

Prévision des défauts sur les lignes de production Valéo

Index:

1. Contexte de l'étude
2. Description du dataset case
 - a. Entrées
 - b. Sortie
 - c. Croisement Entrées/Sortie
3. Import des packages Python et rechargement automatique des packages du projet
4. Chargement des données 'Training'
5. Exploration tabulaire des données:
 - a. Visualisation tabulaire des données - Affichage du type 'head()'
 - b. Rapport sémantique des données - Affichage du type 'info()'
 - c. Données manquantes par type de 'feature'
 - d. Statistique descriptive
 - Nombre de features égales à Nulle (0)
 - Ratio d'observations ayant des features en outlier
 - Comparaison des valeurs statistiques entre le dataset initial et le dataset dépourvu des outliers
 - e. Distribution du jeu des données:
6. Exploration graphique univariable des données:
 - a. Histogramme des features après gestion des valeurs manquantes
 - Histogramme comparant la feature 'OP100_Capuchon_insertion_mesure' avant et après la gestion des valeurs manquantes
 - b. Histogramme des features après gestion des valeurs manquantes + application de 'RobustScaler'
 - c. Histogramme des features après gestion des valeurs manquantes + transformation 'log10'
 - d. Violon et boîte à moustaches des features avec gestion des valeurs manquantes
 - e. Violon et boîte à moustaches des features avec gestion des valeurs manquantes + application de 'RobustScaler'
7. Exploration graphique bivariées 'feature/target' des données:
 - a. Matrice de corrélation et heatmap après gestion des valeurs manquantes
 - b. Matrice de corrélation et heatmap après gestion des valeurs manquantes et rescale
 - c. Nuage de points entre la 'target' et les autres 'features'
8. Feature Engineering/Sélection et choix faits/Hypothèses choisies
9. Analyse de la target:
 - a. Vérification de l'équilibre des données
 - b. Distribution du dataset selon les classes de la target
 - c. Histogramme de distribution du dataset selon les classes de la target
10. Analyse de la target après un oversampling SMOTE

A FAIRE

- a. Régénération SMOTE de la classe minoritaire de la target
 - b. Statistique descriptive du nouveau dataset
 - c. Nouvelle distribution équilibrée du nouveau dataset
 - d. Matrice de corrélation et heatmap du nouveau dataset
 - e. Violon et boîte à moustaches des features du nouveau dataset
 - f. Ratio d'observations ayant des features en outlier du nouveau dataset
11. Modèle à base d'arbres : Balanced Random Forest Classifier **FAIT PARTIELLEMENT => A COMPLETER**
- a. Train / Test / Split + F1 et ROC
 - b. Cross Validation + F1 et ROC
12. Modèle à base de distance : SMOTE et Logistique regression **FAIT PARTIELLEMENT => A COMPLETER**
- a. Train / Test / Split + F1 et ROC
 - b. Cross Validation + F1 et ROC
13. Modèle à base de Réseau de Neuronne ou bien de Stacking (Ensemble learning) **A FAIRE**
14. Conclusions **A FAIRE**
15. Perspectives **A FAIRE**
16. Annexe: Code Python

Récapitulatif du Reste à faire

1. Feature engineering: **Point 8 de l'index ci dessus**
 - Appliquer une transformation log10 pour les features dont les distributions sont asymétriques
 - Enrichir les données avec l'horodatage d'assemblage en les extrayant de l'identifiant technique 'PROC_TRACEINFO'
 - Dans les histogrammes 9.c qui représentent la distribution des "features numériques" sur les 2 classes OK et KO, on constate que la classe minoritaire se retrouve délimité à l'intérieur d'une plage de valeurs pour certains features. (ex: OP070_V_1/2_angle_value, OP110_Vissage_M8_torque_Value). Pour cela, il faudrait vérifier l'impact si on transforme ces features numériques continues en des features catégoriques mettant en avant l'existence de la classe minoritaire KO pour ces catégories.
 - **La mise en place du Feature engineering va induire la régénération du 7, du 8 et 9.**
2. Compléter avec 3 classifieurs de type différents: à base d'arbre / à base de distance / à base de réseau de Neuronne ou bien de Stacking: **Points 11, 12, 13 de l'index ci dessus**
 - Pour chaque classifieur faire: TrainTestSplit / CV / SearchGridCV
 - Pour chaque classifieur: Analyse et interprétation des résultats / Matrice de confusion / F1, Roc / Graphe F1 et ROC
 - **Si le temps le permet** alors faire pour chaque classifieur: Graphe Overfit Underfit / Graphe avec des valeurs différents des hyperparamètres.
3. Expliquer pourquoi une classification déséquilibrée pose un défi pour la modélisation prédictive.

4. Identifier la/les motivations pour avoir une distribution Normale "bell shape". Citer les avantages d'une telle distribution.
5. Mise en forme selon les indications du document "DSSP14_Guidelines_Projet_Professionnel.doc"

Questions auxquelles j'aimerais avoir une réponse:

1. Quand on a une distribution asymétrique pour une feature et qu'on voudrait lui appliquer une transformation logarithmique pour s'approcher d'une distribution normale:
 - Est ce qu'il vaut mieux appliquer la transformation sur les features asymétriques seulement ?
Ou bien on peut l'appliquer sur la totalité des features de la dataframe (ça nous évite de choisir une-par-une les features à transformer)
 - Cette transformation sera appliqué sur le TrainSet;
Est ce qu'il faut l'appliquer aussi sur le TestSet au moment de la prédiction ?
2. Les opérations d'imputations(ex: IterativeImputer) et de scaling (ex: RobustScaling) sont appliquées sur le TrainSet afin d'honorer les pré-requis d'apprentissage de certains algorithmes de machine learning.

Est-ce que la prédiction sur une observation (TestSet) fonctionnera correctement au cas où ***l'observation*** pour laquelle on effectue la prédiction:
 - Manque certains features en "missing values" ?
 - Ou bien si les features de l'observation ne sont pas scalés selon l'attente de l'algorithme en phase d'apprentissage ?
3. Quand on fait des splits de Train/Test par Cross Validation, comme le **11.b** et le **12.b**, la méthode 'cross_validate(..)' retourne autant de fitted classifiers qu'il y a de folds.

C'est à dire: Si le Cross Validation est effectué sur 5 folds, alors la méthode 'cross_validate' retourne 5 fitted classifiers.

Questions:

- Parmi ces fitted classifiers, lequel faut-il choisir afin de l'utiliser ?
 - Est ce qu'on choisit celui dont le roc_auc est le plus élevé ? tel que c'était fait dans le **11.b** et le **12.b**
4. Est ce qu'il faut commenter davantage les graphes et les mesures figurant dans les chapitres 'Exploration des données', chapitres 5,6 et 7 ?
 5. Au niveau de ce document, est ce qu'il faut agrandir les graphes ?

Ou bien, ils seront consultés sur un support électronique (pdf, doc, projection, ...) et par conséquent ils seront agrandis électroniquement ?

6. Quelles sont les parties que je dois développer davantage ? Ou bien être plus concis ?
7. Est-ce qu'il faut garder dans le document, les petits bout de code Python qu'on trouve tout au long des chapitres (hors chapitre 16 – Annexe : code Python) ? Ou bien il faut les supprimer ?

Merci :-)

1 - Contexte de l'étude

L'étude correspond à un 'Challenge Data ENS' qui a pour objectif de prévoir les défauts sur les lignes de production des démarreurs de l'équipementier Valeo. Lors de l'assemblage des démarreurs sur la ligne de production, les différentes valeurs (couples, angles ...) sont mesurées sur les différentes stations de montage.

En fin de ligne, des mesures supplémentaires sont effectuées sur deux bancs de test afin d'isoler les défauts. Par conséquent, les échantillons sont étiquetés "OK" ou "KO". L'objectif est de concevoir un modèle qui pourrait identifier de tels défauts avant l'étape du banc d'essai.

L'étude concerne la classification des données déséquilibrée avec des valeurs de données manquantes. C'est un problème classique dans l'industrie et dans bien d'autres domaines : détection de fraude, détection de spam, domaine médical,

Les classifications déséquilibrées posent un défi pour la modélisation prédictive. La classe minoritaire est plus importante et donc le problème est plus sensible aux erreurs de classification pour la classe minoritaire que pour la classe majoritaire.

A l'heure de la rédaction de ce document, mon meilleur modèle est basé sur le classifieur « Balanced Random Forest », il occupe le 65^{ème} rang sur un nombre total de 116 participants. Mon score (roc_auc) est égal à 0.6344 sur une plage allant de 0.4 jusqu'à 0.76. J'ai déjà identifié certaines pistes d'amélioration que je n'ai pas encore implémentées, notamment au niveau « features engineering ».

2 - Description du data set case

a - Entrées:

Les caractéristiques d'entrée sont des mesures collectées sur différentes stations d'assemblage avec des capteurs connectés à des contrôleurs logiques programmables qui les stockent tous.

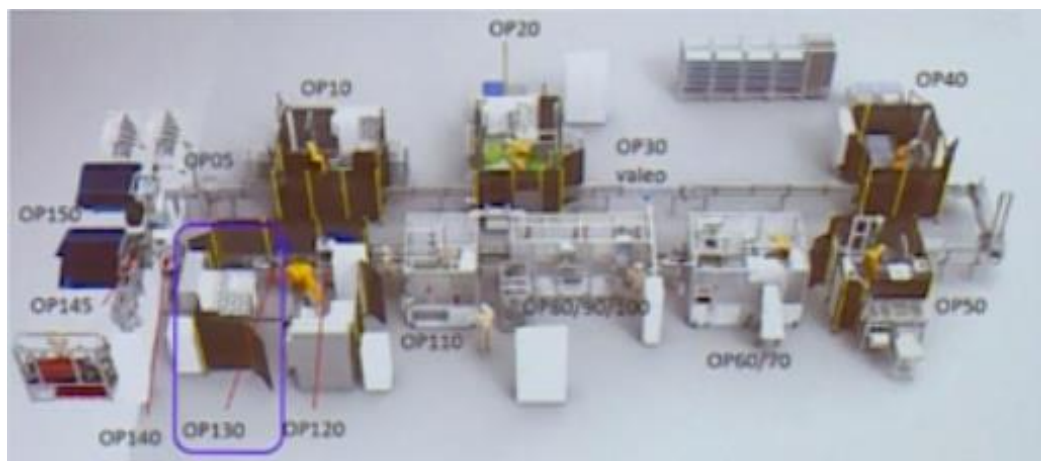
On distingue par exemple:

OP070_V_1_angle_value	V1 Valeur d'angle,
OP070_V_1_torque_value	V1 Valeur de couple,
OP070_V_2_angle_value	V2 Valeur d'angle,
OP070_V_2_torque_value	V2 Valeur de couple,
OP090_StartLinePeakForce_value	Start Line Peak Force value,
OP090_SnapRingMidPointForce_value	Anneau élastique Mid Point Force val,
OP090_SnapRingPeakForce_value	Anneau élastique Peak Force value,
OP090_SnapRingFinalStroke_value	Valeur finale du coup d'arrêt,
OP100_Capuchon_insertion_mesure	Mesure d'insertion capuchon
OP110_Vissage_M8_angle_value	Valeur d'angle Vissage M8,
OP110_Vissage_M8_torque_value	Valeur de couple Vissage M8,
OP120_Rodage_I_mesure_value	Rodage I mesure la valeur,
OP120_Rodage_U_mesure_value	Rodage U mesure value,

b - Sortie:

Il s'agit de la valeur de résultat de l'OP130, banc d'essai: Binar OP130_Resultat_Global_v.

La valeur 0 est affectée aux échantillons OK (réussie) et la valeur 1 est affectée aux échantillons KO (échoué). Il s'agit du résultat combiné de multiples tests électriques, acoustiques et vibro-acoustiques.



L'objectif est de trouver la meilleure prédiction: $\text{Sortie} = f(\text{entrées})$. L'ensemble de données contient 34515 échantillons d'apprentissage et 8001 échantillons de test.

c - Croisement Entrées/Sortie:

Les données de training sont réparties dans 2 fichiers csv:

- [project-root]/data/train/traininginputs.csv
- [project-root]/data/train/trainingoutput.csv

Un identifiant technique 'PROC_TRACEINFO' permet de croiser le fichier d'entrée au fichier de sortie.

C'est un code unique donné attribué au démarreur assemblé.

Exemple: I-B-XA1207672-190701-00494.

- XA1207672 est la référence.
- 190701 est la date: ici le 01 juillet de l'année 2019.
- 00494 est le code unique donné au produit, ce nombre est augmenté de 1 pour chaque nouveau produit.

On dispose aussi des données d'entrée de test: [project-root]/data/test/testinputs.csv

Les données de sortie de test sont générés par l'étude et sont uploader sur la plateforme 'Data Challenge ENS' <https://challengedata.ens.fr/participants/challenges/36/>

3 - Import des packages et rechargement automatique des packages du projets

```
import os
import sys
import logging
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from pandas.plotting import scatter_matrix
import seaborn as sns
%matplotlib inline

from sklearn.model_selection import train_test_split
from imblearn.over_sampling import SMOTE

from imblearn.ensemble import BalancedBaggingClassifier, RUSBoostClassifier,
BalancedRandomForestClassifier
from imblearn.over_sampling import RandomOverSampler, ADASYN, SMOTE, SVMSMOTE, KMeansSMOTE,
BorderlineSMOTE
from imblearn.over_sampling.base import BaseOverSampler
from imblearn.pipeline import Pipeline

from sklearn.ensemble import GradientBoostingClassifier, AdaBoostClassifier, RandomForestClassifier
from sklearn.ensemble._hist_gradient_boosting.gradient_boosting import HistGradientBoostingClassifier
from sklearn.cluster import MiniBatchKMeans
from sklearn.compose import ColumnTransformer
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model._stochastic_gradient import SGDClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC, NuSVC, LinearSVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import cross_validate, StratifiedKFold

import xgboost as xgb

# Import "valeo" module
sys.path.append("../")
from valeo.infrastructure.LogManager import LogManager as lm
# NB: Initializing logger here allows "class loaders of application classes" to benefit from the global
initialization
logger = lm().logger(__name__)
from valeo.infrastructure import Const
from valeo.infrastructure.tools.DfUtil import DfUtil
from valeo.infrastructure.tools.ImgUtil import ImgUtil
from valeo.infrastructure.XY_Loader import XY_Loader
from valeo.infrastructure.XY_metadata import XY_metadata as XY_metadata
from valeo.domain.ValeoModeler import ValeoModeler
from valeo.domain.ValeoPredictor import ValeoPredictor
import valeo.infrastructure.Transformer as transf

# Notebook automatic reload
%load_ext autoreload
%reload_ext autoreload
%aimport valeo.infrastructure.Transformer
%aimport valeo.infrastructure.LogManager
%aimport valeo.infrastructure.Const
%aimport valeo.infrastructure.tools.DfUtil
%aimport valeo.infrastructure.tools.ImgUtil
%aimport valeo.infrastructure.XY_Loader
%aimport valeo.infrastructure.XY_metadata
%aimport valeo.domain.ValeoModeler
%aimport valeo.domain.ValeoPredictor
```


4 - Chargement des données 'Training'

```
data = DfUtil.read_csv([Const.rootDataTrain() , "traininginputs.csv"])
Y_data = DfUtil.read_csv([Const.rootDataTrain(), "trainingoutput.csv"])
```

5 - Exploration et analyse tabulaire des données

a - Visualisation tabulaire des données - Affichage du type 'head()':

```
data.head()
```

	PROC_TRACEINFO	OP070_V_1_angle_value	OP090_SnapRingPeakForce_value	OP070_V_2_angle_value	OP120_Rodage_I_mesure_value
0	I-B-XA1207872-190429-00688	180.4	190.51	173.1	113.64
1	I-B-XA1207872-190828-00973	138.7	147.70	163.5	109.77
2	I-B-XA1207872-190712-03462	180.9	150.87	181.2	109.79
3	I-B-XA1207872-190803-00051	173.5	159.56	151.8	113.25
4	I-B-XA1207872-190508-03248	174.5	172.29	177.5	112.88

OP090_SnapRingFinalStroke_value	OP110_Vissage_M8_torque_value	OP100_Capuchon_insertion_mesure	OP120_Rodage_U_mesure_value
12.04	12.16	NaN	11.97
12.12	12.19	0.39	11.97
11.86	12.24	NaN	11.97
11.82	12.35	0.39	11.97
12.07	12.19	NaN	11.97

OP070_V_1_torque_value	OP090_StartLinePeakForce_value	OP110_Vissage_M8_angle_value	OP090_SnapRingMidPointForce_val	OP070_V_2_torque_value
6.62	26.37	18.8	109.62	6.60
6.41	21.03	18.5	105.48	6.40
6.62	25.81	17.5	100.03	6.61
6.62	24.62	15.6	104.94	6.61
6.62	29.22	33.6	99.19	6.61

Un simple affichage du type 'head()' permet de voir à quoi ressemble les données.

b - Rapport sémantique des données - Affichage du type 'info()':

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 34515 entries, 0 to 34514
```

```
Data columns (total 14 columns):
```

#	Column	Non-Null Count	Dtype
---	-----	-----	-----

0	PROC_TRACEINFO	34515	non-null	object
1	OP070_V_1_angle_value	34515	non-null	float64
2	OP090_SnapRingPeakForce_value	34515	non-null	float64
3	OP070_V_2_angle_value	34515	non-null	float64
4	OP120_Rodage_I_mesure_value	34515	non-null	float64
5	OP090_SnapRingFinalStroke_value	34515	non-null	float64
6	OP110_Vissage_M8_torque_value	34515	non-null	float64
7	OP100_Capuchon_insertion_mesure	15888	non-null	float64
8	OP120_Rodage_U_mesure_value	34515	non-null	float64
9	OP070_V_1_torque_value	34515	non-null	float64
10	OP090_StartLinePeakForce_value	34515	non-null	float64
11	OP110_Vissage_M8_angle_value	34515	non-null	float64
12	OP090_SnapRingMidPointForce_val	34515	non-null	float64
13	OP070_V_2_torque_value	34515	non-null	float64

dtypes: float64(13), object(1)
memory usage: 3.7+ MB

Un affichage sémantique du type 'info()' met en évidence le type des données et le nombre des valeurs manquantes 'missing values'.

On constate que:

- Toutes les features sont numériques et continues, pas de features catégoriques
- Plus de la moitié des valeurs de la feature 7 'OP100_Capuchon_insertion_mesure' sont manquants
=> Cette feature doit être traitée en lui imputant des valeurs. Un imputer de type IterativeImputer(stratégie 'médiane') sera utilisé.
- PROC_TRACEINFO de type object (=> String), c'est l'identifiant de ligne permettant de croiser les 'features' avec la 'target'. Cette feature porte l'horodatage de l'assemblage des démarreurs, la date sous jacente sera extraite et utilisée dans la phase de 'features engineering'

c - Données manquantes par type de 'feature':

