

# Credit card fraud - imbalance dataset with SMOTE

References:

- SMOTE documentation
- Metrics specific to imbalanced learning
- How to Effortlessly Handle Class Imbalance with Python and SMOTE\*
- SMOTE Users Python
- SMOTE for Imbalanced Classification with Python
- Overcoming Class Imbalance using SMOTE Techniques\*\*
- dataset - credit card fraud

```
In [36]: import sys
sys.path.append('./')
sys.path.append('./')

import warnings
warnings.filterwarnings('ignore', category=FutureWarning)
warnings.filterwarnings('ignore', category=UserWarning)
```

```
In [2]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.classification import *
from sklearn.preprocessing import *
from sklearn.metrics import *
```

```
In [3]: df = pd.read_csv('./datasets/credit_dataset.csv')
df.head()
```

target	count	percent
0 Normal	24712	98.32
1 Fraud	422	1.68

A bar chart illustrating the distribution of the 'TARGET' variable. The x-axis is labeled 'TARGET' and has two categories: 'Normal' and 'Fraud'. The y-axis is labeled 'count' and ranges from 0 to 25,000. The 'Normal' category has a count of 24,712, represented by a blue bar. The 'Fraud' category has a count of 422, represented by an orange bar.

TARGET	count
Normal	24712
Fraud	422

Overly unbalance dataset, most likely would make high-accuracy but low-recall models

## Data descriptions.

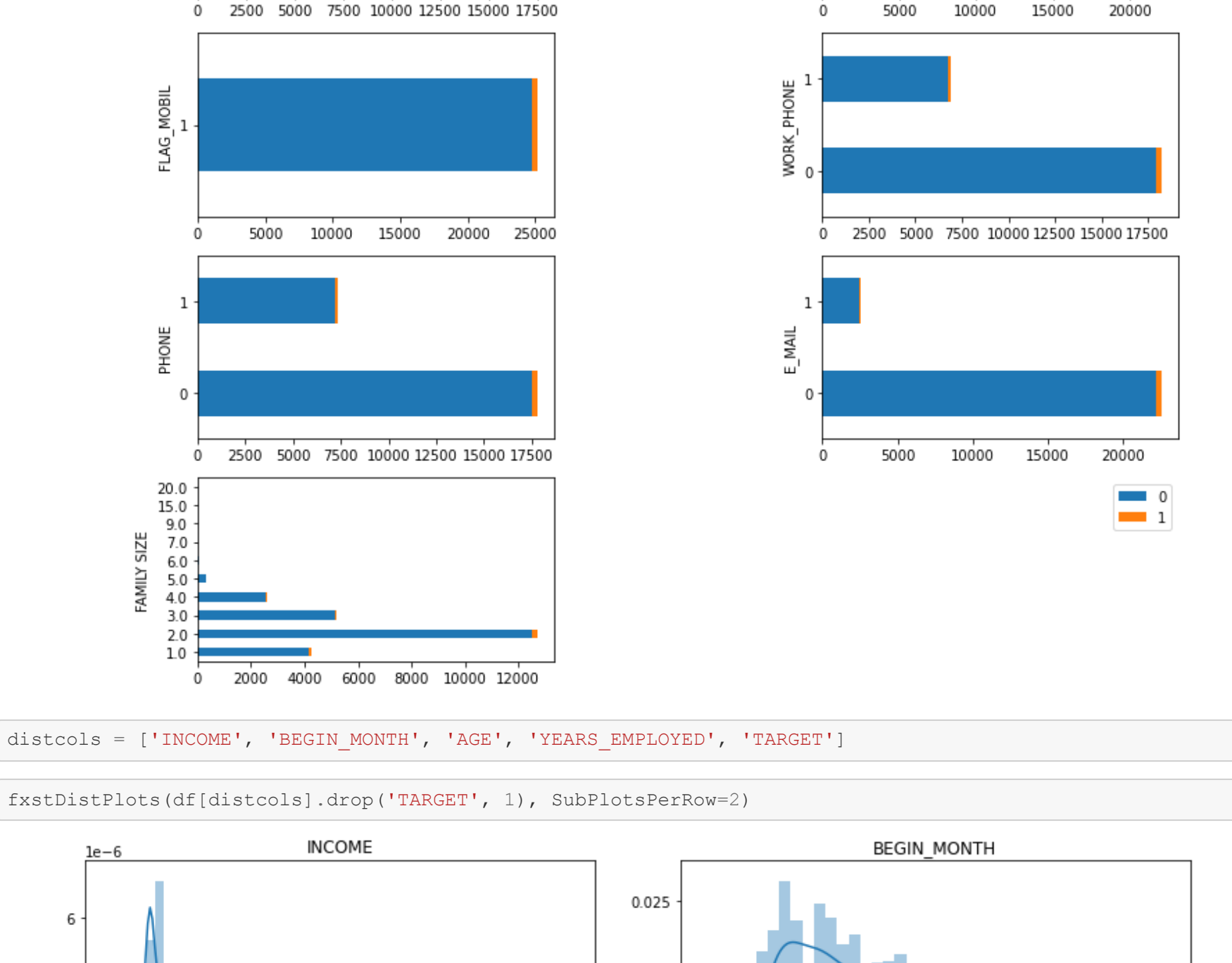
- BEGIN\_MONTH - The month of the extracted data is the starting point, backwards, 0 is the current month, -1 is the previous month, and so on
- REALITY - Is there a property

```
In [4]: df.columns

Out [4]: Index(['Unnamed: 0', 'ID', 'GENDER', 'CAR', 'REALITY', 'NO_OF_CHILD', 'INCOME', 'INCOME_TYPE', 'EDUCATION_TYPE', 'FAMILY_TYPE', 'HOUSE_TYPE', 'FLAG_MOBIL', 'WORK_PHONE', 'PHONE', 'E_MAIL', 'FAMILY_SIZE', 'BEGIN_MONTH', 'AGE', 'YEARS_EMPLOYED', 'TARGET'], dtype='object')
```

