```python
 1 from imblearn.ensemble import BalancedBaggingClassifier, RUSBoostClassifier, BalancedRandomForestClassifier
 2 from imblearn.over_sampling import RandomOverSampler, ADASYN, SMOTE, SVMSMOTE, KMeansSMOTE, BorderlineSMOTE
 3 from imblearn.over_sampling.base import BaseOverSampler
 4 from imblearn.pipeline import Pipeline
 5 from sklearn.cluster import MiniBatchKMeans
 6 from sklearn.compose import ColumnTransformer
 7
 8 from sklearn.ensemble._hist_gradient_boosting.gradient_boosting import HistGradientBoostingClassifier
 9 from sklearn.impute import SimpleImputer
10 # from sklearn.impute._iterative import IterativeImputer
11 from sklearn.experimental import enable_iterative_imputer   # explicitly require this experimental feature
12 from sklearn.impute import IterativeImputer
13 from sklearn.linear_model import LogisticRegression, BayesianRidge
14 from sklearn.preprocessing import Normalizer
15 from sklearn.preprocessing import RobustScaler, MinMaxScaler, label_binarize, StandardScaler
16 from sklearn.svm import SVC
17 import xgboost as xgb
18
19 import pandas as pd
20 import numpy as np
21
22 from valeo.infrastructure.LogManager import LogManager
23 from valeo.infrastructure.tools.DebugPipeline import DebugPipeline
24 from valeo.infrastructure import Const as C
25
26 '''
27 https://github.com/scikit-learn-contrib/imbalanced-learn/tree/master/examples
28 https://towardsdatascience.com/introduction-to-bayesian-linear-regression-e66e60791ea7
29 https://towardsdatascience.com/custom-transformers-and-ml-data-pipelines-with-python-20ea2a7adb65
30 https://jorisvandenbossche.github.io/blog/2018/05/28/scikit-learn-columntransformer/
31 '''
32
33 class ValeoModeler :
34     logger = None
35
36     def __init__(self):
37         logger = LogManager.logger(__name__)
38
39     def build_transformers_pipeline(self, features_dtypes:pd.Series) -> ColumnTransformer:
40         rand_state = 48
41         numerical_features = (features_dtypes == 'int64') | (features_dtypes == 'float64')
42         # categorical_features = ~numerical_features
43         # nan_imputer    = SimpleImputer(strategy='mean', missing_values=np.nan, verbose=False)
44         nan_imputer    = IterativeImputer(estimator=BayesianRidge(), missing_values=np.nan,  max_iter=10,
   initial_strategy = 'median',  add_indicator=True, random_state=rand_state)
45         zeroes_imputer = IterativeImputer(estimator=BayesianRidge(), missing_values=0,  max_iter=10,
   initial_strategy = 'median',  add_indicator=True, random_state=rand_state)
46         scaler         = RobustScaler(with_centering=True, with_scaling=True, quantile_range=(25.0, 75.0))  #
   Normalizer()  # RobustScaler() #StandardScaler() # RobustScaler(with_centering=True, with_scaling=False)  #
   MinMaxScaler()
47         # scaler  = Normalizer(norm='l1')
48         # NB: When using log tranformer: Adopt this transformation -> log(-2) = -1×(log(abs(-2)+1))
49         # dbg = DebugPipeline()
50         num_transformers_pipeline = Pipeline([ #('dbg_0', dbg),
51             ('nan_imputer', nan_imputer),      # ('dbg_1', dbg),
52             ('zeroes_imputer', zeroes_imputer), # ('dbg_2', dbg),
53             ('scaler', scaler),                # ('dbg_3', dbg)
54         ])
55         return ColumnTransformer([('transformers_pipeline',num_transformers_pipeline, numerical_features)],
   remainder='passthrough')
56
57                          # ENS(0.61) without explicit overSampling / test_roc_auc : [0.6719306  0.58851217 0.
   58250362 0.6094371  0.55757417]
58     BBC  = "BBC"        # BalancedBaggingClassifier(base_estimator=HGBR,  sampling_strategy=1.0, replacement=
   False, random_state=48)
59     HGBC = "HGBR"       # HistGradientBoostingClassifier(max_iter = 8 , max_depth=8,learning_rate=0.35,
   l2_regularization=500)
60
61     BRFC = "BRFC"       # BalancedRandomForestClassifier(n_estimators = 50 , max_depth=20)
62     RUSBoost = "RUSBoost" # RUSBoostClassifier(n_estimators = 8 , algorithm='SAMME.R', random_state=42)
63     KNN = "KNN"         # KNeighborsClassifier(3),
64     SVC = "SVC"         # SVC(kernel="rbf", C=0.025, probability=True)
65     NuSVC = "NuSVC"     # NuSVC(probability=True),
66     RFC = "RFC"         # RandomForestClassifier(n_estimators=10, max_depth=10, max_features=10, n_jobs=4))
67     DTC = "DTC"         # DecisionTreeClassifier())  # so bad
68     ADABoost = "ADABoost" # AdaBoostClassifier()
69     GBC  = "GBC"        # GradientBoostingClassifier()
70     LRC  = "LRC"        # LogisticRegression(max_iter=500))  # Best for Recall 1
71     XGBC = "XGBC"       # xgb.XGBClassifier()
72     #  ('classification', GaussianNB())  # 0.5881085402220386
73     #  ('classification', ComplementNB())  # 0.523696690978335
74     #  ('classification', MultinomialNB())  # 0.523696690978335
75     Imbl_Resampler =  "Imbl_Resampler"  # ('imbalancer_resampler', self.build_resampler(sampler_type,
   sampling_strategy='not majority'))
76
77     def build_predictor_pipeline(self, features_dtypes:pd.Series, clfTypes:[str]) -> Pipeline:
78         cls = self.__class__
79         clfs = {
```

```python
80              cls.HGBC : HistGradientBoostingClassifier(max_iter = 100 , max_depth=10,learning_rate=0.10,
      l2_regularization=5),
81              cls.BBC  : BalancedBaggingClassifier(base_estimator=HistGradientBoostingClassifier(),  n_estimators=
      50, sampling_strategy='auto', replacement=False, random_state=48),
82
83              # scale_pos_weight
84              # ESTIM:100 depth:20 [6155 4108]/[41 51] - P:0.0123 - R:0.5543 - roc_auc:0.5770 - f1:0.0240 |
85              #.ESTIM:300 depth:10 [6085 4178]/[37 55] - P:0.0130 - R:0.5978 - roc_auc:0.5954 - f1:0.0254
86              # ESTIM:300 depth:15 [6306 3957]/[37 55] - P:0.0137 - R:0.5978 - roc_auc:0.6061 - f1:0.0268
87              # ESTIM:300 depth:20 [6057 4206]/[33 59] - P:0.0138 - R:0.6413 - roc_auc:0.6157 - f1:0.0271   ***
88              # ESTIM:300 depth:20 class_weight:{0:1, 1:100} [2860 7403]/[22 70] - P:0.0094 - R:0.7609 - roc_auc:0
      .5198 - f1:0.0185
89              # [6127 4136]/[35 57] - P:0.0136 - R:0.6196 - roc_auc:0.6083 - f1:0.0266
90              # [6184 4079]/[37 55] - P:0.0133 - R:0.5978 - roc_auc:0.6002 - f1:0.0260
91              # [6121 4142]/[37 55] - P:0.0131 - R:0.5978 - roc_auc:0.5971 - f1:0.0256
92              # ESTIM:300 depth:30 [6223 4040]/[36 56] - P:0.0137 - R:0.6087 - roc_auc:0.6075 - f1:0.0267
93              # ESTIM:300 depth:40 [6243 4020]/[39 53] - P:0.0130 - R:0.5761 - roc_auc:0.5922 - f1:0.0255
94              #.ESTIM:200 depth:10 [6236 4027]/[39 53] - P:0.0130 - R:0.5761 - roc_auc:0.5919 - f1:0.0254
95              # ESTIM:200 depth:20 [6104 4159]/[34 58] - P:0.0138 - R:0.6304 - roc_auc:0.6126 - f1:0.0269
96              # ESTIM:200 depth:40 [6227 4036]/[37 55] - P:0.0134 - R:0.5978 - roc_auc:0.6023 - f1:0.0263
97              cls.BRFC : BalancedRandomForestClassifier(n_estimators = 300 , max_depth=20, random_state=0),
98
99              cls.RUSBoost : RUSBoostClassifier(n_estimators = 8 , algorithm='SAMME.R', random_state=42),
100             cls.XGBC :  xgb.
101                 XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
102                             colsample_bynode=1, colsample_bytree=1, gamma=0,
103                             learning_rate=0.1, max_delta_step=0, max_depth=10, #max_depth=3,
104                             min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
105                             nthread=None, objective='binary:logistic', random_state=0,
106                             reg_alpha=0, reg_lambda=1, scale_pos_weight=100, seed=42,
107                             silent=None, subsample=1, verbosity=1)
108         }
109         dbg = DebugPipeline()
110         pl= Pipeline([('preprocessor', self.build_transformers_pipeline(features_dtypes)) ,
111                     # ('imbalancer_resampler', self.build_resampler(C.smote_over_sampling,sampling_strategy='
      minority')),  # ('dbg_1', dbg),
112                     ('classifier', clfs[clfTypes[0]])  # ex: bbc : ENS(0.61) without explicit overSampling /
      test_roc_auc : [0.6719306  0.58851217 0.58250362 0.6094371  0.55757417]
113                     ])
114         for i, s in enumerate(pl.steps) :
115             # Ex: 0 -> ('preprocessor', ColumnTransformer( ... +  1 -> ('classifier', BalancedBaggingClassifier(
      base_......
116             print(f"{i} -> {s[0]} / {str(s[1])[:70]}")
117         return pl
118
119
120     '''
121     SMOTe is a technique based on nearest neighbors judged by Euclidean Distance between data points in feature
      space.
122     random_over_sampling : The most naive strategy is to generate new samples by randomly sampling with
      replacement the current available samples.
123     adasyn_over_sampling : Adaptive Synthetic: focuses on generating samples next to the original samples which
      are wrongly classified using a k-Nearest Neighbors classifier
124     smote_over_sampling  : Synth Minority Oversampling Techn: will not make any distinction between easy and
      hard samples to be classified using the nearest neighbors rule
125     ---
126     https://medium.com/towards-artificial-intelligence/application-of-synthetic-minority-over-sampling-technique
      -smote-for-imbalanced-data-sets-509ab55cfdaf
127     https://imbalanced-learn.readthedocs.io/en/stable/over_sampling.html
128     NB:
129     How to apply SMOTE : Shuffling and Splitting the Dataset into Training and Validation Sets and THEN applying
       SMOTe on the Training Dataset.
