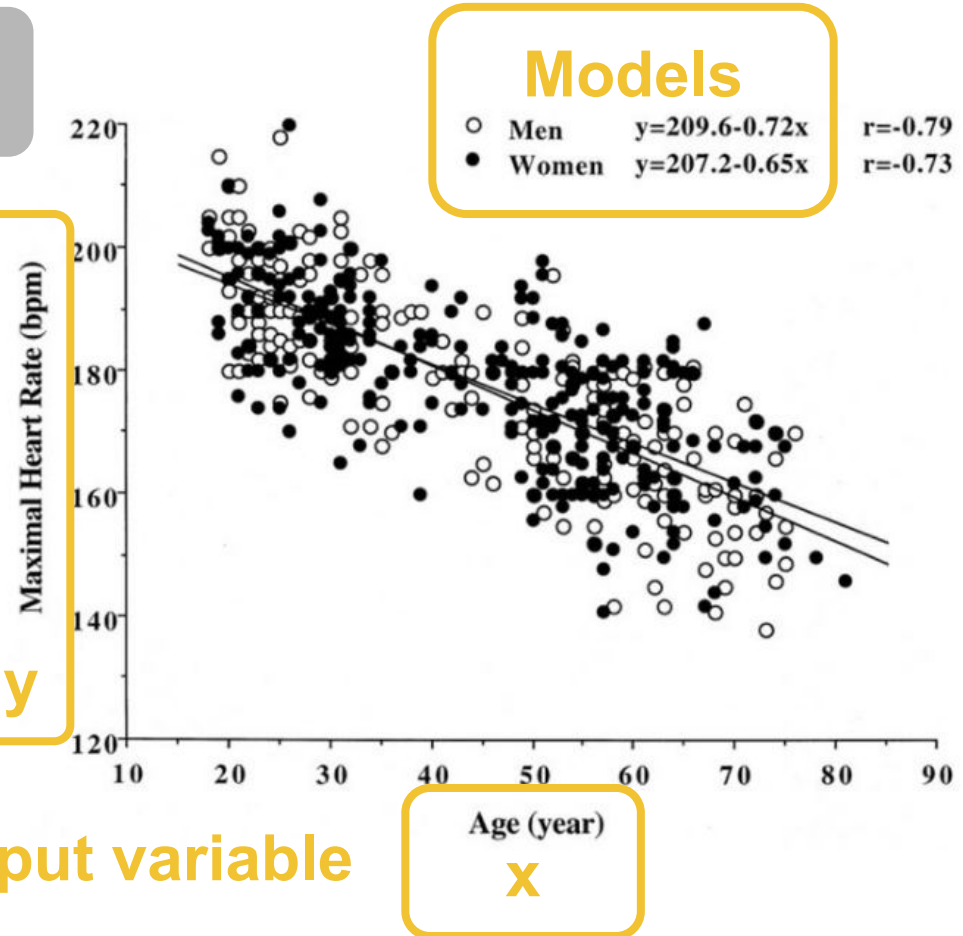


## Example of a model:

### Results

output variable **y**

input variable **x**



# ¿Machine Learning?

# Wikipedia

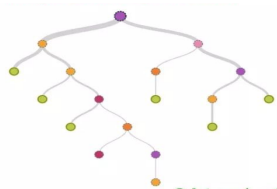


# Modelling Methodologies

## General types of learnings

### Supervised learning

Here is the data set where the right answers [labels] are given for each example. Please, produce more right answers.



### Unsupervised learning

Here is the unlabelled data. Please, find peculiarities, similarities or structures (e.g. clusters) in the data yourself.



### Reinforcement learning

Learn to do something yourself purely by maximising your expected reward.



# Modelling Methodologies

General  
types of  
learnings

## Supervised learning

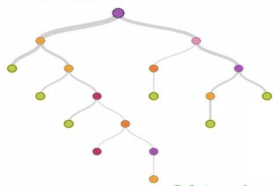
Here is the data set where the right answers [labels] are given for each example. Please, produce more right answers.

### Regression

Predict a continuous valued output.

### Classification

Predict a discrete valued output (eg. a label or a class).



## Unsupervised learning

Here is the unlabelled data. Please, find peculiarities, similarities or structures (e.g. clusters) in the data yourself.

### Clustering

Group similar examples into subsets called clusters.

### Dimensionality reduction

Maybe you don't need all the data. What is the essence of the data?



## Reinforcement learning

Learn to do something yourself purely by maximising your expected reward.