```
data.isna().sum()
```

PROC_TRACEINFO	0
OP070_V_1_angle_value	0
OP090_SnapRingPeakForce_value	0
OP070_V_2_angle_value	0
OP120_Rodage_I_mesure_value	0
OP090_SnapRingFinalStroke_value	0
OP110_Vissage_M8_torque_value	0
OP100_Capuchon_insertion_mesure	18627
OP120_Rodage_U_mesure_value	0
OP070_V_1_torque_value	0
OP090_StartLinePeakForce_value	0
OP110_Vissage_M8_angle_value	0
OP090_SnapRingMidPointForce_val	0
OP070_V_2_torque_value	0

dtype: int64

L'identifiant 'PROC_TRACEINFO' est supprimé **provisoirement** de l'ensemble des features:

```
X_data = data.drop(columns = "PROC_TRACEINFO")
X_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 34515 entries, 0 to 34514
Data columns (total 13 columns):
```

#	Column	Non-Null Count	Dtype
0	OP070_V_1_angle_value	34515 non-null	float64
1	OP090_SnapRingPeakForce_value	34515 non-null	float64
2	OP070_V_2_angle_value	34515 non-null	float64
3	OP120_Rodage_I_mesure_value	34515 non-null	float64
4	OP090_SnapRingFinalStroke_value	34515 non-null	float64
5	OP110_Vissage_M8_torque_value	34515 non-null	float64
6	OP100_Capuchon_insertion_mesure	15888 non-null	float64
7	OP120_Rodage_U_mesure_value	34515 non-null	float64
8	OP070_V_1_torque_value	34515 non-null	float64
9	OP090_StartLinePeakForce_value	34515 non-null	float64
10	OP110_Vissage_M8_angle_value	34515 non-null	float64
11	OP090_SnapRingMidPointForce_val	34515 non-null	float64
12	OP070_V_2_torque_value	34515 non-null	float64

```
dtypes: float64(13)
```

```
memory usage: 3.4 MB
```

d - Statistique descriptive

```
X_data.sort_index(axis=1).describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
OP070_V_1_angle_value	34515.0	159.906922	15.662650	101.80	148.70	158.00	169.30	198.30
OP070_V_1_torque_value	34515.0	6.548403	0.097602	5.67	6.41	6.61	6.62	6.67
OP070_V_2_angle_value	34515.0	159.618236	15.091490	82.00	149.40	158.70	168.90	198.10
OP070_V_2_torque_value	34515.0	6.550867	0.094814	5.74	6.42	6.61	6.61	6.67
OP090_SnapRingFinalStroke_value	34515.0	11.970190	0.169873	0.00	11.85	12.04	12.08	12.19
OP090_SnapRingMidPointForce_val	34515.0	97.700978	6.837714	0.00	94.31	98.50	102.23	127.30
OP090_SnapRingPeakForce_value	34515.0	156.915055	11.271492	0.00	149.21	156.18	164.38	196.92
OP090_StartLinePeakForce_value	34515.0	23.630152	2.546341	0.00	22.28	23.88	25.29	43.41
OP100_Capuchon_insertion_mesure	15888.0	0.388173	0.024425	0.24	0.38	0.39	0.41	0.42
OP110_Vissage_M8_angle_value	34515.0	17.878398	6.785079	6.30	13.50	16.40	20.20	84.60
OP110_Vissage_M8_torque_value	34515.0	12.256785	0.065319	12.03	12.21	12.26	12.30	12.50
OP120_Rodage_I_mesure_value	34515.0	113.350222	3.528522	99.99	111.04	113.16	115.38	177.95
OP120_Rodage_U_mesure_value	34515.0	11.971027	0.003050	11.97	11.97	11.97	11.97	11.99

On constate que :

- OP070_V_2_angle_value : Outlier côté Min => Utiliser un 'robust scaler' pour réduire l'effet Outlier
- OP090_StartLinePeakForce_value, OP090_SnapRingMidPointForce_val, OP090_SnapRingPeakForce_value, OP090_SnapRingFinalStroke_value :

Identification de valeurs nulle, 'min' égal à 0. Normalement ces mesures physiques ne doivent pas être nulle, le fait qu'elles soient nulles laisse penser qu'elles sont nulles à tort et par conséquent il faut les considérer comme des valeurs manquantes et seront traitées dans la phase '**Feature Engineering**'

- OP090_StartLinePeakForce_value : Outlier côté Max => Utiliser un 'robust scaler'
- OP110_Vissage_M8_angle_value : Outlier côté Max => Utiliser un 'robust scaler'
- OP110_Vissage_M8_torque_value : Presque Constant
- OP100_Capuchon_insertion_mesure: Plus de la moitié sans valeurs => Utiliser 'missing values' Imputer
- OP120_Rodage_U_mesure_value : Très petite variance
- OP120_Rodage_I_mesure_value : Outlier cote Max (=> Utiliser un 'robust scaler') + petite variance

Nombre de features égales à Nulle:

```
data.query('OP090_StartLinePeakForce_value == 0 or OP090_SnapRingMidPointForce_val == 0 or  
OP090_SnapRingPeakForce_value == 0 or OP090_SnapRingFinalStroke_value == 0')
```

	PROC_TRACEINFO	OP070_V_1_angle_value	OP090_SnapRingPeakForce_value	OP070_V_2_angle_value	OP120_Rodage_I_mesure_value
549	I-B-XA1207872-190907-01953	137.4	0.0	166.7	105.51
1651	I-B-XA1207872-190821-01367	178.7	0.0	170.4	112.95
22483	I-B-XA1207872-190424-02168	166.4	0.0	171.5	117.26

Seulement 3 observations dont les valeurs des features sont égales à 0. A cela s'ajoute la moitié des valeurs de 'OP100_Capuchon_insertion_mesure' qui sont manquantes.

Ratio d'observations ayant des features en outlier:

```
Q1 = X_data.quantile(0.25)  
Q3 = X_data.quantile(0.75)  
IQR = Q3 - Q1  
#  
outliers = ((X_data < (Q1 - 1.5 * IQR)) |(X_data > (Q3 + 1.5 * IQR))).any(axis=1)  
print(f"Le ratio d'outlier est de {len(X_data[outliers].index)/len(X_data.index)}")
```

Le ratio d'outlier est de 0.24256120527306968

Un ratio élevé => L'éventualité de supprimer les observations n'est pas viable. D'autant plus qu'on ne connaît pas la raison de ces outliers:

- Est ce que c'est une erreur
- Ou bien c'est une vraie donnée dont le pattern est différent

Pour limiter l'effet des outliers:

- Utiliser un modèle résistant aux outliers, comme les arbres
- Transformer les données en utilisant la fonction Log Lors de la visualisation graphique des données on va retrouver des distributions biaisées (skewed)

Comparaison des valeurs statistiques entre le dataset initial et le dataset dépourvu des outliers

```
# 1 - Le dataset dépourvu des outliers  
X_data_out = X_data[~outliers]  
  
# 2 - Créer les 2 dataframes des valeurs statistiques descriptives  
Xt = X_data.sort_index(axis=1).describe().transpose()  
Xt_out = X_data_out.sort_index(axis=1).describe().transpose()  
  
# 3 - Fusionner les afin de pouvoir les comparer  
xt_merged = pd.merge(left=Xt, right=Xt_out, how='inner', left_on=Xt.index, right_on=Xt_out.index,  
suffixes=('', '-o'))  
xt_merged = xt_merged.set_index(['key_0'])  
xt_merged.sort_index(axis=1)
```

	25%	25%-o	50%	50%-o	75%	75%-o	count	count-o	max	max-o	mean	mean-o	min	min-o
key_0														
OP070_V_1_angle_value	148.70	148.80	158.00	158.30	169.30	169.700	34515.0	26143.0	198.30	198.20	159.906922	180.147298	101.80	118.00
OP070_V_1_torque_value	6.41	6.41	6.61	6.61	6.62	6.620	34515.0	26143.0	6.67	6.67	6.548403	6.547524	5.67	6.10
OP070_V_2_angle_value	149.40	149.60	158.70	158.80	168.90	169.100	34515.0	26143.0	198.10	197.90	159.618236	159.847145	82.00	120.20
OP070_V_2_torque_value	6.42	6.42	6.61	6.61	6.61	6.610	34515.0	26143.0	6.67	6.67	6.550887	6.549834	5.74	6.15
OP090_SnapRingFinalStroke_value	11.85	11.85	12.04	12.04	12.08	12.080	34515.0	26143.0	12.19	12.19	11.970190	11.967375	0.00	11.67
OP090_SnapRingMidPointForce_val	94.31	94.91	98.50	98.82	102.23	102.385	34515.0	26143.0	127.30	114.10	97.700978	98.384817	0.00	82.43
OP090_SnapRingPeakForce_value	149.21	149.29	156.18	156.06	164.38	163.990	34515.0	26143.0	198.92	188.87	156.915055	156.769002	0.00	126.52
OP090_StartLinePeakForce_value	22.28	22.43	23.88	23.91	25.29	25.300	34515.0	26143.0	43.41	29.80	23.630152	23.789882	0.00	17.77
OP100_Capuchon_insertion_mesure	0.38	0.38	0.39	0.40	0.41	0.410	15888.0	12028.0	0.42	0.42	0.388173	0.392272	0.24	0.34
OP110_Vissage_M8_angle_value	13.50	13.40	16.40	16.10	20.20	19.400	34515.0	26143.0	84.60	30.20	17.878398	16.754443	6.30	6.30
OP110_Vissage_M8_torque_value	12.21	12.21	12.26	12.26	12.30	12.300	34515.0	26143.0	12.50	12.43	12.256785	12.258490	12.03	12.08
OP120_Rodage_I_mesure_value	111.04	111.17	113.16	113.25	115.38	115.420	34515.0	26143.0	177.95	121.88	113.350222	113.341812	99.99	104.56
OP120_Rodage_U_mesure_value	11.97	11.97	11.97	11.97	11.97	11.970	34515.0	26143.0	11.99	11.97	11.971027	11.970000	11.97	11.97
min	min-o	std	std-o											
101.80	118.00	15.662650	1.561998e+01											
5.67	6.10	0.097602	9.687328e-02											
82.00	120.20	15.091490	1.497825e+01											
5.74	6.15	0.094814	9.459255e-02											
0.00	11.67	0.169873	1.257830e-01											
0.00	82.43	6.837714	5.792290e+00											
0.00	126.52	11.271492	1.104504e+01											
0.00	17.77	2.546341	2.186637e+00											
0.24	0.34	0.024425	1.951393e-02											
6.30	6.30	6.785079	4.503697e+00											
12.03	12.08	0.065319	6.206688e-02											
99.99	104.56	3.528522	3.106980e+00											
11.97	11.97	0.003050	3.636272e-12											

D'une manière générale, les valeurs sont approximativement similaires, sauf pour le 'max' et le 'std' de quelques features:

- OP090_StartLinePeakForce_value, OP110_Vissage_M8_angle_value, OP120_Rodage_I_mesure_value: Le 'max' a chuté considérablement
- OP110_Vissage_M8_angle_value : Variance considérablement plus petite

e - Distribution du jeu des données:

```
starter_count = len(Y_data[Const.Binar_OP130_Resultat_Global_v])
starter_count_ok = Y_data[Const.Binar_OP130_Resultat_Global_v].value_counts()[0]
starter_count_ko = Y_data[Const.Binar_OP130_Resultat_Global_v].value_counts()[1]
#
print(f'Nombre total des démarreurs : {starter_count}')
print(f'Nombre total des démarreurs OK => Nombre de Classes Negatives : {starter_count_ok} soit
{round(starter_count_ok/starter_count * 100,2)} % du dataset')
print(f'Nombre total des démarreurs KO => Nombre de Classes Positives : {starter_count_ko} soit
{round(starter_count_ko/starter_count * 100,2)} % du dataset')
```

Nombre total des démarreurs : 34515

Nombre total des démarreurs OK => Nombre de Classes Negatives : 34210 soit
99.12 % du dataset

Nombre total des démarreurs KO => Nombre de Classes Positives : 305 soit 0.88
% du dataset

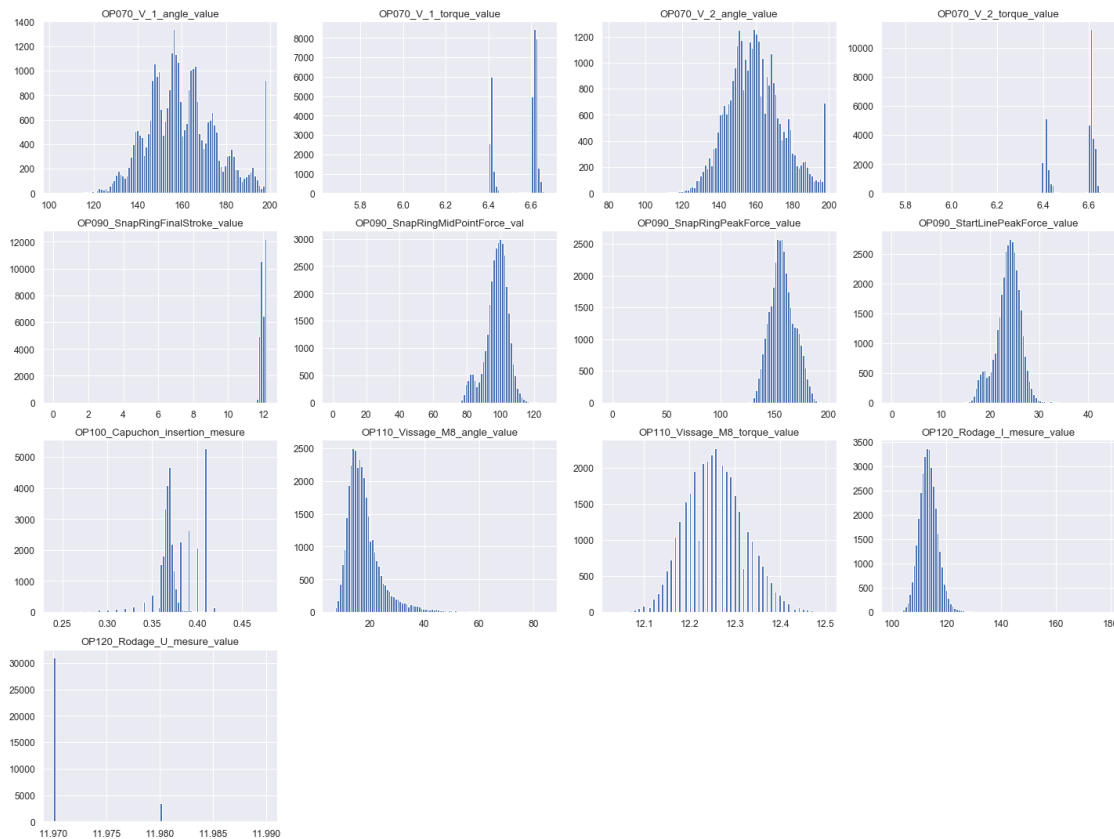
L'étude concerne la classification des données déséquilibrée avec des valeurs de données manquantes. C'est un problème classique dans l'industrie et dans bien d'autres domaines : détection de fraude, détection de spam, domaine médical,

On constate qu'on est sur une classification déséquilibrée dans la répartition ce qui pose un défi pour la modélisation prédictive. La classe minoritaire est plus importante et donc le problème est plus sensible aux erreurs de classification pour la classe minoritaire que pour la classe majoritaire.

6 - Exploration graphique univariable des données

a - Histogramme des features après gestion des valeurs manquantes:

```
tsf = transf.Transformer()
X_data_transformed = tsf.iterative_imputer_transform(X_data)
ImgUtil.save_df_hist_plot(X_data_transformed,"X_data_imputed",figsize=(20,15), bins=100)
plt.show()
```



NB:

tsf.iterative_imputer_transform(X_data) méthode de la classe 'Transformer' du module Python `valeo.infrastructure.Transformer`. Elle applique un imputer du type `IterativeImputer(estimator=BayesianRidge, missing_values, initial_strategy = 'median')`

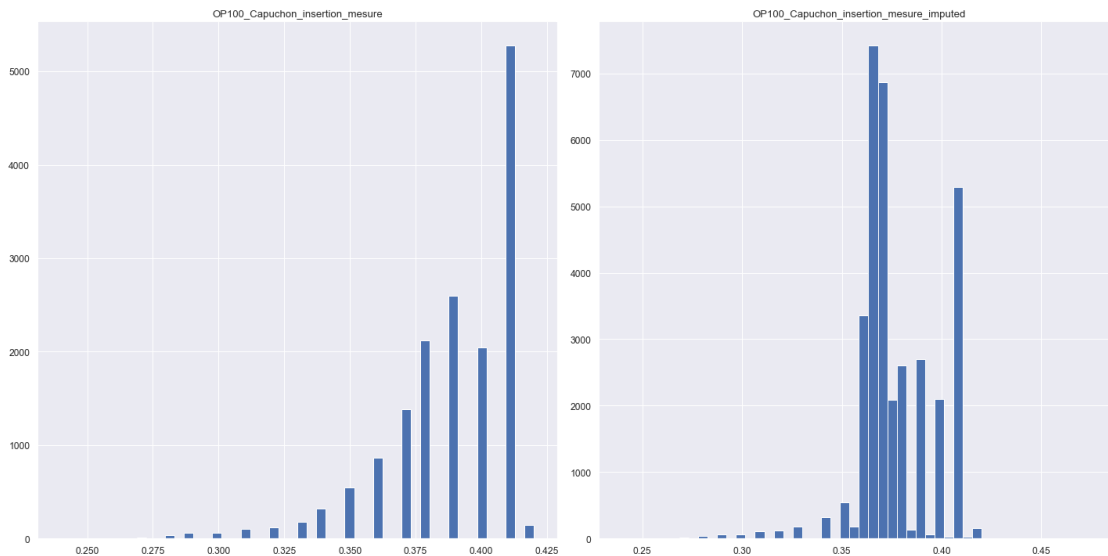
En observant les graphes des différents features, on constate:

- La plupart des distributions sont asymétriques notamment pour: (=> Appliquer une transformation logarithmique dans l'étape F.Engineering)
 - OP070_V_1_angle_value
 - OP090_SnapRingMidPointForce_val
 - OP090_SnapLinePeakForce_value

- OP100_Capuchon_insertion_mesure # feature dont la moitié des mesures n'existait pas
- OP110_Vissage_M8_angle_value
- Une valeur de plafonnement (capping value) pour:
 - OP070_V_1_angle_value
 - OP070_V_2_angle_value
- Les distributions suivantes représentent 2 catégories d'observations indépendantes: (=> Représenter chaque feature par 2 catégories dans l'étape F.Engin.):
 - OP070_V_1_torque_value
 - OP070_V_2_torque_value
 - OP120_Rodage_U_mesure_value
 -