130     '''
131     def build_resampler(self, sampler_type: str, sampling_strategy='auto', k_neighbors=5) -> BaseOverSampler :
132         rand_state = 48
133         if sampler_type.lower() == C.random_over_sampler :
134             return RandomOverSampler(sampling_strategy=sampling_strategy, random_state=rand_state)
135         elif sampler_type.lower() == C.adasyn_over_sampling :
136             return ADASYN(sampling_strategy=sampling_strategy, random_state=rand_state, n_neighbors=k_neighbors)
137         elif sampler_type.lower() == C.smote_over_sampling :
138             return SMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=k_neighbors)
139         # elif sampler_type.lower() == C.smote_nc_over_sampling :      # SMOTE for dataset containing
      continuous and categorical features.
140         #     return SMOTENC(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
      k_neighbors)
141         elif sampler_type.lower() == C.smote_svm_over_sampling :     # Use an SVM algorithm to detect sample to
      use for generating new synthetic samples
142             return SVMSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
      k_neighbors, svm_estimator=SVC())
143         elif sampler_type.lower() == C.smote_kmeans_over_sampling : # Apply a KMeans clustering before to over-
      sample using SMOTE
144             return KMeansSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
      k_neighbors, kmeans_estimator=MiniBatchKMeans(n_clusters=2), cluster_balance_threshold=5)
145         elif sampler_type.lower() == C.smote_bline_over_sampling :  # Borderline samples will be detected and
      used to generate new synthetic samples.
146             return BorderlineSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
      k_neighbors, m_neighbors=3)
147         else :
```

```python
148                raise ValueError(f"Unexpected argument [sampler_type:{sampler_type}] for method '
     compute_sampler_preprocessor'")
149
150
151 # classifiers = [
152 #     KNeighborsClassifier(3),
153 #     SVC(kernel="rbf", C=0.025, probability=True),
154 #     NuSVC(probability=True),
155 #     DecisionTreeClassifier(),
156 #     RandomForestClassifier(),
157 #     AdaBoostClassifier(),
158 #     GradientBoostingClassifier()
159 # ]
160 # for classifier in classifiers:
161 #     pipe = Pipeline(steps=[('preprocessor', preprocessor),
162 #                           ('classifier', classifier)])
163 #     pipe.fit(X_train, y_train)
164 #     print(classifier)
165 #     print("model score: %.3f" % pipe.score(X_test, y_test))
```

```python
1  from valeo.infrastructure.LogManager import LogManager
2
3  import matplotlib.pyplot as plt
4  from valeo.infrastructure import Const as C
5  from sklearn.metrics import roc_auc_score, precision_recall_curve, roc_curve, average_precision_score
6
7  from valeo.infrastructure.tools.ImgUtil import ImgUtil
8
9
10 class MetricPlotter :
11     logger = None
12
13     def __init__(self):
14         MetricPlotter.logger = LogManager.logger(__name__);
15
16     def plot_roc(self, y_test, y_pred):
17         # y_test = label_binarize(y_test.values, classes=[0, 1])  # y_test 'Series'
18         # y_pred = label_binarize(y_pred, classes=[0, 1])         # y_pred  'numpy.ndarray'
19         plt.figure()
20         lw = 2
21         roc = roc_curve(y_test, y_pred)
22         plt.plot(roc[0], roc[1], color='darkorange', lw=lw, label='ROC curve (area = %0.4f)' % roc_auc_score(
   y_test, y_pred))
23         plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
24         plt.xlim([0.0, 1.0])
25         plt.ylim([0.0, 1.05])
26         plt.xlabel('False Positive Rate')
27         plt.ylabel('True Positive Rate')
28         plt.title('Receiver operating characteristic')
29         plt.legend(loc="lower right")
30         ImgUtil.save_fig("ROC_curve")
31         plt.show()
32
33     def plot_precision_recall(self, y_test, y_pred):
34         average_precision = average_precision_score(y_test, y_pred)
35         plt.figure()
36         lw = 2
37         pr = precision_recall_curve(y_test, y_pred)
38         plt.plot(pr[0], pr[1], color='darkorange', lw=lw, label='Precision Recall curve (area = %0.4f)' %
   average_precision)
39         plt.xlim([0.0, 1.05])
40         plt.ylim([0.0, 1.05])
41         plt.xlabel('Recall')
42         plt.ylabel('Precision')
43         plt.title('Precision Recall curve')
44         plt.legend(loc="upper right")
45         ImgUtil.save_fig("PR_curve")
46         plt.show()
47         #
48         # for i in range(0, len(pr[0]) ) :
49         #     print(f"{i}: ({pr[0][i]},{pr[1][i]})")
50
```

```python
1  from imblearn.pipeline import make_pipeline, Pipeline
2  from imblearn.ensemble import BalancedBaggingClassifier
3  from sklearn.linear_model import LogisticRegression
4  from sklearn.metrics import f1_score, auc, roc_auc_score
5  from sklearn.tree import DecisionTreeClassifier
6
7  import pandas as pd
8  from sklearn.model_selection import train_test_split
9
10 from valeo.domain.ValeoPreprocessor import ValeoPreprocessor
11 from valeo.infrastructure import Const as C
12 from valeo.infrastructure.LogManager import LogManager
13
14
15 class ValeoPipeline:
16     logger = None
17
18     def __init__(self):
19         self.preproc = ValeoPreprocessor()
20         ValeoPipeline.logger = LogManager.logger(__name__)
21
22     def pplSmote(self):
23         Pipeline([('column_preprocessor', self.preproc.build_column_preprocessor()) ,
24                   ('smote_resampler', self.preproc.build_resampler(C.smote_over_sampling))])
25
26     # https://towardsdatascience.com/custom-transformers-and-ml-data-pipelines-with-python-20ea2a7adb65
27     def execute(self, X_df:pd.DataFrame, y_df:pd.DataFrame, sampler_type: str):
28         # setting up testing and training sets
29         X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_size=0.25, random_state=48)
30
31         #Create an object of the classifier.
32         bbc = BalancedBaggingClassifier(base_estimator=DecisionTreeClassifier(),
33                                         sampling_strategy='auto',
34                                         replacement=False,
35                                         random_state=0)
36
37         p = Pipeline([('column_preprocessor', self.preproc.build_column_preprocessor()) ,
38                       # ('smote_resampler', self.preproc.build_resampler(sampler_type)),
39                       # ('classification', LogisticRegression())
40                       # ('classifier', bbc)
41                       ])
42
43         # p.fit_resample(X_train, y_train)
44         p.fit_transform(X_train, y_train)
45         # p.fit(X_train, y_train)
46         #
47         # y_predict = p.predict(X_test)
48         # x = f1_score(y_test, y_predict)
49         # y = 0 # auc(y_test, y_predict)
50         # z = 0 # roc_auc_score(y_test, y_predict)
51         # ValeoPipeline.logger.info(f"F1:{x} - auc:{y} - roc_auc:{z}")
52
53
54 # --------------------------------
55 # Exemple Type : PipeLine entier
56 # --------------------------------
57 # >>> pca = PCA()
58 # >>> smt = SMOTE(random_state=42)
59 # >>> knn = KNN()
60 # >>> pipeline = Pipeline([('smt', smt), ('pca', pca), ('knn', knn)])
61
62
63
64 # ----------------------
65 # Exemple_1
66 # ----------------------
67 # from sklearn.compose import ColumnTransformer
68 # from sklearn.pipeline import Pipeline
69 #
70 # # 1 - Define Categorical pipe_line
71 # cat_col = ['sex', 'embarked', 'pclass']
72 # cat_pipeline = Pipeline(steps=[
73 #     ("constant-imputer", SimpleImputer(strategy='constant', fill_value='missing')),
74 #     ("ordinal-encoder", OrdinalEncoder()),
75 # ])
76 #
77 # # 2 - Define Numerical pipe_line
78 # num_cols = ['age', 'parch', 'fare']
79 # num_pipeline = SimpleImputer(
80 #     strategy="mean", add_indicator=True,
81 # )
82 #
83 # # 3 - Define Column Transformer
84 # preprocessor = ColumnTransformer(transformers=[
85 #     ("cat-preprocessor", cat_pipeline, cat_col),
86 #     ("num-preprocessor", num_pipeline, num_cols),
87 # ])
88 #
```

```
 89  # model = Pipeline(steps=[
 90  #     ("preprocessor", preprocessor),
 91  #     ("clf", RandomForestClassifier(n_estimators=100))
 92  # ])
 93  #
 94  # _ = model.fit(X_train, y_train)
 95  #
 96  # (model.named_steps["preprocessor"]
 97  #  .named_transformers_["cat-preprocessor"]
 98  #  .named_steps["ordinal-encoder"].categories_)
 99
100
101
102  # ----------------------
103  # Exemple_2
104  # ----------------------
105  # define the pipelines
106  # cat_pipe = make_pipeline(
107  #     SimpleImputer(strategy='constant', fill_value='missing'),
108  #     OrdinalEncoder(categories=categories)
109  # )
110  # num_pipe = SimpleImputer(strategy='mean')
111  #
112  # preprocessing = ColumnTransformer(
113  #     [('cat_preprocessor', cat_pipe, cat_col),
114  #      ('num_preprocessor', num_pipe, num_cols)]
115  # )
```

```python
1
2  from imblearn.metrics._classification import classification_report_imbalanced
3  # https://imbalanced-learn.readthedocs.io/en/stable/api.html#module-imblearn.pipeline
4  from sklearn.base import BaseEstimator
5  from sklearn.metrics import f1_score, auc, roc_auc_score, confusion_matrix, classification_report, \
6      precision_recall_curve, precision_recall_fscore_support, roc_curve, plot_precision_recall_curve, \
7      average_precision_score, precision_score, recall_score, accuracy_score, balanced_accuracy_score
8  # from sklearn.impute import SimpleImputer
9  from sklearn.model_selection import cross_validate, StratifiedKFold, GridSearchCV, RandomizedSearchCV
10
11 import pandas as pd
12
13 from valeo.domain.MetricPlotter import MetricPlotter
14 from valeo.domain.ValeoModeler import ValeoModeler
15 from valeo.infrastructure.tools.DfUtil import DfUtil
16 from valeo.infrastructure.LogManager import LogManager
17 from valeo.infrastructure import Const as C
18
19 import xgboost as xgb
20
21
22 class ValeoPredictor :
23     logger = None
24
25     def __init__(self):
26         ValeoPredictor.logger = LogManager.logger(__name__)
27         self.modeler = ValeoModeler()
28         self.metricPlt = MetricPlotter()
29
30
31     def fit_cv_grid_search(self, X:pd.DataFrame, y:pd.DataFrame, clfTypes:[str] , n_splits=5) -> ([BaseEstimator
], dict): # (estimator, cv_results)
32         model = self.modeler.build_predictor_pipeline(X.dtypes, clfTypes) # sampler_type)
33         CV = StratifiedKFold(n_splits=n_splits) # , random_state=48, shuffle=True
34         # HGBC
35         # param_grid = {
36         #     'classifier__n_estimators': [3, 5, 10, 20, 50],
37         #     'classifier__base_estimator__l2_regularization': [5, 50, 100, 50],
38         #     'classifier__base_estimator__max_iter' : [100],
39         #     'classifier__base_estimator__max_depth' : [10,50,10]
40         # }
41         # BRFC
42         param_grid = {
43             'classifier__n_estimators': [250,300],
44             'classifier__max_depth': [15,20,25],
45             'classifier__max_features' : ['auto',13]
46         }
47
48         grid = GridSearchCV(model, param_grid=param_grid, n_jobs=-1, cv=CV) # if is_grid else
49         grid.fit(X, y)
50         print(f"Best Estimator: {grid.best_estimator_}")
51         df_results = pd.DataFrame(grid.cv_results_)
52                 # columns_to_keep = ['param_clf__max_depth', 'param_clf__n_estimators', 'mean_test_score', '
std_test_score',]
53                 # df_results = df_results[columns_to_keep]