This can be a goal on its own but is often used as a pre-processing step for other ML tasks



General  
types of  
problems



## Supervised Learning

Learning from past elements

A diagram for Supervised Learning. The title 'Supervised Learning' is in large yellow font. Below it, 'Learning from past elements' is in black font. A dashed yellow arc with small square markers surrounds the text.

### Regression

Predict a number:  
eg PRICE

### Classification

Predict a label:  
eg PURCHASE / NON-PURCHASE

## Unsupervised Learning

Learning by comparison

### Clustering

Finding elements alike:  
eg CUSTOMER SEGMENTATION  
according to their habits

### Dimensionality Reduction

Explaining elements with less  
attributes

## Example: Predicting House Prices

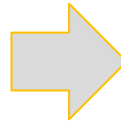
**Problem Statement** We would like to predict the price of a house according to its characteristics.

## Example: Predicting House Prices

Machine Learning requires explaining reality as a table of features.

### Data

100100011101010100110101010101010100101011101101  
01010001001111101010101010101010010101001101010101  
0010011010101010100101010110010110011011110101000  
01010101010101010101010101010101010101100110010001  
110101010011010101010101010010101110110101010001  
00111110101010101010101001010100110101010100100110  
10101010100101010110010110011011111010100001010101  
0101010101010101010101010101010110011010101010010101  
011001011001101111101010000101010101010101010101  
01010101010101100110010001110101010011010101010101  
01010010101110110101010001001111101010101010101010  
010101001101010101001001101010101010010101100101



Id	LotArea	GrLivArea	GarageArea	PoolArea	SalePrice
1	8450	196.0	1710	548	208500
2	9600	0.0	1262	460	181500
3	11250	162.0	1786	608	223500
4	9550	0.0	1717	642	140000
5	14260	350.0	2198	836	250000

## Example: Predicting House Prices

### Features Description

each line  
represents  
a house

Id	LotArea	GrLivArea	GarageArea	PoolArea	SalePrice
1	845	196.0	171	54	208500
2	960	0.0	126	46	181500
3	1125	162.0	178	60	223500
4	955	0.0	171	64	140000
5	1426	350.0	219	83	250000

target

## Example: Predicting House Prices

Machine Learning uses models to explain relationship between features.

In this case, we look for a function that explains the SalePrice:

$$\text{SalePrice} = f(\text{Features})$$

## Example: Predicting House Prices

### Linear Regression

A linear regression tries to find a linear function.

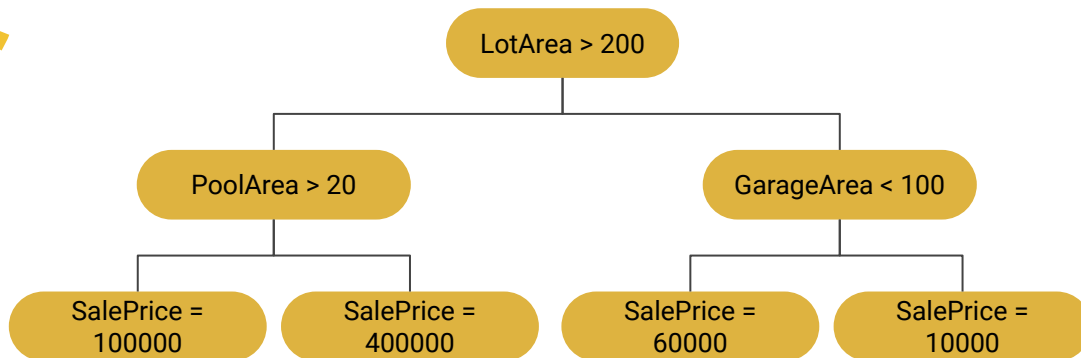
$$\text{SalePrice} = a * \text{LotArea} + b * \text{GrLivArea} + c * \text{GarageArea} + d * \text{PoolArea}$$

We need to find the best a,b,c,d to have the best relationship.

## Example: Predicting House Prices

### Decision Tree

A decision tree tries to find a path to explain the target:



We need to find the best splits for our predictions.



## Example: Predicting House Prices

The Data Scientist works to find the best parameters for the model.

In this case:

- Linear Regression:  $a, b, c$  and  $d$ .
- Decision Tree: number of splits....

How do you find the best parameters?

## Example: Predicting House Prices

There are several metrics that can be chosen for this task. For example:

- Bias: Average of the errors
- MAE: Average of the absolute values of errors
- RMSE: Square root of Average of the square of errors

$$MBE = \frac{1}{N} \sum_{i=1}^N (x_{f,i} - x_{o,i})$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |x_{f,i} - x_{o,i}|$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_{f,i} - x_{o,i})^2}$$

We chose the model parameters that provide the best METRIC.

