

Universidade de São Paulo
Instituto de Física de São Carlos

SCC00277-201-2021 Project 2

Éverton Luís Mendes da Silva (10728171)

Contents

1	Introduction	2
2	Question One	2
2.1	Item a	2
2.2	Item b	3
2.3	Ideal scenario	3
2.4	Scenario with Errors	4
2.5	Random scenario	5
2.6	Inverse scenario	6
3	Questão 2	7
3.1	Item a and b	7
4	Question Three	9
4.1	Feature Engineering	11
4.2	Undersample and Bootstrap	12
4.3	Preprocessing	13
5	Question Four	15
5.1	Hyperparameter dependency	16
5.1.1	Bernoulli Naive Bayes	16
5.1.2	Passive Agressive Classifier	17
5.1.3	Ridge Classifier Negative	17
5.1.4	Ridge Classifier Positive	18
5.1.5	Perceptron	18
5.1.6	Stochastic Gradient Descent Classifier	19
5.1.7	Multi-Layer Perceptron Classifier	20
5.1.8	Quadratic Discriminant Analysis	21
5.2	Best Model	21
5.3	My Kaggle Submission	22
6	Question Five	22
7	Reference	26

1 Introduction

Undoubtedly, with the advent of technology, the use of cards (credit or debit) has increased over the years. Thus, both banks and customers are concerned about the security that this type of operation can provide. With this in mind, this project aims to build models for the prediction of fraud in transactions provided by the kaggle database, IEEE- CIS Fraud Detection. Therefore, in order to predict frauds, 7 types of ML (machine learning) models were trained using data modeling techniques (Pipelines for Preprocessing) and hyperparameter optimization (Bayesian optimization).

2 Question One

2.1 Item a

For the first part of this project, the AUC (area under curve) metric was analyzed for three types of models, each model are listed below:

- Model that randomly ranks (50% chance of saying it is fraud and 50% of saying it's not) (a)
- Model that classifies all cases as fraud (b)
- Model that classifies all cases as non-fraud. (c)

Furthermore, we can find below the AUC measurements for each of the models listed.

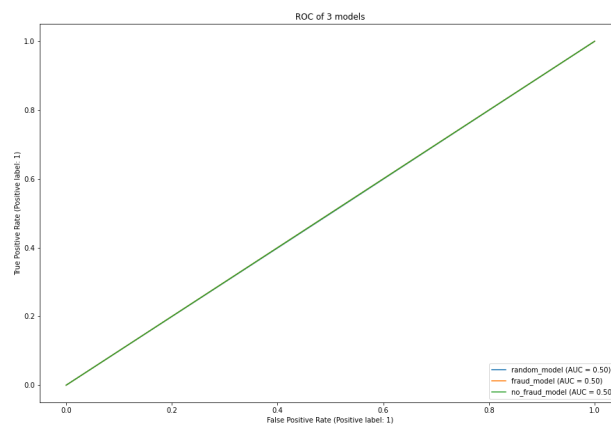


Figure 1: AUC of 3 models

2.2 Item b

To evaluate the graph above, it is necessary to understand a little more about this method of measuring performance in binary classifications, that is, to understand how Area Under the Receiver Operating Characteristics has the ability to distinguish classes. For example, the higher the AUC value, the better the model is at predicting classes 0 as 0 and 1 as 1 (sick patients as sick and healthy as healthy).

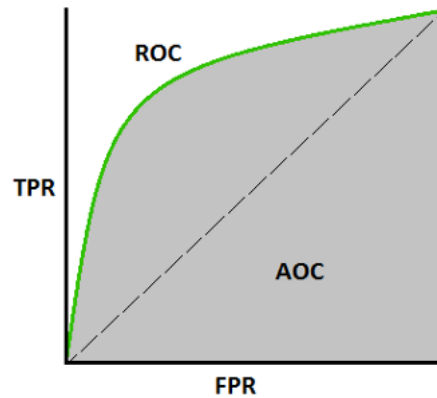


Figure 2: AUC-ROC(Roc Curve)

With this in mind, below we have possible scenarios for the distribution of probabilities for each class.