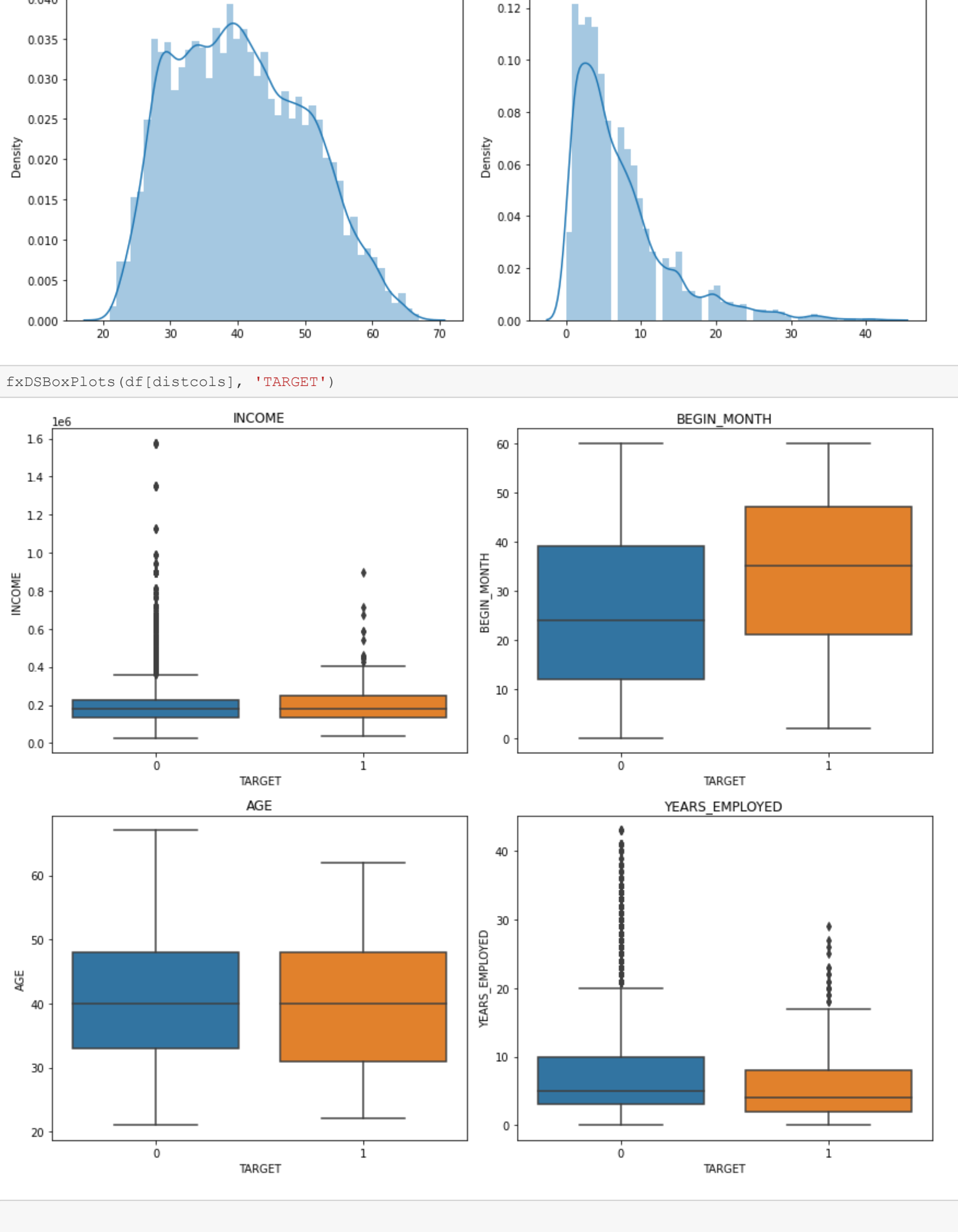
```
In [5]: fxCountPlot(df, 'TARGET', targetlabel=(0: 'Normal', 1: 'Fraud'), plotitem='count')

Out [5]:
```



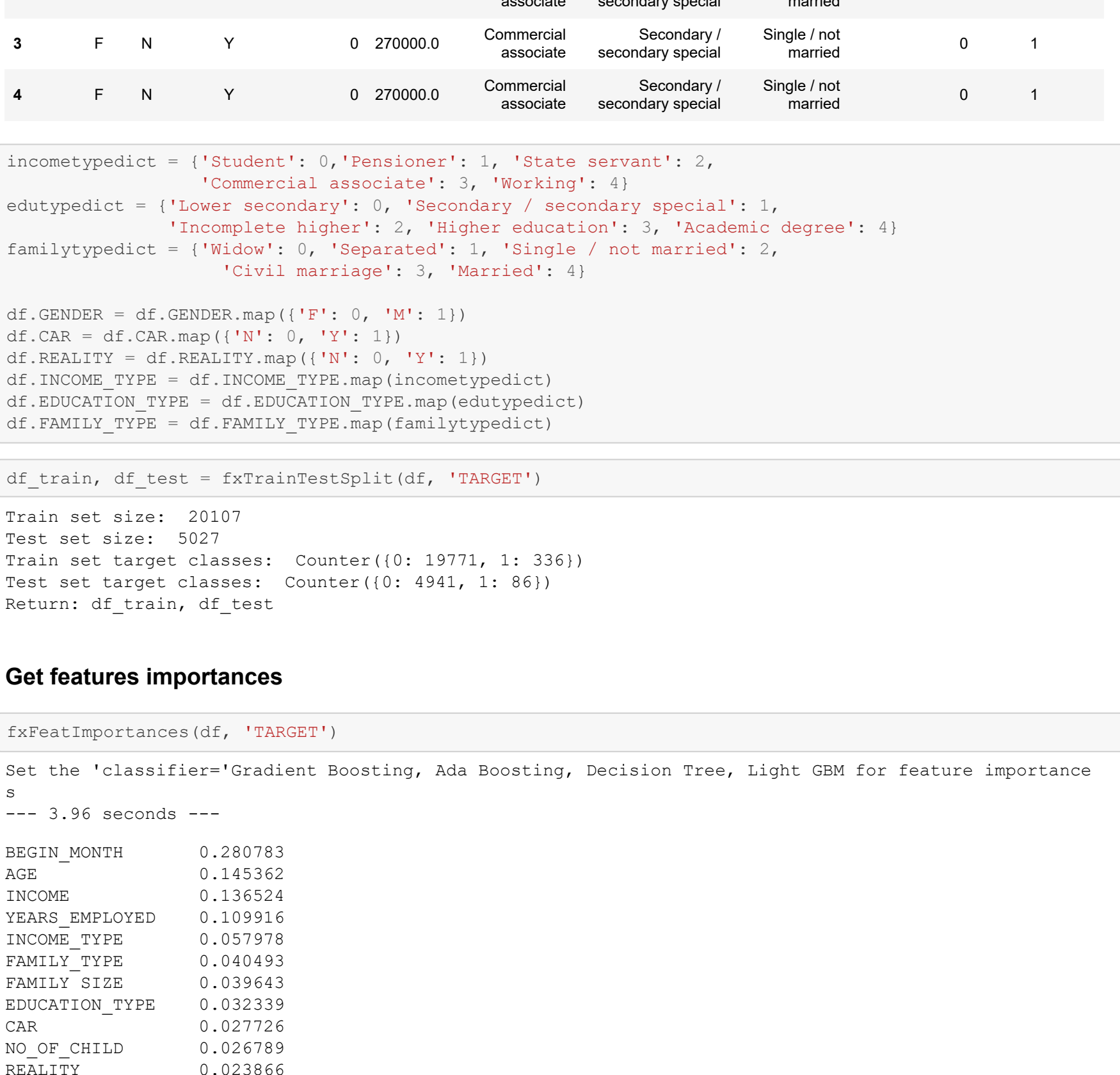
Overly unbalance dataset, most likely would make high-accuracy but low-recall models

```
In [6]: fxFactPlots(df, subcat='TARGET', ignoreCatType=True, subPlotsPerRow=2,
excl=('Unnamed: 0', 'ID', 'INCOME', 'BEGIN_MONTH', 'AGE', 'YEARS_EMPLOYED'))
```