54         DfUtil.write_df_csv( df_results.sort_values(by='mean_test_score', ascending=False), C.ts_pathanme([C.
rootReports(), 'grid_search_cv.csv']) )
55
56     def fit_cv_randomized_search(self, X:pd.DataFrame, y:pd.DataFrame, clfTypes:[str] , n_splits=5) -> ([
BaseEstimator], dict): # (estimator, cv_results)
57         model = self.modeler.build_predictor_pipeline(X.dtypes, clfTypes) # sampler_type)
58         CV = StratifiedKFold(n_splits=n_splits) # , random_state=48, shuffle=True
59         # HGBC
60         # param_grid = {
61         #     'classifier__n_estimators': [3, 5, 10, 20, 50],
62         #     'classifier__base_estimator__l2_regularization': [5, 50, 100, 50],
63         #     'classifier__base_estimator__max_iter' : [100],
64         #     'classifier__base_estimator__max_depth' : [10,50,10]
65         # }
66
67         grid = RandomizedSearchCV(model, param_distributions=param_grid, n_jobs=-1, cv=CV) # if is_grid else
68         grid.fit(X, y)
69         df_results = pd.DataFrame(grid.cv_results_)
70         DfUtil.write_df_csv( df_results.sort_values(by='mean_test_score', ascending=False), C.ts_pathanme([C.
rootReports(), 'grid_search_csv']) )
71
72     def print_model_params_keys(self, model:BaseEstimator):
73         for param in model.get_params().keys():
74             print(param)
75
76     # 1 - Fit without any Cross Validation
77     def fit_and_plot(self, X_train:pd.DataFrame, y_train:pd.DataFrame,  X_test:pd.DataFrame, y_test:pd.DataFrame
, clfTypes:[str]) -> BaseEstimator:
78         # Q1 = X_train.quantile(0.25)
79         # Q3 = X_train.quantile(0.75)
80         # IQR = Q3 - Q1
81         # # print(IQR)
82         # to_remove = ((X_train < (Q1 - 1.5 * IQR)) |(X_train > (Q3 + 1.5 * IQR))).any(axis=1)
```

```python
83          # y_train = y_train.drop(axis=0, index=X_train[to_remove].index)
84          # X_train = X_train[~to_remove]
85          #
86          fitted_model = self.fit(X_train, y_train, clfTypes)
87          # print(f"Type:{type(fitted_model)} - {fitted_model.get_params()}")
88          # self.print_model_params_keys(fitted_model)
89          self.predict_and_plot(fitted_model, X_test, y_test)
90          return fitted_model
91
92      def fit(self, X_train:pd.DataFrame, y_train:pd.DataFrame, clfTypes:[str]) -> BaseEstimator:
93          model = self.modeler.build_predictor_pipeline(X_train.dtypes, clfTypes)
94          return model.fit(X_train, y_train)
95
96      ''' 2 - Fit with Cross Validation
97          NB :
98          a - roc-auc-avo + roc-auc-ovr :
99              https://stackoverflow.com/questions/59453363/what-is-the-difference-of-roc-auc-values-in-sklearn
100             roc_auc is the only one suitable for binary classification. The weighted, ovr and ovo are use for
    multi-class problems
101
102         b - Micro-Average + Macro-Average (for Precision / Recall / F1) :
103             http://rushdishams.blogspot.com/2011/08/micro-and-macro-average-of-precision.html
104             https://datascience.stackexchange.com/questions/15989/micro-average-vs-macro-average-performance-in-
    a-multiclass-classification-settin
105             Ex: Micro-P = (TP1 + TP2) / ( TP1 + FP1 + TP2 + F2)
106                 Macro-P = (P1 + P2) / 2
107             Suitability:
108             . Macro-average method can be used when you want to know how the system performs overall across the
    sets of data
109             . Micro-average method can be a useful measure when your dataset varies in size.
110
111         c - How can we report 'confusion matrix' while using 'cross_validate' ?
112             https://stackoverflow.com/questions/40057049/using-confusion-matrix-as-scoring-metric-in-cross-
    validation-in-scikit-learn
113             c1. Either use 'cross_val_predict' and deduce confusion-matrix:
114                 y_pred = cross_val_predict(clf, x, y, cv=10)
115                 conf_mat = confusion_matrix(y_test, y_pred)
116                 BUT BEWARE: Passing these predictions into an evaluation metric may not be a valid way to
    measure generalization performance.
117                     Results can differ from cross_validate and cross_val_score unless all tests sets
    have equal size and the metric decomposes over samples.
118             c2. If you want to obtain confusion matrices for multiple evaluation runs (such as cross validation
    ) you have to do this by hand:
119                 conf_matrix_list_of_arrays = []
120                 kf = cross_validation.KFold(len(y), n_folds=5)
121                 for train_index, test_index in kf:
122                     X_train, X_test = X[train_index], X[test_index]   # Panda-Column index 'train_index' are of
    type 'numpy array'
123                     y_train, y_test = y[train_index], y[test_index]
124                     #
125                     model.fit(X_train, y_train)
126                     conf_matrix = confusion_matrix(y_test, model.predict(X_test))
127                     conf_matrix_list_of_arrays .append(conf_matrix)
128
129                 On the end you can calculate your mean of list of numpy arrays (confusion matrices) with:
130                 mean_of_conf_matrix_arrays = np.mean(conf_matrix_list_of_arrays, axis=0
    )
131     '''
132     def fit_cv(self, X:pd.DataFrame, y:pd.DataFrame, clfTypes:[str], n_splits=5) -> ([BaseEstimator], dict):
    # (estimator, cv_results)
133         model = self.modeler.build_predictor_pipeline(X.dtypes, clfTypes)
134         CV = StratifiedKFold(n_splits=n_splits) # , random_state=48, shuffle=True
135         cv_results =  cross_validate(model, X, y, cv=CV, scoring=('f1', 'f1_micro', 'f1_macro', 'f1_weighted', '
    recall', 'precision', 'average_precision', 'roc_auc'), return_train_score=True, return_estimator=True)
136         fitted_estimators = []
137         for key in cv_results.keys() :
138             if str(key) !=  "estimator" :
139                 print(f"{key} : {cv_results[key]}")
140             fitted_estimators.append(cv_results[key])
141         return fitted_estimators, cv_results
142
143     '''
144         - Print metrics
145         - Print report
146         - Plot ROC : TP vs FP
147         - Plot AUC : Precison vs Recall
148     '''
149     def predict_and_plot(self, fitted_model: BaseEstimator, X_test:pd.DataFrame, y_test:pd.DataFrame):
150         y_pred = fitted_model.predict(X_test)
151         #
152         print(f"- Model score: {fitted_model.score(X_test, y_test)}")
153         print(f"- Accuracy score: {accuracy_score(y_test, y_pred)}")
154         print(f"- Balanced accuracy score: {balanced_accuracy_score(y_test, y_pred)} / The balanced accuracy to
    deal with imbalanced datasets. It is defined as the average of recall obtained on each class.")
155         # print(f"- auc : {auc(y_test, y_pred)}")  # ValueError: x is neither increasing nor decreasing : [0 0
    0 ... 0 0 0]
156         print(f"- Average_precision_score: {average_precision_score(y_test, y_pred)}")
157         print(f"- Precision_score: {precision_score(y_test, y_pred)}")
```
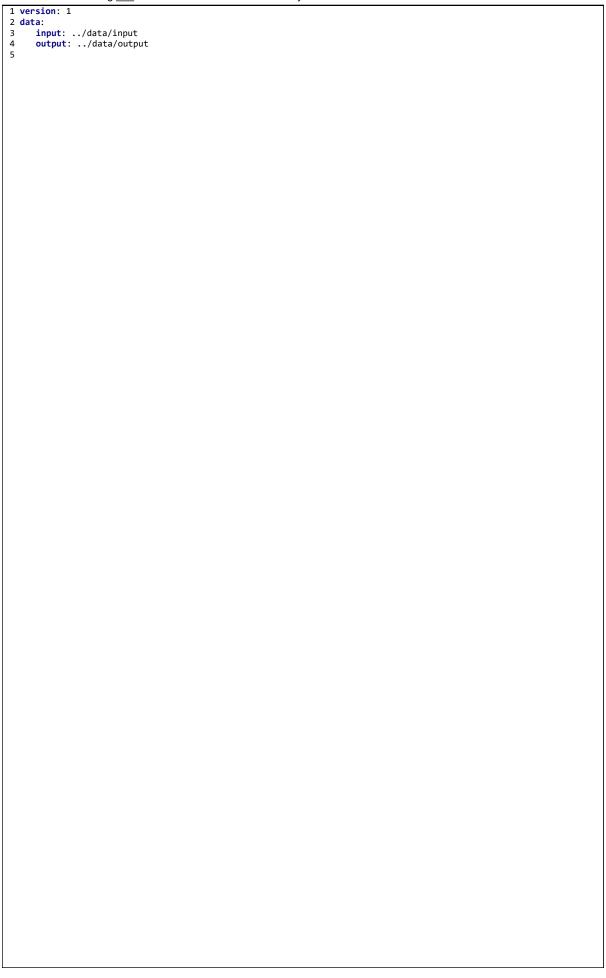
```python
158         print(f"- Recall score: {recall_score(y_test, y_pred)}")
159         print(f"- Roc_auc_score: {roc_auc_score(y_test, y_pred)}")
160         print(f"- F1 score: {f1_score(y_test, y_pred)}")
161         m = confusion_matrix(y_test, y_pred)
162         print(f"- {m[0]}/{m[1]} - P:{precision_score(y_test, y_pred):0.4f} - R:{recall_score(y_test, y_pred):0.
