

# **Project of Application of Big Data**

**version**

**Vianney Gonnot, Louis Gailhac, Elodie Gueuret**

novembre 23, 2021



# Contents

<b>Welcome to Application of big data's documentation!</b>	<b>1</b>
Part 1 :	1
Part 2 :	2
Part 3 :	2
<b>Indices and tables</b>	<b>3</b>



# Welcome to Application of big data's documentation!

**Application of bigdata** (/Our project/) is a python project, that train us to apply tools and concepts seen in course. It pulls data from the *DataSet ofHome Credit Risk Classification* <<https://www.kaggle.com/c/home-credit-default-risk/overview>>.

To run our program correctly, you will need to run the different python scripts in the following order :

1. Cleaning\_Dataset.ipynb
2. Features\_Engineering.ipynb
3. Training\_model.ipynb
4. Shap.ipynb

## Part 1 :

In the first part, we build a machine learning project using jupyter notebook, github, a conda environnement and sphinx. We tried to separate the different workflow into different scripts, one for the data preparation, one for the data preparation, one for the feature engineering, one for the models training and a last one for the prediction.

### Data preparation :

We first clean the dataset from all the NAN values.

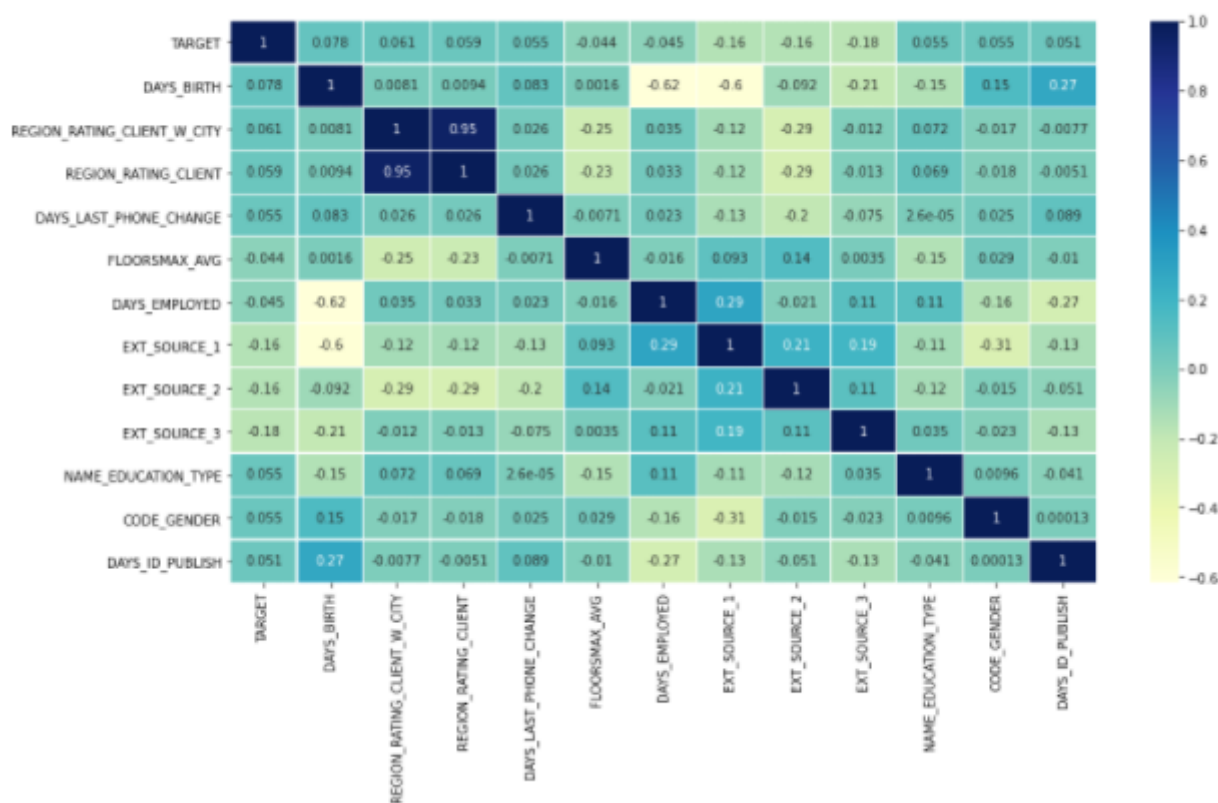
- `init()`, will return the cleaned dataset

### Feature engineering :

We have done a correlation matrix, and from that we have kept the most correlated features and deleted the least correlated ones.