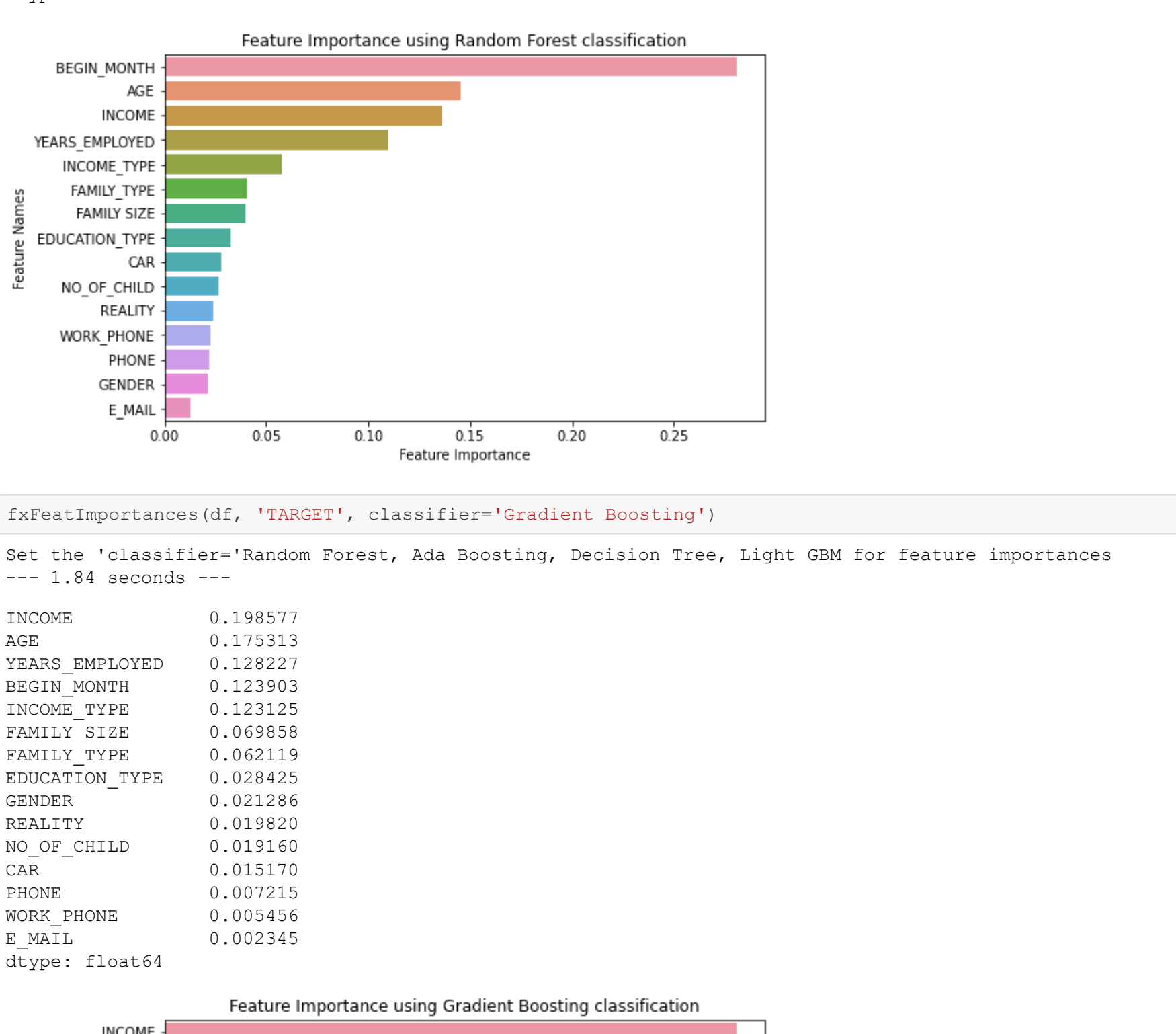


```
In [7]: distcols = ['INCOME', 'BEGIN_MONTH', 'AGE', 'YEARS_EMPLOYED', 'TARGET']
```

```
In [8]: fxFactPlots(df[distcols].drop('TARGET', 1), subPlotsPerRow=2)
```



```
In [9]: fxBoxPlots(df[distcols], 'TARGET')
```



```
In [ ]:
```

## Data preparation for machine learning

```
In [10]: droptems = ['Unnamed: 0', 'ID', 'FLAG_MOBIL', 'HOUSE_TYPE']
df.drop(droptems, axis=1, inplace=True)
df.head()
```

```
Out [10]:
```

