# CREDIT RISK CLASSIFICATION

PROJET APPLICATION OF BIG DATA

**Réalisé par :**

- **AMOIN Famien Le Roi (DE1)**
- **NGUEYAP NGUIKAM Godelieve (DE2)**

## INTRODUCTION

The project we have developed aims to predict the probability of loan repayment for a lending company. To do this, we used the XGBoost model and implemented a web interface using Flask to allow users to make real-time predictions.

## REALISATION

### PRELIMINARY

GIT for team collaboration, code & model versioning

- ● We use a template cookie cutter.
- ● We use a conda environment t for all our libraries.
- ● Use a documentation library Sphinx, then we updated it. it is located in the directory

*Credit_Risk_Classification\docs\build\html\* launch *index.html*

### DATA

The data used for this project comes from the Kaggle's lending app dataset. The dataset is divided into two parts: a training part and a test part. The training part includes more than 300,000 observations and the test part includes more than 100,000 observations.

### DATA PREPARATION

Before training the XGBoost model, we prepared the data using the data preparation function. This function performed the following tasks:

- Deleting rows with missing values
- Replace missing values with mean value.
- Encoding of categorical variables

## FEATURE ENGINEERING:

After preparing the data, we created new features using the Feature Engineering function. This function made it possible to create new variables which improved the performance of the XGBoost model and select the most important variables thanks to One HOT Encoding

## XGBOOST MODEL

The XGBoost model used for this project was trained using the model training feature. The model was trained using the default hyperparameters such as learning rate, regularization, base value, etc. Hyperparameters can be modified by the user through the web interface.

## MODEL EVALUATION

The model was evaluated using the following metrics: precision, recall, F1 score, and accuracy score. The output of prediction is located in *credit_Risk_Classification\output\prediction\prediction.csv*

## MLFLOW PROJECT

- MLProject file is located in *credit_Risk_Classification\MLproject*
- MLflow outputs is located in *credit_Risk_Classification\mlruns\0*
- MLflow ui

## REST SERVER IMPLEMENTATION

We implemented a REST Server using Flask that allows users to make real-time predictions by providing the training and test data paths along with the desired hyperparameters.

- Form



- route for prediction

## SHAP LIBRARY FOR EXPLANATIONS

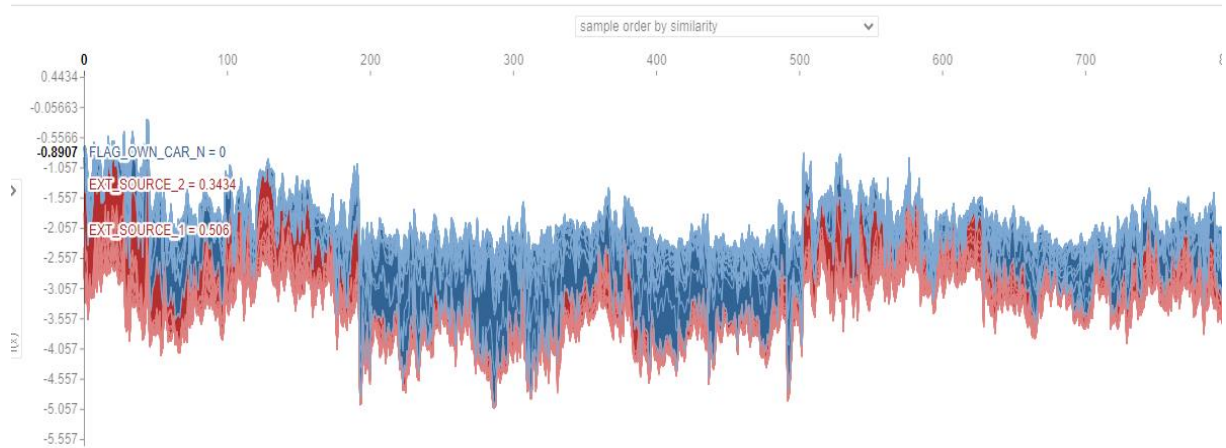SHAP (SHapley Additive exPlanations) is a game theoretic approach to explain the output of any machine learning model. the explanations are in the repository:  *Credit_Risk_Classification\output\figures*

● Visualize explanations for a specific point of your data set: ***explanationSpecificPoint.html*** ( the image is saved in .html because it is interactive)
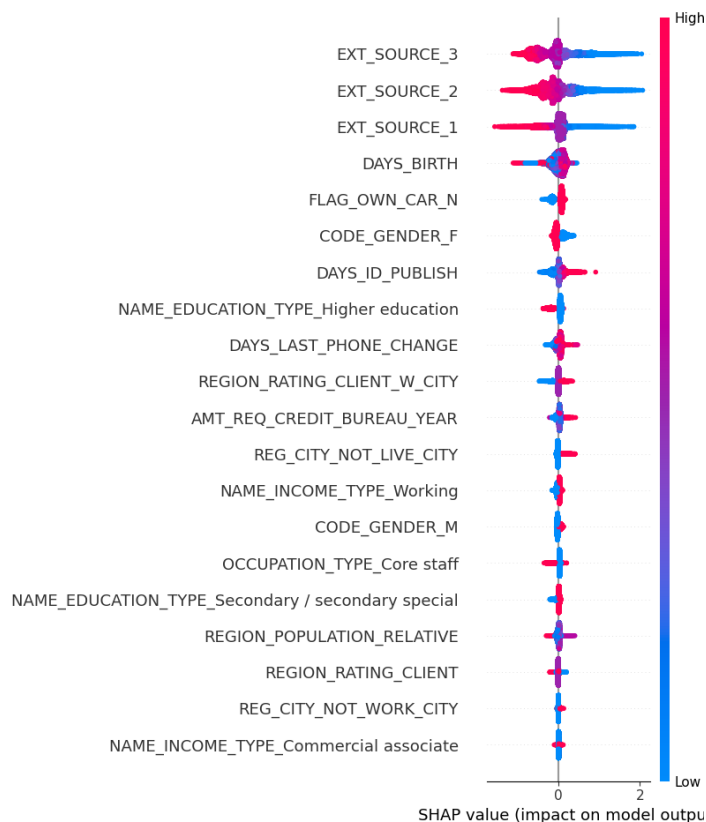


● Visualize explanations for all points of your data set at once: ***explanationAllPoint.html*** (

the image is saved in .html because it is interactive)



● Visualize a summary plot for each class on the whole dataset: **SummaryEachClass.png**



## CONCLUSION

In conclusion, we have successfully built a machine learning project using MLflow and a Flask REST API. We were able to track experiments and log parameters, code versions, and performance metrics, making it easier to reproduce our results. The Flask REST API allowed us to deploy our model and make predictions on new data. And SHAP to explain the prediction. Overall, this project highlights the importance of utilizing tools such as MLflow to streamline and simplify the machine learning development process.