In this case, let's use MAE:

$$\frac{|RealPrice1 - PredictedPrice1| + |RealPrice2 - PredictedPrice2| + \dots}{Number\ of\ houses}$$

## Example: Predicting House Prices

Once we have the best parameters for a type of model, we can rank the models with the best metric they had

- In this case: MAE:
  - Linear Regression: 15
  - Decision Tree: 10

And we ultimately choose the model with the best metric.

## Example: Predicting chewing gum buyers

**Problem Statement** We would like to predict the customers who buy chewing gum.

## Example: Predicting chewing gum buyers

### Features Description

each line  
represents  
a customer

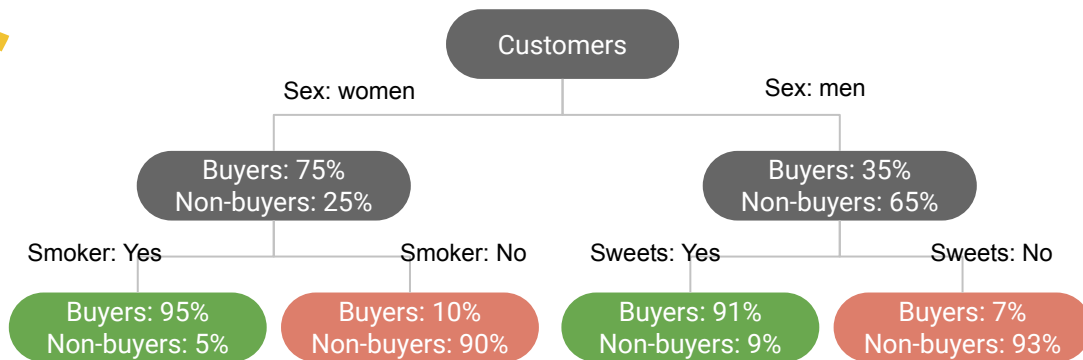
Id	Sex	Age	Smoker	Sweets	GumBuyer
1	Woman	19	Yes	No	Yes
2	Woman	45	No	No	No
3	Man	53	Yes	No	Yes
4	Woman	21	No	Yes	Yes
5	Man	35	No	Yes	No

target

## Example: Predicting chewing gum buyers

### Decision Tree

A decision tree tries to find a path to explain the target:



We need to find the best parameters for the model: number of splits, minimum number of member per group...

## Classification metrics: Confusion matrix

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive (TP)	False Negative (FN) <b>Type II Error</b>
	Negative	False Positive (FP) <b>Type I Error</b>	True Negative (TN)

## Classification metrics: Confusion matrix

		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) <b>Type II Error</b>	<b>Sensitivity</b> $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) <b>Type I Error</b>	True Negative (TN)	<b>Specificity</b> $\frac{TN}{(TN + FP)}$
		<b>Precision</b> $\frac{TP}{(TP + FP)}$	<b>Negative Predictive Value</b> $\frac{TN}{(TN + FN)}$	<b>Accuracy</b> $\frac{TP + TN}{(TP + TN + FP + FN)}$



## Example: Salmon and sea bass

### Problem Statement

I would like to discern  
between salmon and  
sea bass.



## Example: Salmon and sea bass

A fish-packing plant wants to automate the process of sorting incoming fish according to species.

As a pilot project, it is decided to try to separate sea bass from salmon using ML.

We know they are different but we don't know what features make them different.

Let's go.

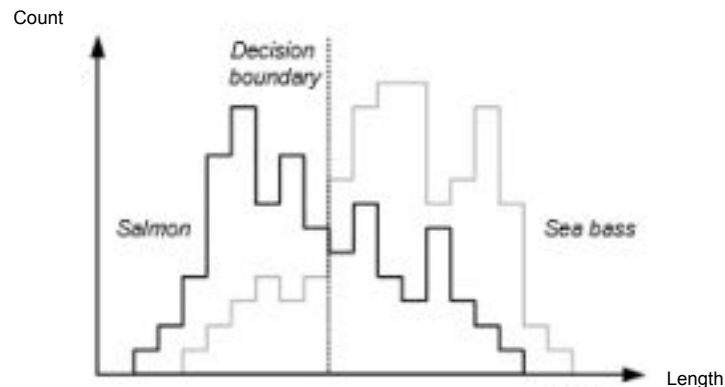
## Example: Salmon and sea bass

Some features to explore:

- Length
- Scale density
- Width
- Position of the mouth
- ...