2.3 Ideal scenario

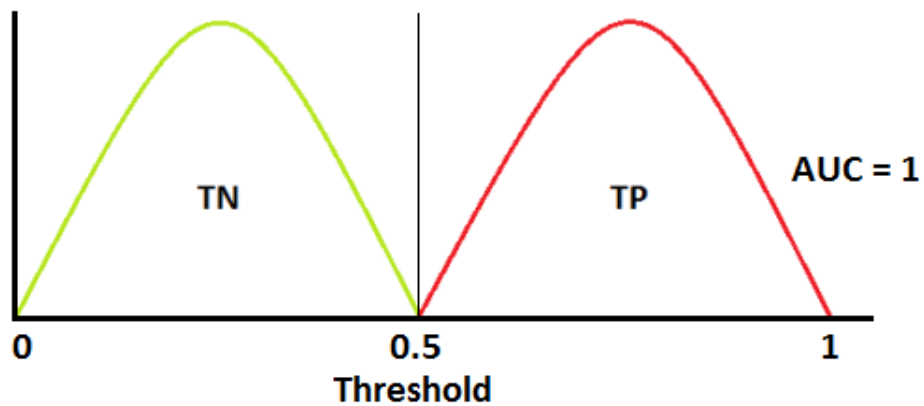


Figure 3: Probabilities

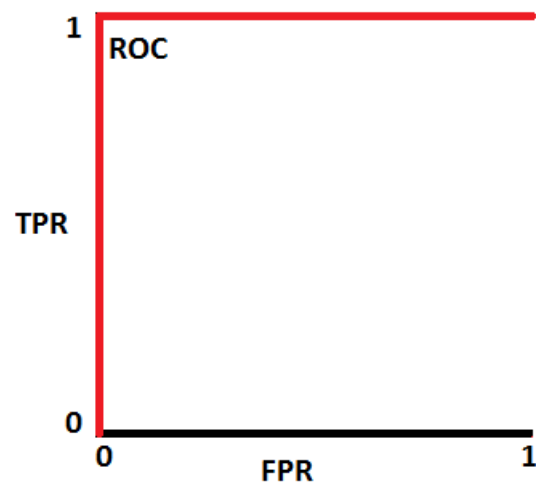


Figure 4: Roc Curve

2.4 Scenario with Errors

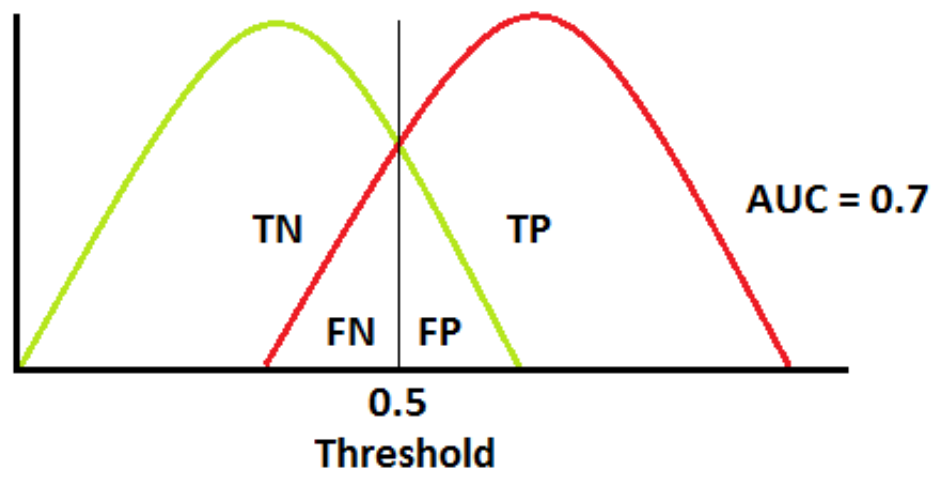


Figure 5: Probabilities

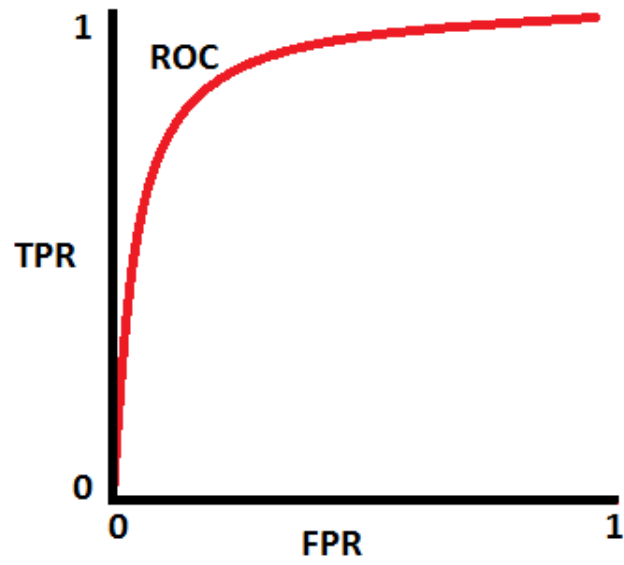


Figure 6: Roc Curve

2.5 Random scenario

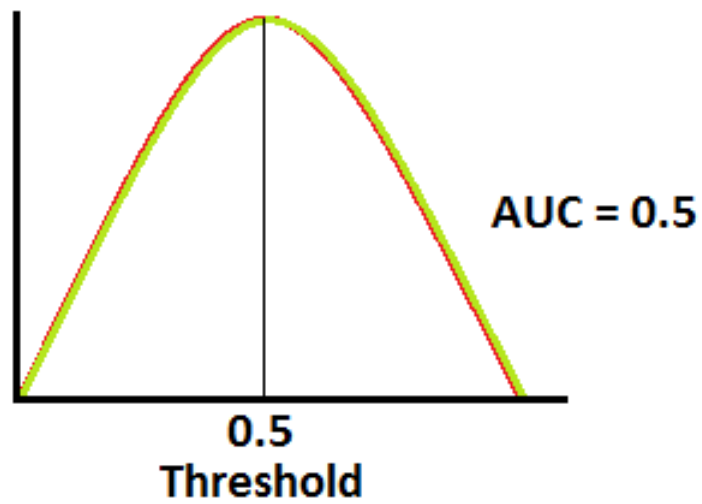


Figure 7: Probabilities

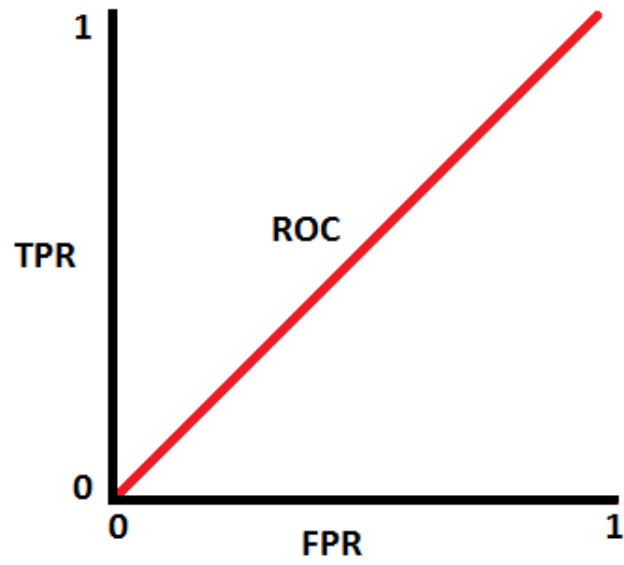


Figure 8: Roc Curve

2.6 Inverse scenario

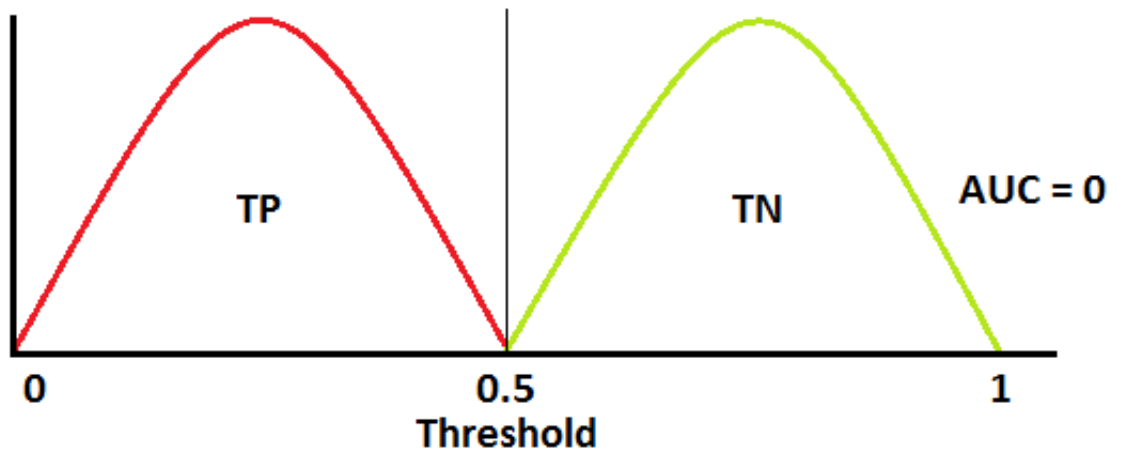


Figure 9: Probabilities

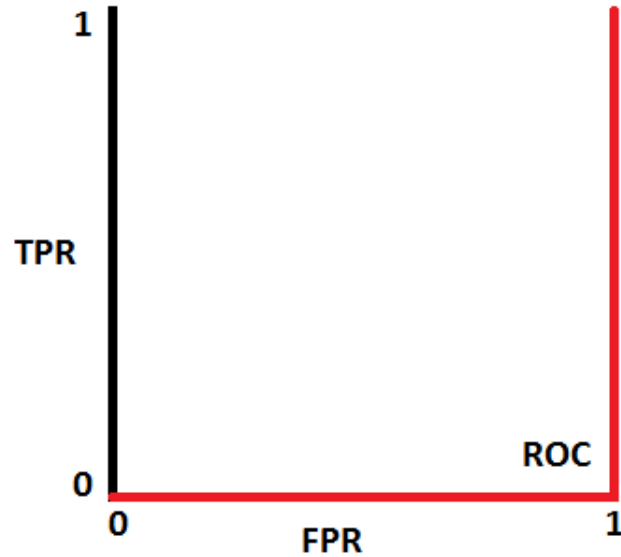


Figure 10: Roc Curve

In short, after presenting the possible scenarios, it is evident that the three types of models are random. That is, they predict the classes with the greatest possible chance of error.

3 Quest o 2

3.1 Item a and b

In this second part of the project, new explanatory variables were created for the training of models. The basic idea for these new variables is the reliability of the domains, type of device or browser used. That is, eight new variants were created that assume the value between 0 and 1 according to the unique values found. For example, the column "P_emaildomain" has several types of domains, and each one of them has a certain number of people who use it. Thus, the new reliability column takes into account the number of people using that domain with the total number of transactions performed, that is, if 800 people use the "yahoo" domain in a total of 1000 transactions, this domain acquires a reliability of 0.8 .

$$\left\{ \begin{array}{ll} \frac{N_Domains_{x_1}}{Total\ of\ Transactions} & \text{if } N_Domains_{x_1} > Reliability\ Limit * Total\ of\ Transactions \\ 0 & \text{if } N_Domains_{x_1} < Reliability\ Limit * Total\ of\ Transactions \end{array} \right.$$

In our case, we consider that there is a reliability limit (between 0 and 1), where domains with a very low number of people are obligatorily zero.

With this in mind, below we have a list of the variables created by this method along with the corresponding code.

- P_emaildomain, id_30, id_31, id_33, id_34, DeviceInfo, card6, R_emaildomain

```

1 import numpy as np
2 import pandas as pd
3
4
5 def domain_reliability(df, columns, p_trustfull_above=0.01):
6     '''creates a column with the domain reliability of each
7     sample
8     Args:
9         df, DataFrame
10        columns, list of columns to analyse realiability
11        p_trustfull_above, percentage of trusted domains
12    Return:
13        domain_values, dict with the numbers of samples in each
14        domain
15    '''
16    columns = [column for column in columns if column in df.
17               columns]
18
19    for column in columns:
20        domain_values = df[column].value_counts()
21        n_all_domains = domain_values.sum()
22        trustfull_domains = domain_values[domain_values >
23                                          p_trustfull_above*
24                                          n_all_domains].to_dict()
25
26        df[str(column)+'_reliability'] = [trustfull_domains[
27                                          domain]/n_all_domains
28                                          if domain in