   4f} - roc_auc:{roc_auc_score(y_test, y_pred):0.4f} - f1:{f1_score(y_test, y_pred):0.4f}")
163         print(f"- {confusion_matrix(y_test, y_pred)}")
164         print(f"- classification_report_imbalanced:\n{classification_report_imbalanced(y_test, y_pred)}")
165         print(f"- classification_report:\n{classification_report(y_test, y_pred)}")
166         print(f"- precision_recall_curve: {precision_recall_curve(y_test, y_pred)}")
167         print(f"- precision_recall_fscore_support: {precision_recall_fscore_support(y_test, y_pred)}")
168         print(f"- roc_curve: {roc_curve(y_test, y_pred)}")
169         #
170         self.metricPlt.plot_roc(y_test, y_pred)
171         self.metricPlt.plot_precision_recall(y_test, y_pred)
172         # self.plot_roc(y_test, y_pred)
173         # self.plot_precision_recall(y_test, y_pred)
174
175
176
177 # https://medium.com/towards-artificial-intelligence/application-of-synthetic-minority-over-sampling-technique-
   smote-for-imbalanced-data-sets-509ab55cfdaf
178 # from sklearn.ensemble import GradientBoostingClassifier
179 # from sklearn.model_selection import GridSearchCV
180 # parameters = {'n_estimators':[100,150,200,250,300,350,400,450,500],
181 #               'max_depth':[3,4,5]}
182 # clf= GradientBoostingClassifier()
183 # grid_search = GridSearchCV(param_grid = parameters, estimator = clf,
184 #                            verbose = 3)
185 # grid_search_2 = grid_search.fit(X_train,y_train)
186
187 # GOOGLE ON: classifier over sampled imbalanced dataset
188 # https://sci2s.ugr.es/imbalanced  : Tres interessant****
189 # https://www.datacamp.com/community/tutorials/diving-deep-imbalanced-data  : Tres interessant****
190 #---
191 # https://journalofbigdata.springeropen.com/articles/10.1186/s40537-017-0108-1  : Tres interessant****
192 # xperimental evaluation
193 # The projected technique works on binary-class/multi-class imbalanced Big Data sets in the organization to
   recommended LVH. Four basic classifiers viz. Random Forest (-P 100-I 100-num-slots 1-K 0-M 1.0-V 0.001-S 1),
   Naïve Bayes, AdaBoostM1 (-P 100-S 1-I 10-W weka.classifiers.trees.DecisionStump) and MultiLayer Perceptron (-L 0
   .3-M 0.2-N 500-V 0-S 0-E 20-H a) are applied to over_sampled data sets using dissimilar values of cross-
   validation and KNN. Lastly the results, based on the F-measure and AUC values are used to compare between
   benchmarking (SMOTE/Borderline-SMOTE/ADASYN/SPIDER2/SMOTEBoost/MWMOTE) and planned technique (UCPMOT). Tables 3
   , 4, 5, 6, 7, 8, 9, 10, 11, 12 and 13 describe the results in detail. The analysis of results validates the
   superiority of UCPMOT for enhancing the classification.
194
195 # GOOGLE ON: scikit learn imbalanced dataset resampling type cross validation shuffle
196 # https://www.kaggle.com/rafjaa/resampling-strategies-for-imbalanced-datasets
197 # https://machinelearningmastery.com/cross-validation-for-imbalanced-classification/
198 # CV = ShuffleSplit(n_splits=10, test_size=0.25, random_state=48)
199 # https://www.alfredo.motta.name/cross-validation-done-wrong/
200
201 # https://www.dataschool.io/simple-guide-to-confusion-matrix-terminology/
202 #     https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6
203
204 # http://www.cs.nthu.edu.tw/~shwu/courses/ml/labs/08_CV_Ensembling/08_CV_Ensembling.html
205 # https://github.com/arrayslayer/ML-Project
206
207 # https://www.kaggle.com/shiqbal/first-data-exploration/notebook  applied on Porto Seguro's
208
209
210 ''' TODO :
211
212 -> Faire Ressortir l'importance et la contribution de chaque feature:
213    https://scikit-learn.org/stable/auto_examples/inspection/plot_permutation_importance.html
214
215 -> Essayer ces scénarios de modélisation:
216    scen1: Oversampling + LogisticR
217    scen2: BaggingClassifier + Histo
218
219 '''
```

```python
1  from imblearn.over_sampling import RandomOverSampler, SMOTE, SMOTENC, SVMSMOTE, KMeansSMOTE, BorderlineSMOTE,
   ADASYN
2  from imblearn.over_sampling.base import BaseOverSampler
3  # from sklearn.experimental import enable_iterative_imputer   # explicitly require this experimental feature
4  # from sklearn.impute import IterativeImputer
5  from sklearn.impute import SimpleImputer  # now you can import normally from sklearn.impute
6  from sklearn.linear_model import BayesianRidge
7  from sklearn.model_selection import train_test_split
8  from sklearn.preprocessing import RobustScaler, StandardScaler, MinMaxScaler
9
10 from sklearn.compose import ColumnTransformer, make_column_transformer
11 from imblearn.pipeline import make_pipeline, Pipeline
12
13 import pandas as pd
14 import numpy as np
15
16 from valeo.infrastructure import Const as C
17 from valeo.infrastructure.LogManager import LogManager
18
19 from valeo.infrastructure.tools.DebugPipeline import DebugPipeline
20
21
22 class ValeoPreprocessor:
23     logger = None
24
25     def __init__(self):
26         ValeoPreprocessor.logger = LogManager.logger(__name__)
27
28     # def build_iterative_preprocessor(self)  -> ColumnTransformer:
29     #     imp_cols = [C.OP100_Capuchon_insertion_mesure]
30     #     it_imp = IterativeImputer(estimator=BayesianRidge(),  max_iter=10, initial_strategy = 'median',
   add_indicator=True, random_state=48)
31     #     return ColumnTransformer([('iterative_imputer', it_imp, imp_cols)] )
32
33     def build_column_preprocessor(self) -> Pipeline: # ColumnTransformer:
34         rand_state = 48
35         # 1 - IterativeImputer : models each feature with missing values as a function of other features, and
   uses that estimate for imputation
36         # imputer_pipe = Pipeline( IterativeImputer(estimator=BayesianRidge(), missing_values=[np.nan, 0],
   max_iter=10, initial_strategy = 'median') )
37
38         # imputer_nan_values = IterativeImputer(estimator=BayesianRidge(), missing_values=np.nan,  max_iter=10,
   initial_strategy = 'median',  add_indicator=True, random_state=rand_state)
39         # imputer_nan_values = IterativeImputer(estimator=BayesianRidge(), missing_values=np.nan,  max_iter=10,
   initial_strategy = 'median',  add_indicator=False, random_state=rand_state)
40         imputed_nan_cols = [C.OP100_Capuchon_insertion_mesure]  # columns having too much missing values
41         #
42         # imputer_zeroes_values = IterativeImputer(estimator=BayesianRidge(), missing_values=0,  max_iter=10,
   initial_strategy = 'median',  add_indicator=False, random_state=rand_state)
43         imputed_zeroes_cols = [C.OP100_Capuchon_insertion_mesure, C.OP090_StartLinePeakForce_value, C.
   OP090_SnapRingMidPointForce_val,   # columns equals to 0 for a few rows
44                                C.OP090_SnapRingPeakForce_value, C.OP090_SnapRingFinalStroke_value]
45
46         # 2 - Scale features using statistics that are robust to outliers.
47         scaler      = StandardScaler() #  MinMaxScaler() # StandardScaler() # RobustScaler(with_centering=True,
   with_scaling=False)
48         scaled_cols = [C.OP070_V_1_angle_value, C.OP070_V_1_torque_value,
49                        C.OP070_V_2_angle_value, C.OP070_V_2_torque_value,
50                        C.OP090_StartLinePeakForce_value, C.OP090_SnapRingMidPointForce_val,
51                        C.OP090_SnapRingPeakForce_value,  C.OP090_SnapRingFinalStroke_value,
52                        C.OP100_Capuchon_insertion_mesure,
53                        C.OP110_Vissage_M8_angle_value, C.OP110_Vissage_M8_torque_value,
54                        C.OP120_Rodage_I_mesure_value,  C.OP120_Rodage_U_mesure_value]
55
56         dbg = DebugPipeline()
57         return  Pipeline([ ('dbg_0', dbg, scaled_cols),
58                            ('imputer_preprocessor_nan', SimpleImputer(strategy='mean', missing_values=np.
   nan, verbose=True), imputed_nan_cols),
59                            ('dbg_1', dbg, scaled_cols),
60                            ('imputer_preprocessor_zeroes', SimpleImputer(strategy='mean', missing_values=
   0.0, verbose=True), imputed_zeroes_cols),
61                            ('dbg_2', dbg, scaled_cols),
62                            ('scaler_preprocessor', scaler, scaled_cols),
63                            ('dbg_3', dbg, scaled_cols),
64                            # #
65                            # ('imputer_preprocessor_nan_bis', SimpleImputer(strategy='mean'), scaled_cols
   ),
66                            # ('dbg_1_bis', DebugPline(),scaled_cols)
67                            ])
68                    # , remainder='passthrough')
69
70
71     def execute(self, X_train:pd.DataFrame):
72         imputed_nan_cols = [C.OP100_Capuchon_insertion_mesure]
73         imputed_zeroes_cols = [C.OP100_Capuchon_insertion_mesure, C.OP090_StartLinePeakForce_value, C.