Here is the correlation matrix :

Correlation matrix :



- `matrice_corr(df_train,df_test)`, is a void function that show us the correlation matrix

## Part 2 :

- **setup\_train(df\_train,df\_test)**, will return four values (X\_train, X\_test, y\_train and y\_test)

### Models training and predict :

We had to train three models: XGboost, Random Forest and Gradient Boosting. The XGboost model, is done with the optimized distributed gradient boosting library, XGboost. The Random Forest model, consists of many decision trees. The Gradient Boosting model, is an ensemble of weak prediction models(decision trees).

- **XGBC\_model(X\_train,X\_test,y\_train,y\_test,learning\_rate,max\_depth,scale\_pos\_weight)**, The XGBOOST model is a supervised learning algorithm whose principle is to combine the results of a set of models. The idea is simple: instead of using a single model, the algorithm will use several which will then be combined to obtain a single result.
- **RF\_model(X\_train,X\_test,y\_train,y\_test)**, The random forest algorithm performs parallel learning on multiple randomly constructed decision trees trained on different subsets of data.
- **GB\_model(X\_train,X\_test,y\_train,y\_test)**, Gradient boosting is a machine learning technique used in regression and classification tasks, among others. It gives a prediction model in the form of an ensemble of weak prediction models, which are typically decision trees. When a decision tree is the weak learner, the resulting algorithm is called gradient-boosted trees; it usually outperforms random forest.

Those functions train the different models.

All three model, succeed in predicting if a client could get a loan. Most had each around 0.91 of accuracy.

RF accuracy score:

0.9173535353535354

GB accuracy score:

0.9185959595959596

XGBOOST accuracy score:

0.9186363636363636

## Part 2 :

In this part, we got introduced to MLFLOW. We decided to track the parameters of the XGboost model. It helped us to choose the best parameter, to have better result, with our model.

Here we can have a look at MLFlow:

Run ID	Run Name	Location	Driver	Status	Version	Model URI	Accuracy	Learning rate	Max depth	Scale pos weight
1	10 minutes ago	23.5%	-	start	1	0.000000	0.011	0.7	35	0.48
2	10 minutes ago	6.7%	-	start	2	0.000000	0.010	0.4	14	0.15
3	10 minutes ago	3.7%	-	start	3	0.000000	0.010	0.6	2	0.26
4	10 minutes ago	46.7%	-	start	4	0.000000	0.010	0.2	82	0.12
5	10 minutes ago	12.3%	-	start	5	0.000000	0.014	0.8	15	0.48
6	10 minutes ago	14.8%	-	start	6	0.000000	0.010	0.25	20	0.7
7	10 minutes ago	27.8%	-	start	7	0.000000	0.010	0.2	45	0.2
8	10 minutes ago	33.3%	-	start	8	0.000000	0.008	0.8	50	0.5
9	10 minutes ago	19.2%	-	start	9	0.000000	0.010	0.3	25	0.7
10	10 minutes ago	11.8%	-	start	10	0.000000	0.010	0.1	15	0.8
11	10 minutes ago	16.5%	-	start	11	0.000000	0.010	0.1	20	0.1
12	10 minutes ago	15.2%	-	start	12	0.000000	0.010	0.1	35	0.1
13	10 hours ago	15.0%	-	start	13	0.000000	0.010	0.1	35	0.1

## Part 3 :

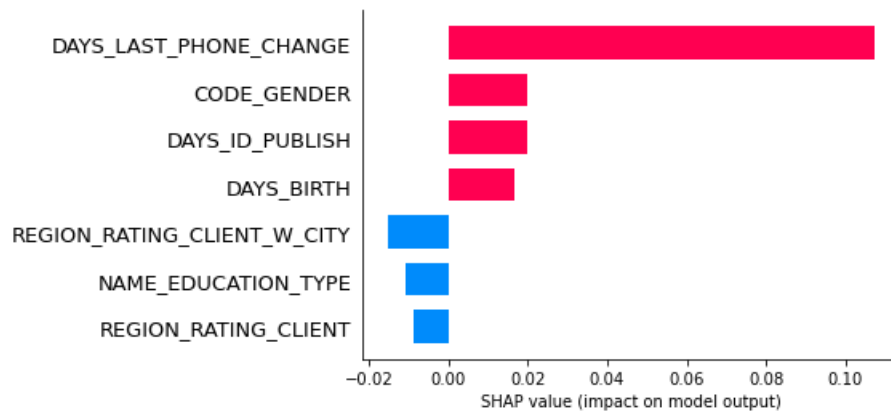
Finally, we used SHAP Library on our XGboost model to understand it. We can visualize three graph by running the following function.