29                                          trustfull_domains else 0 for domain in df[column]]
30
31    df.drop(columns=columns, inplace=True)

```

```

26
27     return domain_values

```

Listing 1: Part of data_treatment.py

4 Question Three

For this study, 7 types of ML models were used, below we can find which they are and their respective hyperparameter search spaces, which took into account the AUC(ROC) metric. This search was performed by the Bayesian optimization found in the `optimize(BayesSearchCV)` scikit.

- QuadraticDiscriminantAnalysis
- BernoulliNB
- PassiveAggressiveClassifier
- RidgeClassifier
- Perceptron
- SGDClassifier
- MLPClassifier

```

1 from sklearn.naive_bayes import BernoulliNB
2 from sklearn.linear_model import RidgeClassifier
3 from sklearn.linear_model import Perceptron
4 from sklearn.linear_model import SGDClassifier
5 from sklearn.discriminant_analysis import
   QuadraticDiscriminantAnalysis
6 from sklearn.linear_model import PassiveAggressiveClassifier
7 from sklearn.neural_network import MLPClassifier
8
9 from skopt.space import Real, Categorical, Integer
10
11
12 mlp_clf = {
13     'model': Categorical([MLPClassifier()]),
14     'model__activation': Categorical(['identity', 'logistic',
   'tanh', 'relu']),
15     'model__solver': Categorical(['sgd', 'adam']),
16     'model__alpha': Real(0.0000001, 0.001, 'uniform'),
17     'model__learning_rate': Categorical(['constant', '
   invscaling', 'adaptive']),

```

```

18     'model__learning_rate_init': Real(0.0000001, 0.001, 'uniform'),
19     'model__power_t': Real(0.0005, 0.5, 'uniform'),
20     'model__max_iter': Integer(100, 5000, 'uniform'),
21     'model__momentum': Real(0.1, 0.99, 'uniform'),
22     'model__beta_1': Real(0.1, 0.99, 'uniform'),
23     'model__beta_2': Real(0.1, 0.99, 'uniform'),
24     'model__epsilon': Real(0.0000001, 0.001, 'uniform'),
25 }
26
27
28 bernoulli_clf={
29     'model': Categorical([BernoulliNB()]),
30     'model__alpha': Real(0.01, 0.99, 'uniform')
31 }
32
33 QDA_clf={
34     'model': Categorical([QuadraticDiscriminantAnalysis()]),
35     'model__tol': Real(0.0000001, 0.01, 'uniform')
36 }
37
38 PassiveAggressive_clf={
39     'model': Categorical([PassiveAggressiveClassifier()]),
40     'model__max_iter': Integer(1000, 10000, 'uniform'),
41     'model__tol': Real(0.0000001, 0.01, 'uniform'),
42     'model__C': Real(0.01, 0.99, 'uniform'),
43     'model__loss': Categorical(['hinge', 'squared_hinge'])
44 }
45
46
47 ridge_clf_positive = {
48     'model': Categorical([RidgeClassifier(positive=True)]),
49     'model__tol': Real(0.0000001, 0.001, 'uniform'),
50 }
51
52 ridge_clf_false = {
53     'model': Categorical([RidgeClassifier(positive=False)]),
54     'model__solver': Categorical(['svd', 'cholesky', 'sparse_cg', 'sag', 'saga']),
55     'model__tol': Real(0.0000001, 0.001, 'uniform')
56 }
57
58 perceptron_clf = {
59     'model': Categorical([Perceptron(fit_intercept=False)]),
60     'model__penalty': Categorical(['l2', 'l1', 'elasticnet']),
61     'model__alpha': Real(0.00000001, 0.001, 'uniform'),
62     'model__l1_ratio': Real(0.01, 0.99, 'uniform'),
63     'model__max_iter': Integer(1000, 10000, 'uniform'),