   OP090_SnapRingMidPointForce_val,   # columns equals to 0 for a few rows
74                                C.OP090_SnapRingPeakForce_value, C.OP090_SnapRingFinalStroke_value]
75         scaled_cols = [C.OP070_V_1_angle_value, C.OP070_V_1_torque_value,
```

```python
 76                          C.OP070_V_2_angle_value, C.OP070_V_2_torque_value,
 77                          C.OP090_StartLinePeakForce_value, C.OP090_SnapRingMidPointForce_val,
 78                          C.OP090_SnapRingPeakForce_value,  C.OP090_SnapRingFinalStroke_value,
 79                          C.OP100_Capuchon_insertion_mesure,
 80                          C.OP110_Vissage_M8_angle_value, C.OP110_Vissage_M8_torque_value,
 81                          C.OP120_Rodage_I_mesure_value,  C.OP120_Rodage_U_mesure_value]
 82          # DebugPipeline.counter += 10
 83          d = DebugPipeline()
 84          #
 85          d.fit_transform(X_train[scaled_cols])
 86          s = SimpleImputer(strategy='mean', missing_values=np.nan, verbose=True)
 87          X_train[imputed_nan_cols] = s.fit_transform(X_train[imputed_nan_cols])
 88          d.fit_transform(X_train[scaled_cols])
 89          #
 90          s = SimpleImputer(strategy='mean', missing_values=0.0, verbose=True)
 91          X_train[imputed_zeroes_cols] = s.fit_transform(X_train[imputed_zeroes_cols])
 92          d.fit_transform(X_train[scaled_cols])
 93          #
 94          scaler      = MinMaxScaler() # StandardScaler() # RobustScaler(with_centering=True, with_scaling=False)
 95          X_train[scaled_cols]  = scaler.fit_transform(X_train[scaled_cols])
 96          d.fit_transform(X_train[scaled_cols])
 97
 98
 99
100      '''
101      def build_column_preprocessor(self) -> ColumnTransformer:
102          rand_state = 48
103          # 1 - IterativeImputer : models each feature with missing values as a function of other features, and
    uses that estimate for imputation
104          imputer_pipe = make_pipeline( IterativeImputer(estimator=BayesianRidge(), missing_values=np.nan,
    max_iter=10, initial_strategy = 'median',  add_indicator=True, random_state=rand_state) )
105          # imputer_pipe = Pipeline( IterativeImputer(estimator=BayesianRidge(), missing_values=[np.nan, 0],
    max_iter=10, initial_strategy = 'median') )
106          imputed_cols = [C.OP100_Capuchon_insertion_mesure]  # columns having too much missing values
107                      # C.OP090_StartLinePeakForce_value, C.OP090_SnapRingMidPointForce_val,  # columns equals
     to 0 for a few rows
108                      # C.OP090_SnapRingPeakForce_value, C.OP090_SnapRingFinalStroke_value]
109
110          # 2 - Scale features using statistics that are robust to outliers.
111          scaler_pipe = make_pipeline(RobustScaler(with_centering=True, with_scaling=False))
112          scaled_cols = [C.OP070_V_1_angle_value, C.OP070_V_1_torque_value,
113                      C.OP070_V_2_angle_value, C.OP070_V_2_torque_value,
114                      C.OP090_StartLinePeakForce_value, C.OP090_SnapRingMidPointForce_val,
115                      C.OP090_SnapRingPeakForce_value,  C.OP090_SnapRingFinalStroke_value,
116                      C.OP100_Capuchon_insertion_mesure,
117                      C.OP110_Vissage_M8_angle_value, C.OP110_Vissage_M8_torque_value,
118                      C.OP120_Rodage_I_mesure_value,  C.OP120_Rodage_U_mesure_value,]
119
120          return  ColumnTransformer([('imputer_preprocessor', imputer_pipe, imputed_cols),
121                                      ('scaler_preprocessor', scaler_pipe, scaled_cols)] )
122      '''
123
124      '''
125     SMOTe is a technique based on nearest neighbors judged by Euclidean Distance between data points in feature
    space.
126     random_over_sampling : The most naive strategy is to generate new samples by randomly sampling with
    replacement the current available samples.
127     adasyn_over_sampling : Adaptive Synthetic: focuses on generating samples next to the original samples which
    are wrongly classified using a k-Nearest Neighbors classifier
128     smote_over_sampling  : Synth Minority Oversampling Techn: will not make any distinction between easy and
    hard samples to be classified using the nearest neighbors rule
129     ---
130     https://medium.com/towards-artificial-intelligence/application-of-synthetic-minority-over-sampling-technique
    -smote-for-imbalanced-data-sets-509ab55cfdaf
131     https://imbalanced-learn.readthedocs.io/en/stable/over_sampling.html
132     NB:
133     How to apply SMOTE : Shuffling and Splitting the Dataset into Training and Validation Sets and THEN applying
     SMOTe on the Training Dataset.
134     '''
135      def build_resampler(self, sampler_type: str, sampling_strategy='auto', k_neighbors=5) -> BaseOverSampler :
136          rand_state = 48
137          if sampler_type.lower() == C.random_over_sampler :
138              return RandomOverSampler(sampling_strategy=sampling_strategy, random_state=rand_state)
139          elif sampler_type.lower() == C.adasyn_over_sampling :
140              return ADASYN(sampling_strategy=sampling_strategy, random_state=rand_state, n_neighbors=k_neighbors)
141          elif sampler_type.lower() == C.smote_over_sampling :
142              return SMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=k_neighbors)
143          # elif sampler_type.lower() == C.smote_nc_over_sampling :      # SMOTE for dataset containing
    continuous and categorical features.
144          #     return SMOTENC(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
    k_neighbors)
145          elif sampler_type.lower() == C.smote_svm_over_sampling :    # Use an SVM algorithm to detect sample to
    use for generating new synthetic samples
146              return SVMSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
    k_neighbors)
147          elif sampler_type.lower() == C.smote_kmeans_over_sampling : # Apply a KMeans clustering before to over-
    sample using SMOTE
148              return KMeansSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
```

```python
148 k_neighbors)
149         elif sampler_type.lower() == C.smote_bline_over_sampling :  # Borderline samples will be detected and
    used to generate new synthetic samples.
150             return BorderlineSMOTE(sampling_strategy=sampling_strategy, random_state=rand_state, k_neighbors=
    k_neighbors)
151         else :
152             raise ValueError(f"Unexpected argument [sampler_type:{sampler_type}] for method '
    compute_sampler_preprocessor'")
153
154
155
156 # --------------------------------
157 # Exemple Type : PipeLine entier
158 # --------------------------------
159 # >>> pca = PCA()
160 # >>> smt = SMOTE(random_state=42)
161 # >>> knn = KNN()
162 # >>> pipeline = Pipeline([('smt', smt), ('pca', pca), ('knn', knn)])
163
164
165
166 # ----------------------
167 # Exemple_1
168 # ----------------------
169 # from sklearn.compose import ColumnTransformer
170 # from sklearn.pipeline import Pipeline
171 #
172 # # 1 - Define Categorical pipe_line
173 # cat_col = ['sex', 'embarked', 'pclass']
174 # cat_pipeline = Pipeline(steps=[
175 #     ("constant-imputer", SimpleImputer(strategy='constant', fill_value='missing')),
176 #     ("ordinal-encoder", OrdinalEncoder()),
177 # ])
178 #
179 # # 2 - Define Numerical pipe_line
180 # num_cols = ['age', 'parch', 'fare']
181 # num_pipeline = SimpleImputer(
182 #     strategy="mean", add_indicator=True,
183 # )
184 #
185 # # 3 - Define Column Transformer
186 # preprocessor = ColumnTransformer(transformers=[
187 #     ("cat-preprocessor", cat_pipeline, cat_col),
188 #     ("num-preprocessor", num_pipeline, num_cols),
189 # ])
190 #
191 # model = Pipeline(steps=[
192 #     ("preprocessor", preprocessor),
193 #     ("clf", RandomForestClassifier(n_estimators=100))
194 # ])
195 #
196 # _ = model.fit(X_train, y_train)
197 #
198 # (model.named_steps["preprocessor"]
199 #   .named_transformers_["cat-preprocessor"]
200 #   .named_steps["ordinal-encoder"].categories_)
201
202
203
204 # ----------------------
205 # Exemple_2
206 # ----------------------
207 # define the pipelines
208 # cat_pipe = make_pipeline(
209 #     SimpleImputer(strategy='constant', fill_value='missing'),
210 #     OrdinalEncoder(categories=categories)
211 # )
212 # num_pipe = SimpleImputer(strategy='mean')
213 #
214 # preprocessing = ColumnTransformer(
215 #     [('cat_preprocessor', cat_pipe, cat_col),
216 #      ('num_preprocessor', num_pipe, num_cols)]
217 # )
```

```yaml
1 version: 1
2 data:
3     input: ../data/input
4     output: ../data/output
5
```

```yaml
 1 version: 1
 2 formatters:
 3   standard:
 4     format: '%(asctime)s - %(levelname)s - %(name)s - %(message)s'
 5 #    format: '%(asctime)s - %(levelname)s - %(module)s - %(message)s'
 6   verbose:
 7     format: '%(asctime)s - %(levelname)s <PID %(process)d:%(processName)s> %(module)s.%(funcName)s(): %(message)s'
 8 handlers:
 9   console:
10     class: logging.StreamHandler
11     level: DEBUG
12     formatter: standard
13     stream: ext://sys.stdout
14 #    propagate: yes
15   valeo_log:
16     class: logging.FileHandler
17     filename: C:/EXED/Training/___VALEO/log/valeo.log
18     mode: w
19     level: DEBUG
20     formatter: standard
21 #    propagate: yes
22
23 #  errors:
24 #    class: logging.FileHandler
25 #    filename: mplog-errors.log
26 #    mode: w
27 #    level: ERROR
28 #    formatter: detailed
29 #loggers:
30 #  simpleExample:
31 #    level: DEBUG
32 #    handlers: [console]
33 #    propagate: no
34 root:
35   level: DEBUG
36   handlers: [console, valeo_log]
37   formatter: standard
```

```python
1  # ENVIRONMENT keys used to refer to configuration files
2  ENV_KEY_CONFIG_FILE_PATHNAME = '__VALEO__APP_CONFIG_FILE_PATHNAME' # ex: SET __VALEO__APP_CONFIG_FILE_PATHNAME
   =...../valeo.yaml'
3  ENV_KEY_LOG_FILE_PATHNAME    = '__VALEO__APP_LOG_FILE_PATHNAME'    # ex: SET __VALEO__APP_LOG_FILE_PATHNAME
   =...../logging.yaml'
4  #
5  # Symbolic name of configuration files
6  APP_DEFAULT_CONFIG_FILE = 'valeo.yaml'
7  APP_DEFAULT_LOG_FILE    = 'logging.yaml'
8  #
9  # Valeo Dataset columns names
10 PROC_TRACEINFO          = 'PROC_TRACEINFO'
11 OP070_V_1_angle_value  = 'OP070_V_1_angle_value'
12 OP070_V_1_torque_value = 'OP070_V_1_torque_value'
13 OP070_V_2_angle_value  = 'OP070_V_2_angle_value'
14 OP070_V_2_torque_value = 'OP070_V_2_torque_value'
15 OP090_StartLinePeakForce_value = 'OP090_StartLinePeakForce_value'
16 OP090_SnapRingMidPointForce_val = 'OP090_SnapRingMidPointForce_val'
17 OP090_SnapRingPeakForce_value    = 'OP090_SnapRingPeakForce_value'
18 OP090_SnapRingFinalStroke_value  = 'OP090_SnapRingFinalStroke_value'
19 OP100_Capuchon_insertion_mesure  = 'OP100_Capuchon_insertion_mesure'
20 OP110_Vissage_M8_angle_value  = 'OP110_Vissage_M8_angle_value'
21 OP110_Vissage_M8_torque_value = 'OP110_Vissage_M8_torque_value'
22 OP120_Rodage_I_mesure_value   = 'OP120_Rodage_I_mesure_value'
23 OP120_Rodage_U_mesure_value   = 'OP120_Rodage_U_mesure_value'
24 Binar_OP130_Resultat_Global_v = 'Binar OP130_Resultat_Global_v'
25 #
26 # Imbalanced resampling type
27
28 random_over_sampler = 'random_over_sampler' # The most naive strategy is to generate new samples by randomly
   sampling with replacement the current available samples.
