# Using Gradient Boosting Tree algorithms on Porto-Seguro's safe drive competition

Binh Nguyen

Université Paris–Saclay & UEVE

January 16, 2018

# Overview

# Introduction

- Hosted on Kaggle, by Porto Seguro - a car insurance company.
- Objective: build a model that predicts the probability that a driver will initiate an auto insurance claim in the next year.
- Done on Python & Google Cloud Platform

# Problem settings

- Supervised learning
- Binary classification.
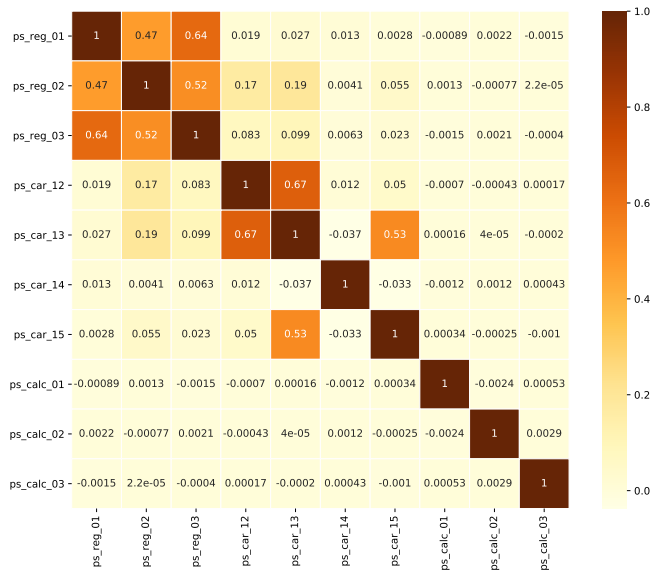- Performance metrics: Normalized Gini Index.

# Data Analysis

## General info

- `train.csv` & `test.csv`
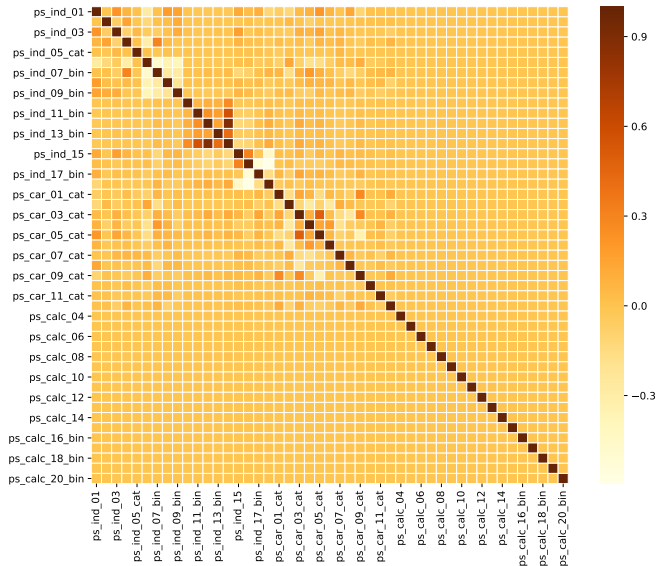- For train: $n = 595212$; $p = 59$. For test: $n = 892816$