```

```

64     'model__tol': Real(0.0000001, 0.001, 'uniform'),
65 }
66
67
68 sgd_clf = {
69     'model': Categorical([SGDClassifier(fit_intercept=False)
70 ]),
71     'model__loss': Categorical(['hinge', 'log', '
modified_huber', 'squared_hinge', 'perceptron', '
squared_error', 'huber', 'epsilon_insensitive', '
squared_epsilon_insensitive']),
72     'model__alpha': Real(0.00000001, 0.001, 'uniform'),
73     'model__max_iter': Integer(1000, 10000, 'uniform'),
74     'model__epsilon': Real(0.0000001, 0.001, 'uniform'),
75     'model__power_t': Real(0.01, 0.99, 'uniform'),
76     'model__eta0': Real(0.01, 0.99, 'uniform'),
77     'model__warm_start': Categorical([True, False]),
78     'model__tol': Real(0.0000001, 0.001, 'uniform'),
79     'model__penalty': Categorical(['l2', 'l1', 'elasticnet'])
80 ,
81     'model__l1_ratio': Real(0.01, 0.99, 'uniform'),
82     'model__learning_rate': Categorical(['constant', 'optimal
', 'invscaling', 'adaptive']),
83 }

```

Listing 2: hyperparameters for each model

4.1 Feature Engineering

Before training the models, it was necessary to perform data engineering to create new variables and perform normalizations. First, we removed columns with more than 90% the Nan values and columns that had a number of unique values close to the number of samples. So below, I have the code that performs these parts.

```

1 def drop_Nan_columns(df, drop_limit=0.7):
2     ''' drop columns with my Nan values
3     Args:
4         df, DataFrame
5         drop_limit, percentage limit of Nan values in a column
6     Return:
7         Nan_dropped, list of features dropped
8     '''
9     Nan_dropped = [feature for feature in df.columns if df[
feature]
10                     [df[feature].isna() == True].shape[0] >
drop_limit*df.shape[0]]
11     df.drop(columns=Nan_dropped, inplace=True)

```

```

12
13     return Nan_dropped
14
15
16 def unique_upperBound_columns(df, p_upper_bound=0.85,
17 drop_nunique=False):
18     '''get columns with unique values above upper bound
19     Args:
20         df, DataFrame
21         upper_bound, percentage of samples above upper bound
22         drop_nunique, drop the nunique_features above upper
23         bound if True
24     Return:
25         nunique_features, dict of samples of above upper bound
26     '''
27     nunique_features = {feature: df[feature].nunique() for
28 feature in df.columns if df[feature].nunique() >
29 p_upper_bound*df.shape[0]}
30     if drop_nunique:
31         df.drop(columns=nunique_features.keys(), inplace=True)
32
33     return nunique_features

```

Listing 3: feature engineering

4.2 Undersample and Bootstrap

Also, looking at the dataset, you can see that the classes are unbalanced (about 500,000 non-cheats to approximately 20,000 cheats). Undeniably, this imbalance creates an absurd bias that was taken away with an undersampling of the data, with a bootstrap afterwards to have a greater generality of the reduced class.

```

1 import numpy as np
2 import pandas as pd
3
4
5 def undersample_bootstrap(inputs: pd.DataFrame, targets: pd.
6 DataFrame, bootstrap_size=1000):
7     '''undersample a DataFrame with features(inputs) and
8     targets
9     Args:
10         inputs, DataFrame with the features
11         targets, DataFrame with the targets
12     Return:
13         undersampled_data,
14         undersampled_targets,

```

```

13     '''
14     min_sample = min(targets.value_counts().tolist())
15
16     undersampled_data = pd.DataFrame(columns=inputs.columns)
17     undersampled_targets = pd.DataFrame()
18
19     # undersample
20     for class_type in targets.value_counts().index:
21
22         indices_class = np.where(targets == class_type)[0]
23
24         indices_attr_sample = np.random.choice(
25             a=indices_class, size=min_sample, replace=False)
26         undersampled_data = undersampled_data.append(
27             inputs.iloc[indices_attr_sample])
28         undersampled_targets = undersampled_targets.append(
29             targets.iloc[indices_attr_sample].tolist())
30
31     # bootstrap
32     for class_type in targets.value_counts().index:
33
34         indices_class = np.where(targets == class_type)[0]
35
36         indices_attr_sample = np.random.choice(
37             a=indices_class, size=bootstrap_size, replace=
38 True)
39         undersampled_data = undersampled_data.append(
40             inputs.iloc[indices_attr_sample])
41         undersampled_targets = undersampled_targets.append(
42             targets.iloc[indices_attr_sample].tolist())
43
44     return undersampled_data, undersampled_targets

```

Listing 4: undersample and bootstrap

4.3 Preprocessing

Finally, these data went to the processing part in order to replace the Nan values and modify categorical variables with OneHotEncoder. These procedures were done using numeric and categorical pipelines, as shown in the code below.

```

1 from sklearn.pipeline import Pipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.compose import ColumnTransformer
4 from sklearn.impute import SimpleImputer
5 from sklearn.preprocessing import OneHotEncoder,
   OrdinalEncoder

```

```

6
7
8 def preproc_normalize(X_train=None, X_test=None, y_train=None
, y_test=None, scaler=None, scaler_trigger=False,
X_test_trigger=True):
9     '''normilize the data with MinMaxScaler
10     Args:
11         X_train, X_test, y_train, y_test
12     Return:
13         X_train, X_test, y_train, y_test
14     '''
15
16     if scaler_trigger == False:
17         scaler = StandardScaler()
18         scaler.fit(X_train)
19
20     X_train = scaler.transform(X_train)
21     if X_test_trigger == True:
22         X_test = scaler.transform(X_test)
23
24     return X_train, X_test, y_train, y_test, scaler
25
26
27 def preprocess(X_train=None, X_test=None, y_train=None,
y_test=None, categorical_features=None, numerical_features
=None, normalize_fn=None, scaler_fn=None, scaler_trigger=
False, X_test_trigger=True):
28     '''replace Nan values of categorical and numerical
features. Moreover, transform categorical features in
numerical data
29     Args:
30         X_train, X_test, y_train, y_test
31         categorical_features, list with the name of the
categorical columns
32         numerical_features, list with the name of the numerical
columns
33     Return:
34         X_train, X_test, y_train, y_test
35     '''
36
37     numerical_pipeline = Pipeline(steps=[
38         ('imputer', SimpleImputer(strategy='mean'))])
39
40     categorical_pipeline = Pipeline(steps=[
41         ('imputer', SimpleImputer(strategy='most_frequent')),
42         ('onehot', OneHotEncoder())])
43
44     transformation = ColumnTransformer(
45         transformers=[