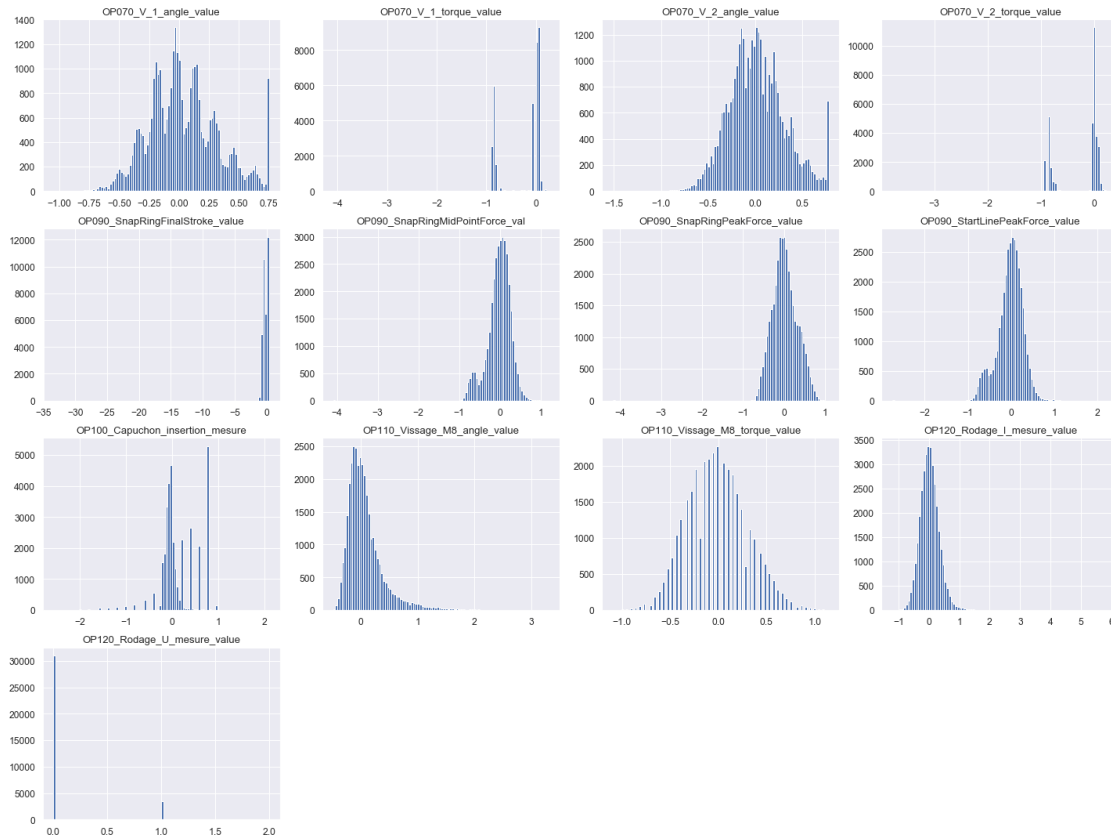
Histogramme comparant la feature 'OP100_Capuchon_insertion_mesure' avant et après la gestion des valeurs manquantes:

```
dff_ = pd.DataFrame(X_data[Const.OP100_Capuchon_insertion_mesure])
dff_[Const.OP100_Capuchon_insertion_mesure + "_imputed"] =
X_data_transformed[Const.OP100_Capuchon_insertion_mesure]
ImgUtil.save_df_hist_plot(dff_,Const.OP100_Capuchon_insertion_mesure, figsize=(20,10))
plt.show()
```



b - Histogramme des features avec gestion des valeurs manquantes + Application de 'RobustScaler':

```
X_data_transformed_scaled = tsf.robust_scaler_transform(X_data_transformed)
ImgUtil.save_df_hist_plot(X_data_transformed_scaled,"X_data_imputed_robust_scaled",bins=100)
plt.show()
```



Pourquoi appliquer un scaling:

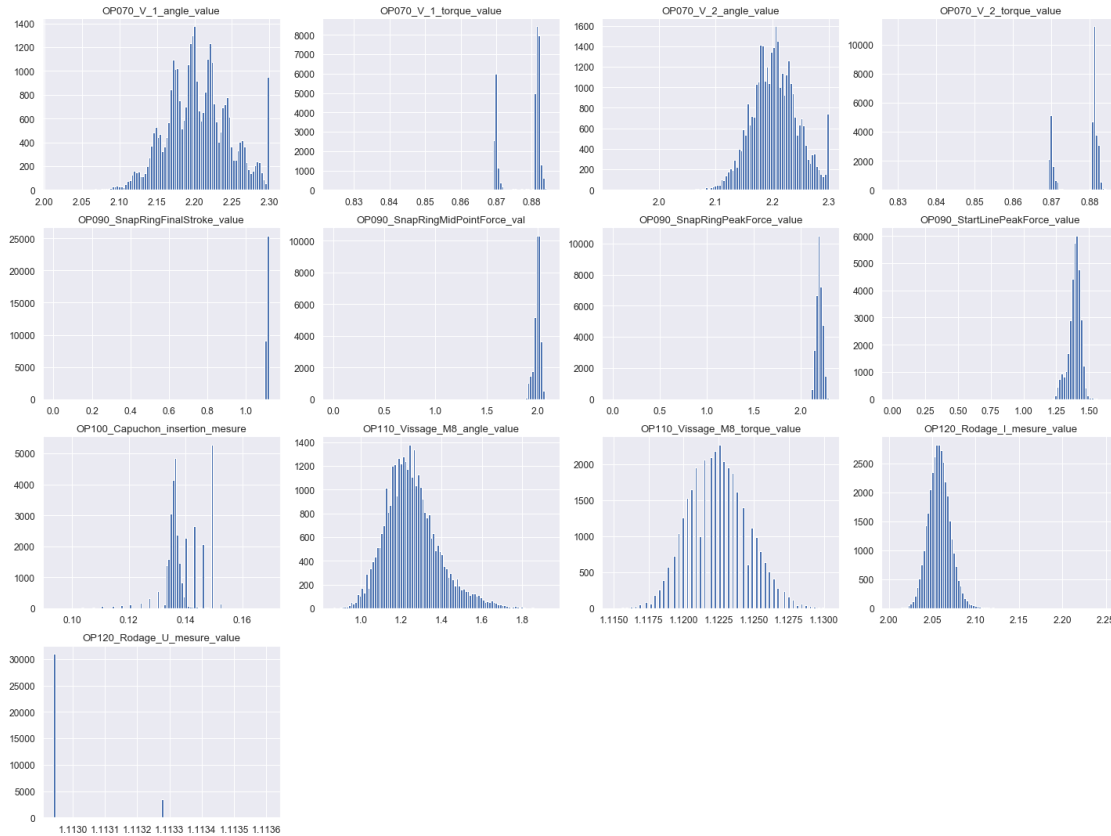
Les algorithmes d'apprentissage automatique prennent en compte uniquement la magnitude des mesures, mais pas les unités de ces mesures. Par la suite, une caractéristique exprimée en une magnitude (nombre) très élevée, peut affecter la prévision beaucoup plus qu'une caractéristique tout aussi importante.

Notez que tous les algorithmes se comportent pas de cette façon et par la suite l'application du scaling n'est pas un pré-requis pour tout les algorithmes.

Les algorithmes à base d'arbres et de Naive Bayes ne nécessitent pas de mise à l'échelle des fonctionnalités, car ils fonctionnent. Les algorithmes qui exploitent des distances ou des similitudes (par exemple sous forme de produit scalaire) entre des échantillons de données, tels que k-NN et SVM, nécessitent souvent une mise à l'échelle des fonctionnalités.

c - Histogramme des features après gestion des valeurs manquantes + Transformation 'log10':

```
tsf = transf.Transformer()
X_data_offset_1 = X_data_transformed + 1
X_data_transformed_log10 = X_data_offset_1.applymap(np.log10)
ImgUtil.save_df_hist_plot(X_data_transformed_log10,"X_data_imputed_log10", bins=100)
plt.show()
```



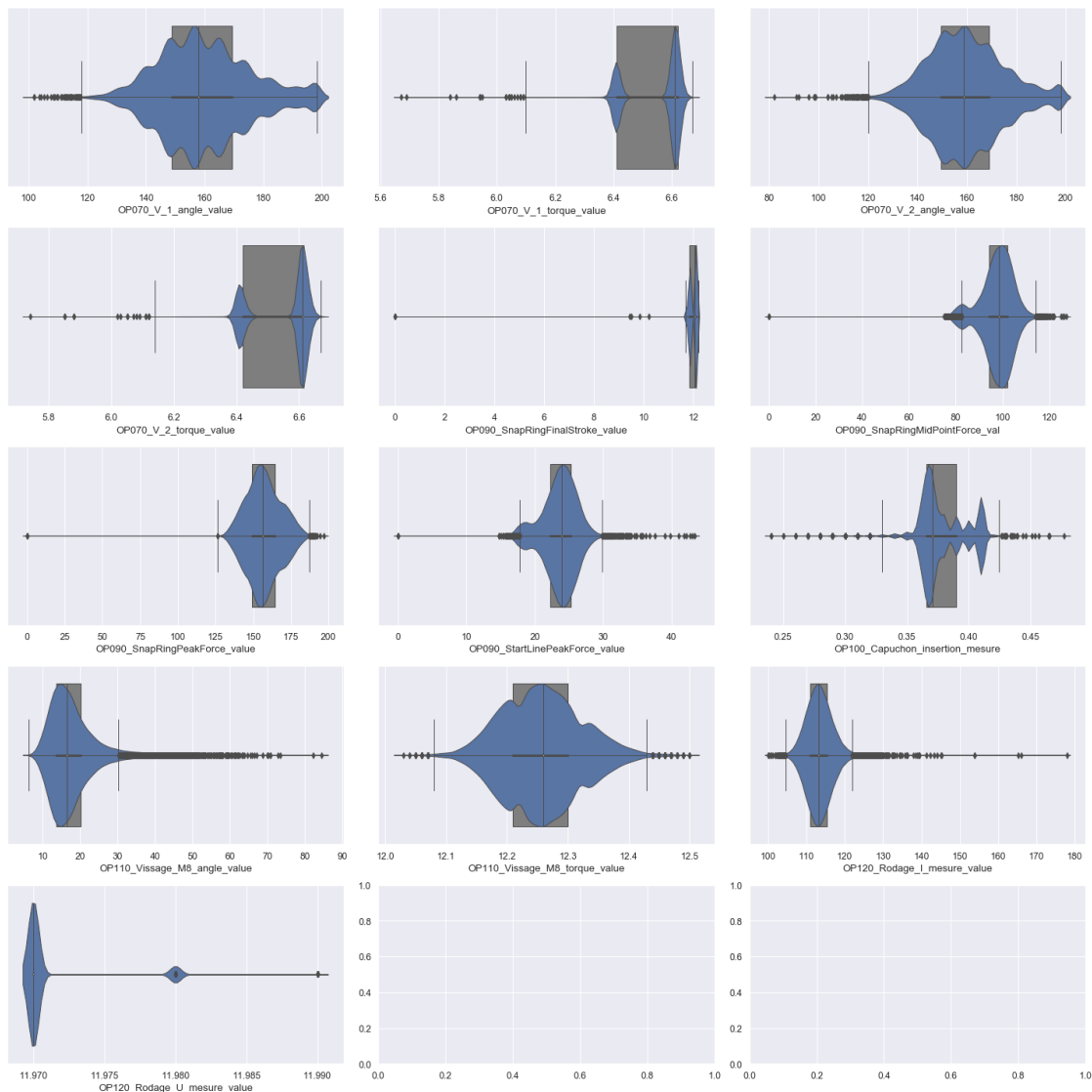
On constate que les distributions ci dessus correspondant aux features transformées par log10 **ressemblent plus** à des distributions Normales.

d - 'Violon' et 'boîte à moustaches' des features avec gestion des valeurs manquantes:

Chacun des graphes suivants correspond à la superposition de 2 graphes: Celui d'une 'boîte à moustaches' et d'un 'violon'.

La 'boîte à moustache' représente clairement Q1, Q3, médiane, moustaches, min, max, outliers. Alors que le 'violon' montre bien la distribution des données à l'intérieur.

```
ImgUtil.save_df_violin_plot(X_data_transformed, 'X_data_distribution', 3)
```

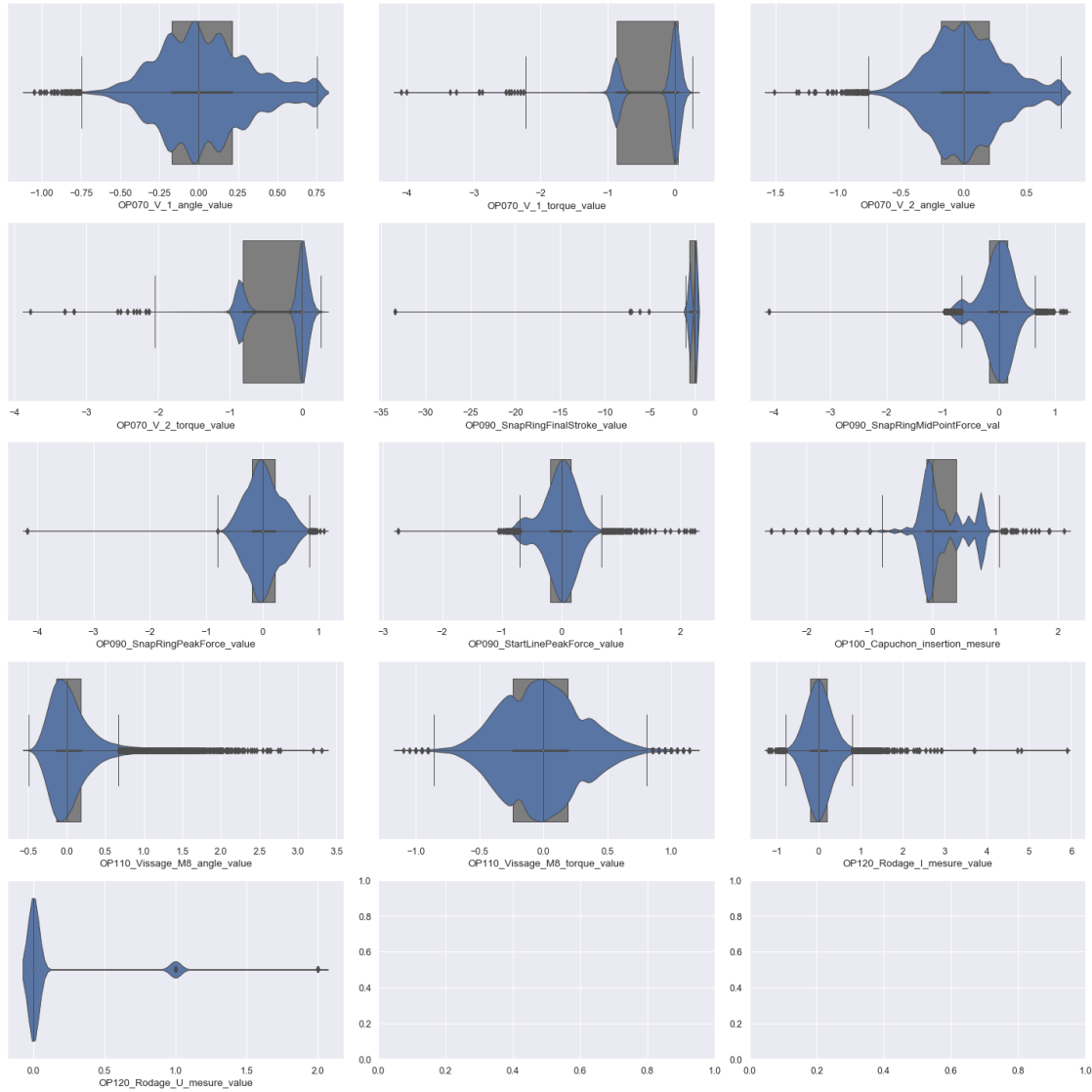


Représenter la distribution d'une feature par un 'violin plot' superposé à un 'box plot' permet de:

- Visualiser la taille de la distribution d'une feature en fonction de sa valeur
- Afficher Q1, Q3 et la médiane
- Afficher les moustaches inférieur et supérieur ainsi que les outliers

e - Violon et boîte à moustaches des features avec gestion des valeurs manquantes + application de 'RobustScaler':

```
ImgUtil.save_df_violin_plot(X_data_transformed_scaled, 'X_data_scaled_distribution_scaled', 3)
```



7 - Exploration graphique bivariées 'feature/target' des données

a - Matrice de corrélation (pearson) et heatmap après gestion des valeurs manquantes :

La corrélation des données est un moyen de comprendre la relation entre plusieurs features/target dans un ensemble de données.