GENDER	CAR	REALITY	NO_OF_CHILD	INCOME	INCOME_TYPE	EDUCATION_TYPE	FAMILY_TYPE	WORK_PHONE	PHONE	E_MAIL
0	M	Y	Y	0	112500.0	Working	Secondary / secondary special	Married	0	0
1	F	N	Y	0	270000.0	Commercial associate	Secondary / secondary special	Single / not married	0	1
2	F	N	Y	0	270000.0	Commercial associate	Secondary / secondary special	Single / not married	0	1
3	F	N	Y	0	270000.0	Commercial associate	Secondary / secondary special	Single / not married	0	1
4	F	N	Y	0	270000.0	Commercial associate	Secondary / secondary special	Single / not married	0	1

```
In [11]: incometypedict = {'Student': 0, 'Pensioner': 1, 'State servant': 2,
                        'Commercial associate': 3, 'Working': 4}
edutypedict = {'Lower secondary': 0, 'Secondary / secondary special': 1,
               'Incomplete higher': 2, 'Higher education': 3, 'Academic degree': 4}
familydict = {'Widow': 0, 'Separated': 1, 'Single / not married': 2,
              'Civil marriage': 3, 'Married': 4}

df.GENDER = df.GENDER.map({'F': 0, 'M': 1})
df.CAR = df.CAR.map({'N': 0, 'Y': 1})
df.REALITY = df.REALITY.map({'N': 0, 'Y': 1})
df.INCOME_TYPE = df.INCOME_TYPE.map(incometypedict)
df.EDUCATION_TYPE = df.EDUCATION_TYPE.map(edutypedict)
df.FAMILY_TYPE = df.FAMILY_TYPE.map(familydict)
```

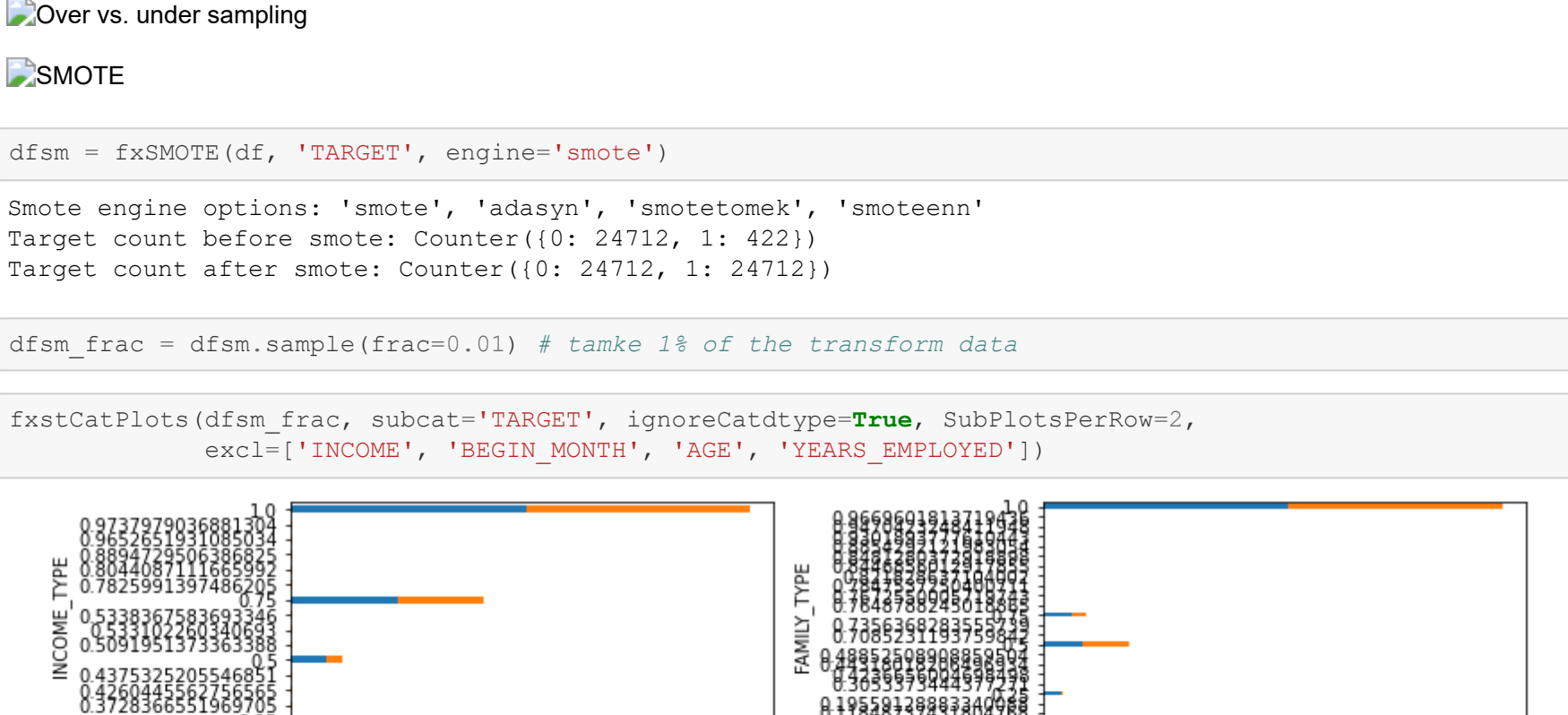
```
In [12]: df_train, df_test = fxTrainTestSplit(df, 'TARGET')
```

Train set size: 20107  
Test set size: 5027  
Train set target classes: Counter((0: 19771, 1: 336))  
Test set target classes: Counter((0: 4941, 1: 86))  
Return: df\_train, df\_test