The first one is to visualize the explanations for a specific value, we choose to select the 100th value, and we observed that the day of the last time the person changed his phone had a lot of influence on the result of the prediction.

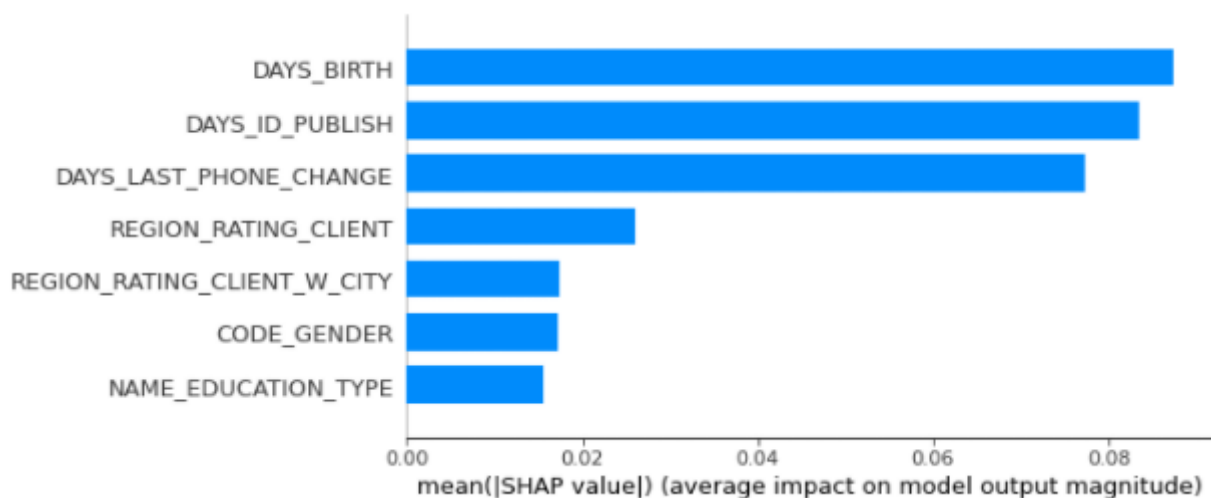
The second, is a the same as the one just seen but for all the values of the dataset. We can see the day of birth is the most influent between all the features.

The last one, is a summary plot for each class of the dataset.

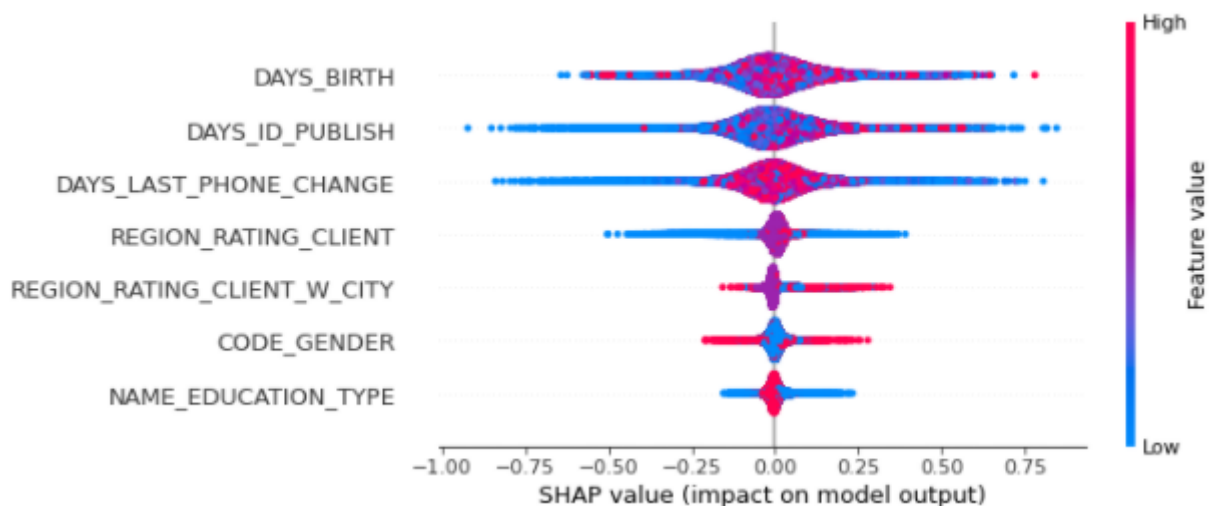
Here is the graph for a specific value :



Here is the graph for all values :



Here a summary plot for each class on the whole dataset :



- `get_explainer(xg_clf,X_train_test)`, is a void function that print the three graphs

## Indices and tables

- `search`