1 - Charger les features et la target en les croisant:

```
XY_data_with_id = pd.merge(left=data, right=Y_data, how='inner', left_on=Const.PROC_TRACEINFO,
right_on=Const.PROC_TRACEINFO)
XY_data_with_id.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 34515 entries, 0 to 34514
```

```
Data columns (total 15 columns):
```

#	Column	Non-Null Count	Dtype
0	PROC_TRACEINFO	34515 non-null	object
1	OP070_V_1_angle_value	34515 non-null	float64
2	OP090_SnapRingPeakForce_value	34515 non-null	float64
3	OP070_V_2_angle_value	34515 non-null	float64
4	OP120_Rodage_I_mesure_value	34515 non-null	float64
5	OP090_SnapRingFinalStroke_value	34515 non-null	float64
6	OP110_Vissage_M8_torque_value	34515 non-null	float64
7	OP100_Capuchon_insertion_mesure	15888 non-null	float64
8	OP120_Rodage_U_mesure_value	34515 non-null	float64
9	OP070_V_1_torque_value	34515 non-null	float64
10	OP090_StartLinePeakForce_value	34515 non-null	float64
11	OP110_Vissage_M8_angle_value	34515 non-null	float64
12	OP090_SnapRingMidPointForce_val	34515 non-null	float64
13	OP070_V_2_torque_value	34515 non-null	float64
14	Binar OP130_Resultat_Global_v	34515 non-null	int64

```
dtypes: float64(13), int64(1), object(1)
```

```
memory usage: 4.2+ MB
```

2 - Rajout des missing values afin d'avoir une meilleure représentation

```
XY_data = XY_data_with_id.drop(columns = Const.PROC_TRACEINFO)
```

```
XY_data_transformed = tsf.iterative_imputer_transform(XY_data)
```

3 - Corrélation entre la target "Binar OP130_Resultat_Global_v" et les autres attributs

```
corr_matrix = XY_data_transformed.corr()
```

```
corr_matrix[Const.Binar_OP130_Resultat_Global_v].sort_values(ascending=False)
```

Binar OP130_Resultat_Global_v	1.000000
OP100_Capuchon_insertion_mesure	0.040366
OP090_SnapRingFinalStroke_value	0.015148
OP090_SnapRingMidPointForce_val	0.014273
OP090_StartLinePeakForce_value	0.010720
OP110_Vissage_M8_angle_value	0.005470
OP120_Rodage_I_mesure_value	0.003763

```
OP110_Vissage_M8_torque_value      -0.002984
OP070_V_2_angle_value              -0.006342
OP090_SnapRingPeakForce_value      -0.007290
OP120_Rodage_U_mesure_value         -0.010492
OP070_V_1_angle_value              -0.012793
OP070_V_1_torque_value              -0.037438
OP070_V_2_torque_value              -0.039752
Name: Binar OP130_Resultat_Global_v, dtype: float64
```

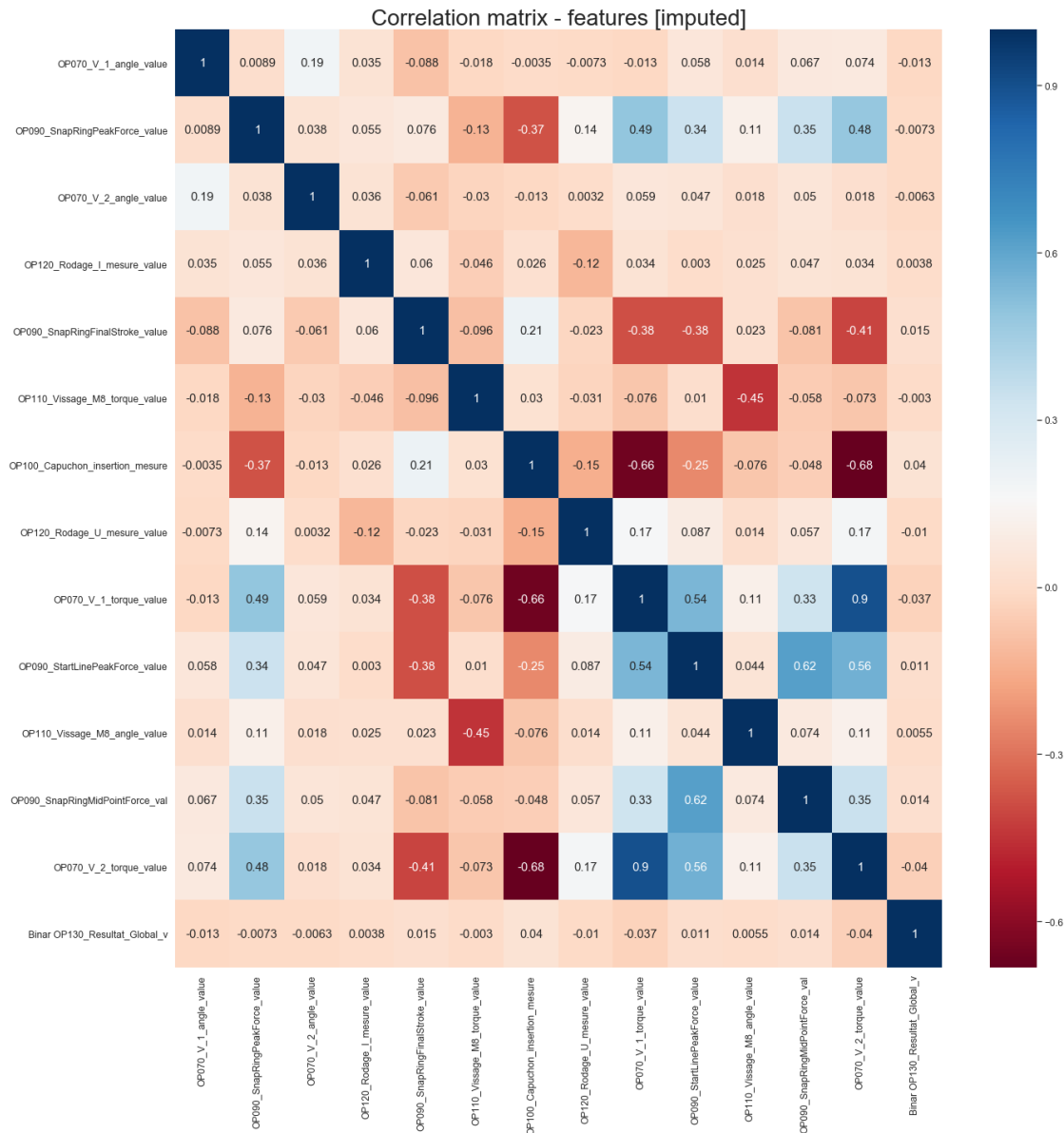
Le coefficient de correlation varie entre -1 et 1 :

- Proche de 1 => il y a une correlation forte positive.
- Proche de -1 => il y a une correlation forte négative.
- Proche de 0 => Il n'y a pas de **correlation linéaire**

Le coefficient de correlation mesure uniquement les correlations linéaires

4 - Dessiner la Heatmap

```
title = 'Correlation matrix - features [{0}]'
ImgUtil.save_df_heatmap_plot(corr_matrix,title.format('imputed'))
```



En observant la matrice de correlation, on constate:

- L'inexistence d'aucune correlation forte entre la target 'Binar OP130_Resultat_Global_v' et n'importe quel feature.
- L'existence de correlations positives (0.54, 0.49, 0.48, ..) et negatives (-0.68, -0.45, -0.38, ...) parmi les autres features

b - Matrice de correlation et heatmap avec gestion des valeurs manquantes et rescale:

1 - Appliquer la transformation 'Robust Scaler'

```
XY_data_transformed_scaled =  
tsf.robust_scaler_transform(XY_data_transformed.drop(columns=Const.Binar_OP130_Resultat_Global_v, axis=1))
```

2 - Rajouter la target à la dataframe

```
XY_data_transformed_scaled[Const.Binar_OP130_Resultat_Global_v] =  
XY_data_transformed[Const.Binar_OP130_Resultat_Global_v]
```

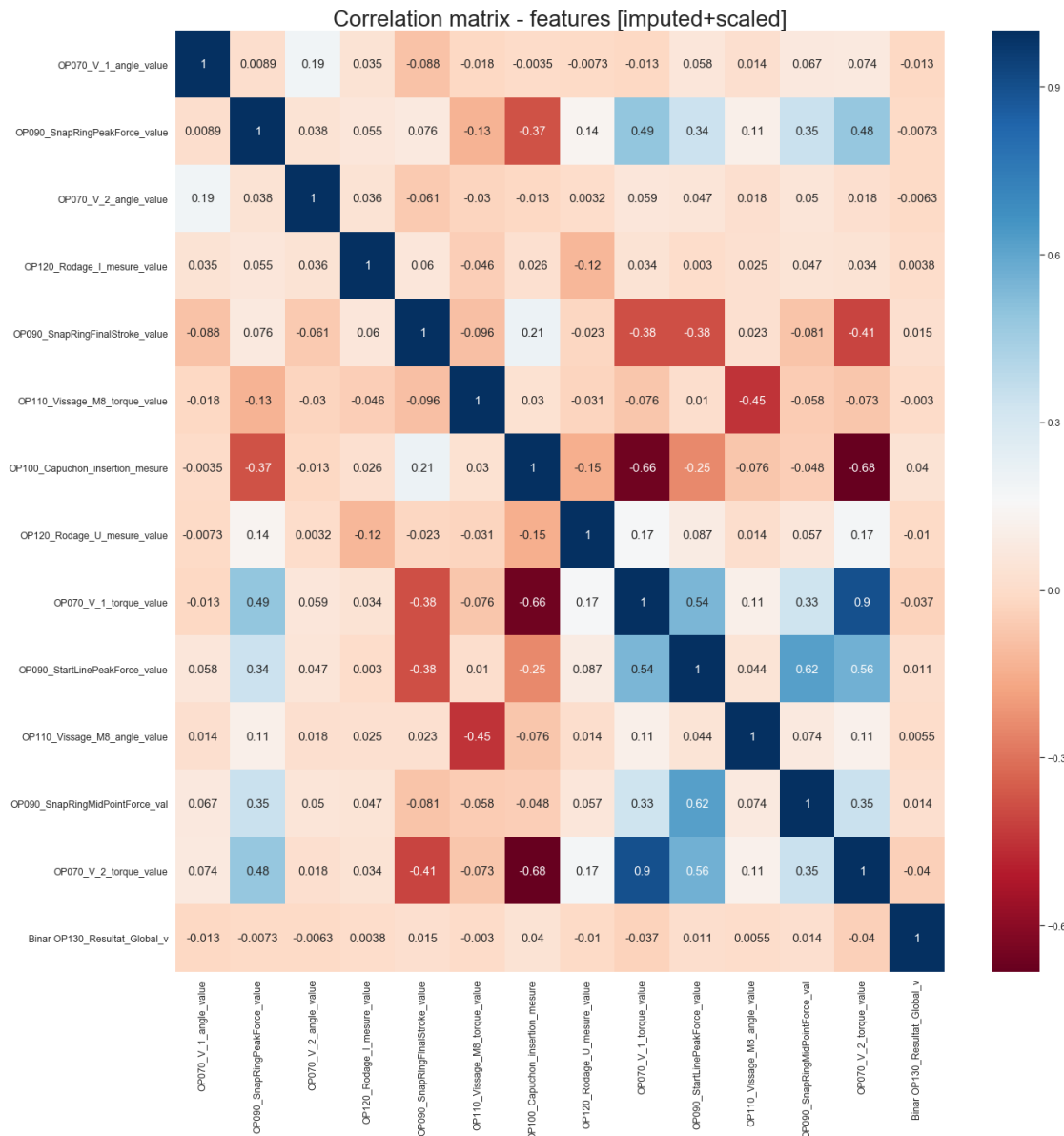
3 - Correlation entre la target "Binar OP130_Resultat_Global_v" et les autres attributs

```
corr_matrix_scaled = XY_data_transformed_scaled.corr()  
corr_matrix_scaled[Const.Binar_OP130_Resultat_Global_v].sort_values(ascending=False)
```

```
Binar_OP130_Resultat_Global_v      1.000000  
OP100_Capuchon_insertion_mesure    0.040366  
OP090_SnapRingFinalStroke_value    0.015148  
OP090_SnapRingMidPointForce_val    0.014273  
OP090_StartLinePeakForce_value     0.010720  
OP110_Vissage_M8_angle_value       0.005470  
OP120_Rodage_I_mesure_value         0.003763  
OP110_Vissage_M8_torque_value      -0.002984  
OP070_V_2_angle_value              -0.006342  
OP090_SnapRingPeakForce_value      -0.007290  
OP120_Rodage_U_mesure_value        -0.010492  
OP070_V_1_angle_value              -0.012793  
OP070_V_1_torque_value              -0.037438  
OP070_V_2_torque_value              -0.039752  
Name: Binar_OP130_Resultat_Global_v, dtype: float64
```

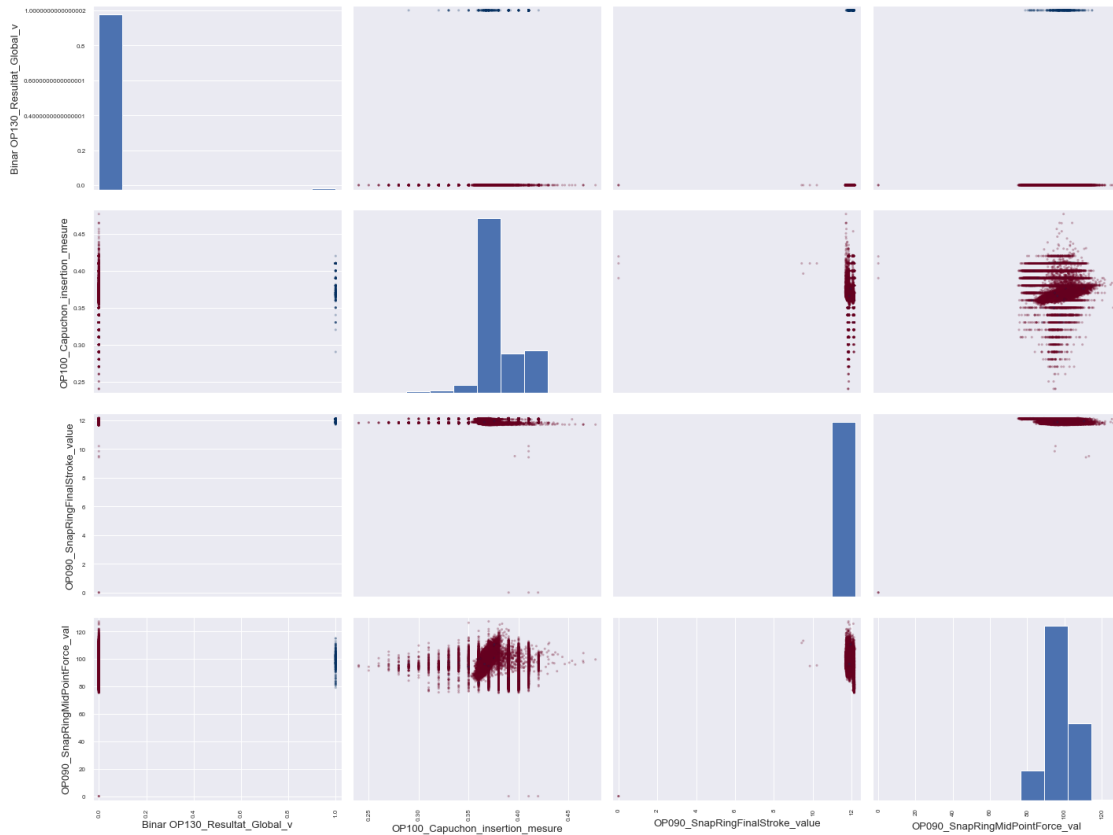
4 - Dessiner la Heatmap

```
ImgUtil.save_df_heatmap_plot(corr_matrix,title.format('imputed+scaled'))
```



c - Nuage de points entre la target 'Binar OP130 Resultat Global v' et les autres features:

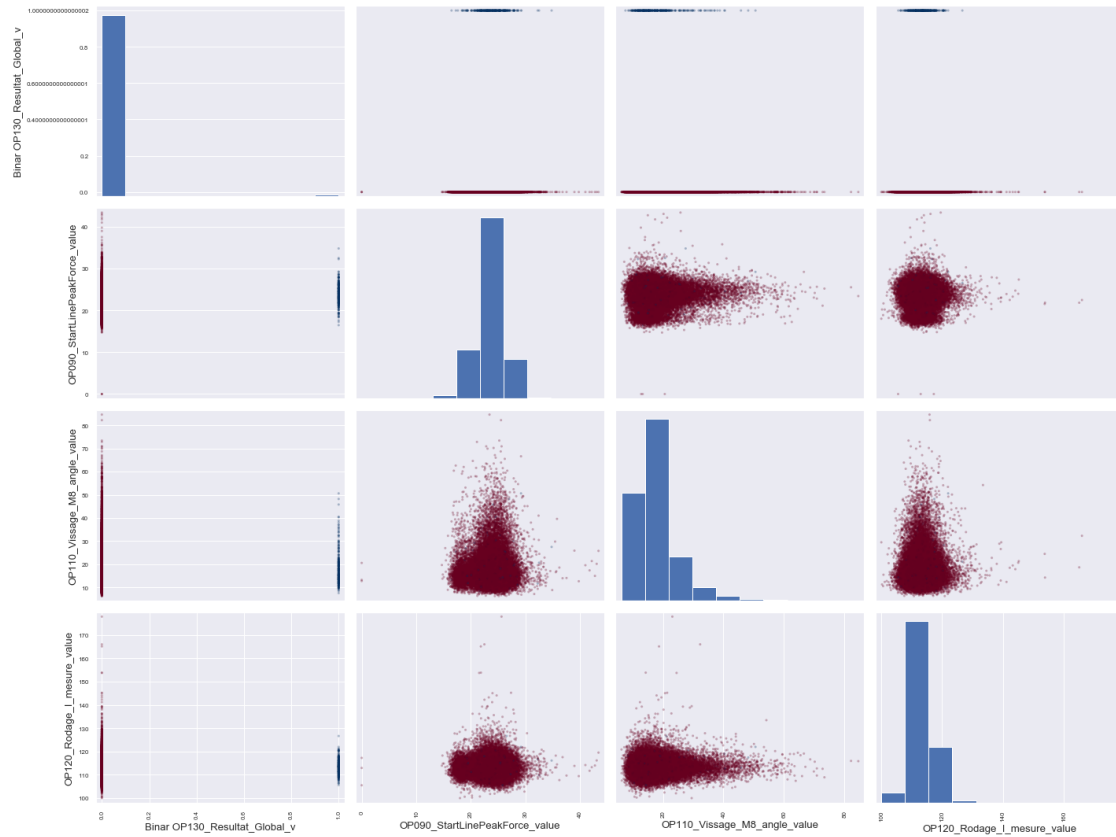
```
features = [Const.Binar_OP130_Resultat_Global_v,
            Const.OP100_Capuchon_insertion_mesure, Const.OP090_SnapRingFinalStroke_value,
            Const.OP090_SnapRingMidPointForce_val]
ImgUtil.save_df_scatter_matrix_plot(XY_data_transformed[features], "XY_data_imputed_corr_pos_1",
cfield=Const.Binar_OP130_Resultat_Global_v)
```



La diagonale allant du coin-gauche-haut au coin-droite-bas représente des barres droites d'histogramme, ces graphes représentent le nombre d'observations d'une feature (ou de la target) en fonction des différentes valeurs que cette feature peut prendre. Le nuage rouge représente les démarreurs étiquetés OK(O) et le bleu représente les KO(1)

NB: On constate que le graphe correspondant à la target (Binar_OP130_Resultat_Global_v) représente une distribution fortement déséquilibrée entre les 2 valeurs '0' et '1' que peut prendre la target.

