*Contact details:*

*Yannick Léo*

*Associate Partner & Data Science Director, Emerton Data*

[yannick.leo@emerton-data.com](mailto:yannick.leo@emerton-data.com)

EMERTON DATA

# Data Science for Business: project exercises and instructions

## To deliver before March 12

# Data Science for Business: dates and timeline

| Dates | Timeline |
|-------|----------|
| Tuesday, 3th January, 17:15-20:30 | Sharing of the regression exercise |
| Thursday, 5th January, 17:15-20:30 | Sharing of the regression exercise |
| Thursday, 26th January, 17:15-20:30 | Sharing of the shap extension exercise |
| Sunday, 12th March, 23:59 | Projects delivery including 1. regression with 2. shap extension and 3. clustering |
| Sunday, 19th March, 23:59 | Feedback regarding the projects |
| Thursday, 24th March, 10:00-12:00 | Exam |

# Exercises and instructions

## Exercises

1. Build ML model for regression (auto insurance problem)

2. Perform an interpretability study of the built model in the previous exercise

3. Implement the clustering exercises

## Delivery instructions

1. Create a google colab with Python environment

2. Access to dataset (see next page)

3. Develop your methodology on a notebook

4. Export colab notebook in .ipynb format

5. The complete restitution is the March 12 and should contain all the exercices

## Important to get all points

1. Follow the steps presented in the lessons

2. A cell not ran is to considered

3. Work to get a good result (high mape)

4. Easy to read notebook with few comments and logical steps

5. Follow the delivery instructions

# Agenda

# Exercise: build full ML methodology to predict Auto Insurance pricing

## Problem statement

In the context of price optimization, we need to inject in the demand model a variable that reflects competitiveness. Fortunately we have at disposal a market index price on a set of quotes.  We don't know the details on this index (it could be the best price offer, or any kinds of averaged price offer), but we know it is useful to segment elasticity.

# Presentation of the dataset and information regarding access

→ Id: categorical, quotes identifier number

→ Y: response numerical variable, 'Response_Market_Index'

→ X: 38 variables (17 are categorical, 21 are numerical), they are related to:
  → Behavioral information (BEH prefix): 3 variables
  → Claims information (CLA prefix): 5 variables
  → Geocoded information (GEO): 2 variables
  → Personal information (PER): 13 variables
  → Product information (PRO): 1 variable
  → Vehicle information (VEH): 14 variables

**FILES DESCRIPTION AND ACCESS**

→ Files in the folder data_market_prediction:
  → data_train_insurance.csv: (; separated) dataset to be used, contains Id,Y,X
  → data_type.csv: contains the list of variables and related
→ Access to files:
  → https://drive.google.com/drive/folders/19g17pqo2DWT1lcOdkwX3lSgj-tyATWNx?usp=sharing

# Description of metrics for assessment

**CRITERIA**

→ Mean Absolute Percentage Error – MAPE
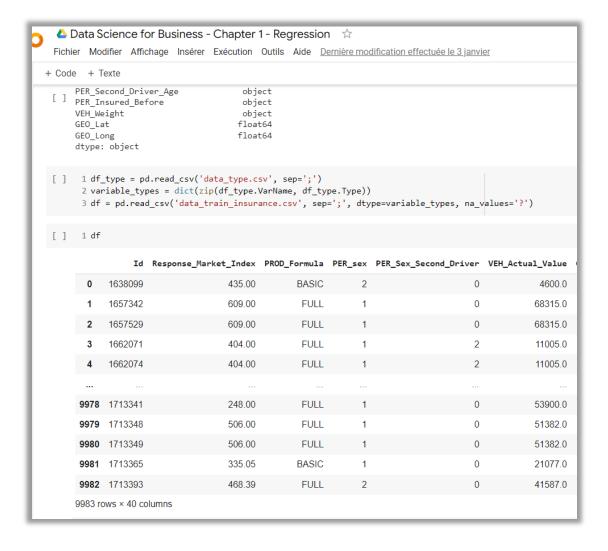
$$MAPE = \frac{1}{N}\sum_{i=1}^{N}\frac{|y_i - \widehat{y_i}|}{y_i}$$

→ Inverse percentile of APE at 20% level

$$IP_{20\%} = card\left\{i \in 1,\dots,N: \frac{|y_i - \widehat{y_i}|}{y_i} < 20\%\right\}/N$$

