# My 1st Machine Learning model

Atelier « Maker »
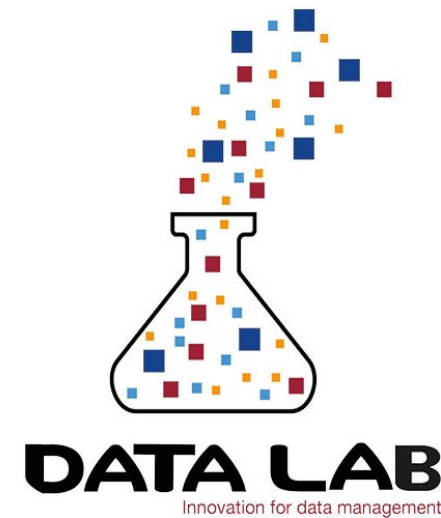
# Who we are

**Nicolas FROT**
Data Squad

**Florian BERGAMASCO**
EP/EXPLO/GTS/IGR
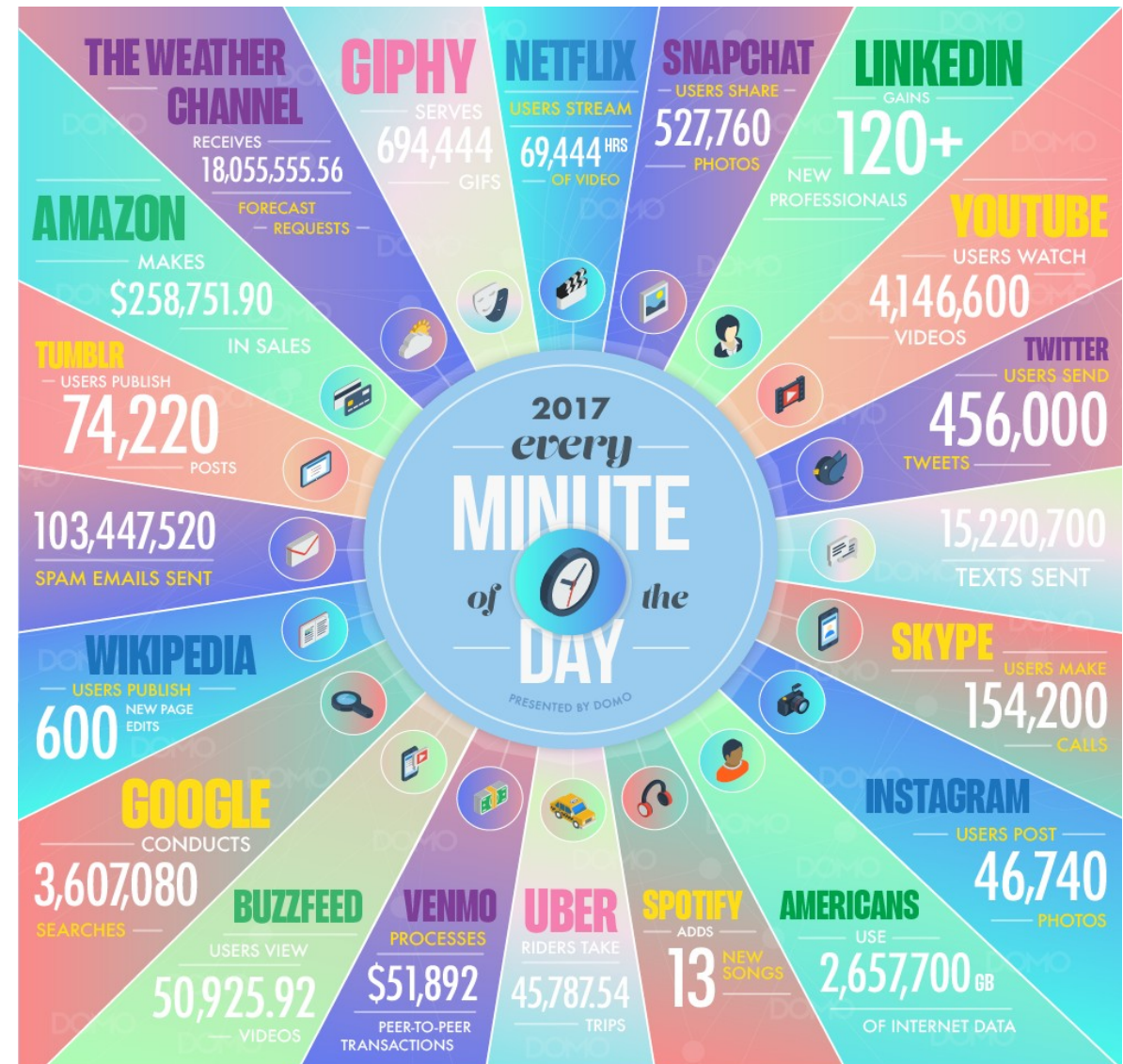
# What is Big Data?
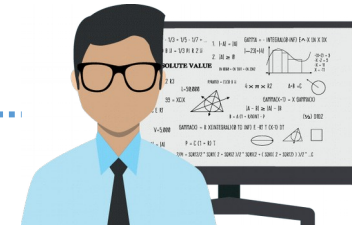