## Get features importances

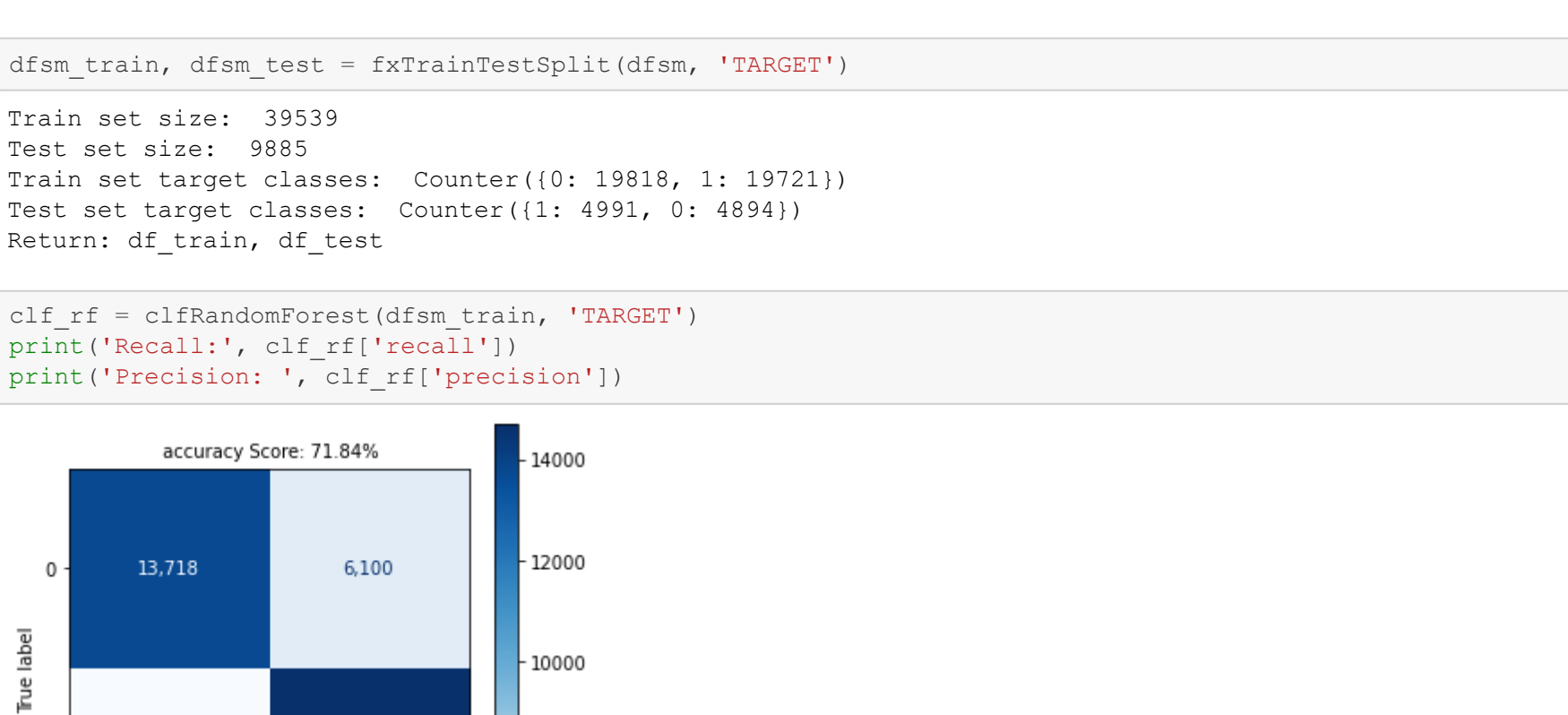
```
In [13]: fxFeatImportances(df, 'TARGET')
```

Set the 'classifier'='Gradient Boosting, Ada Boosting, Decision Tree, Light GBM for feature importance  
--- 3.96 seconds ---



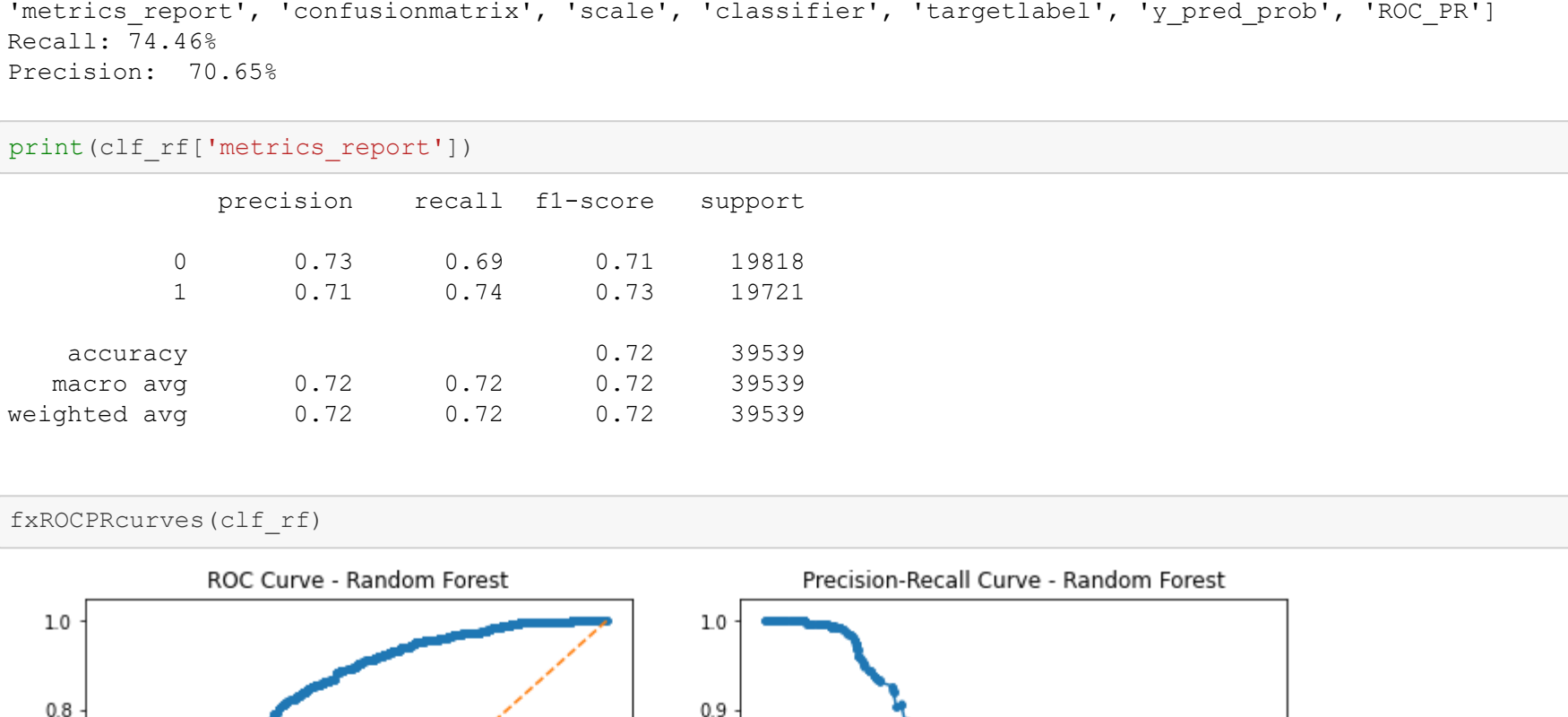
```
In [14]: fxFeatImportances(df, 'TARGET', classifier='Gradient Boosting')
```