## Example: Salmon and sea bass

Start by checking length of each fish:



They have different distributions. Notice that salmon tends to be shorter than sea bass.

## Example: Salmon and sea bass

Find the best length  $L$  threshold:

fish length  $< L$



classify as a salmon

fish length  $> L$



classify as a sea bass

After searching through all possible thresholds  $L$ , the best  $L = 9$ , and still 30% of fish is misclassified.

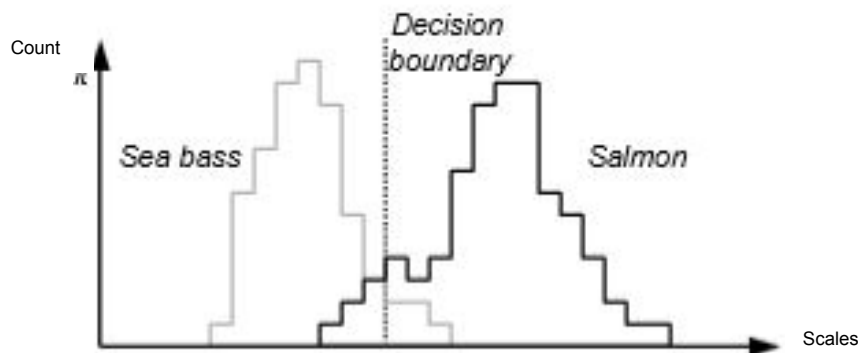
We cannot reliably separate sea bass from salmon by length alone!

Example: Salmon and sea bass

I have to continue  
looking for...

## Example: Salmon and sea bass

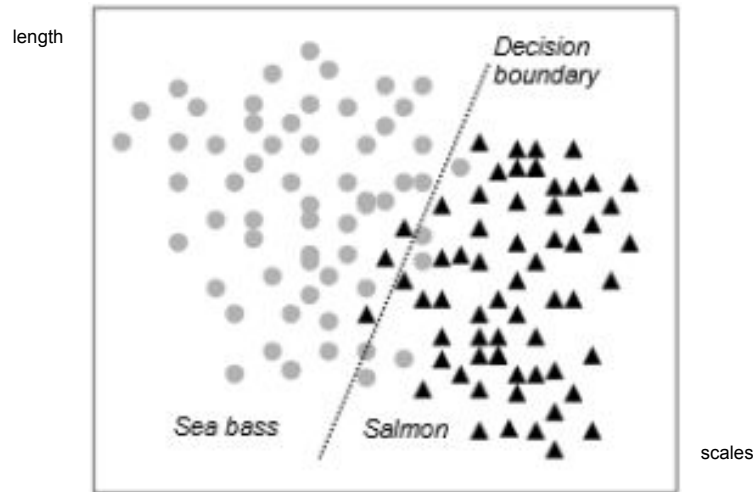
Next feature I can consider:  
Scale density



The two classes are much better separated!

## Example: Salmon and sea bass

Lets to combine the two features:



We have reduced the error rate below 5%.

Should we be satisfied with this result?



Example: Salmon and sea bass

And... what could we  
do now?

## Example: Salmon and sea bass

### Options we have:

- Consider additional features:
  - Which ones?
  - Some features may be redundant (e.g., if eye color perfectly correlates with width, then we gain no information by adding eye color as feature.)
  - It may be costly to attain more features
  - Too many features may hurt the performance
- Use a more complex model

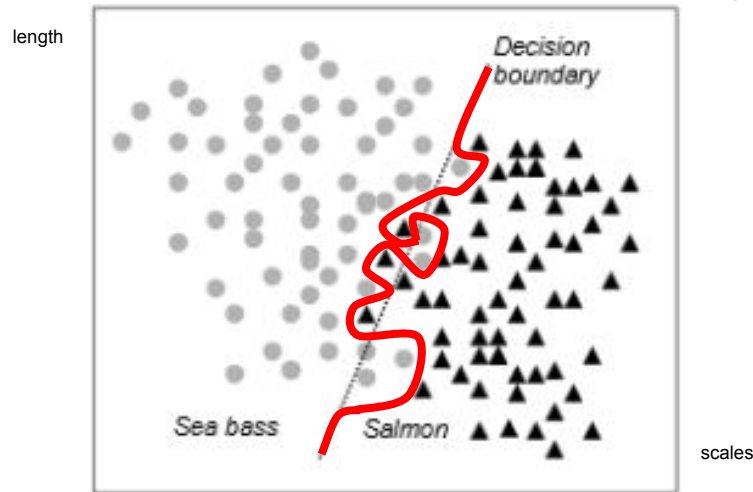
## Example: Salmon and sea bass

### Options we have:

- Consider additional features:
  - Which ones?
  - Some features may be redundant (e.g., if eye color perfectly correlates with width, then we gain no information by adding eye color as feature.)
  - It may be costly to attain more features
  - Too many features may hurt the performance
- Use a more complex model

## Example: Salmon and sea bass

With a more complex model all data are perfectly separated :



Should we be satisfied now??