29 adasyn_over_sampling = 'adasyn_over_sampling' # Adaptive Synthetic: focuses on generating samples next to the
   original samples which are wrongly classified using a k-Nearest Neighbors classifier
30 smote_over_sampling  = 'smote_over_sampling'  # Synth Minority Oversampling Techn: will not make any distinction
    between easy and hard samples to be classified using the nearest neighbors rule
31 smote_nc_over_sampling   = 'smote_nc_over_sampling'
32 smote_svm_over_sampling  = 'smote_svm_over_sampling'
33 smote_kmeans_over_sampling = 'smote_kmeans_over_sampling'
34 smote_bline_over_sampling  = 'smote_bline_over_sampling'
35
36
37 import os
38 from datetime import datetime
39 # timestamp : none / suffix / prefix
40 ts_none = 0
41 ts_sfix = 1
42 ts_pfix = 2
43
44 def rootProject() -> str :
45     return  os.path.join(os.path.abspath(os.path.dirname(__file__)), '..', '..', '..')  # this_folder = D:/
   Training.git/trunk/___VALEO/src/valeo/infrastructure
46
47 def rootSrc() -> str :
48     return  os.path.join(rootProject(),  'src' )
49
50 def rootData() -> str :
51     return  os.path.join(rootProject(),  'data' )
52
53 def rootDataTrain() -> str :
54     return  os.path.join(rootData(),  'train' )
55
56 def rootDataTest() -> str :
57     return  os.path.join(rootData(),  'test' )
58
59 def rootImages() -> str :
60     return  os.path.join(rootProject(),  'images' )
61
62 def rootReports() -> str :
63     return  os.path.join(rootProject(),  'reports' )
64
65 def rootResources() -> str :
66     return  os.path.join(rootProject(), 'src', 'valeo', 'resources')
67
68 def ts_pathanme(pathAsStrList : [], ts_type=ts_sfix) -> str:
69     if not isinstance(pathAsStrList,list) :
70         pathAsStrList = [pathAsStrList]
71     fname_with_ext = os.path.splitext(pathAsStrList[-1])
72     return os.path.join(pathAsStrList[0], '' if len(pathAsStrList) <= 2 else str(*pathAsStrList[1:-1] ),
73             f"{fname_with_ext[0]}{datetime.now().strftime('_%Y_%m_%d-%H.%M.%S')}{fname_with_ext[1]}" if
   ts_type == ts_sfix else \
74             (f"{datetime.now().strftime('%Y_%m_%d-%H.%M.%S_')}{pathAsStrList[-1]}" if ts_type == ts_pfix  else
   pathAsStrList[-1]) )
75
```

```python
1  import os
2
3  import pandas as pd
4  import numpy as np
5  from sklearn.model_selection import ShuffleSplit
6
7  import valeo.infrastructure.XY_metadata as XY_metadata
8  from valeo.infrastructure import Const
9  from valeo.infrastructure.LogManager import LogManager
10 from valeo.infrastructure.tools.DfUtil import DfUtil
11
12
13 class XY_Loader:
14     logger = None
15
16     def __init__(self):
17         XY_Loader.logger = LogManager.logger(__name__)
18
19     def get_cv(X, y):
20         cv = ShuffleSplit(n_splits=8, test_size=0.5, random_state=57)
21         return cv.split(X)
22
23
24     def load_XY_df(self, mt: XY_metadata, delete_XY_join_cols=True) -> ():
25         # X_df = pd.read_csv(mt.X_pathname, na_values='')  # NaN
26         X_df = pd.read_csv(mt.X_pathname)  # NaN
27         # print(X_df[Const.OP100_Capuchon_insertion_mesure].head(20))
28         # X_df[[Const.OP100_Capuchon_insertion_mesure]] = X_df[[Const.OP100_Capuchon_insertion_mesure]].fillna(0
   .0)
29         # print(X_df[Const.OP100_Capuchon_insertion_mesure].head(20))
30
31         # 1 - Check whether Y is in separate file or in the same as X
32         if mt.is_XY_in_separate_file() :
33             Y_df = pd.read_csv(mt.Y_pathname)
34             XY_df = pd.merge(left=X_df, right=Y_df, how='inner', left_on=mt.X_join, right_on=mt.Y_join, suffixes
   =('',''))
35         else :
36             Y_df = None
37             XY_df = X_df
38
39         # 2 - When not reading a Test dataset (it means there is a Target dataset) THEN Let X_df group only
   features and Y_df only target
40         if mt.is_training_set() :
41             Y_df = XY_df[mt.target_col_name]
42             X_df = XY_df.drop(mt.target_col_name, axis=1)
43
44         # 3 - Check whether we should remove joining columns
45         if delete_XY_join_cols :
46             X_df = X_df.drop(mt.X_join, axis=1)
47             try :
48                 X_df = X_df.drop(mt.Y_join, axis=1)
49             except :
50                 pass
51
52         #
53         # XY_Loader.logger.debug(f'X_df.columns: {X_df.columns}')
54         # if Y_df is not None :
55         #     XY_Loader.logger.debug(f'type(Y_df):{type(Y_df)}\nY_df: {Y_df}')
56         return X_df, Y_df
57
58     def load_XY_values(self, mt: XY_metadata, delete_XY_join_cols=True) -> ():
59         X_df, Y_df = self.load_XY_df(mt, delete_XY_join_cols)
60         return X_df.values if X_df is not None else None, \
61                Y_df.values if Y_df is not None else None
62
```

```python
 1  # https://docs.python.org/3/library/logging.html#logrecord-attributes + Useful Handlers
 2  # https://docs.python-guide.org/writing/logging/
 3  # https://github.com/Delgan/loguru
 4  # https://kingspp.github.io/design/2017/11/06/the-head-and-tail-of-logging.html
 5  # https://stackoverflow.com/questions/4690600/python-exception-message-capturing
 6  import logging.config
 7  import os
 8
 9  from valeo.infrastructure.tools.ConfigLoader import ConfigLoader
10  import valeo.infrastructure.Const as Const
11
12  class LogManager():
13
14      # NB: The ctor() initializes the logging configuration
15      def __init__(self):
16          self.log_config = LogLoader().load()
17
18      @classmethod
19      def logger(cls,logname):
20          # l = logging.getLogger(logname)
21          # l.
22          return logging.getLogger(logname)
23
24
25  class LogLoader(ConfigLoader):
26      """
27      Load the logging configuration file
28      """
29      def load(self) -> dict:
30          try :
31              dict = super().load(os.path.join(Const.rootResources(), Const.APP_DEFAULT_LOG_FILE), Const.
    ENV_KEY_LOG_FILE_PATHNAME)
32              logging.config.dictConfig(dict)
33              return dict
34          except Exception as ex:
35              logging.basicConfig(level=logging.INFO)
36              logging.warning(f'Error while loading logging configuration file:\n' \
37                              f'\t- APP_RESOURCE_PATH = {Const.rootResources()}\n' \
38                              f'\t- APP_DEFAULT_LOG_FILE = {Const.APP_DEFAULT_LOG_FILE}\n' \
39                              f'\t- ENV_KEY_LOG_FILE_PATHNAME = {Const.ENV_KEY_LOG_FILE_PATHNAME}')
40              logging.exception(ex)
41              return None
42
```

```python
1  # explicitly require this experimental feature
2  from sklearn.experimental import enable_iterative_imputer
3  # now you can import normally from sklearn.impute
4  from sklearn.impute import IterativeImputer
5  from sklearn.linear_model import BayesianRidge
6
7  import pandas as pd
8  import numpy as np
9  from sklearn.preprocessing import RobustScaler
10
11 from valeo.infrastructure.LogManager import LogManager
12
13
14 class Transformer() :
15     logger = LogManager.logger(__name__)
16
17     # def __init__(self):
18     #     lm = LogManager()
19     #     self.logger = lm.logger(__name__)
20
21     '''
22     A strategy for imputing missing values by modeling each feature with missing values as a function of other
   features in a round-robin fashion.
23     Multivariate imputer that estimates each feature from all the others.
24     https://scikit-learn.org/stable/modules/impute.html#iterative-imputer
25
26     Arg:
27     ----
28     estimator : The estimator to use at each step of the round-robin imputation.