## Interesting findings

- Number of $n$ in test is much larger than train.
- Features are encoded under generic names (for example `ps_car_01`, `ps_reg_02`) that give no specific indication on the meaning.
- There are 3 types: continuous features, categorical features and binary features.
- Values of -1 indicate that the feature was missing from the observation, or `NaN` value as normally seen in many datasets.
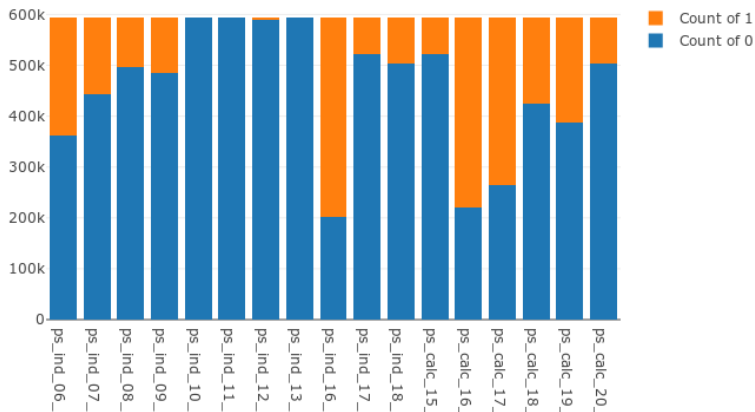
# Correlation matrix of continuous features (`float64`)
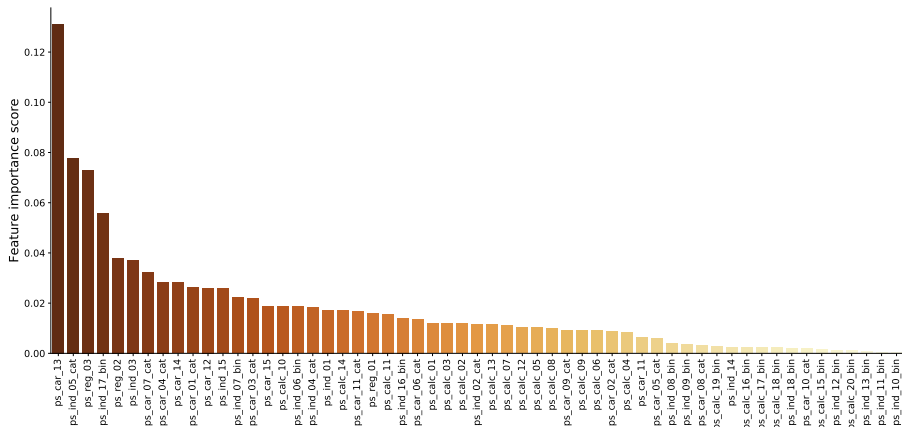
# Correlation matrix of continuous features (`int32`)

# Variable important ranking using Random Forest Classifier

# Feature Engineering

## Eliminating unimportant features

Top 35 variables which have highest feature important score are kept.

## Dealing with missing values

One could possibly either set `NaN` to some extreme value (e.g. $10^6$) or use the value that represented the property of the particular variable ($\max(X_i)$, $\mathrm{median}(X_i)$). In this case, we choose to keep the default value $-1$ to represent `NaN`.

## Likelihood encoding

The central idea behind likelihood (or target) encoding technique is to estimate likelihood target variable given a particular value of a categorical variable:

$$Pr(y = 1 \mid x = \texttt{catValue}) \ = \frac{\sum \mathbb{1}_{y=1 \mid x=\texttt{catValue}}(y)}{n_{\text{observation}}}$$

**Algorithm 1** Target encoding categorical variables

**Input:** `categoricalArray`, `labelArray`, $k$, $l$

1: `data` ← **join** (`categoricalArray`, `labelArray`)
2: `targetEncodedArray` ← Empty Array
3: `categoricalValue` ← unique values of `categoricalArray`
4: **split** `data` into $k$ parts
5: **for each** `fold[i]` in $k$ folds **do**
6:     `valueEncoded` ← 0
7:     `fold[-i]` ← **join** remaining $k - 1$ folds except $i$
8:     **split** `fold[-i]` into $l$ parts
9:     **for each** `fold[j]` in new $l$ parts **do**
10:         **for each** `value` in `categoricalValue` **do**
11:             `valueEncoded` ← `valueEncoded` + $\dfrac{Pr(label \mid value)}{l}$
12:         **end for**
13:     **end for**
14:     `valueEncodedFold[i]` ← **join** all `valueEncoded`
15: **end for**
16: `targetEncodedArray` ← **join** $k$ parts of `valueEncodedFold[i]`
17: **return** `targetEncodedArray`

# Implementation

## Model Selection

- Gradient Boosting Tree: XGBoost[a] & LightGBM[b]
- Extremely fast, can also be used with GPU