## Example: Salmon and sea bass

### **We must consider:**

Which decisions will the classifier take on novel patterns, i.e. fishes not yet seen?

Will the classifier suggest the correct actions?


This is the issue of  
**GENERALIZATION**

## Measurement of the model

We take the available features table and we split in 2 sets:

1 set to calculate parameters: **Training Set**

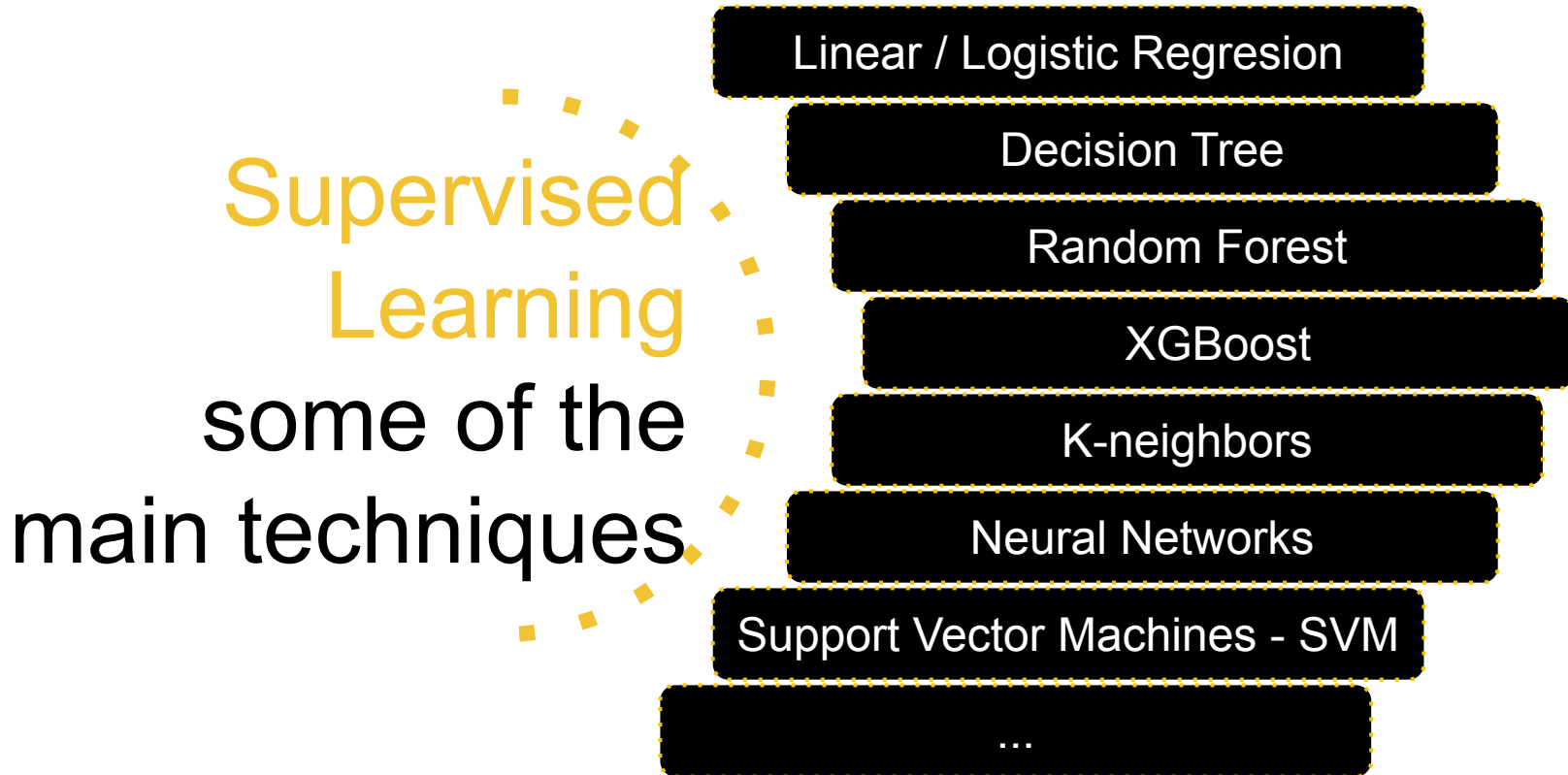
1 set to calculate metrics: **Test Set**



	Length	DensityScale	Width	...
1	10000	37	1385	...
2	12350	71	1613	...
3	9871	25	2100	...
4	15431	22	1985	...
5	5200	48	1791	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...