```

```

46         ('numerical transformation', numerical_pipeline,
numerical_features),
47         ('categorical transformation',
48          categorical_pipeline, categorical_features),
49     ])
50
51     X_train = transformation.fit_transform(X_train)
52     if X_test_trigger == True:
53         X_test = transformation.transform(X_test)
54
55     if scaler_trigger == False and X_test_trigger == False:
56         X_train, X_test, y_train, y_test, scaler =
normalize_fn(
57             X_train=X_train, y_train=y_train, X_test_trigger=
False)
58
59     elif scaler_trigger == True and X_test_trigger == False:
60         X_train, X_test, y_train, y_test, scaler =
normalize_fn(
61             X_train=X_train, y_train=y_train, X_test_trigger=
False, scaler=scaler_fn, scaler_trigger=True)
62
63     elif scaler_trigger == False and X_test_trigger == True:
64         X_train, X_test, y_train, y_test, scaler =
normalize_fn(
65             X_train, X_test, y_train, y_test)
66     else:
67         X_train, X_test, y_train, y_test, scaler =
normalize_fn(
68             X_train=X_train, X_test=X_test, y_train=y_train,
y_test=y_test, scaler=scaler_fn, X_test_trigger=False,
scaler_trigger=True)
69
70     return X_train, X_test, y_train, y_test, scaler

```

Listing 5: replace Nan

5 Question Four

The purpose of this section is to discuss the results obtained by the Bayesian optimization of several models and hyperparameters, in order to choose the one with the best Score for the AUC metric. With this in mind, in the image below we can see a convergence plot for the models trained in each period of study.

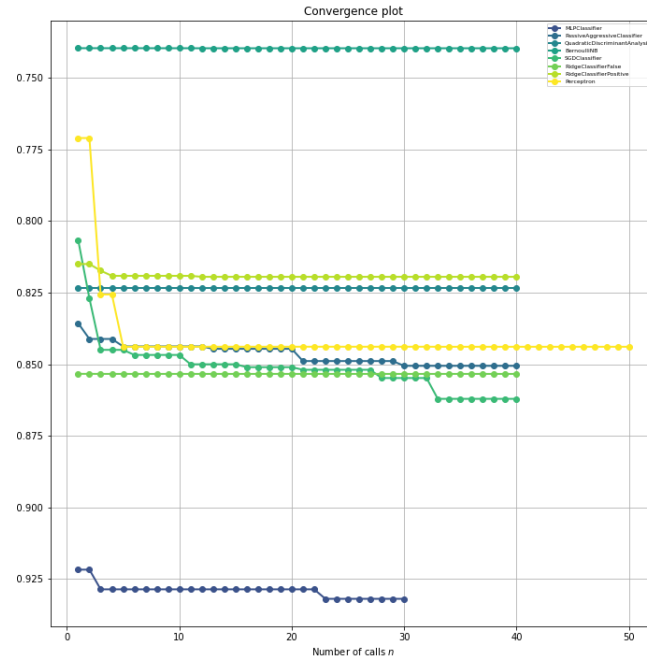


Figure 11: Convergence of models

5.1 Hyperparameter dependency

At first, we can see the model MLPClassifier has the best score for the metric. However, to have a deeper analysis we need to see how each of the hyperparameters of each model depend on each other. As a result, below we have the dependency graphs for each of the optimized models.

