

# Summary

Best Gini: 0.517997

Stand Dev: 0.06114332752194878  
Total Submissions: 182

Steps	Explanation	What worked	What didn't work	Wishlist
EDA	<b>Quick EDA revealed:</b> <ul style="list-style-type: none"><li>- small train set (864 observations); very skewed classes: only 10% fraud cases, 90% non-fraud (legitimate) transactions</li><li>- large number of variables (no feature name or info provided); features: binary, categorical and float</li><li>- missing values: only few variables with missing values, no duplicate rows</li><li>- if_var_68 and if_var_69 most likely transaction amounts (given value distribution); test-dataset contained an additional variable (contract date)</li><li>- distribution: some variables skewed or log distribution</li></ul>			n/a
Preprocessing/ Feature Transformation	Minimal Preprocessing was necessary given invariance to scale of Random Forest as well as relative cleanliness of data (only few missing values)	<b>1. Replacing 0's:</b> <ul style="list-style-type: none"><li>a. binary variables: with either 0 or 1, depending on column mean (1 if &gt;0.5, else 0)</li><li>b. categorical variables: mode</li><li>c. float variables: mean</li></ul> <b>2. Removing var with no information content</b> (only 1 unique value): "ib_var_12"	<b>1. Removing Outliers</b> – this was expected, since Fraud cases can often represent itself as outliers and by removing those we would not train our model on those frauds <b>2. Scaling Features</b> - explicitly tried for skewed variables if_var_68 and if_var_69, however this did not impact results. With Random Forest as final model, scaling was not necessary	<b>1. Binning Categorical Variables</b> - although no variable names were given this could have been performed, e.g. using scikit learns "KBinsDiscretize" <b>2. Creating Dummies with Categorical Variables</b> - relevant for distance-based algorithms (e.g. KNN) to ensure that a value of 4 (Category4) is not read as 4x as big as the value of 1 (category1)
Sampling	<b>Hypothesis:</b> Improve the classifier by providing a balanced dataset.  Sampling performed at two steps in the process:  1. Sampling before calculating correlation, IV and PSI: this was performed in order to understand correlations and information value of all variables to then conduct feature selection. For feature selection, key features were selected from the unsampled X_train and X_test and entered training unsampled (see, 2. Sampling during model training).  2. Sampling during model training: To prevent data leakage, over-/undersampling should be performed during cross-validation, not before. Therefore, sampling was conducted through a pipeline during training and cross-validation	<b>1. Oversampling &amp; Undersampling vs. Only Oversampling:</b> mixed performance of conducting both Oversampling (using SMOTE) and Undersampling (using RandomUnderSampling), before hypertuning the model performed better with only oversampling (to a ratio of 1:1), after hyperparameter tuning using both Oversampling & Undersampling, each with 0.5 sampling strategy (resulting in 2:1 ratio of legitimate to fraud cases) performed best (also tested vs. other sampling ratios) <b>2. Oversampling with SMOTENC:</b> although improvement was very minor, a slightly better performance was achieved when using synthetic oversampling technique for nominal and categorical features (treats categorical features differently) <b>3. Using Sampling during Cross-validation:</b> following common literature, to prevent data leakage although no actual performance improvement was observed vs. sampling before model training	<b>1. SMOTE alternatives:</b> BorderlineSMOTE, KNNSMOTE and ADASYN SMOTE did not provide better performance (Gini between 0.49-0.5x; although tested only during a very small sample of trials and not systematically due to time constraints) <b>2. class_weight="balanced"</b> of scikit-learns RandomForestClassifier as alternative to imbalanced-learn's SMOTE, this did however significantly decrease performance	<b>1. Trying SAMPLING alternatives using smote-variants, which provides a larger number of different smote methods,</b> e.g. weighting fraud cases differently depending on some criteria, however smote-variants was incompatible with my conda set-up. <b>2. Research methods to augment the entire dataset</b> (both fraud and non-fraud cases) since so little training data was available, possibly leading to the severe overfitting
Feature Selection:	<b>Hypothesis:</b> improving the classifier by selecting only those features which would best classify fraud and non-fraud cases). However, given RandomForest built in Feature Selection, no significant improvement expected.  Features were assessed using correlation (Spearman), IV as well as PSI for a sampled (=balanced) dataset. Based on this, key features were selected.	<b>1. Filter-based method: feature selection using correlation (Spearman)</b> with target variable - very quick'n'dirty method <b>2. Keeping majority of the features:</b> Classifier performed best when dropping only a very limited number of features (15) and keeping the majority (66 features)	<b>1. Dropping a large share of features</b> -> Gini dropped to values below 0.45, indicates that most features contained relevant information and/or that Random Forest did a better job in Feature selection than the filter-based method (using correlation)	<b>1. PCA</b> - although this would only work for the numerical (non-binary, non-categorical) features (16 in total), so this would need to be applied alongside other feature selection methods <b>2. GAN</b> - Genetic Algorithm <b>3. Further Filter-Methods:</b> such as Chi- <sup>2</sup> Test <b>4. Wrapper-based methods:</b> recursive feature elimination