## Summary: ML is about features, models and metrics

- ML requires explaining reality as a table of **features**
- We then use **models** to explain relationships between features for the given ML task
- Each model is evaluated according to a **metric**



And many others... The list is growing thanks to the active research.



## Hands on Exercises

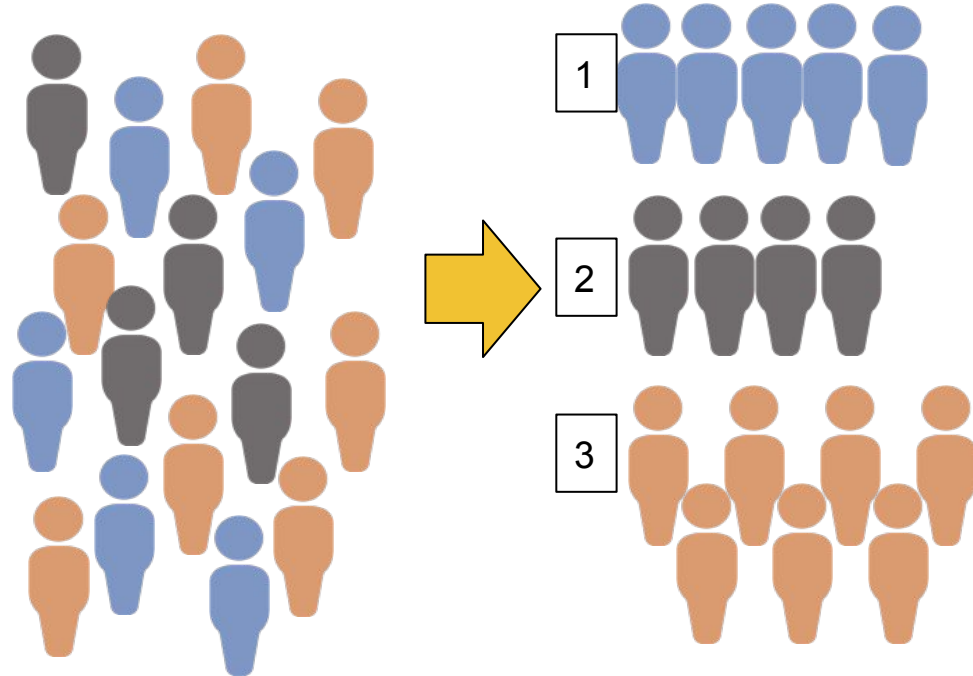


Let's try to do next two exercises:

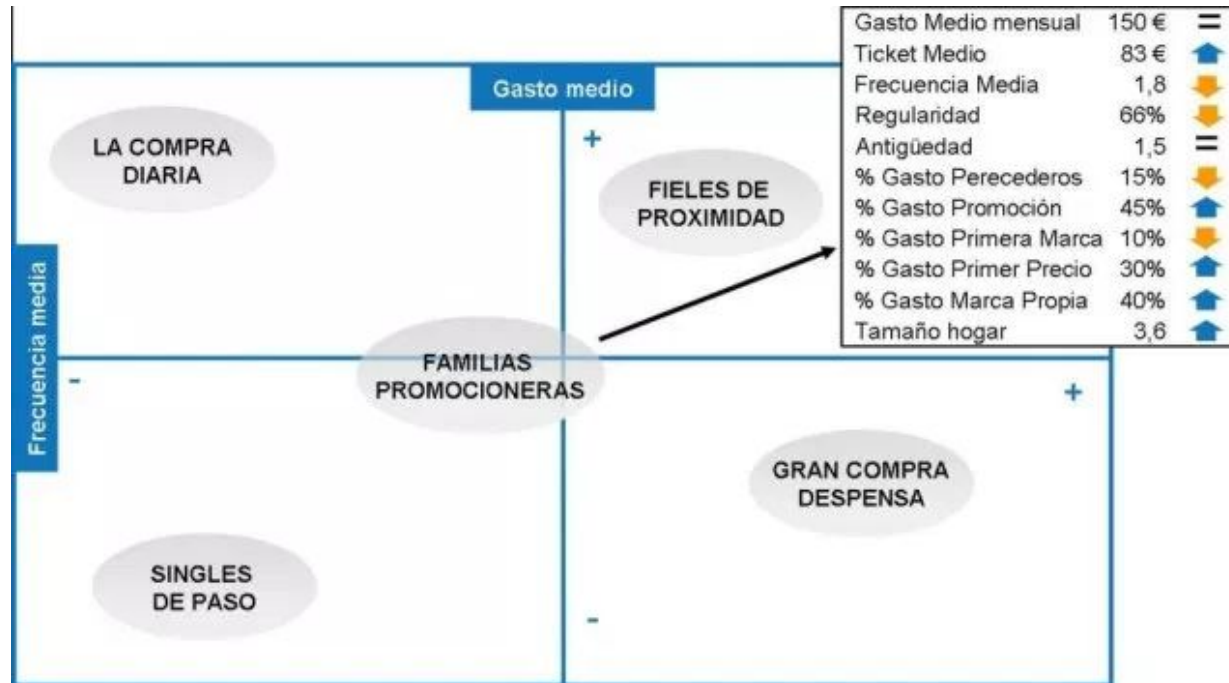
1. Exercise\_1: Compare models' predictions and decide which is the best one with different metrics.
2. Exercise\_2: Build the confusion matrix and calculate main metrics about it.
  - a. Which is the best model in general?
  - b. Which identify better the 1 class?

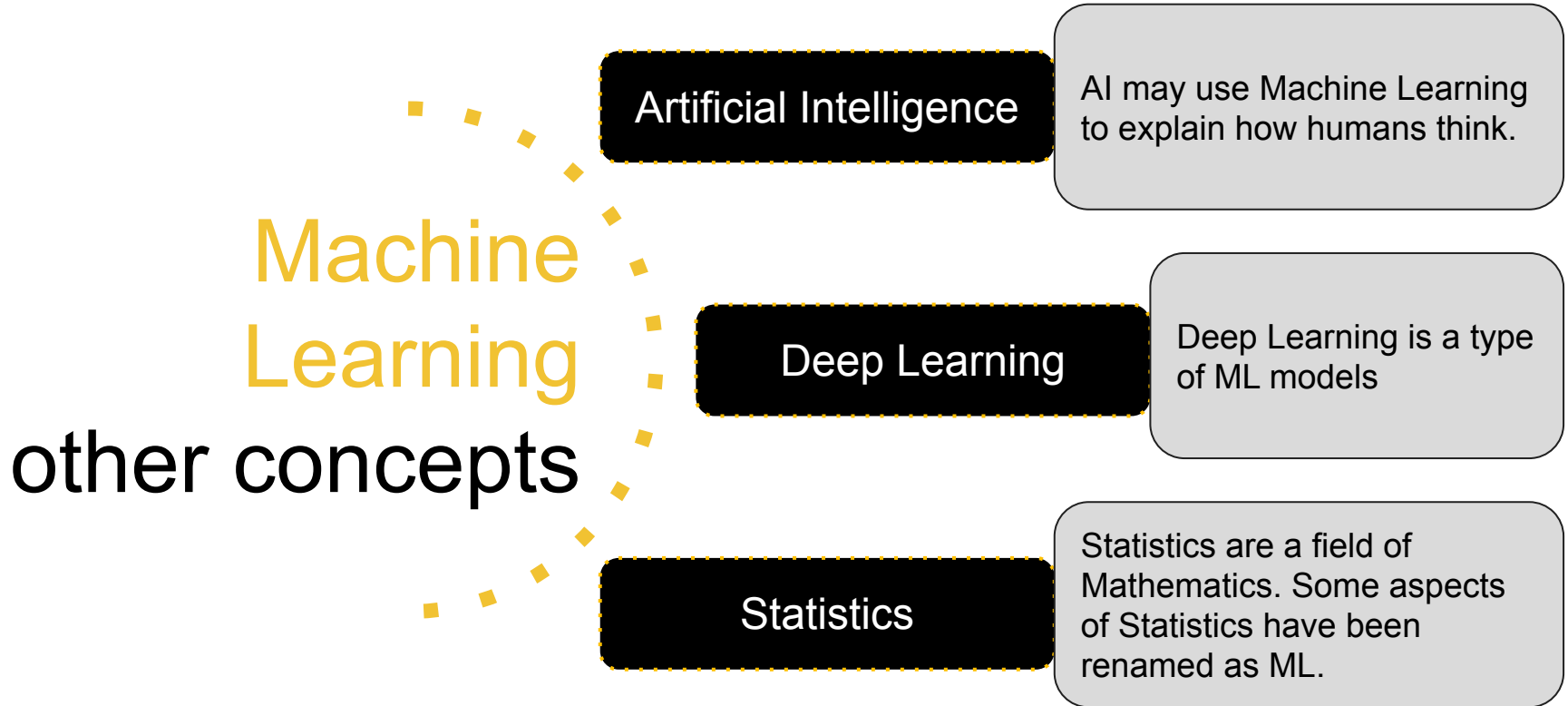
Unsupervised  
Learning  
the main type of  
problems.