Set the 'classifier'='Random Forest, Ada Boosting, Decision Tree, Light GBM for feature importances  
--- 1.84 seconds ---



```
In [15]: fxFeatImportances(df, 'TARGET', classifier='Ada Boosting')
```

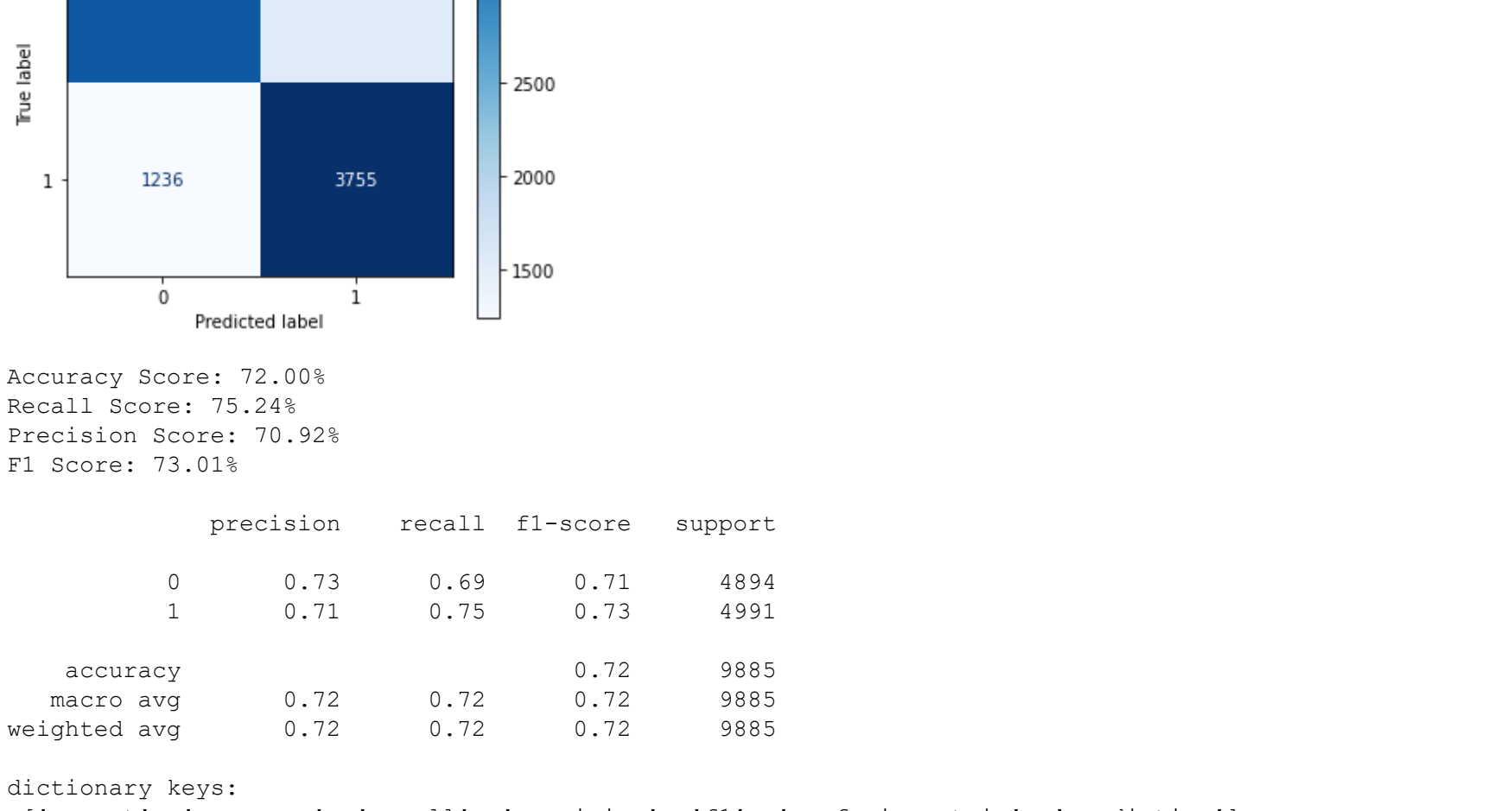
Set the 'classifier'='Random Forest, Gradient Boosting, Decision Tree, Light GBM for feature importances  
--- 0.76 seconds ---



```
In [16]: # drop less important features
df.drop(['CAR', 'EDUCATION_TYPE', 'GENDER', 'NO_OF_CHILD', 'REALITY',
        'E_MAIL', 'PHONE', 'WORK_PHONE', 'AGE', 'YEARS_EMPLOYED'], axis=1, inplace=True)
```