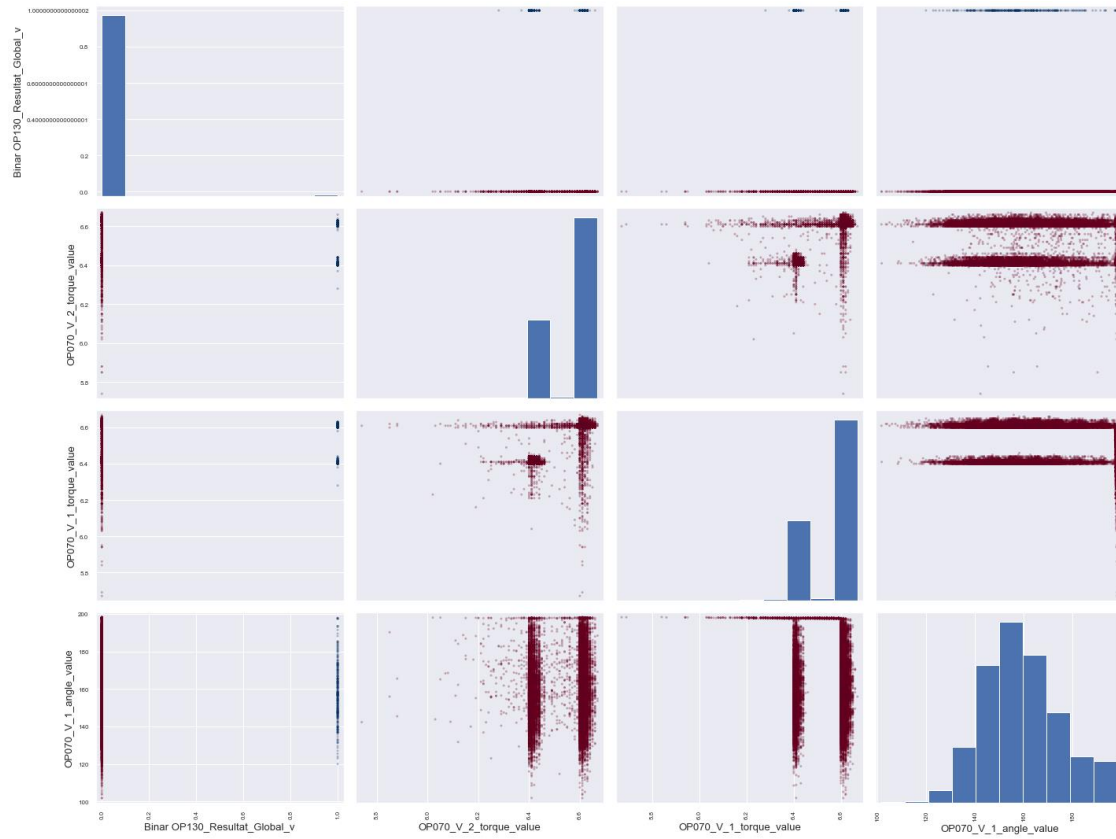
```
features = [Const.Binar_OP130_Resultat_Global_v,
            Const.OP090_StartLinePeakForce_value, Const.OP110_Vissage_M8_angle_value,
            Const.OP120_Rodage_I_mesure_value]
ImgUtil.save_df_scatter_matrix_plot(XY_data_transformed[features], "XY_data_imputed_corr_pos_2",
cfield=Const.Binar_OP130_Resultat_Global_v)
```



```

features = [Const.Binar_OP130_Resultat_Global_v,
            Const.OP070_V_2_torque_value, Const.OP070_V_1_torque_value,
            Const.OP070_V_1_angle_value]
ImgUtil.save_df_scatter_matrix_plot(XY_data_transformed[features], "XY_data_imputed_neg_1",
cfield=Const.Binar_OP130_Resultat_Global_v)

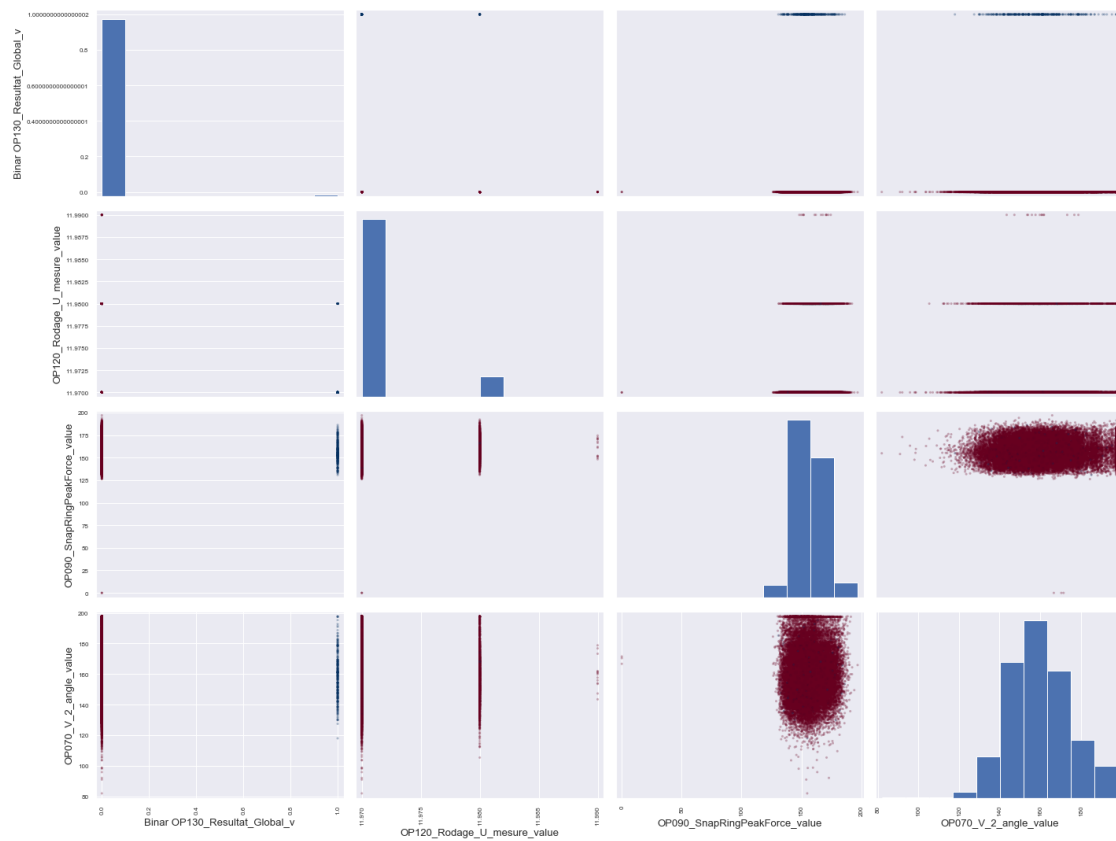
```



```

features = [Const.Binar_OP130_Resultat_Global_v,
            Const.OP120_Rodage_U_mesure_value, Const.OP090_SnapRingPeakForce_value,
            Const.OP070_V_2_angle_value]
ImgUtil.save_df_scatter_matrix_plot(XY_data_transformed[features], "XY_data_imputed_neg_2",
cfield=Const.Binar_OP130_Resultat_Global_v)

```



D'après les graphes, on constate qu'il n'y aucune relation lineaire !!

8 - Feature Engineering/Sélection et choix faits/Hypothèses choisies : TODO

9 - Analyse de la target

a - Vérification de l'équilibre des données:

```
starter_count = len(Y_data[Const.Binar_OP130_Resultat_Global_v])
starter_count_ok = Y_data[Const.Binar_OP130_Resultat_Global_v].value_counts()[0]
starter_count_ko = Y_data[Const.Binar_OP130_Resultat_Global_v].value_counts()[1]
#
print(f'Nombre total des démarreurs : {starter_count}')
print(f'Nombre total des démarreurs OK => Nombre de Classes Negatives : {starter_count_ok} soit
{round(starter_count_ok/starter_count * 100,2)} % du dataset')
print(f'Nombre total des démarreurs KO => Nombre de Classes Positives : {starter_count_ko} soit
{round(starter_count_ko/starter_count * 100,2)} % du dataset')
```

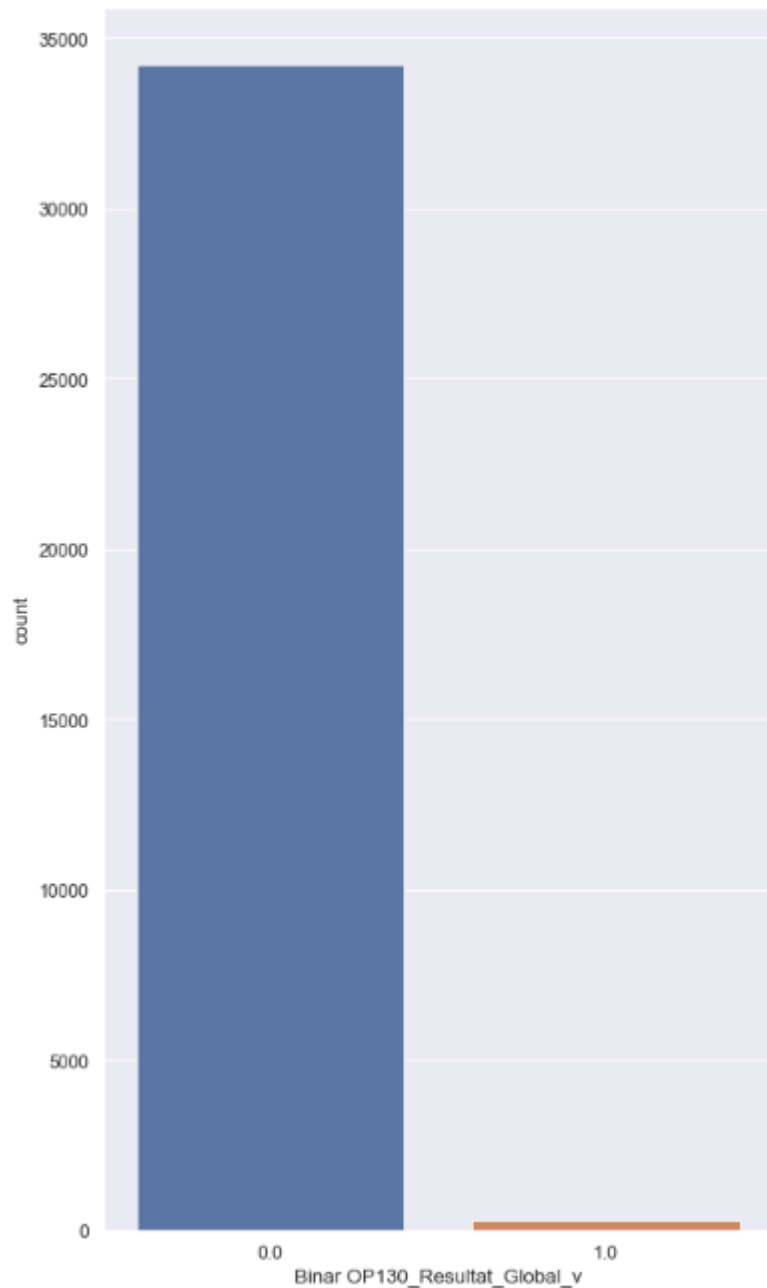
Nombre total des démarreurs : 34515

Nombre total des démarreurs OK => Nombre de Classes Negatives : 34210 soit
99.12 % du dataset

Nombre total des démarreurs KO => Nombre de Classes Positives : 305 soit 0.88
% du dataset

b - Distribution du dataset selon les classes de la target:

```
plt.figure(figsize=(8, 15))  
sns.countplot(Const.Binar_OP130_Resultat_Global_v, data=XY_data_transformed)
```

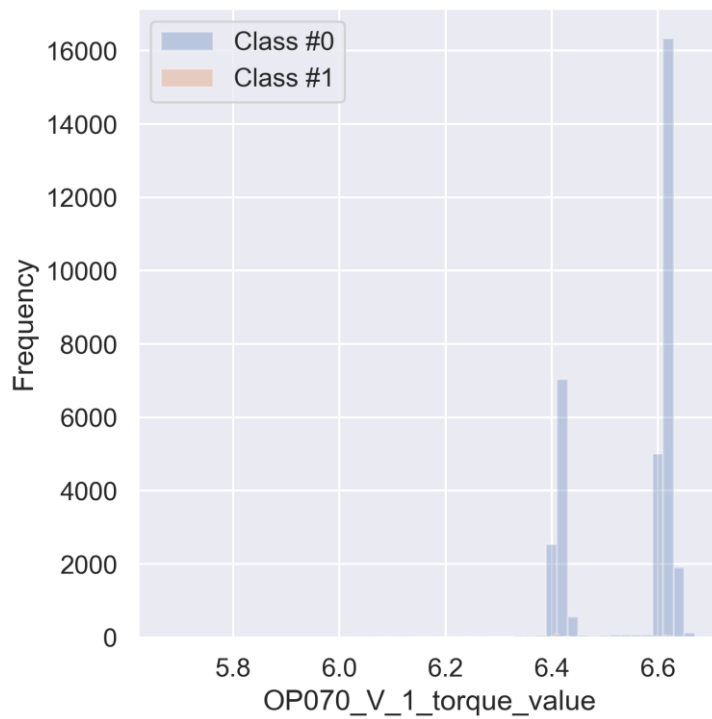
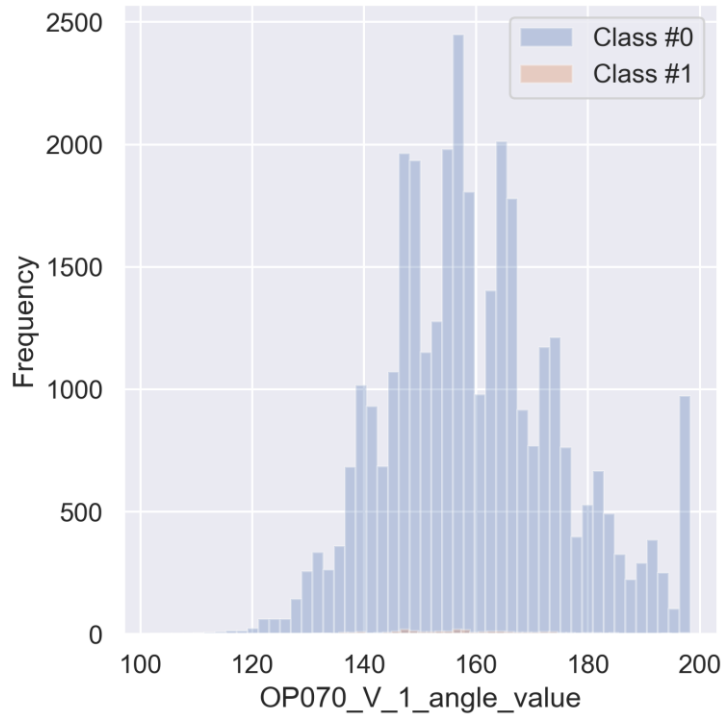


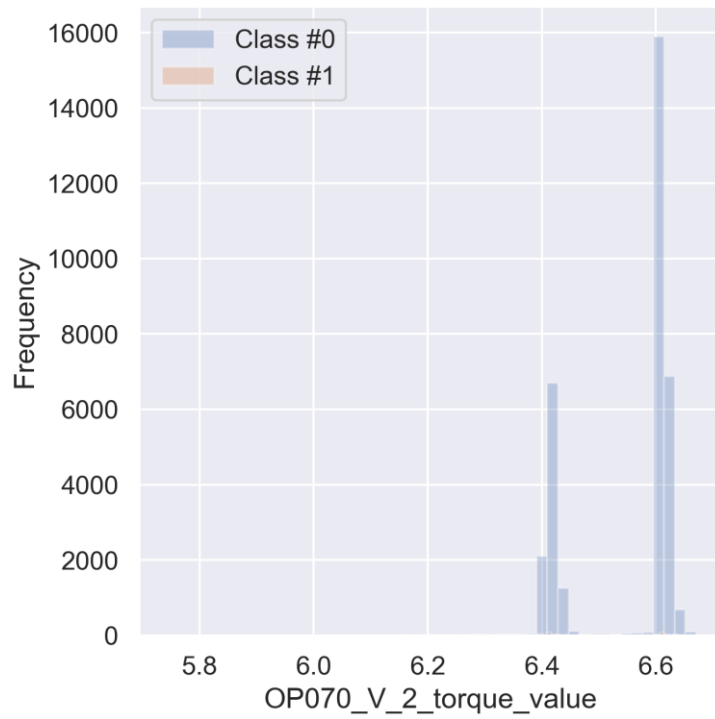
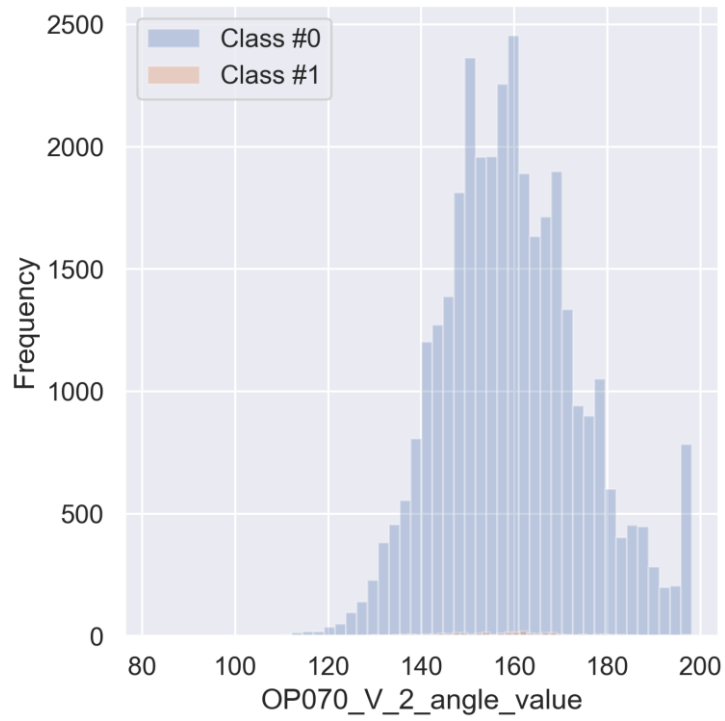
On constate que le jeu de données est fortement déséquilibré.\ Presque totalité des démarreurs (99.12%) ne sont pas défectueux lors de la sortie de la ligne de production.