## Clustering



## Example of clustering: department store





## ML Canvas

## Translating a Business Problem to ML

Any Business Problem can be analyzed through the following axes:

1. Business Value and Measure of Success.
2. Data Sources
3. ML task + metric
4. Features Extraction + Model Creation
5. Industrialization and Maintenance



# Modelling Methodologies

Opportunity:

Estimated Value:

Estimated Cost:

## Value Proposition

Business Description  
Resulting Action  
Measure of Success (KPI)  
ML Initiative Cost

## Machine Learning

ML task  
Methods of evaluation

## Data

Sources  
Frequency  
Legal  
History

Acceptable Quality  
Main Features

## Ranking of Models

List of the models and their metric

## Industrialization

Description  
Implementation Cost  
Maintenance Cost  
Model Specific Cost

By:

Iteration:

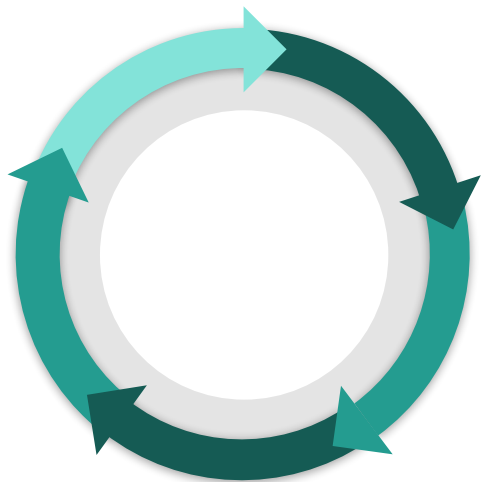
Date:

Once the first version of the canvas is created, Data Scientist perform the following actions:

- [illegible]



## The whole process is iterative



Product Managers and Data Scientists can revise their initial plans based on discoveries made during the different steps:

- Insufficient Data
- Metric not good enough for the task
- ...

# Modelling Methodologies

Opportunity: House Price Estimation

Estimated Value:

Estimated Cost:

## Value Proposition

Business Description  
Resulting Action  
Measure of Success (KPI)  
ML Initiative Cost

Business Description:  
Negotiate better house prices

Resulting Action:  
Daily report of the new houses with offered price and estimated price.

Measure of Success (KPI):  
Money saved

ML Initiative Cost:  
1 month study

## Machine Learning

ML task  
Methods of evaluation

Regression  
Method of evaluation: MAE

## Data

Sources  
Frequency  
Legal  
History

Data Sources:  
Idealista, Kaggle Data, ...

Frequency:  
Weekly

Legal:  
Public Data

History:  
1 year

Acceptable Quality  
Main Features

Acceptable Quality:  
Yes

Main Features:  
Area of house,  
Bathrooms, ...

## Ranking of Models

List of the models and their metric

Decision Tree  
MAE:

Linear Regression:  
MAE:

## Industrialization

Description  
Implementation Cost  
Maintenance Cost  
Model Specific Cost

Weekly job that retrains and produces new model.

Every morning, we run the predictions against a set of newly available houses.

Yearly evaluation to check changes

By:

Iteration:

Date:

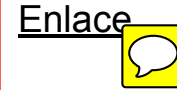
## Hands on Exercises



Let's try to do that exercise for those 4 use cases:

1. Weather Forecast
2. Movie Recommendation
3. Email Spam Filters
4. Pollution Prediction

## Some actual examples of Machine Learning



**IBM SlamTracker**



## Data science competitions:



*“Cuando eres capaz de  
ver lo sutil, es fácil  
ganar”*

*Sun Tzu*

# Thank you!

[irenetorresvalle@gmail.com](mailto:irenetorresvalle@gmail.com)