5.1.1 Bernoulli Naive Bayes

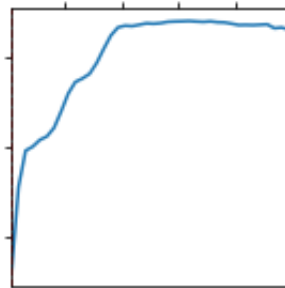


Figure 12: Dependence Plot

5.1.2 Passive Agressive Classifier

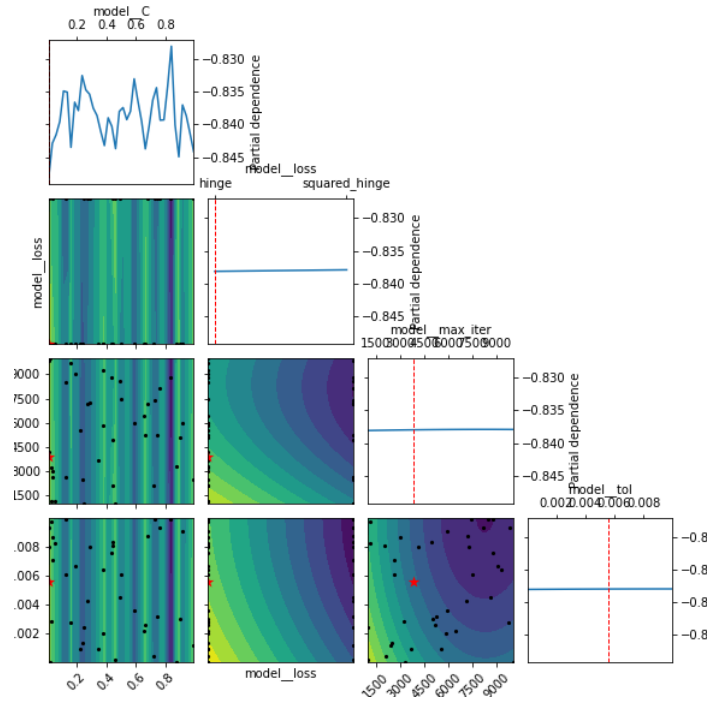


Figure 13: Dependence Plot

5.1.3 Ridge Classifier Negative

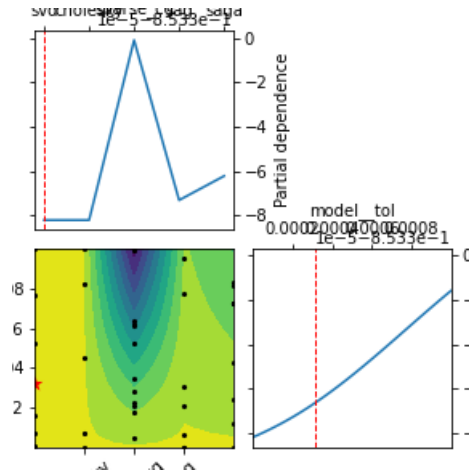


Figure 14: Dependence Plot

5.1.4 Ridge Classifier Positive

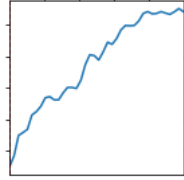


Figure 15: Dependence Plot

5.1.5 Perceptron

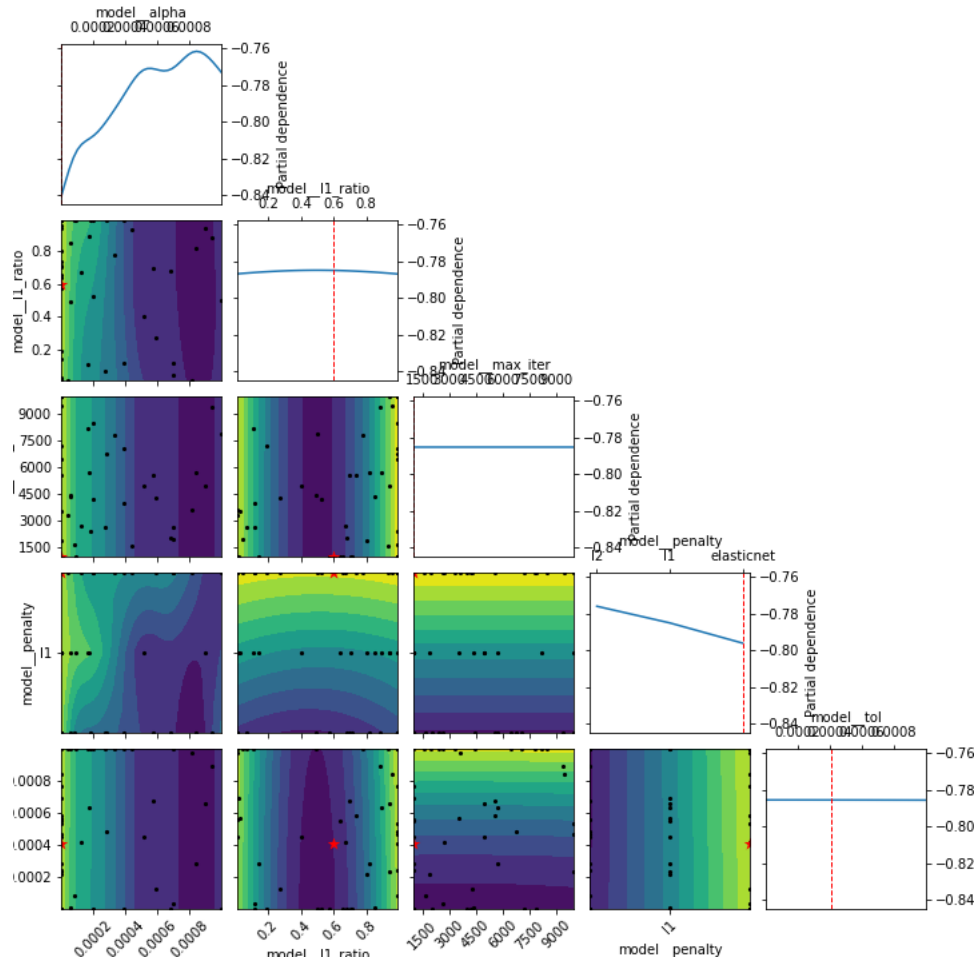


Figure 16: Dependence Plot

5.1.6 Stochastic Gradient Descent Classifier