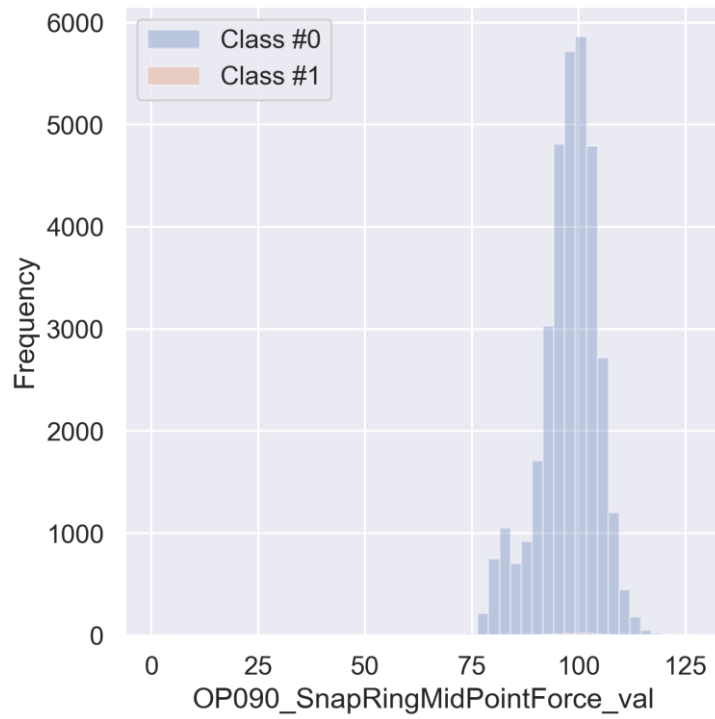
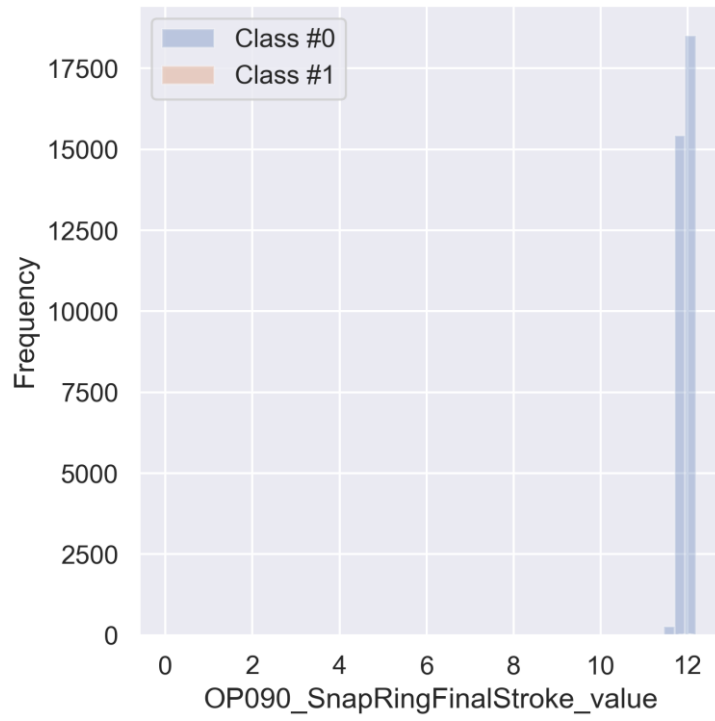
En utilisant cette base de données comme base pour les modèles prédictifs et pour les analyses, on pourrait obtenir beaucoup d'erreurs par des algorithmes inadaptés car ils 'supposeront' que les 'demarreurs' ne sont pas défectueux.\ On cherche un modèle capable de détecter les patterns qui prédisent les défauts sur les lignes de production du démarreur.

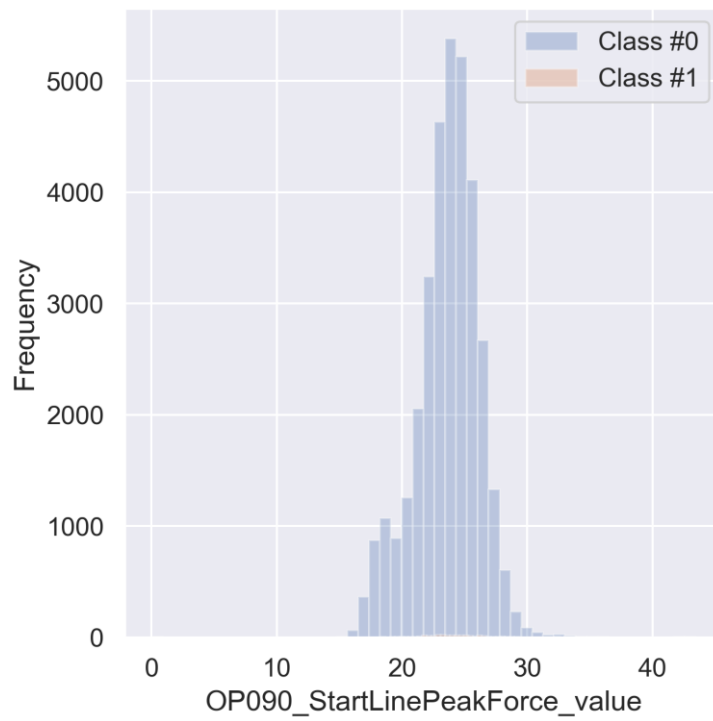
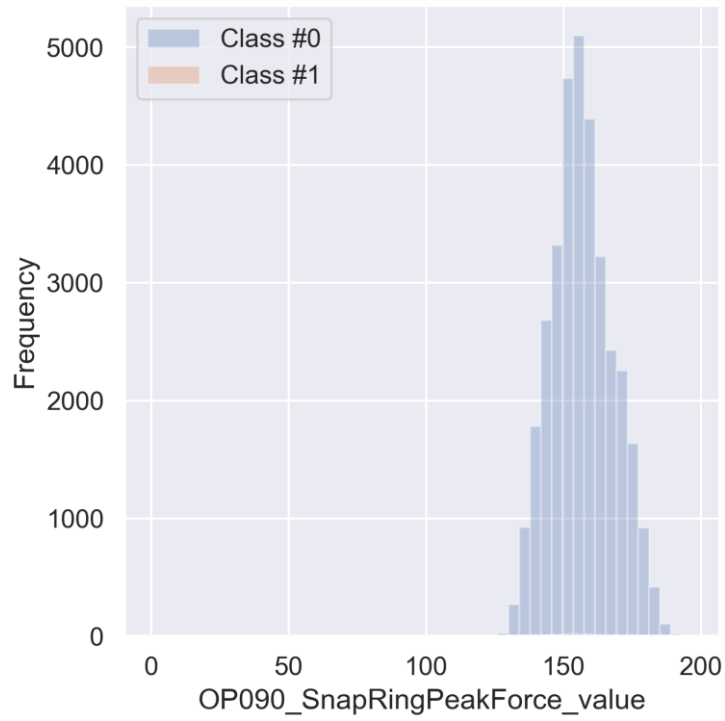
c - Histogramme de distribution du jeu de données selon les classes de la target:

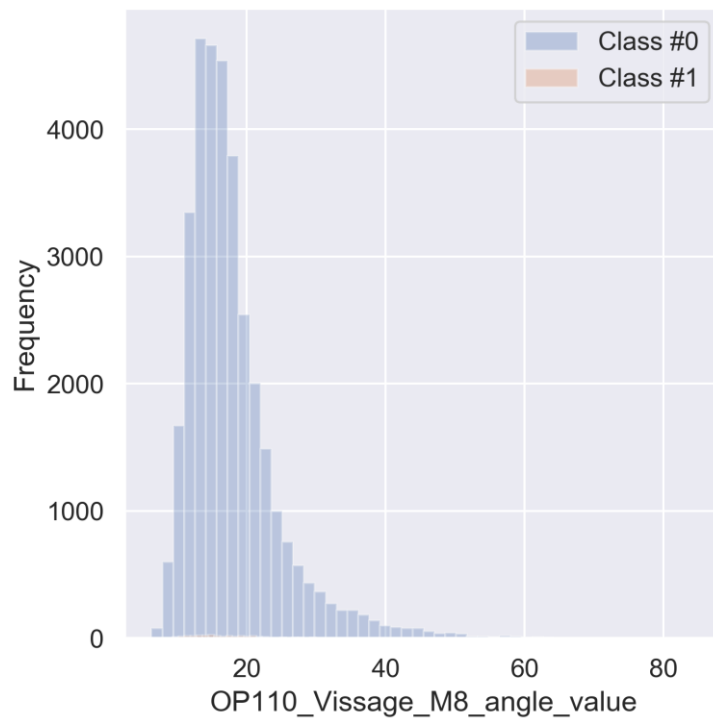
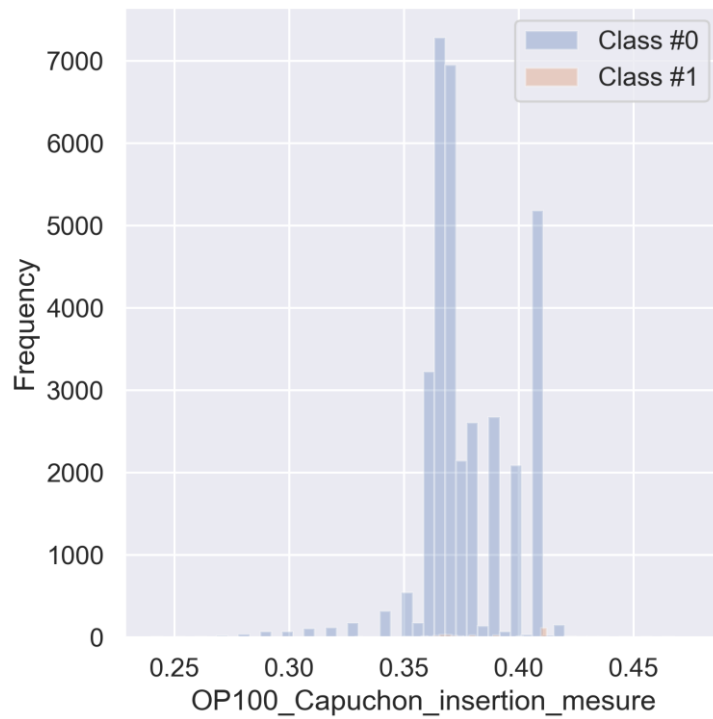
```
ImgUtil.save_df_XY_hist_plot(XY_data_transformed, "XY_imputed",  
y_target_name=Const.Binar_OP130_Resultat_Global_v)
```

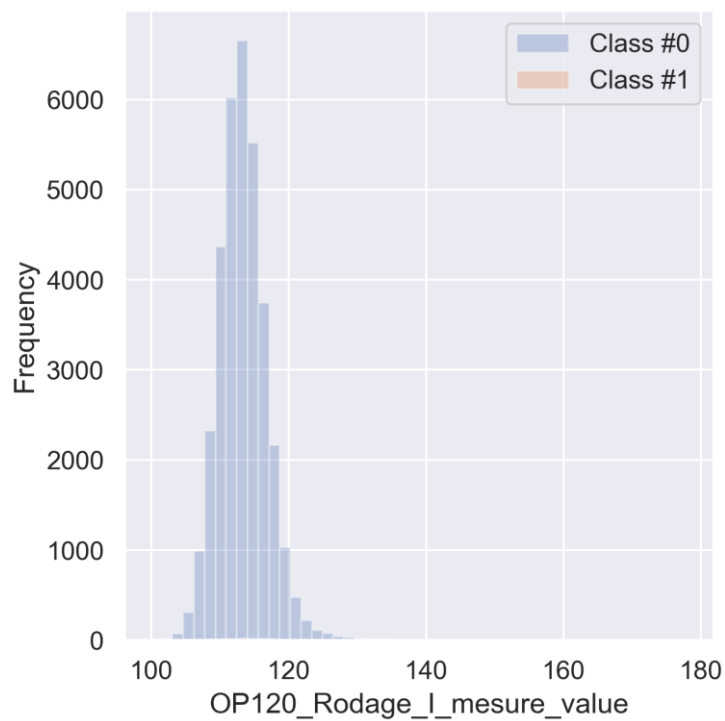
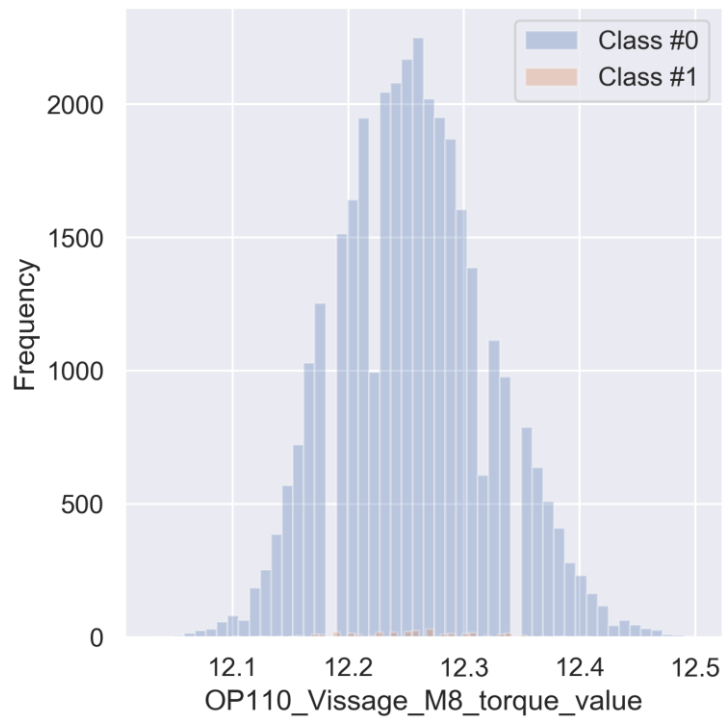


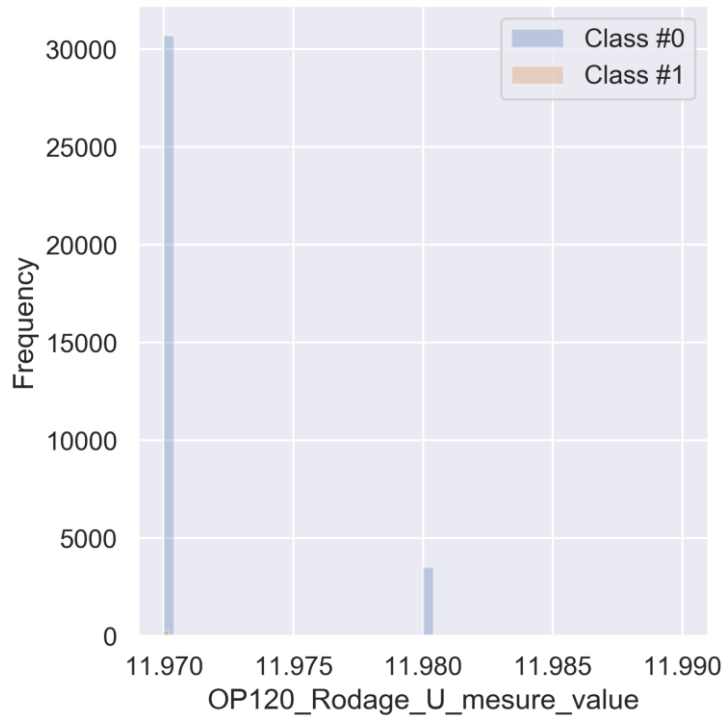












On constate que la classe minoritaire se retrouve délimité à l'intérieur d'une plage de valeurs pour certains features. (ex: OP070_V_1/2_angle_value, OP110_Vissage_M8_torque_Value). Pour cela, il faudrait vérifier l'impact si on transforme ces features numériques continues en des features catégoriques mettant en avant l'existence de la classe minoritaire KO pour ces catégories.

10 - Analyse de la target après un oversampling SMOTE

a - Regénération SMOTE de la classe minoritaire de la target:

```
sm = SMOTE(sampling_strategy='minority', random_state=7)
#
oversampled_X, oversampled_Y = sm.fit_sample(XY_data_transformed.drop(Const.Binar_OP130_Resultat_Global_v,
axis=1),
```

```
XY_data_transformed[Const.Binar_OP130_Resultat_Global_v])
oversampled_XY = pd.concat([pd.DataFrame(oversampled_X), pd.DataFrame(oversampled_Y)], axis=1)
oversampled_XY.columns = XY_data_transformed_scaled.columns
```

```
oversampled_XY.head()
```

	OP070_V_1_angle_value	OP090_SnapRingPeakForce_value	OP070_V_2_angle_value	OP120_Rodage_I_mesure_value	OP090_SnapRingFinalStroke_value
0	180.4	190.51	173.1	113.64	12.04
1	138.7	147.70	163.5	109.77	12.12
2	180.9	150.87	181.2	109.79	11.88
3	173.5	159.56	151.8	113.25	11.82
4	174.5	172.29	177.5	112.88	12.07

	OP110_Vissage_M8_torque_value	OP100_Capuchon_insertion_mesure	OP120_Rodage_U_mesure_value	OP070_V_1_torque_value	OP090_StartLinePeakForce_value
	12.16	0.373146	11.97	6.62	26.37
	12.19	0.390000	11.97	6.41	21.03
	12.24	0.370876	11.97	6.62	25.81
	12.35	0.390000	11.97	6.62	24.62
	12.19	0.368966	11.97	6.62	29.22

OP110_Vissage_M8_angle_value	OP090_SnapRingMidPointForce_val	OP070_V_2_torque_value	OP130_Resultat_Global_v	Binar
18.8	109.62	6.60		0.0
18.5	105.48	6.40		0.0
17.5	100.03	6.61		0.0
15.6	104.94	6.61		0.0
33.6	99.19	6.61		0.0

b - Statistique descriptive du nouveau dataset:

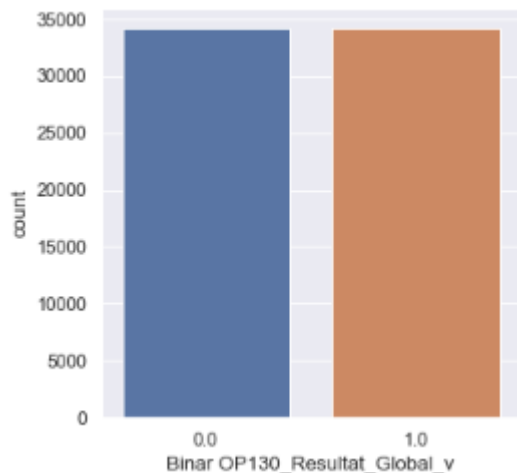
```
oversampled_XY.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
OP070_V_1_angle_value	68420.0	158.781088	14.897773	101.80	148.200000	157.200000	167.800000	198.300000
OP090_SnapRingPeakForce_value	68420.0	156.387713	10.917060	0.00	149.547469	155.290000	163.090000	196.920000
OP070_V_2_angle_value	68420.0	159.251503	14.508312	82.00	149.000000	158.900000	168.300000	198.100000
OP120_Rodage_I_mesure_value	68420.0	113.425963	3.141876	99.99	111.370000	113.354241	115.303464	177.950000
OP090_SnapRingFinalStroke_value	68420.0	11.983641	0.137183	0.00	11.890000	12.036788	12.077263	12.190000
OP110_Vissage_M8_torque_value	68420.0	12.256063	0.058422	12.03	12.214422	12.256656	12.291956	12.500000
OP100_Capuchon_insertion_mesure	68420.0	0.382027	0.019601	0.24	0.367719	0.376890	0.400000	0.476894
OP120_Rodage_U_mesure_value	68420.0	11.970863	0.002622	11.97	11.970000	11.970000	11.970000	11.990000
OP070_V_1_torque_value	68420.0	6.528569	0.096079	5.67	6.410000	6.600000	6.610000	6.670000
OP090_StartLinePeakForce_value	68420.0	23.771639	2.352532	0.00	22.449532	23.890000	25.280000	43.410000
OP110_Vissage_M8_angle_value	68420.0	17.876205	6.390749	6.30	13.797719	16.400445	20.100000	84.600000
OP090_SnapRingMidPointForce_val	68420.0	98.209668	6.177013	0.00	95.300000	98.816196	102.120000	127.300000
OP070_V_2_torque_value	68420.0	6.530620	0.094590	5.74	6.416717	6.600000	6.610000	6.670000
Binar OP130_Resultat_Global_v	68420.0	0.500000	0.500004	0.00	0.000000	0.500000	1.000000	1.000000

c - Nouvelle distribution équilibrée du nouveau dataset:

```
plt.figure(figsize=(5, 5))
```

```
sns.countplot(Const.Binar_OP130_Resultat_Global_v, data=oversampled_XY)
```



d - Matrice de corrélation et heatmap du nouveau dataset:

3 - Correlation entre la target "Binar OP130_Resultat_Global_v" et les autres attributs

```
corr_matrix_oversampled = oversampled_XY.corr()
```

```
corr_matrix_oversampled[Const.Binar_OP130_Resultat_Global_v].sort_values(ascending=False)
```

```
Binar OP130_Resultat_Global_v    1.000000
OP100_Capuchon_insertion_mesure  0.229819
OP090_SnapRingFinalStroke_value  0.099822
```

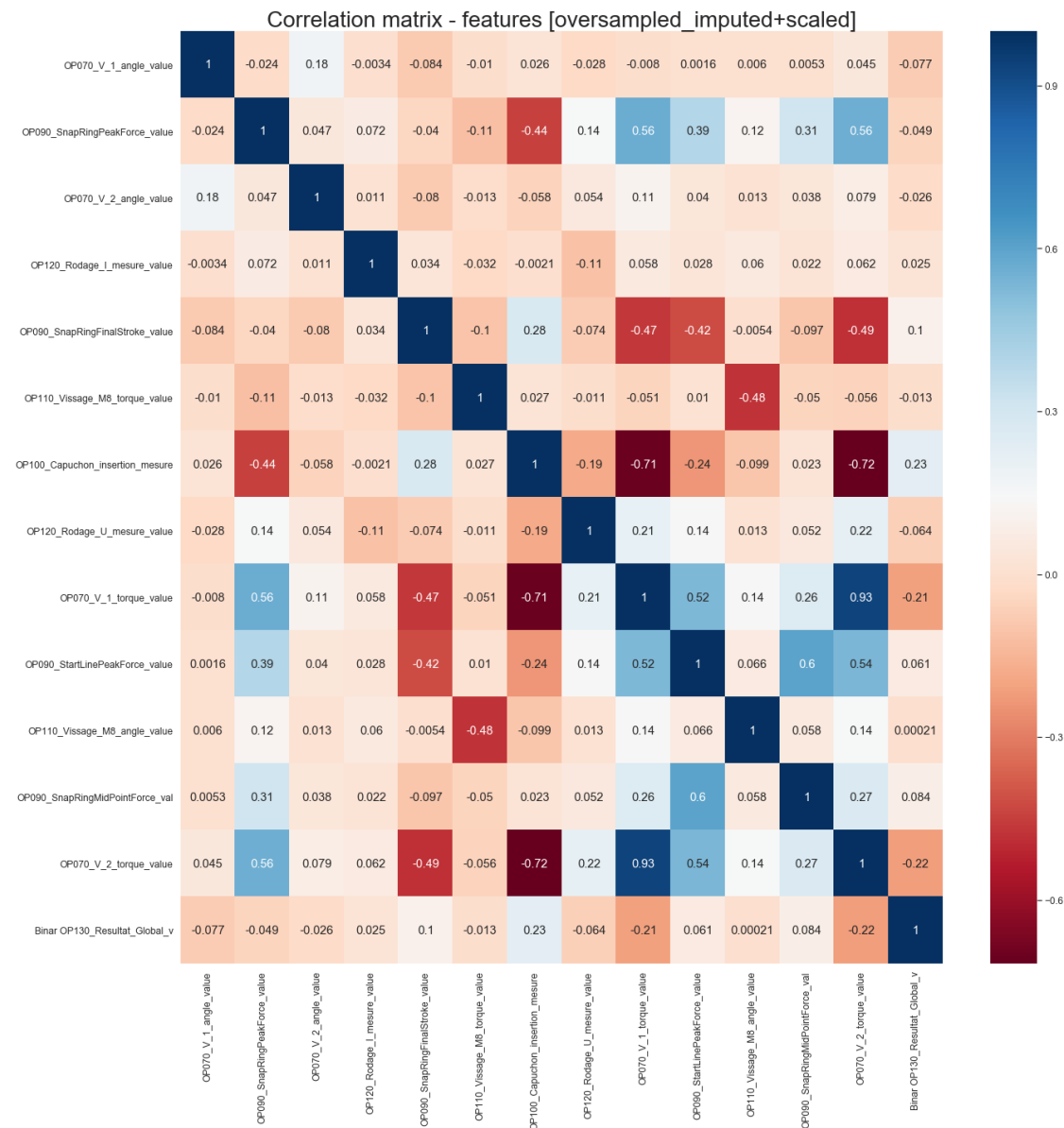
```

OP090_SnapRingMidPointForce_val    0.083845
OP090_StartLinePeakForce_value     0.061239
OP120_Rodage_I_mesure_value        0.024506
OP110_Vissage_M8_angle_value       0.000205
OP110_Vissage_M8_torque_value      -0.012686
OP070_V_2_angle_value              -0.025901
OP090_SnapRingPeakForce_value      -0.049015
OP120_Rodage_U_mesure_value        -0.064017
OP070_V_1_angle_value              -0.076841
OP070_V_1_torque_value             -0.210033
OP070_V_2_torque_value             -0.217814
Name: Binar OP130_Resultat_Global_v, dtype: float64

```

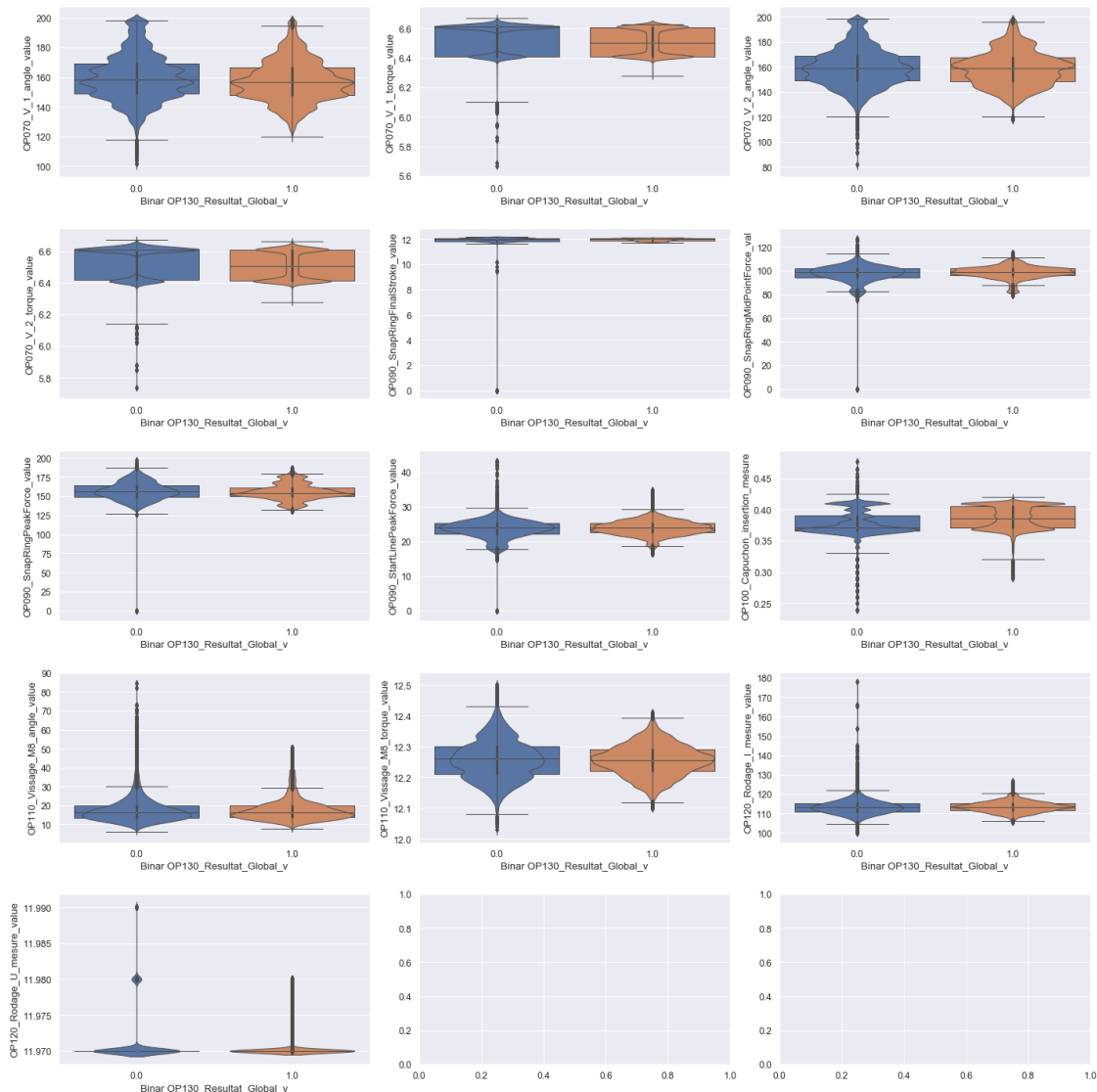
4 - Dessiner la Heatmap

```
ImgUtil.save_df_heatmap_plot(corr_matrix_oversampled,title.format('oversampled_imputed+scaled'))
```



e - Violon et boîte à moustaches des features du nouveau dataset:

```
ImgUtil.save_df_XY_violin_plot(oversampled_XY, Const.Binar_OP130_Resultat_Global_v,
'XY_oversampled_data_distribution', 3)
```



f - Ratio d'observations ayant des features en outlier du nouveau dataset:

```
Q1 = oversampled_X.quantile(0.25)
Q3 = oversampled_X.quantile(0.75)
IQR = Q3 - Q1
#
outliers = ((oversampled_X < (Q1 - 1.5 * IQR)) |(oversampled_X > (Q3 + 1.5 * IQR))).any(axis=1)
print(f"Le ratio d'outlier est de {len(oversampled_X[outliers].index)/len(oversampled_X.index)}")
```

Le ratio d'outlier est de 0.2762496346097632

Le nombre d'outlier est considérable, à peu près 25% des données => On ne peut pas supprimer les observations correspondantes.

11 - Modèle à base d'arbre : Balanced Random Forest Classifier:

Définissons un ensemble de clés de classifieur afin d'y accéder plus facilement

```
HGBC = HistGradientBoostingClassifier(max_iter = 100 , max_depth=10, learning_rate=0.10,
l2_regularization=5)
BBC = BalancedBaggingClassifier(base_estimator=HistGradientBoostingClassifier(), n_estimators=300,
sampling_strategy='auto', replacement=False, random_state=48)
BRFC = BalancedRandomForestClassifier(n_estimators = 300 , max_depth=20, random_state=0)
BRFC_ = BalancedRandomForestClassifier(n_estimators = 300 , max_depth=20, random_state=0,
replacement=True)

BRFC_W = BalancedRandomForestClassifier(n_estimators = 300 , max_depth=20, random_state=0,
class_weight={0:1, 1:1})
RUSBoost = RUSBoostClassifier(n_estimators = 8 , algorithm='SAMME.R', random_state=42)
XGBC = xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=1, gamma=0,
learning_rate=0.1, max_delta_step=0, max_depth=10, #max_depth=3,
min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
nthread=None, objective='binary:logistic', random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=100, seed=42,
silent=None, subsample=1, verbosity=1)

KNN = KNeighborsClassifier(3)
RFC = RandomForestClassifier(n_estimators=300, max_depth=10, max_features=10, n_jobs=4, class_weight=
{0:1,1:100})
DTC = DecisionTreeClassifier()
ADABOOST = AdaBoostClassifier()
GBC = GradientBoostingClassifier()
LRC = LogisticRegression(max_iter=500)
# SVC = SVC(kernel="rbf", C=0.025, probability=True)
# GNB = GaussianNB()
# NuSVC = NuSVC(probability=True),
# LinearSVC = LinearSVC(C=0.1, class_weight={'1':100})
# SGDClassifier = SGDClassifier(class_weight='balanced')
```

Chargement du jeu de données training - Commun à tout les modèles

```
# 1 - Rechargement des données
mt_train = XY_metadata([Const.rootDataTrain(), 'traininginputs.csv'], [Const.rootDataTrain(),
'trainingoutput.csv'],
[Const.PROC_TRACEINFO], [Const.PROC_TRACEINFO], Const.Binar_OP130_Resultat_Global_v)
xy_loader = XY_Loader();
X_df, y_df = xy_loader.load_XY_df(mt_train)
```

a – Balanced Random Forest Classifier - Train / Test / Split + F1 et ROC:

```
# 2 - Split Training et Validation
X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_size=0.3, random_state=48,
stratify=y_df)

# 3 - Imputer et Scaler + classifieur
modeler = ValeoModeler()
pred = ValeoPredictor()
pl = Pipeline([('preprocessor', modeler.build_transformers_pipeline(X_train.dtypes)), # --> Imputer +
Scaler
('classifier', BRFC) # --> Balanced Random Forest Classifier
])

# 4 - Fit, train, predict and plot ROC and F1
pl.fit(X_train, y_train)
pred.predict_and_plot(pl, X_test, y_test)

# 5 - Test using ENS data
X_ens = DfUtil.read_csv([Const.rootDataTest() , "testinputs.csv"])
y_ens = pl.predict(X_ens.drop(columns=[Const.PROC_TRACEINFO]))
DfUtil.write_y_csv(X_ens[Const.PROC_TRACEINFO], y_ens, Const.Binar_OP130_Resultat_Global_v,
[Const.rootDataTest() , "testoutput.csv"])
```

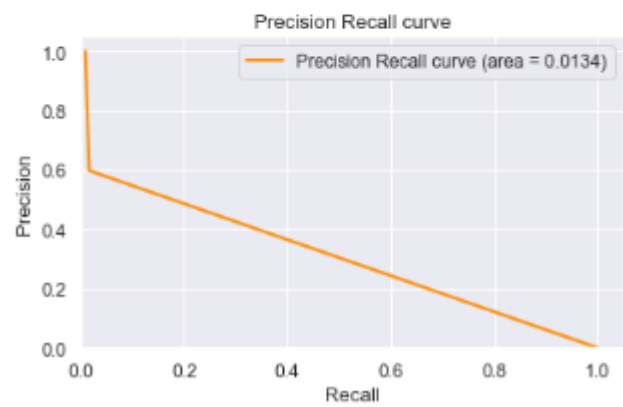
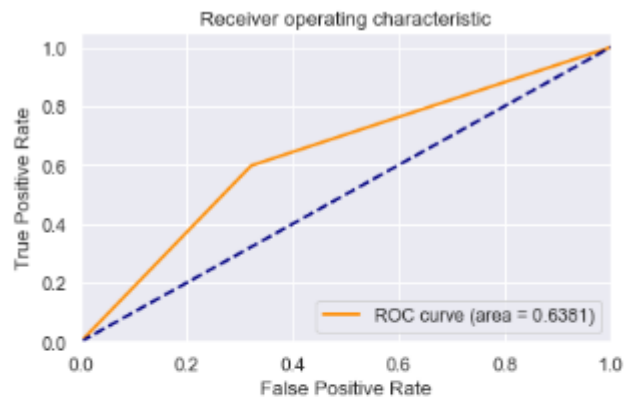
- Model score: 0.6776436504104297
- Accuracy score: 0.6776436504104297
- Balanced accuracy score: 0.6380926205999602 / The balanced accuracy to deal with imbalanced datasets. It is defined as the average of recall obtained on each class.
- Average_precision_score: 0.013370660450719168
- Precision_score: 0.016388557806912993
- Recall score: 0.5978260869565217
- Roc_auc_score: 0.6380926205999602
- F1 score: 0.031902552204176336
- [[6962 3301]/[37 55] - P:0.0164 - R:0.5978 - roc_auc:0.6381 - f1:0.0319
- [[6962 3301]
- [37 55]]

- Classification_report_imbalanced:

		pre	rec	spe	f1	geo	iba
sup							
	0	0.99	0.68	0.60	0.81	0.64	0.41
10263							
	1	0.02	0.60	0.68	0.03	0.64	0.40
92							
avg / total		0.99	0.68	0.60	0.80	0.64	0.41
10355							

- Classification_report:

		precision	recall	f1-score	support
	0	0.99	0.68	0.81	10263
	1	0.02	0.60	0.03	92
accuracy				0.68	10355
macro avg		0.51	0.64	0.42	10355
weighted avg		0.99	0.68	0.80	10355



NB: LA SURFACE F1 CALCULÉE n'est pas correcte ($\text{base} * \text{hauteur} / 2$) => A corriger

b - Balanced Random Forest Classifier - Cross Validation + F1 et ROC:

```
X_train, y_train = X_df, y_df

# 2 - Initialize a CV Split
CV = StratifiedKFold(n_splits=8) # , random_state=48, shuffle=True

# 3 - Imputer et Scaler + classifier
modeler = ValeoModeler()
pred = ValeoPredictor()
BRFC = BalancedRandomForestClassifier(n_estimators = 300 , max_depth=20, random_state=0)
pl = Pipeline([('preprocessor', modeler.build_transformers_pipeline(X_train.dtypes)), # --> Imputer +
Scaler
                ('classifier', BRFC) # --> Balanced Random Forest Classifier
            ])

# 4 - Cross Validate
cv_results = cross_validate(pl, X_train, y_train, cv=CV, scoring=('f1', 'f1_micro', 'f1_macro',
'f1_weighted', 'recall', 'precision', 'average_precision', 'roc_auc'), return_train_score=True,
return_estimator=True)
fitted_estimators = []
for key in cv_results.keys() :
    if str(key) != "estimator" :
        print(f"{key} : {cv_results[key]}")
        fitted_estimators.append(cv_results[key])

fitted_model = cv_results["estimator"][np.argmax(cv_results["test_roc_auc"])]
pred.predict_and_plot(fitted_model,X_test, y_test)