## Feed into Machine Learning model (without SMOTE)

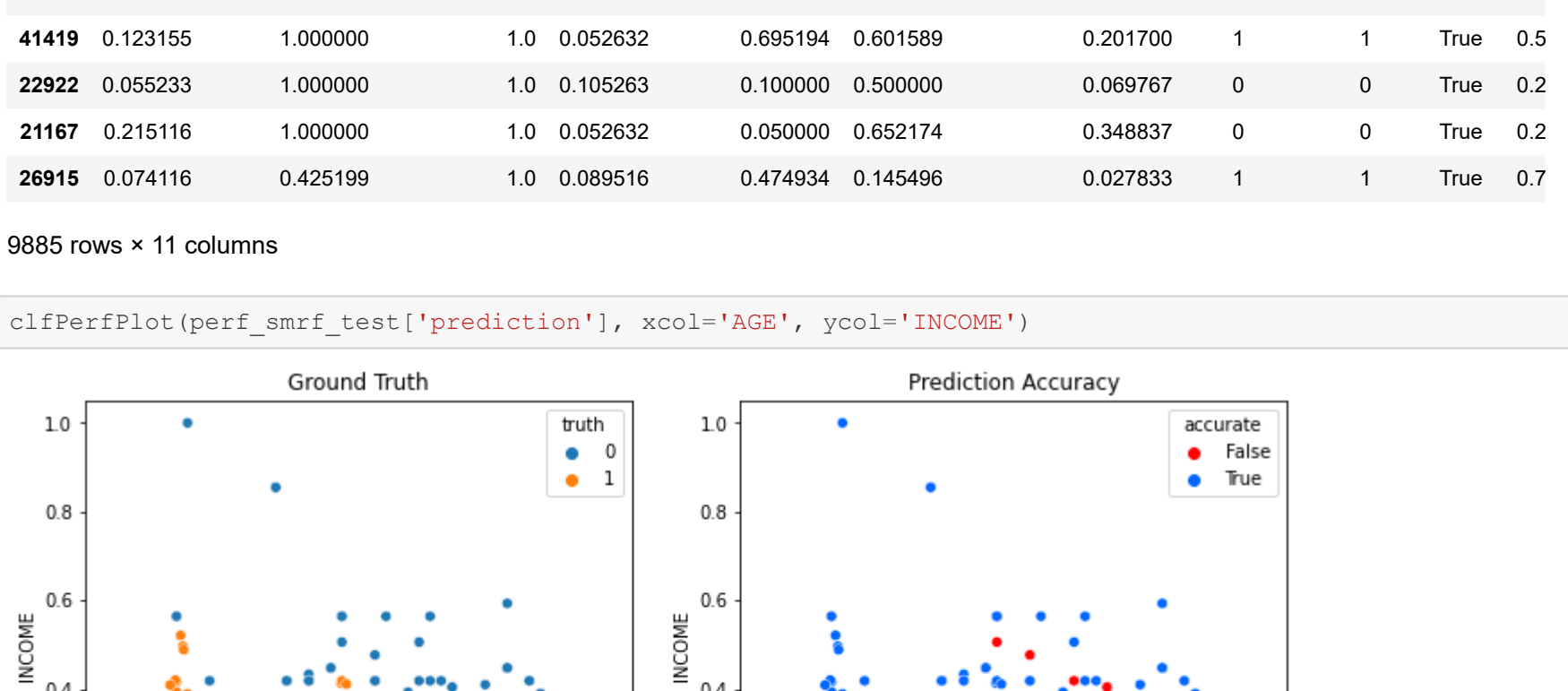
```
In [37]: clf_rf = clfRandomForest(df_train, 'TARGET')
print('Recall:', clf_rf['recall'])
print('Precision:', clf_rf['precision'])
```



Accuracy: 98%  
Recall: 0.6%  
Precision: 100.0%  
F1: 98.58%  
dictionary keys:  
['type', 'prediction', 'accuracy', 'recall', 'precision', 'f1', 'df\_proba', 'params', 'importances', 'metrics\_report', 'confusionmatrix', 'scale', 'classifier', 'targetlabel', 'y\_pred\_prob', 'ROC\_PR']  
Recall: 0.6%  
Precision: 100.0%

- All are being predicted as 0, the machine learning basically does not learn anything

```
In [18]: fxROCPRCurves(clf_rf, plotsize=(12, 5))
```



## Using SMOTE to balance the dataset

### Synthetic Minority Oversampling Technique

Over vs. under sampling

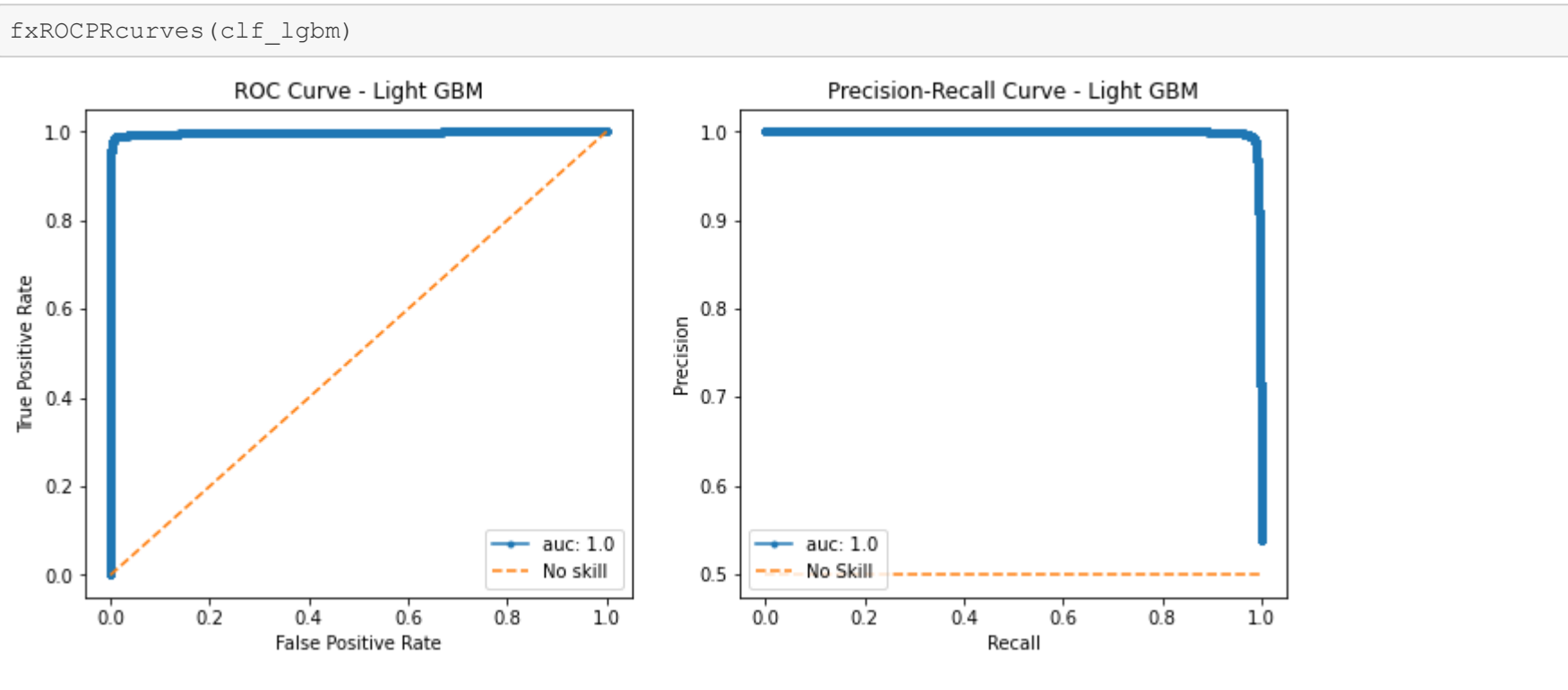
SMOTE

```
In [19]: dfsm = fxSMOTE(df, 'TARGET', engine='smote')
```