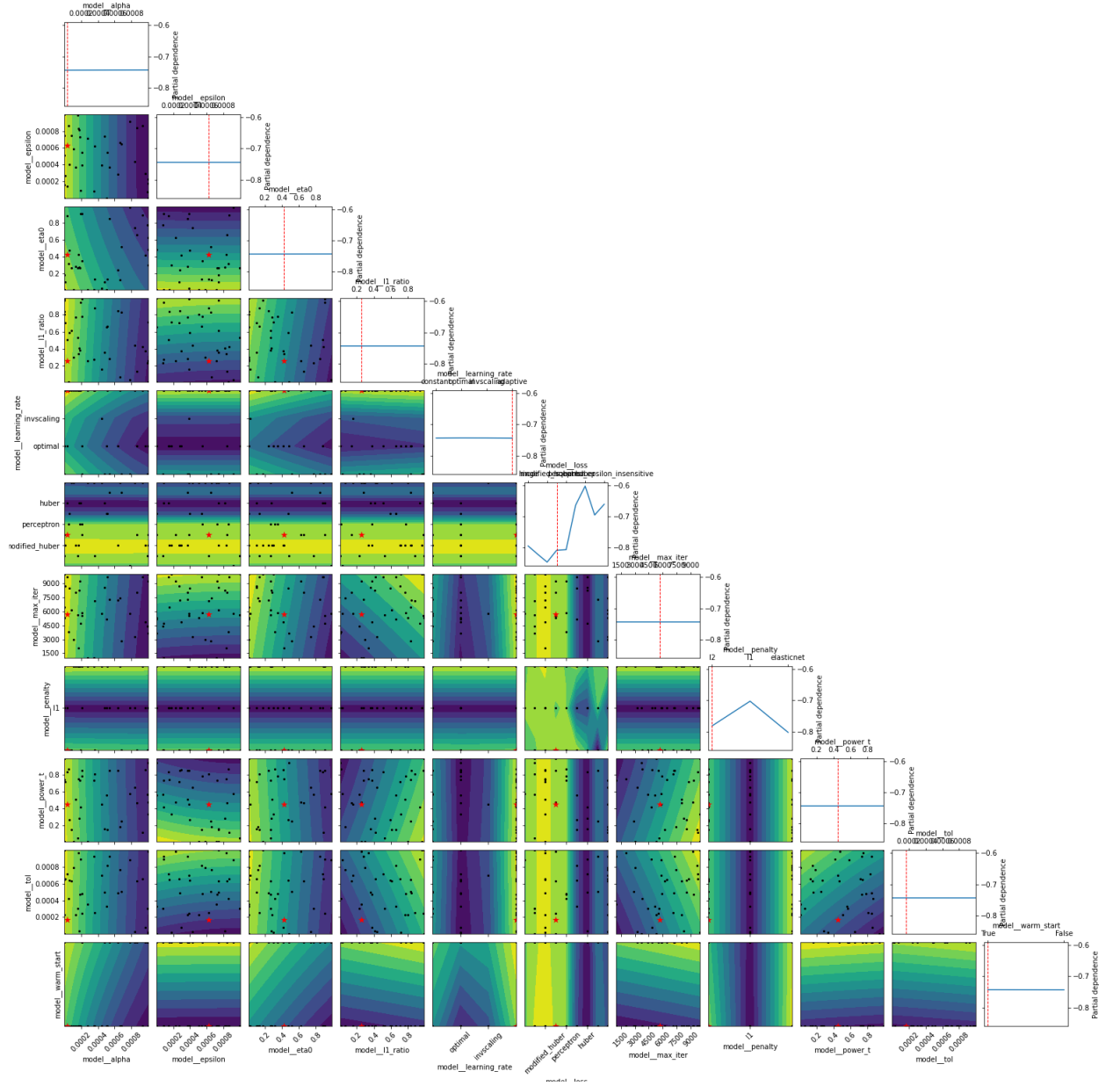


Figure 17: Dependence Plot

5.1.7 Multi-Layer Perceptron Classifier

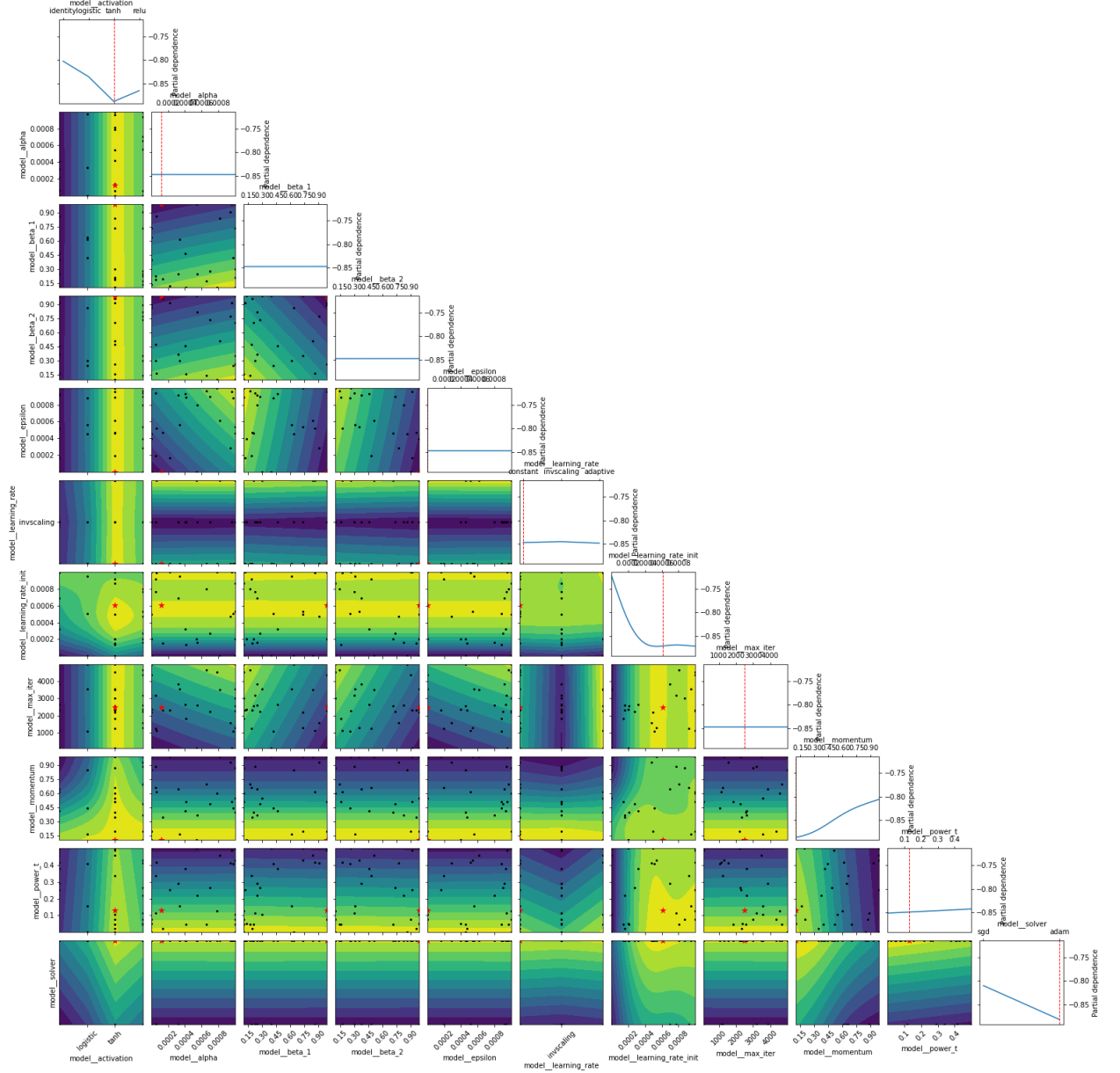


Figure 18: Dependence Plot

5.1.8 Quadratic Discriminant Analysis

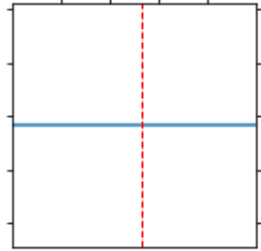


Figure 19: Dependence Plot

5.2 Best Model

In short, the best model of this project achieved a good ROC curve for the parameters of TPR, FPR, TNR, FNR. To reiterate this inference, we have below the images generated for the ROC curve and confusion matrix of the best model, respectively.