29
30     Returns:
31     --------
32     A transformed Dataframe containing all the missing values.
33     NB: The arguement Dataframe is NOT modified => It stills intact
34     https://towardsdatascience.com/introduction-to-bayesian-linear-regression-e66e60791ea7
35     '''
36     def iterative_imputer_transform(self, df_to_transform : pd.DataFrame,  estimator=BayesianRidge(),
   missing_values=np.nan,  max_iter=10, initial_strategy = 'median') -> pd.DataFrame :
37         cols = df_to_transform.columns
38         imputer = IterativeImputer(estimator=estimator, missing_values=missing_values,  max_iter=max_iter,
   initial_strategy=initial_strategy,  add_indicator=False) # It models each feature with missing values as a
   function of other features, and uses that estimate for imputation
39         df_transformed = pd.DataFrame(imputer.fit_transform(df_to_transform))
40         # df_transformed.columns = df_transformed.columns[:-1]
41         df_transformed.columns = cols
42         return df_transformed
43
44     def robust_scaler_transform(self, df_to_transform : pd.DataFrame, with_centering=True, with_scaling=True,
   quantile_range=(5.0, 95.0)):
45         cols = df_to_transform.columns
46         scaler  =  RobustScaler(with_centering=with_centering, with_scaling=with_scaling, quantile_range=
   quantile_range)
47         df_transformed = pd.DataFrame(scaler.fit_transform(df_to_transform))
48         df_transformed.columns = cols
49         return df_transformed
```

```python
1  import os
2
3
4  class XY_metadata :
5
6      def __init__(self,  X_pathname :[], Y_pathname :[], X_join:[], Y_join:[], target_col_name:str):
7          self.X_pathname = os.path.join(X_pathname[0], *X_pathname[1:])
8          self.Y_pathname = None if Y_pathname is None else os.path.join(Y_pathname[0], *Y_pathname[1:])
9          self.X_join = X_join
10         self.Y_join = Y_join
11         self.target_col_name = target_col_name
12
13     def is_training_set(self) -> bool :
14         return True if self.target_col_name is not None else False
15
16     def is_XY_in_separate_file(self) -> bool:
17         return True if self.Y_pathname is not None else False
```

```python
1 from sklearn.impute import SimpleImputer as _SimpleImputer
2
3 from valeo.infrastructure.tools.DfInDfOut import DfInDfOut
4
5
6 class SimpleImputer(_SimpleImputer, DfInDfOut):
7
8     def transform(self, X):
9         Xt = super().transform(X)
10        return super().check_output(Xt, ensure_index=X, ensure_columns=X)
```

```python
1  from sklearn.preprocessing import StandardScaler as _StandardScaler
2
3  from valeo.infrastructure.tools.DfInDfOut import DfInDfOut
4
5
6  class StandardScaler(_StandardScaler, DfInDfOut):
7
8      def transform(self, X):
9          Xt = super().transform(X)
10         return super().check_output(Xt, ensure_index=X, ensure_columns=X)
```

File - C:\EXED\Training\___VALEO\src\valeo\infrastructure\AppConfigManager.py

```python
 1 import valeo.infrastructure.Const as const
 2 from valeo.infrastructure.tools.ConfigLoader import ConfigLoader
 3
 4 class AppConfigManager():
 5
 6     def __init__(self):
 7         cl = AppConfigLoader()
 8         self.app_config = cl.load()
 9
10     def getValue(self, nested_dict:{}, keys:[]) -> str :
11         return nested_dict[keys[0]] if len(keys) ==  1 else self.getValue(nested_dict[keys[0]] , keys[1:])
12
13
14
15 class AppConfigLoader(ConfigLoader) :
16
17     def load(self) -> dict:
18         return super().load(f'{const.rootResources}{const.APP_DEFAULT_CONFIG_FILE}', const.
   ENV_KEY_CONFIG_FILE_PATHNAME)
```

```python
 1 import valeo.infrastructure.Const as const
 2 from valeo.infrastructure.tools.ConfigLoader import ConfigLoader
 3
 4 class AppConfigManager():
 5
 6     def __init__(self):
 7         cl = AppConfigLoader()
 8         self.app_config = cl.load()
 9
10     def getValue(self, nested_dict:{}, keys:[]) -> str :
11         return nested_dict[keys[0]] if len(keys) ==  1 else self.getValue(nested_dict[keys[0]] , keys[1:])
```

```python
1  from datetime import datetime
2
3  from pandas import Series
4  from sklearn.base import BaseEstimator
5
6  from valeo.infrastructure.LogManager import LogManager
7  from valeo.infrastructure import Const as C
8
9  import os
10 import pandas as pd
11 import numpy as np
12
13 class DfUtil() :
14     logger = LogManager.logger(__name__)
15
16
17     # https://stackabuse.com/pythons-classmethod-and-staticmethod-explained/
18     @classmethod
19     def read_csv(cls, pathAsStrList : []) -> pd.DataFrame:
20         try:
21             return pd.read_csv(os.path.join(pathAsStrList[0], *pathAsStrList[1:]) )
22         except Exception as ex :
23             cls.logger.exception("Error while load data from %s", "/".join(pathAsStrList))
24
25     @classmethod
26     def write_y_csv(cls, X_id:Series, y_target: np.ndarray, y_col_name:str, pathAsStrList : [], ts_type=C.
   ts_sfix):
27         DfUtil.write_df_csv( pd.DataFrame(data={X_id.name:X_id, y_col_name:y_target}),  pathAsStrList, ts_type=
   ts_type)
28
29     @classmethod
30     def write_df_csv(cls, df:pd.DataFrame, pathAsStrList : [], ts_type=C.ts_sfix):
31         try :
32             df.to_csv( C.ts_pathanme(pathAsStrList,ts_type), index = False)
33         except Exception as ex:
34             cls.logger = LogManager.logger("DfUtil")
35             cls.logger.exception(f"Error while writing 'df' to CSV '{pathAsStrList}'")
36
37     @classmethod
38     def df_imputer(cls, dfToImpute:pd.DataFrame, imputer:BaseEstimator):
39         '''This method encodes non-null data and replace it in the original data'''
40         # Retains only non-null values. dropna: Remove [rows(default) OR columns] when missing values
41         nonulls = np.array(dfToImpute.dropna())
42         # Reshapes the data for encoding
43         impute_reshape = nonulls.reshape(-1,1)
44         #     #encode date
45         #     impute_ordinal = imputer.fit_transform(impute_reshape)
46         # Assign back encoded values to non-null values
47         dfToImpute.loc[dfToImpute.notnull()] = np.squeeze(imputer.fit_transform(impute_reshape))  # np.squeeze:
   Remove single-dimensional entries from the shape of an array.
48         return dfToImpute
49
50     @classmethod
51     def outlier_ratio(cls, df:pd.DataFrame) -> float:
52         Q1 = df.quantile(0.25)
53         Q3 = df.quantile(0.75)
54         IQR = Q3 - Q1
55         #
56         outliers = ((df < (Q1 - 1.5 * IQR)) |(df > (Q3 + 1.5 * IQR))).any(axis=1)
57         return len(df[outliers].index)/len(df.index)
58
```

File - C:\EXED\Training\___VALEO\src\valeo\infrastructure\tools\ImgUtil.py

```python
1  from valeo.infrastructure import Const as C
2  from valeo.infrastructure.LogManager import LogManager
3  import os
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  import pandas as pd
7
8  class ImgUtil() :
9      logger = LogManager.logger(__name__)
10
11     # https://stackabuse.com/pythons-classmethod-and-staticmethod-explained/
12     @classmethod
13     def save_fig(cls, fig_id:str , tight_layout=True, fig_extension="png", resolution=300, ts_type=C.ts_sfix):
14         path = C.ts_pathanme([C.rootImages() , fig_id + "." + fig_extension], ts_type=ts_type)
15         # cls.logger.debug(f"Saving figure '{fig_id}'")
16         if tight_layout:
17             plt.tight_layout()
18         # Save "the current figure plot" that is set by "df.hist(...))". @ReferTo: pyplot.py / def gcf()
19         plt.savefig(path, format=fig_extension, dpi=resolution)
20
21     @classmethod
22     def save_df_hist_plot(cls, df:pd.DataFrame, fig_id:str , bins=50, figsize=(20,15), tight_layout=True,
23                         fig_extension="png", resolution=300, ts_type=C.ts_sfix):
24         cls.logger.debug(f"Generating 'hist' plot: bins={bins} - figsize={figsize}")
25         df.hist(bins=bins, figsize=figsize)
26         cls.save_fig(fig_id=f"{fig_id}_histogram_{figsize[0]}x{figsize[1]}", tight_layout=tight_layout,
   fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
27
28     @classmethod
29     def save_df_XY_hist_plot(cls, df_XY:pd.DataFrame, fig_id:str, bins=50, figsize=(5, 5), y_target_name=None,
   tight_layout=True,
30                             fig_extension="png", resolution=300, ts_type=C.ts_sfix):
31         cls.logger.debug(f"Generating 'XY_hist' plot: bins={bins} - figsize={figsize}")
32         df_X = df_XY.drop(columns=y_target_name, axis=1)
33         y    = df_XY[y_target_name]
34         fig, ax = plt.subplots(figsize=figsize)
35         for i, col in enumerate(sorted(df_X.columns)) :
36             for clazz in y.unique() :
37                 df_X[y==clazz][col].plot.hist(bins=bins, figsize=figsize, alpha=0.3, label=f'Class #{int(clazz)}
   ')
38             plt.legend()
39             plt.xlabel(col)
40             ImgUtil.save_fig(fig_id=f"{fig_id}_{col}_histogram_{figsize[0]}x{figsize[1]}", tight_layout=
   tight_layout, fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
41             ax.clear()
42     # NB:
43     # df.hist: => Plot 1 Histo per dfColumn
44     # df.plot.hist: => Plot all df-referenced-Columns on same Histo
45
46     @classmethod
47     def save_df_scatter_matrix_plot(cls, df:pd.DataFrame, fig_id:str , figsize=(20,15), cfield=None, tight_layout
   =True,
48                                     fig_extension="png", resolution=300, ts_type=C.ts_sfix):
49         cls.logger.debug(f"Generating 'scatter matrix' plot: figsize:{figsize}")
50         if cfield == None :
51             pd.plotting.scatter_matrix(df, figsize=figsize)
52         else :
53             pd.plotting.scatter_matrix(df, figsize=figsize,  alpha=0.3, c=df[cfield].values, cmap='RdBu')
54         cls.save_fig(fig_id=f"{fig_id}_scatter_matrix_{figsize[0]}x{figsize[1]}", tight_layout=tight_layout,
   fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
55
56     @classmethod
57     def save_df_heatmap_plot(cls, df:pd.DataFrame, fig_id:str , figsize=(20,20), cmap='RdBu', annot=True ,
   annot_kws={'size':15}, fig_extension="png", resolution=300, ts_type=C.ts_sfix):
58         cls.logger.debug(f"Generating 'heatmap' plot: figsize:{figsize}")
59         fig, ax = plt.subplots(figsize=figsize)
60         sns.set(font_scale=1.1)
61         sns.heatmap(df, cmap=cmap, annot=annot , annot_kws=annot_kws, ax=ax)
62         ax.set_title(fig_id, fontsize=28)
63         cls.save_fig(fig_id=f"{fig_id.replace(' ','_')}_heatmap_{figsize[0]}x{figsize[1]}", tight_layout=True,
   fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
64
65     @classmethod
66     def save_df_violin_plot(cls, df:pd.DataFrame, fig_id:str, grid_elmt_x:int, figsize=(20,20), fig_extension="
   png", resolution=300, ts_type=C.ts_sfix):
67         cls.logger.debug(f"Generating 'violin' plot: figsize:{figsize}")
68         grid_elmt_y = len(df.columns) // grid_elmt_x if (len(df.columns) % grid_elmt_x) == 0 else (len(df.