**EVALUATION METRICS**

→ Mean Absolute Percentage Error – MAPE

$$MAPE = \frac{1}{N}\sum_{i=1}^{N}\frac{|y_i - \widehat{y_i}|}{y_i}$$

→ Inverse percentile of APE at 20% level

$$IP_{20\%} = card\left\{i \in 1,\dots,N: \frac{|y_i - \widehat{y_i}|}{y_i} < 20\%\right\}/N$$

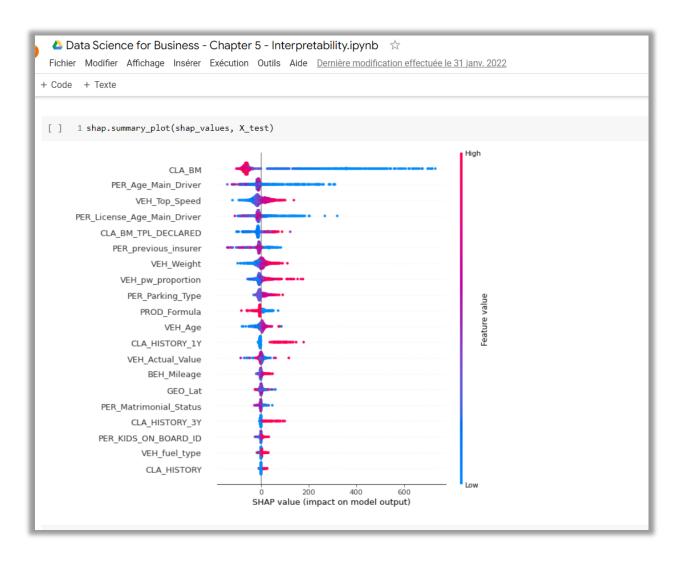# Tips to open the files in the regression notebook

# Agenda

① Regression exercise

② **Extension with the interpretability**

③ Clustering exercise

# Exercise: perform an interpretability study should be implemented during the regression problem as an extension
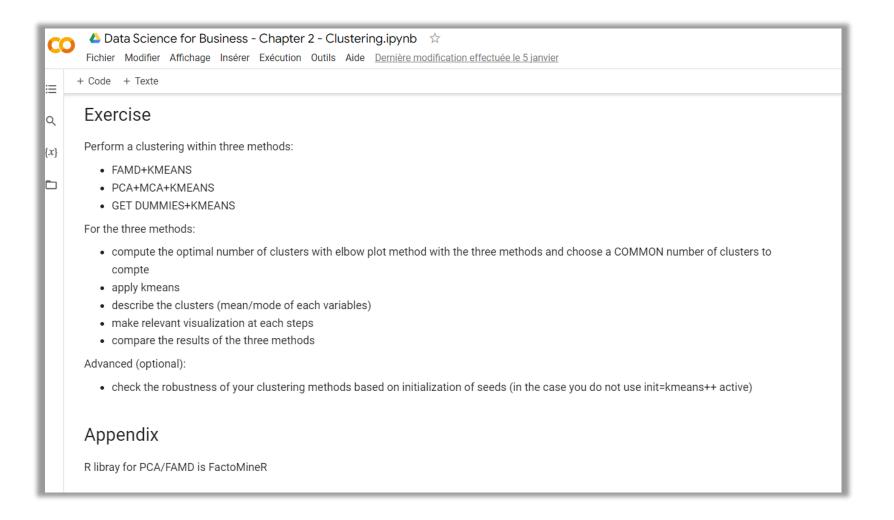
# Agenda

① Regression exercise

② Extension with the interpretability

③ **Clustering exercise**

# Exercise: build three clustering methodologies on credit bank data

**Problem statement:**

➡ Perform a clustering within three methods:
- Get dummies+KMEANS
- FAMD+KMEANS
- PCA+MCA+KMEANS

➡ For the three methods:
- Compute the optimal number of clusters with elbow plot method with the three methods and choose a COMMON number of clusters to compute the cluster
- Apply kmeans other the three methods
- Describe the clusters (mean/mode of each variables)
- Make relevant visualization at each steps
- Compare the results of the three methods

➡ Advanced (optional):
- Check the robustness of your clustering methods based on initialization of seeds (in the case you do not use init=kmeans++ active)

# Exercise on clustering is positioned at the end of the clustering notebook

# EMERTON DATA