Smote engine options: 'smote', 'adasyn', 'smotetomek', 'smoteenn'  
Target count before smote: Counter((0: 24712, 1: 422))  
Target count after smote: Counter((0: 24712, 1: 422))

```
In [20]: dfsm_frac = dfsm.sample(frac=0.01) # take 1% of the transform data
```

```
In [21]: fxFactPlots(dfsm_frac, subcat='TARGET', ignoreCatType=True, subPlotsPerRow=2,
excl=('INCOME', 'BEGIN_MONTH', 'AGE', 'YEARS_EMPLOYED'))
```



## Retrain the model with transformed data

```
In [22]: dfsm_train, dfsm_test = fxTrainTestSplit(dfsm, 'TARGET')
```

Train set size: 39539  
Test set size: 9855  
Train set target classes: Counter((0: 19818, 1: 19721))  
Test set target classes: Counter((0: 4991, 1: 4844))  
Return: df\_train, df\_test

```
In [23]: clf_rf = clfRandomForest(dfsm_train, 'TARGET')
print('Recall:', clf_rf['recall'])
print('Precision:', clf_rf['precision'])
```



Accuracy: 72%  
Recall: 74.46%  
Precision: 70.65%  
F1: 72.51%  
dictionary keys:  
['type', 'prediction', 'accuracy', 'recall', 'precision', 'f1', 'df\_proba', 'params', 'importances', 'metrics\_report', 'confusionmatrix', 'scale', 'classifier', 'targetlabel', 'y\_pred\_prob', 'ROC\_PR']  
Recall: 74.46%  
Precision: 70.65%

```
In [24]: print(clf_rf['metrics_report'])
```

	precision	recall	f1-score	support
0	0.73	0.69	0.71	19818
1	0.71	0.74	0.73	19721
accuracy			0.72	39539
macro avg	0.72	0.72	0.72	39539
weighted avg	0.72	0.72	0.72	39539

```
In [25]: fxROCPRCurves(clf_rf)
```



```
In [26]: # fxLearningCurve(dfsm_train, 'TARGET', clf_rf['classifier']) # take long time
```

## Perform on test data

```
In [27]: perf_smrf_test = clfTestsdPerformance(dfsm_test, 'TARGET', clf_rf['classifier'])
print('Recall:', clf_rf['recall'])
print('Precision:', clf_rf['precision'])
```



Accuracy Score: 72.00%  
Recall Score: 75.24%  
Precision Score: 70.92%  
F1 Score: 73.01%

	precision	recall	f1-score	support
0	0.73	0.69	0.71	4894
1	0.71	0.75	0.73	4991
accuracy			0.72	9885
macro avg	0.72	0.72	0.72	9885
weighted avg	0.72	0.72	0.72	9885

dictionary keys:  
['report', 'accuracy', 'recall', 'precision', 'f1', 'confusionmatrix', 'prediction']  
Recall: 74.46%  
Precision: 70.65%

```
In [28]: perf_smrf_test['prediction']
```

```
Out [28]:
```

INCOME	INCOME_TYPE	FAMILY_TYPE	FAMILY_SIZE	BEGIN_MONTH	AGE	YEARS_EMPLOYED	truth	prediction	accurate	like	
31542	0.142442	0.750000	1.0	0.052632	0.757860	0.456522	0.046512	1	1	True	0.6
25601	0.049771	1.000000	0.5	0.004419	0.814271	0.640523	0.071820	1	1	True	0.6
1536	0.063953	1.000000	0.0	0.157895	0.200000	0.282609	0.302326	0	0	True	0.3
1373	0.098837	1.000000	0.0	0.000000	0.200000	0.717391	0.116279	0	0	True	0.4
19409	0.098837	1.000000	0.5	0.000000	0.516667	0.217391	0.232558	0	0	True	0.4
...	...	...	...	...	...	...	...	...	...	...	...
32887	0.185510	0.677990	1.0	0.105263	0.209601	0.586957	0.126138	1	0	False	0.4
44119	0.123155	1.000000	1.0	0.052632	0.695194	0.601589	0.201700	1	1	True	0.5
22322	0.055233	1.000000	1.0	0.105263	0.100000	0.500000	0.089767	0	0	True	0.2
21617	0.215116	1.000000	1.0	0.052632	0.650000	0.652174	0.348837	0	0	True	0.2
26915	0.074116	0.425199	1.0	0.089516	0.474934	0.145496	0.027933	1	1	True	0.7

9885 rows \* 11 columns

```
In [29]: clfPerfPlot(perf_smrf_test['prediction'], xcol='AGE', ycol='INCOME')
```



## Using Light GBM classifier

```
In [30]: clf_lgbm = clfLightGBM(dfsm_train, 'TARGET')
print('Recall:', clf_lgbm['recall'])
print('Precision:', clf_lgbm['precision'])
```



Accuracy: 97.65%  
Recall: 99.53%  
Precision: 99.53%  
F1: 98.58%  
dictionary keys:  
['type', 'prediction', 'accuracy', 'recall', 'precision', 'f1', 'df\_proba', 'params', 'importances', 'metrics\_report', 'confusionmatrix', 'scale', 'classifier', 'targetlabel', 'y\_pred\_prob', 'ROC\_PR']  
Recall: 97.65%  
Precision: 99.53%

```
In [31]: fxROCPRCurves(clf_lgbm)
```



```
In [32]: perf_sm_lgbm_test = clfTestsdPerformance(dfsm_test, 'TARGET', clf_lgbm['classifier'])
```



Accuracy Score: 98.00%  
Recall Score: 97.32%  
Precision Score: 99.24%  
F1 Score: 98.27%

	precision	recall	f1-score	support
0	0.97	0.99	0.98	4894
1	0.99	0.97	0.98	4991
accuracy			0.98	9885
macro avg	0.98	0.98	0.98	9885
weighted avg	0.98	0.98	0.98	9885

dictionary keys:  
['report', 'accuracy', 'recall', 'precision', 'f1', 'confusionmatrix', 'prediction']



```
In [33]: clfPerfPlot(perf_smlgbm_test['prediction'], xcol='AGE', ycol='INCOME')
```