1956: 5 Mo, $50k

2018: 256 Go, $30



THE WEATHER CHANNEL RECEIVES 18,055,555.56 FORECAST REQUESTS

GIPHY SERVES 694,444 GIFS

NETFLIX USERS STREAM 69,444 HRS OF VIDEO

SNAPCHAT USERS SHARE 527,760 PHOTOS

LINKEDIN GAINS 120+ NEW PROFESSIONALS

AMAZON MAKES $258,751.90 IN SALES

TUMBLR USERS PUBLISH 74,220 POSTS

103,447,520 SPAM EMAILS SENT

WIKIPEDIA USERS PUBLISH 600 NEW PAGE EDITS

GOOGLE CONDUCTS 3,607,080 SEARCHES

2017 every MINUTE of the DAY PRESENTED BY DOMO

YOUTUBE USERS WATCH 4,146,600 VIDEOS

TWITTER USERS SEND 456,000 TWEETS

15,220,700 TEXTS SENT

SKYPE USERS MAKE 154,200 CALLS

INSTAGRAM USERS POST 46,740 PHOTOS

BUZZFEED USERS VIEW 50,925.92 VIDEOS

VENMO PROCESSES $51,892 PEER-TO-PEER TRANSACTIONS

UBER RIDERS TAKE 45,787.54 TRIPS

SPOTIFY ADDS 13 NEW SONGS

AMERICANS USE 2,657,700 GB OF INTERNET DATA

# What are the roles in a Big Data organization?

Data Scientist

Machine Learning Engineer

Data Manager

Data engineer

Data Architect

Data Analyst

# What is Machine Learning?

Machine learning algorithms build a mathematical model of sample data, known as "training data", in order to make predictions or decisions without being explicitly programmed to perform the task.

Machine learning is used in:
- Email filtering
- Image classification
- Fraud detection
- Etc…

# Which use cases we saw at TOTAL?

- ➢ Use Case 1: xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx xxxxxxxxx

- ➢ Use Case 2: xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

- ➢ Use Case 3: xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

- ➢ Use Case 4: xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx xxxxxxxxxxx

# Some basic vocabulary

Pb: predict the quantity of apples sold in a supermarket for a given day

Variable / feature

Target variable

observation

| | Temp ext (°C) | Day of week | ... | Price/ kg for apples |
|---|---|---|---|---|
| 01/01/2017 | -10 | 3 | ... | 2 |
| 02/01/2017 | -8 | 4 | ... | 2,03 |
| 03/01/2017 | 5 | 5 | ... | 2,04 |
| 04/01/2017 | 6 | 6 | ... | 2,50 |
| 05/01/2017 | 2 | 7 | ... | 2,50 |

| Apples sold (kg) |
|---|
| 34 |
| 37 |
| 67 |
| 64 |
| 33 |
| ... |
| 87 |

Dataset

# What are the branches of machine learning?

# Classification

Pb: (Spotify) Will the users buy our premium offer?

| | Nb streams per day | Seniority | Buy after trial? |
|---|---|---|---|
| User 1 | 12 | 1 | |
| User 2 | 56 | 24 | |
| User 3 | 467 | 13 | |
| ... | ... | ... | ... |
| User n | 32 | 4 | |

After the model has been fitted to the training set, let's apply the prediction on a new observation:

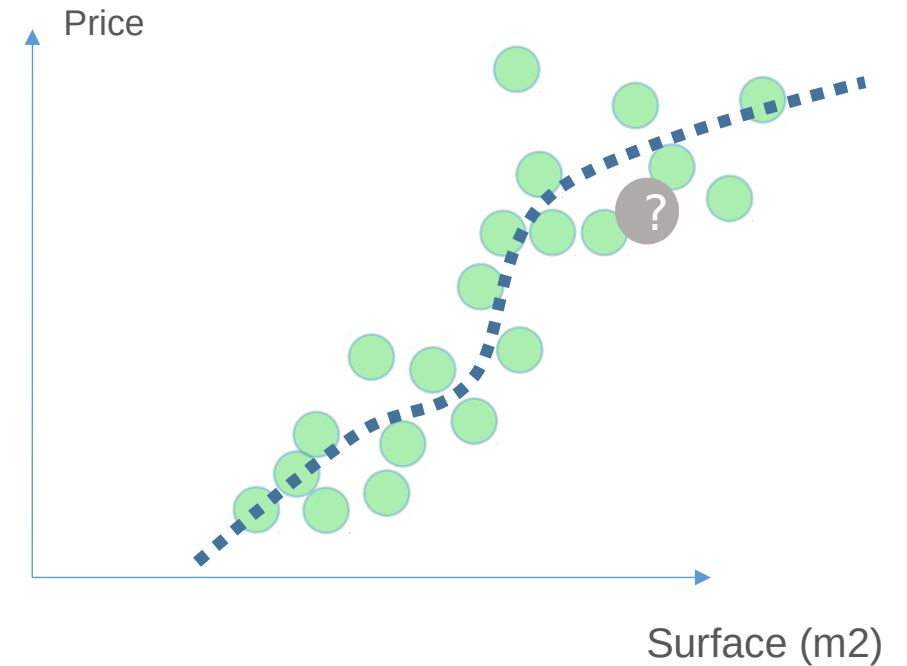| User n+1 | | | ? |
|---|---|---|---|

The model tries to find the best border that splits the positive and negative observations

# Regression

Pb: (MeilleursAgents) How to estimate the price of an appartment?

| | Surface (m2) | Price (k€) |
|---|---|---|
| **Apt 1** | 12 | 200 k |
| **Apt 2** | 56 | 450 k |
| **Apt 3** | 130 | 1200 k |
| **...** | ... | ... |
| **Apt n** | 32 | 300 k |

After the model has been fitted to the training set, let's apply the prediction on a new observation:

| | | ? |
|---|---|---|
| **Apt n+1** | | |

Price

Surface (m2)

# Clustering

Pb: (Netflix) Can I group similar users by behaviour on the app?

|  | Nb movies (/month) | Nb connecti ons |
|---|---|---|
| **User 1** | 12 | 1 |
| **User 2** | 32 | 24 |
| **User 3** | 46 | 13 |
| **…** | … | … |
| **User n** | 32 | 44 |

After the model has been fitted to the training set, let's apply the prediction on a new observation:

| User n+1 |  |  | ? |
|---|---|---|---|

Nb connections

Cluster 1

Cluster 2

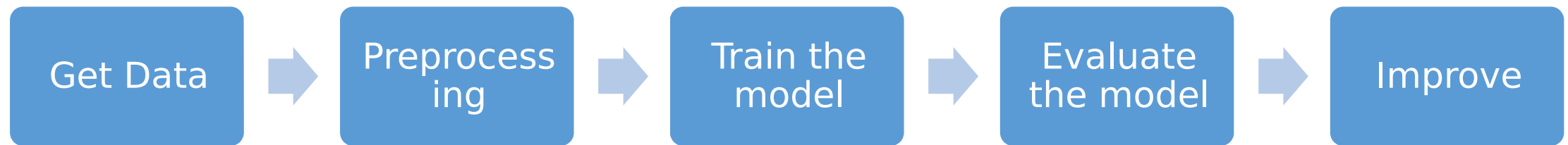Cluster 3

?