Model Choice:	<b>Hypothesis:</b> Ensemble methods provide superior prediction power for this usecase since they combine multiple 'individual' (diverse) models together for a prediction, thus are more robust and better in fighting the "local minima problem"  To obtain a baseline, several models were tried (without any hyperparameters) to compete against each other (Ensemble and non-Ensemble Methods): <ul style="list-style-type: none"><li>- Logistic Regression</li><li>- Random Forest</li><li>- Decision Tree</li><li>- ExtraTrees Classifier,</li></ul> out of which <b>Random Forest performed best</b> . This was also confirmed in class, so Random Forest was the main model chosen for the classification.	<b>1. Ensemble Methods -&gt; Random Forest:</b> ensemble of decision trees, trained via bagging method and thus averaging predictions (fighting the local minima problem). Additionally, it introduces more randomness through searching for the best feature among a random subset of features vs among all features, resulting in more bias and more robustness.	<b>1. ExtraTreesClassifier</b> -> did not improve scores <b>2. Stacking:</b> Hard Voting implemented (combinations of Random Forest and ExtraTreesClassifier as well as LogRegression with different configurations tested), however due to time constraints not well-developed and optimized, thus it performed similar/not better than only Random Forest (Gini: 0.48-0.5x)	<b>1. Improve/Tune Voting Classifier</b> <b>2. Testing BaggingClassifier</b> (tbd if it can also do hard voting or only soft voting) <b>3. Testing Stacking</b> (using e.g. DESlib or own method), but tbd due to small training set <b>4. Testing further classification models:</b> <ul style="list-style-type: none"><li>- KNN</li><li>- Naive Bayes</li><li>- Neural Nets</li><li>- Boosting Methods (XGBoost, GradientBoost, AdaBoost)</li></ul> Not tested due to time constraints and due to general opinion/information shared, that Random Forest (or other bagging methods) would perform best, however it could have been interesting to test if, for example, the Voting algorithm worked better when using bagging methods (for robustness and randomness) with boosting methods, despite boosting methods tending to overfit (and lacking randomness)
Hyperparameter Fine-Tuning	<b>Hypothesis:</b> Improving Model performance (precision) with the correct hyperparameters  Different combinations of hyperparameters for Random Forest were extensively tested using general understanding of the hyperparamters on model behaviour plus RandomSearch first, then GridSearch.	Best results were obtained with: <ul style="list-style-type: none"><li>1. N_Estimators of 1,100 trees -&gt; averaging results over 1,100 trees</li><li>2. Min_samples_split: 3</li><li>3. Max_depth: 15</li><li>4. Max_features: log2</li><li>5. Random-Seed: alternating the random-seed as the very last step on the tuned model provided randomly better (and worse :) ) scores</li><li>-&gt; all other parameters were left on default values</li><li>6. RandomSearch and Gridsearch</li></ul>	<b>1. Bootstrapping:</b> best results were obtained with bootstrapping=False, which means that for each tree all observations were used. This was most likely due to the small training data, with bootstrapping=True only a sample would have been used (via max_sample), resulting in even less training data at a time.	<b>1. max_features:</b> try float instead of default auto, sqrt, log2
Cross Validation	<b>Hypothesis:</b> Using cross-validation to average training results and preventing overfitting  Using Cross-validation improved model-performance; it was used in a pipeline together with the above mentioned sampling strategy so the model would be trained on a balanced training set.	<b>1. Cross-validation:</b> improved model performance, using the following hyperparameters:  cv = StratifiedKfold(n_splits=10,shuffle=True)	The following cross-validation hyperparameters/methods didn't impact model results (no change in Gini!): <b>1. Testing different number of splits</b> during cross-validation didn't impact model results <b>2. Testing StratifiedFold vs. RepeatedStratifiedKfold</b> didn't impact model results	
General/ Evaluation	<b>Conclusion:</b> The model completely overfit during training which could be seen in the evaluation results: Gini-Index=1.0 on the test set vs. 0.5-ish as actual result and an ROC_AUC score of 0.5-0.6 (indicating that the classifier is just slightly better than random guessing). Interestingly, the model still returned best results when the hyperparameters were tuned to low regularization values that would have overfit most models: very high max_depth (15), very little min sample_split (1), etc.). Most likely this was due to the model having only very little training data to learn from.			<b>1. Data Augmentation</b> <b>2. Implementing a cost-sensitive cost-function</b> <b>3. Other Libraries</b> than Scikit Learn (e.g. fastai, weka) or ML Tools, e.g. Dataiku