   columns) // grid_elmt_x) + 1
69         #
70         fig, axs = plt.subplots(grid_elmt_y, grid_elmt_x, figsize=figsize)
71         for i, col in enumerate(sorted(df.columns)) :
72             sns.violinplot(x=df[col], linewidth=1, ax=axs[i//grid_elmt_x, i%grid_elmt_x])
73             # sns.stripplot( x=df[col], color="orange", jitter=0.2, linewidth=1, ax=axs[i//3,i%3])
74             sns.boxplot( x=df[col], linewidth=1, ax=axs[i//grid_elmt_x, i%grid_elmt_x], saturation=0 )
75             # axs.set_title(fig_id, fontsize=28)
76         cls.save_fig(fig_id=f"{fig_id.replace(' ','_')}_violin_{figsize[0]}x{figsize[1]}", tight_layout=True,
   fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
77
```

```python
78
79  def save_df_XY_violin_plot(df_XY:pd.DataFrame, y_target_name:str, fig_id:str, grid_elmt_x:int, figsize=(20,20
    ), fig_extension="png", resolution=300, ts_type=Const.ts_sfix):
80      df = df_XY.drop(columns=y_target_name, axis=1)
81      grid_elmt_y = len(df.columns) // grid_elmt_x if (len(df.columns) % grid_elmt_x) == 0 else (len(df.columns
    ) // grid_elmt_x) + 1
82      #
83      fig, axs = plt.subplots(grid_elmt_y, grid_elmt_x, figsize=figsize)
84      for i, col in enumerate(sorted(df.columns)) :
85          sns.violinplot(x=y_target_name, y=col, data=df_XY, linewidth=1, ax=axs[i//grid_elmt_x, i%grid_elmt_x])
86          sns.boxplot    (x=y_target_name, y=col, data=df_XY, linewidth=1, ax=axs[i//grid_elmt_x, i%grid_elmt_x] )
87      ImgUtil.save_fig(fig_id=f"{fig_id.replace(' ','_')}_violin_{figsize[0]}x{figsize[1]}", tight_layout=True,
    fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
88
89
90      @classmethod
91      # SWARM PLOT did not work correctly
92      def save_df_swarm_plot(cls, df:pd.DataFrame, fig_id:str, grid_elmt_x:int, figsize=(20,20), cfield=None,
    fig_extension="png", resolution=300, ts_type=C.ts_sfix):
93          cls.logger.debug(f"Generating 'swarm' plot: figsize:{figsize}")
94          grid_elmt_y = len(df.columns) // grid_elmt_x if (len(df.columns) % grid_elmt_x) == 0 else (len(df.
    columns) // grid_elmt_x) + 1
95          #
96          fig, axs = plt.subplots(grid_elmt_y, grid_elmt_x, figsize=figsize)
97          for i, col in enumerate(sorted(df.columns)) :
98              sns.swarmplot(x=df[col], linewidth=1, ax=axs[i//grid_elmt_x, i%grid_elmt_x], hue=df[cfield].values)
99          cls.save_fig(fig_id=f"{fig_id.replace(' ','_')}_swarm_{figsize[0]}x{figsize[1]}", tight_layout=True,
    fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
100
101     # def save_df_swarm_plot(df_XY:pd.DataFrame, fig_id:str, figsize=(5,5), y_target_name=None, fig_extension="
    png", resolution=300, ts_type=Const.ts_sfix):
102     #     df_X = df_XY.drop(columns=y_target_name, axis=1)
103     #     y    = df_XY[y_target_name]
104     #     fig, ax = plt.subplots(figsize=figsize)
105     #     for i, col in enumerate(sorted(df_X.columns)) :
106     #         for clazz in y.unique() :
107     #             sns.swarmplot(x=col, hue=y_target_name, data=df_XY[y==clazz])
108     #             # sns.swarmplot(x='data', y='feature',  hue='label', data=df)
109     #         plt.legend()
110     #         plt.xlabel(col)
111     #         ImgUtil.save_fig(fig_id=f"{fig_id}_{col}_swarm_{figsize[0]}x{figsize[1]}", tight_layout=
    tight_layout, fig_extension=fig_extension, resolution=resolution, ts_type=ts_type)
112     #         ax.clear()
113
114
115 # NB: Generic way to plot whathever :
116 # fig, ax = plt.subplots(figsize=(20,20))
117 # sns.heatmap(corr_matrix, cmap='RdBu', annot=True , annot_kws={'size':15}, ax=ax)
118 # ax.set_title("Valeo starter production correlation measures", fontsize=14)
119 # plt.show()
120
121
```

```python
import pandas as pd

class DfInDfOut:

    # https://github.com/scikit-learn/scikit-learn/issues/5523 : Pandas in, Pandas out
    def check_output(self, X, ensure_index=None, ensure_columns=None):
        """
        Joins X with ensure_index's index or ensure_columns's columns when avaialble
        """
        if ensure_index is not None:
            if ensure_columns is not None:
                if type(ensure_index) is pd.DataFrame and type(ensure_columns) is pd.DataFrame:
                    X = pd.DataFrame(X, index=ensure_index.index, columns=ensure_columns.columns)
            else:
                if type(ensure_index) is pd.DataFrame:
                    X = pd.DataFrame(X, index=ensure_index.index)
        return X
```

```python
1
2  import os
3  import yaml
4  import logging
5
6  # from infrastructure.LogManager import LogManager
7
8  class YamlLoader :
9      """
10     Load a yaml  configuration file
11     """
12     logger = None
13
14     def __init__(self):
15         if YamlLoader.logger is None :
16             YamlLoader.logger =  logging.getLogger(__name__)
17             logging.basicConfig(level=logging.INFO)
18
19     def load(self, file_pathname:str) -> dict :
20         if os.path.exists(file_pathname):
21             with open(file_pathname, 'rt') as f:
22                 try:
23                     dict =  yaml.safe_load(f.read())
24                     # YamlLoader.logger.info(f'Loading file "{file_pathname}":\n\t{dict}')
25                     YamlLoader.logger.info(f'Loading file "{file_pathname}":\n{dict}')
26                     return dict
27                 except Exception as ex:
28                     YamlLoader.logger.exception(f'Error while loading file "{file_pathname}"')
29         else:
30             YamlLoader.logger.error(f'Error while loading file "{file_pathname}"')
31
32         return None
33
```

File - C:\EXED\Training\___VALEO\src\valeo\infrastructure\tools\ConfigLoader.py

```python
1
2 import os
3 import logging
4
5 from valeo.infrastructure.tools.YamlLoader import YamlLoader
6
7 class ConfigLoader(YamlLoader) :
8     logger = None
9     """
10    Load an external or a package embedded configuration file.
11    Check first if the environment variable {APP_CONFIG_PATHNAME}
12    """
13
14    def __init__(self):
15        super().__init__()
16        ConfigLoader.logger = logging.getLogger(__name__)
17
18    def load(self, file_pathname:str, env_key_as_config_pathname:str) -> dict :
19        try :
20            path_as_key = os.getenv(env_key_as_config_pathname, None)
21            return super().load(path_as_key if path_as_key else file_pathname )
22        except Exception as ex :
23            ConfigLoader.logger.exception(f'Error while loading file "{file_pathname}"')
24            raise ex
25            # self.logger.error(ex, exc_info=True)
26        return None
27
```

```python
import os
from datetime import datetime

from sklearn.base import BaseEstimator, TransformerMixin
import numpy as np

from valeo.infrastructure import Const as C


class DebugPipeline(BaseEstimator, TransformerMixin):
    OFFSET  = 10
    counter = -OFFSET

    def __init__(self):
        DebugPipeline.counter = ( (DebugPipeline.counter + DebugPipeline.OFFSET) // DebugPipeline.OFFSET) * DebugPipeline.OFFSET

    def transform(self, X, y=None):
        # %f : print micro seconds
        # np.savetxt(os.path.join(C.rootProject(), 'log', 'dbgPipeline_' + datetime.now().strftime("%Y_%m_%d-%H.%M.%S_") + str(DebugPipeline.counter)) + '.txt', X, delimiter=',')
        DebugPipeline.counter += 1

        return X

    def fit(self, X, y=None, **fit_params):
        if y is not None :
            # np.savetxt(os.path.join(C.rootProject(), 'log', 'dbgPipeline_' + datetime.now().strftime("%Y_%m_%d-%H.%M.%S_") + str(DebugPipeline.counter)) + '_Y.txt', y, delimiter=',')
            # np.savetxt(os.path.join(C.rootProject(), 'log', 'dbgPipeline_' + datetime.now().strftime("%Y_%m_%d-%H.%M.%S_") + str(DebugPipeline.counter)) + '_X.txt', X, delimiter=',')
            DebugPipeline.counter += 1
        return self
```