---

[a]https://github.com/dmlc/xgboost
[b]https://github.com/Microsoft/LightGBM

## Ensembling technique

- 5-fold cross-validation
- Using `sklearn.model_selection.GridSearchCV` to find best `learning_rate`, `max_depth` parameters
- Prediction data from these 2 models was fitted by a Logistic Classifier acted as the ensembler to produce final prediction
- Google Cloud cluster with 16 CPU cores and 12GB RAM takes over 2 hours

---

**Algorithm 2** Ensemble classifier (XGBoost + LightGBM + Logistic Regression)

---

**Input:** trainData, testData, $k$

1: objective $\leftarrow$ Normalized Gini Index
2: **split** trainData, testData to $k$ folds
3: **for each** fold[i] in $k$ folds **do**
4:     **run** Algorithm 1
5:     **fit** XGBoost on fold[-i] using fold[i] to maximize objective
6:     **fit** lightGBM on fold[-i] using fold[i] to maximize objective
7:     predictXGB $\leftarrow$ XGBoost(testData)
8:     predictLGBM $\leftarrow$ lightGBM(testData)
9: **end for**
10: predictXGB $\leftarrow \dfrac{\text{predictXGB}}{k}$
11: predictLGBM $\leftarrow \dfrac{\text{predictLGBM}}{k}$
12: ensemblePrediction $\leftarrow$ logisticRegression(predictXGB, predictLGBM)
13: **return** ensemblePrediction

# Results

| Model | Public Leaderboard | Final Leaderboard |
|-------|--------------------|--------------------|
| XGBoost only | 0.27477 | 0.27756 |
| XGBoost + target encode | 0.28204 | 0.28598 |
| Ensemble + target encode | **0.28365** | 0.28763 |
| **Best result** | 0.28332 | **0.28849** |

| 1800 | ▼ 960 | Regis | | 0.28764 | 26 | 2mo |
| 1801 | ▼ 257 | AntiLippasaar | | 0.28764 | 23 | 2mo |
| 1802 | ▲ 544 | scut_430 | | 0.28763 | 3 | 3mo |
| 1803 | ▼ 2 | idc man | | 0.28763 | 68 | 3mo |
| 1804 | ▲ 606 | LightR | | 0.28763 | 1 | 3mo |
| 1805 | ▼ 29 | Rafael Ladeira | | 0.28762 | 9 | 3mo |

✔ Your submission description has been saved.

Sort by  Private Score  ▼

**All**   Successful   Selected

| Submission and Description | Private Score | Public Score | Use for Final Score |
|---|---|---|---|
| **submission-lgb-only-impact-encoding-joinbin-dropcalconly-10fold-17...** <br> 3 months ago by BTNG <br><br> 10 fold cv | 0.28849 | 0.28332 | ☐ |
| **submission-171022-004132.csv** <br> 3 months ago by BTNG <br><br> Upload Submission File submission-171022-004132 stack of 3 <br> including gradient boosting tree with max_depth=5 all other default. <br> No feature elimination. | 0.28838 | 0.28225 | ☐ |
| **submission-lgb-only-impact-encoding-5fold-notps13square-171025-1...** <br> 3 months ago by BTNG <br><br> submission-lgb-only-impact-encoding-5fold-notps13square- <br> 171025-193323.csv not create ps13square, not drop ps 13 | 0.28820 | 0.28321 | ☐ |
| **submission-lgb-only-impact-encoding-joinbin-dropcalconly-171024-1...** <br> 3 months ago by BTNG <br><br> submission-lgb-only-impact-encoding-joinbin-dropcalconly- <br> 171024-171759.csv - create new impact join bin 06-09, maxdepth 6 on <br> lgb learning rate 0.01 | 0.28806 | 0.28287 | ☐ |
| **stacking1-171019-201739.csv** <br> 3 months ago by BTNG <br><br> stacking1-171019-201739. | 0.28806 | 0.28214 | ☐ |

# The End.