Nb movies

The model tries to find the best border that splits the positive and negative observations

# Workflow of a ML project

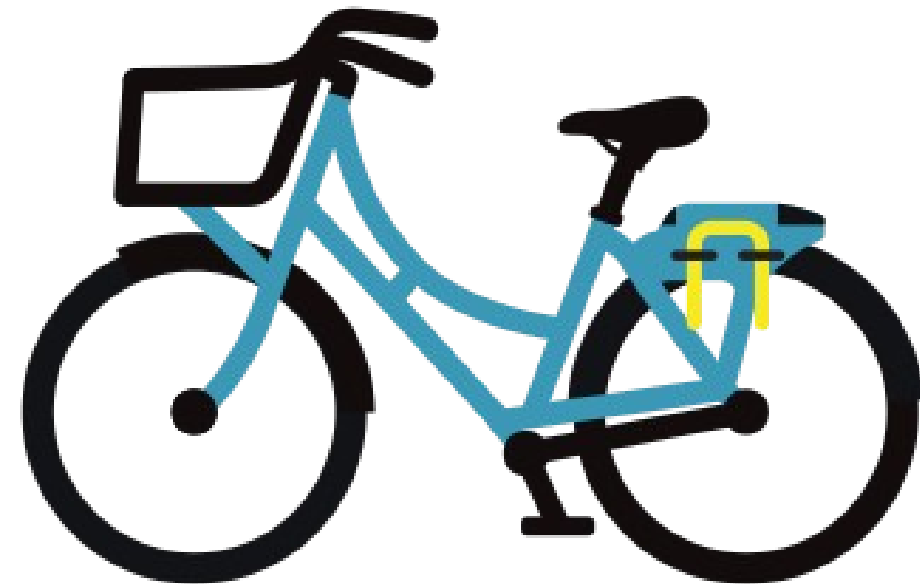Get Data → Preprocessing → Train the model → Evaluate the model → Improve

# What are we going to do today?

## The Challenge

**Predict the number of shared bikes rented every hour in San Diego given meteorological information**

➢ **Features**:
  ➢ Temperature
  ➢ Time
  ➢ Humidity
  ➢ Wheather
  ➢ Weekday
  ➢ Is_holiday
  ➢ Etc…

➢ **Data**: records from 2016/2017

➢ **Tools**: using Python (via Jupyter Notebook)

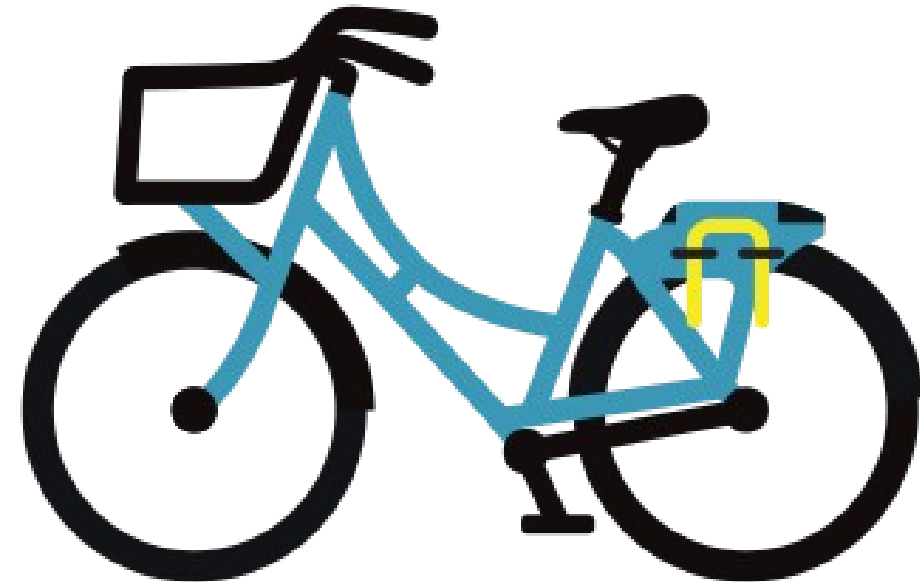# What are we going to do today?

The steps

Theory & Hands on:

➢ **Step 0:**

Introduction to Python and Jupyter

➢ **Step 1:**

Build a first basic model

➢ **Step 2:**

Improve your model: preprocessing

➢ **Step 3:**

Improve your model: model choice and optimise the hyperparameters

# Step 0: Introduction to Python and Jupyter notebooks

# Step 1: Build a 1st basic model

| Get Data | → | Preprocessing | → | Train the model | → | Evaluate the model | → | Improve |

- Load data

Do the minimum:

- Remove rows with NA values

- Dummify categorical values

- Split Dataframe in train and test

Do the minimum:

- Choose 1 model only

Do the minimum:

- Choose an evaluation metric

- Compute the score of your model

# Step 1: Build a 1st basic model

Split the Dataset in order to evaluate your model

**Dataset**

features     Target

TRAIN

TEST

A

The model trains

Make predictions
with the test set

How to compute the
performance of a model ?