# 5 - Test using ENS data
X_ens = DfUtil.read_csv([Const.rootDataTest() , "testinputs.csv"])
y_ens = fitted_model.predict(X_ens.drop(columns=[Const.PROC_TRACEINFO]))
DfUtil.write_y_csv(X_ens[Const.PROC_TRACEINFO], y_ens, Const.Binar_OP130_Resultat_Global_v,
[Const.rootDataTest() , "testoutput.csv"])
```

```
fit_time : [4.20418215 3.88397908 3.76627254 3.901299    3.95574546 3.66990876
4.14086699 4.21298003]
score_time : [0.34399605 0.31799507 0.32755017 0.2992568  0.3776269
0.30384541
0.35573673 0.38599515]
test_f1 : [0.02617801 0.03069054 0.03129445 0.03215434 0.03556772 0.02931379
0.03125    0.02610587]
train_f1 : [0.05155933 0.04991121 0.05410353 0.0492302  0.05065933 0.0512476
0.05085714 0.05507426]
test_f1_micro : [0.65515643 0.64866744 0.68435689 0.65113584 0.67315716
0.6622624
0.66944831 0.68868799]
train_f1_micro : [0.6747351  0.6634106  0.69201987 0.65852124 0.66865336
0.67265985
0.67001093 0.69663256]
test_f1_macro : [0.40832978 0.408071    0.42137812 0.40968667 0.41940261
0.41243997
0.41598833 0.42041947]
train_f1_macro : [0.42763407 0.42269408 0.43508467 0.420559    0.42498069
0.4267282
0.42557285 0.43719279]
test_f1_weighted : [0.78375072 0.77880466 0.80441046 0.78056799 0.79647547
```

```

0.7888166
  0.7939487  0.80778642]
train_f1_weighted : [0.79705901 0.78888536 0.80935449 0.78532214 0.79268348
0.79556972
  0.79366302 0.81255488]
test_recall : [0.52631579 0.63157895 0.56410256 0.65789474 0.68421053
0.57894737
  0.60526316 0.47368421]
train_recall : [1. 1. 1. 1. 1. 1. 1. 1.]
test_precision : [0.01342282 0.01572739 0.01609364 0.01647989 0.01825843
0.01503759
  0.01603905 0.01342282]
train_precision : [0.02646184 0.02559433 0.02780391 0.02523629 0.02598793
0.02629765
  0.02609206 0.02831689]
test_average_precision : [0.01460394 0.01852067 0.0140107  0.02250251
0.03696366 0.01684649
  0.01927391 0.01715631]
train_average_precision : [0.36489053 0.38781819 0.34316751 0.43090551
0.37367985 0.36632035
  0.38315343 0.36073824]
test_roc_auc : [0.61249892 0.64770252 0.63006404 0.70079637 0.73938383
0.65292206
  0.66054724 0.62940648]
train_roc_auc : [0.95913257 0.96622369 0.95572754 0.96683853 0.96490919
0.96163445
  0.96284347 0.96020734]
- Model score: 0.6728958423873678
- Accuracy score: 0.6728958423873678
- Balanced accuracy score: 0.8106188392810079 / The balanced accuracy to deal
with imbalanced datasets. It is defined as the average of recall obtained on
each class.
- Average_precision_score: 0.024277023825104087
- Precision_score: 0.025075659316904454
- Recall score: 0.9508196721311475
- Roc_auc_score: 0.8106188392810079
- F1 score: 0.04886267902274642
- [4587 2255]/[ 3 58] - P:0.0251 - R:0.9508 - roc_auc:0.8106 - f1:0.0489
- [[4587 2255]
  [ 3 58]]

```

```

- Classification_report_imbalanced:

```

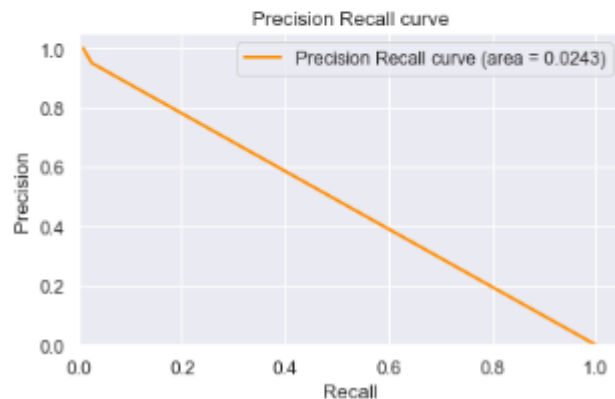
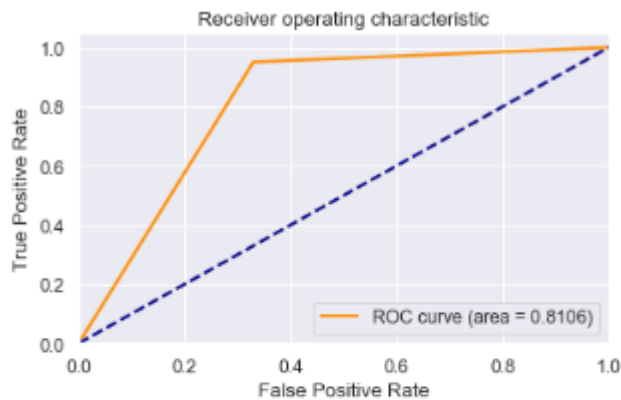
		pre	rec	spe	f1	geo	iba
sup							
	0	1.00	0.67	0.95	0.80	0.80	0.62
6842	1	0.03	0.95	0.67	0.05	0.80	0.66
61							
avg / total		0.99	0.67	0.95	0.80	0.80	0.62

6903

- classification_report:

	precision	recall	f1-score	support
0	1.00	0.67	0.80	6842
1	0.03	0.95	0.05	61
accuracy			0.67	6903
macro avg	0.51	0.81	0.43	6903
weighted avg	0.99	0.67	0.80	6903

```
- precision_recall_curve: (array([0.00883674, 0.02507566, 1.          ]),  
array([1.          , 0.95081967, 0.          ]), array([0, 1], dtype=int64))  
- precision_recall_fscore_support: (array([0.99934641, 0.02507566]),  
array([0.67041801, 0.95081967]), array([0.80248425, 0.04886268]),  
array([6842, 61], dtype=int64))  
- roc_curve: (array([0.          , 0.32958199, 1.          ]), array([0.          ,  
0.95081967, 1.          ]), array([2, 1, 0], dtype=int64))
```



NB: LA SURFACE F1 CALCULÉE n'est pas correcte (base x hauteur / 2) => A corriger

12 - Modèle à base de distance : Logistique regression avec SMOTE oversampling

a – Logistique regression avec SMOTE - Train / Test / Split + F1 et ROC:

```
X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_size=0.3, random_state=48,
stratify=y_df)

# 3 - Imputer et Scaler + SMOTE + Logistic Regression
modeler = ValeoModeler()
pred = ValeoPredictor()
pl = Pipeline(['preprocessor', modeler.build_transformers_pipeline(X_train.dtypes)],      # -->
Imputer + Scaler
              ('imbalancer_resampler', SMOTE(sampling_strategy='minority', random_state=7)),  # -->
SMOTE oversampling
              ('classifier', LRC) # --> Logistic Regression Classifier
              ])
pl.fit(X_train, y_train)
pred.predict_and_plot(pl,X_test, y_test)
```

```
# 5 - Test using ENS data
X_ens = DfUtil.read_csv([Const.rootDataTest() , "testinputs.csv"])
y_ens = pl.predict(X_ens.drop(columns=[Const.PROC_TRACEINFO]))
DfUtil.write_y_csv(X_ens[Const.PROC_TRACEINFO], y_ens, Const.Binar_OP130_Resultat_Global_v,
[Const.rootDataTest() , "testoutput.csv"])

- Model score: 0.6560115886045389
- Accuracy score: 0.6560115886045389
- Balanced accuracy score: 0.6271796321950103 / The balanced accuracy to deal
with imbalanced datasets. It is defined as the average of recall obtained on
each class.
```

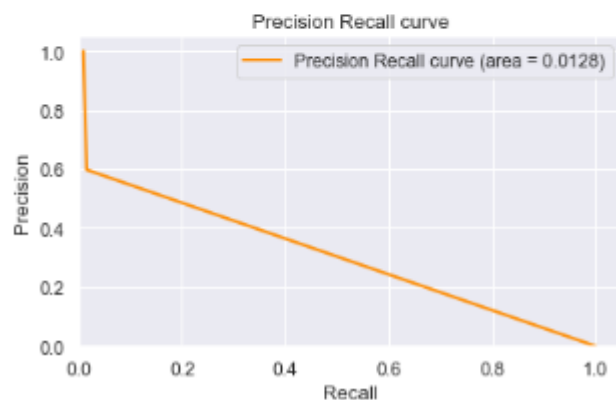
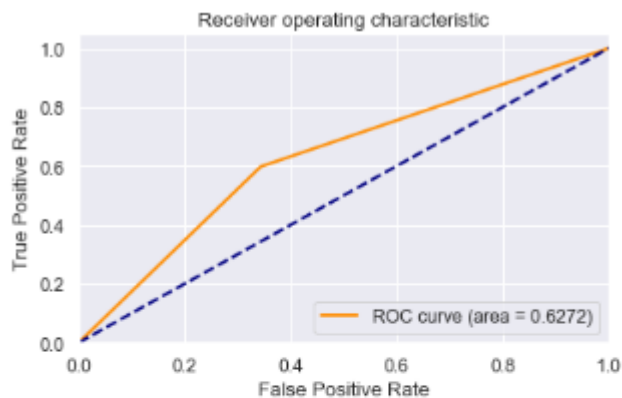
```
- Average_precision_score: 0.012757632055707119
- Precision_score: 0.015363128491620111
- Recall score: 0.5978260869565217
- Roc_auc_score: 0.6271796321950104
- F1 score: 0.02995642701525054
- [6738 3525]/[37 55] - P:0.0154 - R:0.5978 - roc_auc:0.6272 - f1:0.0300
- [[6738 3525]
  [ 37 55]]
```

```
- classification_report_imbalanced:
```

		pre	rec	spe	f1	geo	iba
sup							
	0	0.99	0.66	0.60	0.79	0.63	0.39
10263							
	1	0.02	0.60	0.66	0.03	0.63	0.39
92							
avg / total		0.99	0.66	0.60	0.78	0.63	0.39
10355							

```
- classification_report:
      precision    recall  f1-score   support
```

	0	0.99	0.66	0.79	10263
	1	0.02	0.60	0.03	92
accuracy				0.66	10355
macro avg		0.50	0.63	0.41	10355
weighted avg		0.99	0.66	0.78	10355



NB: LA SURFACE F1 CALCULÉE n'est pas correcte (base x hauteur / 2) => A corriger

b - Logistique regression avec SMOTE - Cross Validation + F1 et ROC:

```
X_train, y_train = X_df, y_df
```

```
# 2 - Initialize a CV Split
```

```
CV = StratifiedKFold(n_splits=8) # , random_state=48, shuffle=True
```

```
# 3 - Imputer et Scaler + classifier
```

```
modeler = ValeoModeler()
```

```
pred = ValeoPredictor()
```

```
pl = Pipeline([('preprocessor', modeler.build_transformers_pipeline(X_train.dtypes)), # -->
```

```
Imputer + Scaler
```

```
                ('imbalancer_resampler', SMOTE(sampling_strategy='minority', random_state=7)), # -->
```

```
SMOTE oversampling
```

```
                ('classifier', LRC) # --> Logistic Regression Classifier
```

```
            ])
```

```
# 4 - Cross Validate
```

```

cv_results = cross_validate(pl, X_train, y_train, cv=CV, scoring=('f1', 'f1_micro', 'f1_macro',
'f1_weighted', 'recall', 'precision', 'average_precision', 'roc_auc'), return_train_score=True,
return_estimator=True)
fitted_estimators = []
for key in cv_results.keys() :
    if str(key) != "estimator" :
        print(f"{key} : {cv_results[key]}")
        fitted_estimators.append(cv_results[key])

fitted_model = cv_results["estimator"][np.argmax(cv_results["test_roc_auc"])]
pred.predict_and_plot(fitted_model,X_test, y_test)

# 5 - Test using ENS data
X_ens = DfUtil.read_csv([Const.rootDataTest() , "testinputs.csv"])
y_ens = fitted_model.predict(X_ens.drop(columns=[Const.PROC_TRACEINFO]))
DfUtil.write_y_csv(X_ens[Const.PROC_TRACEINFO], y_ens, Const.Binar_OP130_Resultat_Global_v,
[Const.rootDataTest() , "testoutput.csv"])

```

```

fit_time : [2.36323762 1.91440773 2.06579304 2.04484057 1.85639429 1.75007534
1.96227741 2.28214383]
score_time : [0.03199959 0.02400422 0.0418551 0.02399921 0.02499723
0.03502941
0.03199911 0.03299737]
test_f1 : [0.02473958 0.02678028 0.0290429 0.03076923 0.0341556 0.03088803
0.0297542 0.03050398]
train_f1 : [0.03089371 0.03005464 0.03097345 0.03165416 0.03028654 0.03054707
0.03066471 0.03193988]
test_f1_micro : [0.65283893 0.62943221 0.65909618 0.64951321 0.64603616
0.65090403
0.65229485 0.66110338]
train_f1_micro : [0.65099338 0.64735099 0.65192053 0.64547532 0.64590576
0.6532234
0.65040893 0.65878613]
test_f1_macro : [0.40678761 0.39896335 0.4111483 0.40842366 0.40873382
0.40899787
0.40897481 0.41258345]
train_f1_macro : [0.4090338 0.40727742 0.40941633 0.40733584 0.40684741
0.40969681
0.40870743 0.41241667]
test_f1_weighted : [0.78210663 0.76459116 0.78634657 0.77942494 0.77671308
0.78044653
0.78151467 0.7879318 ]
train_f1_weighted : [0.78048758 0.7778301 0.78119261 0.7763749 0.7767501
0.7821426
0.78006578 0.78616605]
test_recall : [0.5 0.57894737 0.56410256 0.63157895 0.71052632
0.63157895
0.60526316 0.60526316]
train_recall : [0.62921348 0.61797753 0.63157895 0.65543071 0.62546816
0.61797753
0.62546816 0.63670412]
test_precision : [0.01268358 0.01370717 0.01490515 0.01576873 0.01749838
0.01583113
0.01525199 0.01564626]

```

```

train_precision : [0.01583561 0.01540185 0.01587602 0.01621872 0.015519
0.01566059
0.01571765 0.01638081]
test_average_precision : [0.01615652 0.01554629 0.01585269 0.03865314
0.02021889 0.01731332
0.05091361 0.0200504 ]
train_average_precision : [0.01932375 0.01975579 0.02047861 0.01765319
0.01911704 0.02003393
0.01881078 0.01992643]
test_roc_auc : [0.60311581 0.65022827 0.63546689 0.70105485 0.71515435
0.65763626
0.67102801 0.64764167]
train_roc_auc : [0.68933515 0.679299 0.68218565 0.67620075 0.67052259
0.67673789
0.68010722 0.68292153]
- Model score: 0.6446161274746499
- Accuracy score: 0.6446161274746499
- Balanced accuracy score: 0.6645193370867913 / The balanced accuracy to deal
with imbalanced datasets. It is defined as the average of recall obtained on
each class.
- Average_precision_score: 0.014416439513109575
- Precision_score: 0.016962843295638127
- Recall score: 0.6847826086956522
- Roc_auc_score: 0.6645193370867913
- F1 score: 0.03310562270099843
- [[6612 3651]/[29 63] - P:0.0170 - R:0.6848 - roc_auc:0.6645 - f1:0.0331
- [[6612 3651]
[ 29 63]]

```

```

- classification_report_imbalanced:

```

		pre	rec	spe	f1	geo	iba
sup							
	0	1.00	0.64	0.68	0.78	0.66	0.44
10263							
	1	0.02	0.68	0.64	0.03	0.66	0.44
92							
avg / total		0.99	0.64	0.68	0.78	0.66	0.44
10355							

```

- classification_report:

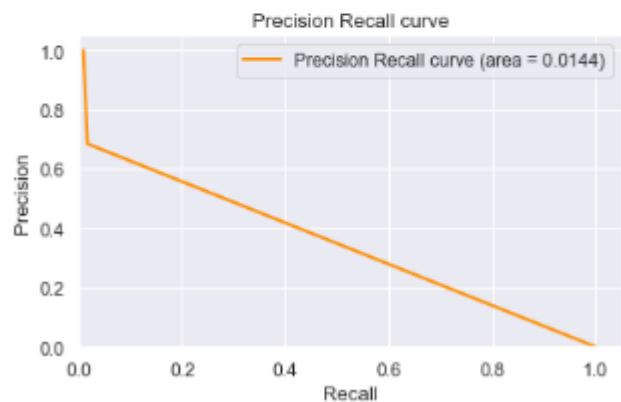
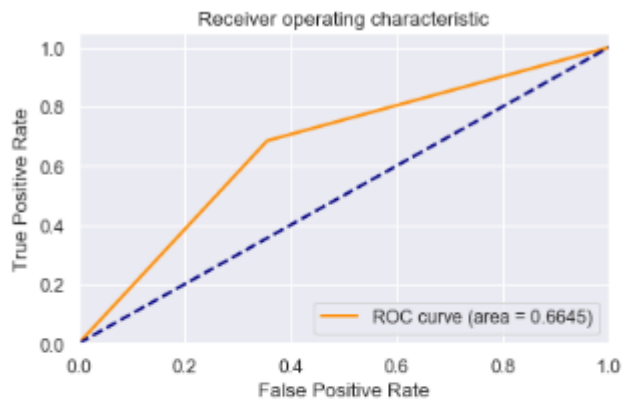
```

		precision	recall	f1-score	support
	0	1.00	0.64	0.78	10263
	1	0.02	0.68	0.03	92
accuracy				0.64	10355
macro avg		0.51	0.66	0.41	10355
weighted avg		0.99	0.64	0.78	10355

```

- precision_recall_curve: (array([0.0088846 , 0.01696284, 1.          ]),
array([1.          , 0.68478261, 0.          ]), array([0, 1], dtype=int64))
- precision_recall_fscore_support: (array([0.99563319, 0.01696284]),
array([0.64425607, 0.68478261]), array([0.78230005, 0.03310562]),
array([10263,    92], dtype=int64))
- roc_curve: (array([0.          , 0.35574393, 1.          ]), array([0.          ,
0.68478261, 1.          ]), array([2, 1, 0], dtype=int64))

```



NB: LA SURFACE F1 CALCULÉE n'est pas correcte (base x hauteur / 2) => A corriger

13 - Modèle à base de réseau de neurones ou de stacking : TODO

14 - Conclusion : TODO

15 - Perspectives : TODO

16 – Annexe : Code Python