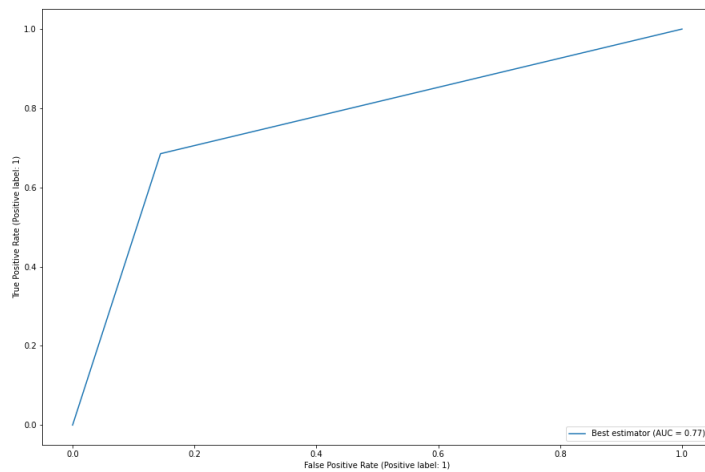


Figure 20: Roc curve

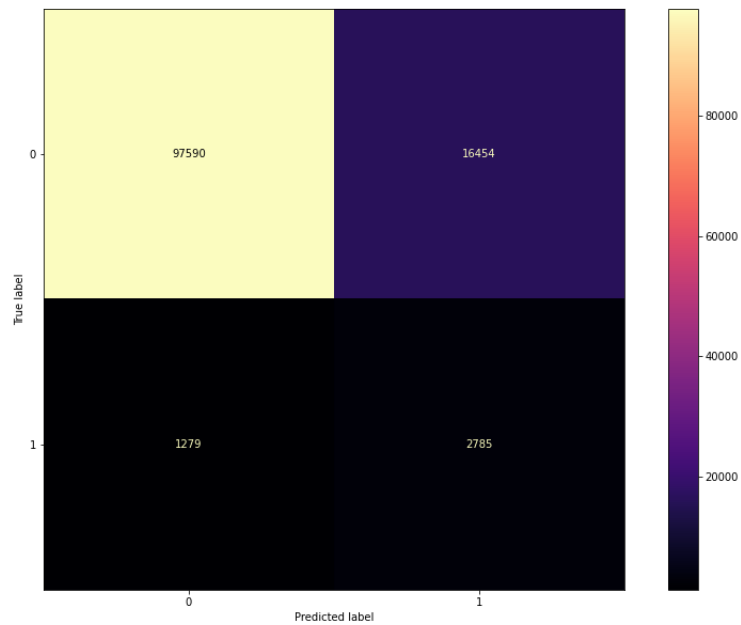


Figure 21: Confusion Matrix

5.3 My Kaggle Submission

kaggle_submission (1).csv 3 days ago by evertonmendes73 add submission details	0.723634	0.771538	<input type="checkbox"/>
--	----------	----------	--------------------------

Figure 22: kaggle submission

6 Question Five

In summary, having the best model, we can obtain the probabilities of each transaction being a fraud. In sklearn, we have the `predict_proba` function that gives us the necessary information to use with decision trees and binomial distributions to find a cutoff point from which transactions will be

considered as fraud and barred/blocked. Taking into account the false positives that interfere with the customer's convenience, making him reflect on going to a competitor in the market.

In this way, we can build a decision tree with the confusion matrix of the best model, obtained in the previous section, and the probability of each transaction being a fraud (through the histogram of the training file, "train_transaction.csv").

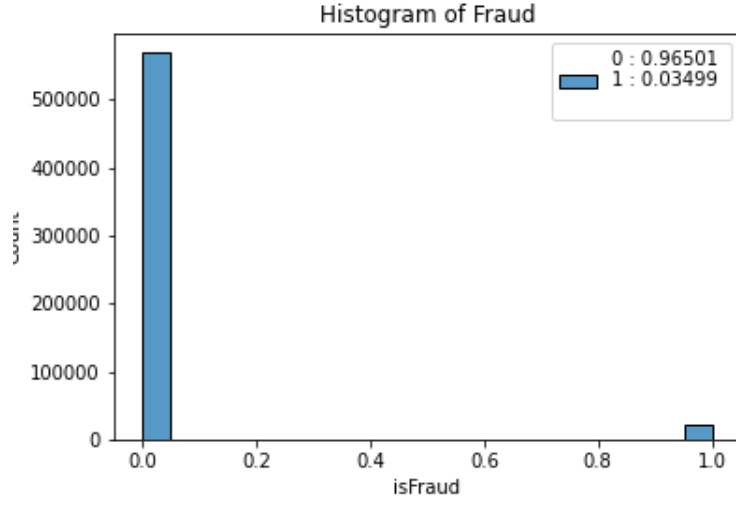
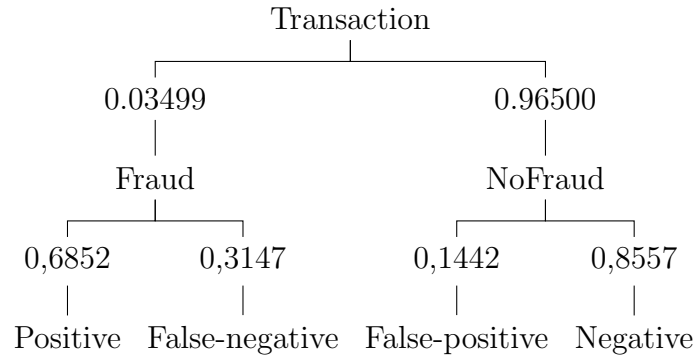


Figure 22: Histogram of "train_transaction.csv"



$$\text{Fraud} \rightarrow \begin{cases} \text{Positive} = 0.02397 \\ \text{False - negative} = 0.01101 \end{cases}$$

$$\text{NoFraud} \rightarrow \begin{cases} \text{False - Positive} = 0.13915 \\ \text{Negative} = 0.82575 \end{cases}$$

Taking the above probability paths into consideration, it is possible to find the probability of correctly hitting the cheat when I predict a cheat, and the probability of missing a cheat when I predict NoFraud.

$$\text{Right_Frauds} = \frac{0.02397}{0.13915 + 0.02397} = 14,6947\% \quad (1)$$

$$Missing_Frauds = \frac{0.01101}{0.01101 + 0.82575} = 1,3157\% \quad (2)$$

With this in mind, let's consider binomial distributions and the concept of value at risk to obtain a cutoff point for the customer's card. The distributions assume a 30-day card history and the respective probability found above, Right_Frauds and Missing_Frauds.

$$P(n, k, p) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} \quad (3)$$

To show how the cutting bridge will be calculated, we will present some variables of the problem.

- number of transactions = n_t
- probability of Right_Frauds = $p_r = 14,6947\%$
- probability of Missing_Frauds = $p_m = 1,3157\%$
- number of Frauds = $k_r[]_i$
- number of Missing Frauds = $k_m[]_i$
- value of each transaction = $t_v[]_i$
- predict_proba Fraud = $pb_r[]_i$
- prdict_proba No Fraud = $pb_m[]_i$

First, we need to initialize the Fraud and Non-Fraud vectors, $k_r[]_i$ and $k_m[]_i$ respectively. The initialization of these vectors depends solely on the output of the predict_proba function. As an example, let's consider only the $k_r[]_i$ algorithm, however it is also valid for the $k_m[]_i$ considering the necessary changes. The initialization follows the next step, we have a K_{aux} that starts with zero and at each interaction of predict_proba the helper's value is incremented if the Fraud probability is greater than not Fraud, inserting the helper's value in position i of the vector $k_r[]_i$.

Finally, we can build the Risk equations, as shown below:

$$R_m = \sum_{i=0}^{n_t} t_{vi} P(n_t, k_{mi}, p_m) pb_{mi} \quad (4)$$

$$R_r = \sum_{i=0}^{n_t} t_{vi} P(n_t, k_{ri}, p_r) pb_{ri} \quad (5)$$

$$if \frac{R_m + R_f}{\sum_i t_{vi}} > 5\% \quad \text{denied card}$$

7 Reference

- [1] Links with the images of Graph and Adjancecy List
- [2] da Silva, Éverton Luís Mendes. Codes from this project
- [3] Images to explain AUC-ROC metric
- [4] da Silva, Éverton Luis Mendes. Codes from this project in Drive