710

POOR     GOOD

**SCORE**

Mae, rmse, etc…

# Step 1: Build a 1st basic model

# Step 2: Improve your model : Preprocessing

| Get Data | → | Preprocessing | → | Train the model | → | Evaluate the model | → | Improve |

- Load data

- Dummify categorical values
- Split Dataframe in train and test
- Impute missing values
- Add some feature engineering

Do the minimum:

- Choose 1 model only

Do the minimum:

- Choose an evaluation metric
- Compute the score of your model
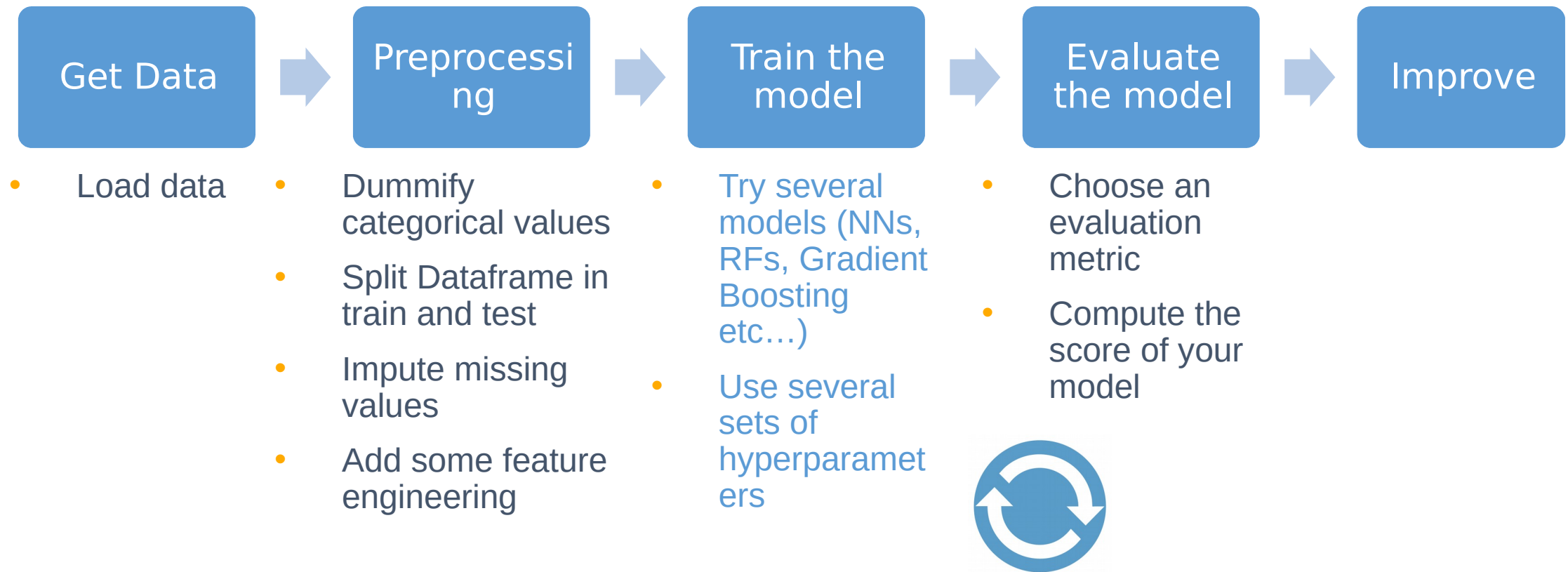
# Step 2: Improve your model : Preprocessing

# Step 3: Improve your model : Models and hyperparameters optimisation

**Get Data** → **Preprocessing** → **Train the model** → **Evaluate the model** → **Improve**

- Load data

- Dummify categorical values
- Split Dataframe in train and test
- Impute missing values
- Add some feature engineering

- Try several models (NNs, RFs, Gradient Boosting etc…)
- Use several sets of hyperparameters

- Choose an evaluation metric
- Compute the score of your model

# Step 3: Improve your model : Models and hyperparameters optimisation