



# Porto Seguro's Safe Driver Prediction

Zexi Yuan, Jianqiao Liu, Wenlin Ou



# Introduction

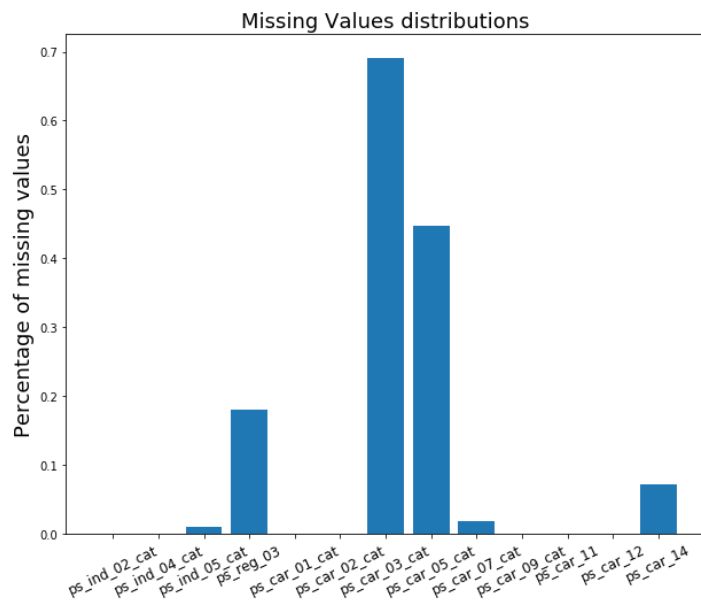
- Problem background
  - Dataset Description
- Approaches
  - Algorithms implemented
  - Evaluation Criteria
- Experiment Results
  - Data Preprocessing
  - Tuning Parameters
- Conclusions
- Q & A



## Part I: Problem Background

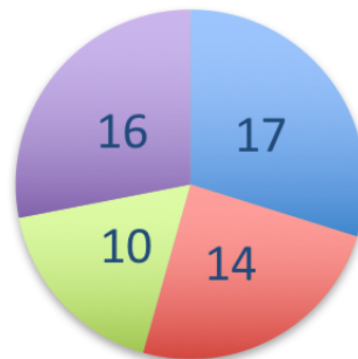
Porto Seguro is one of Brazil's largest auto and homeowner insurance companies, hence setting reasonable insurance prices is beneficial for both the company and customers. We are given a dataset that contains drivers' driving history and whether they have filed a claim in the past year. Our goal is to predict the probability that a driver will file a claim during the next year.

# Dataset



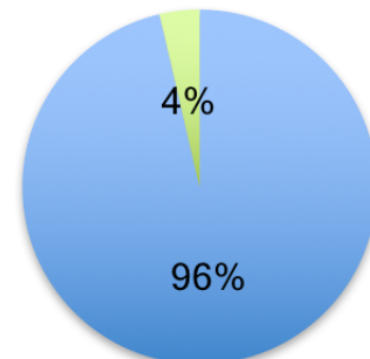
## Dataset:

57 Features



Binary Categorical  
Continuous Ordinal

595212 Instances



No Claim Claim



## Part II: Approaches

- Logistic Regression
- Random Forest
- XGB



## Model Evaluation: AUC & F-Score

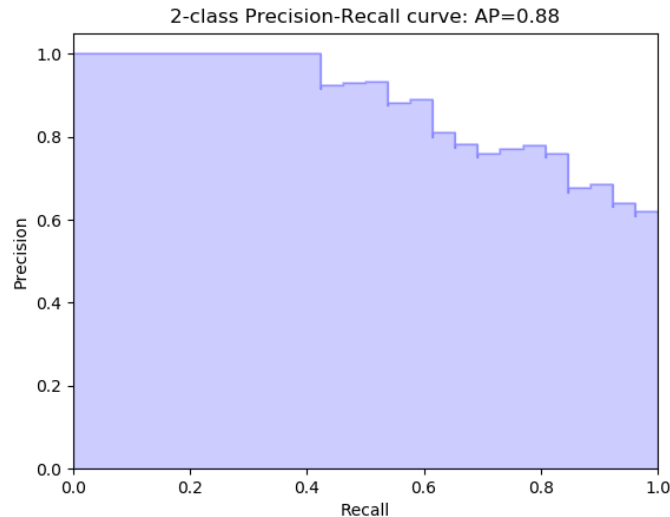
- AUC: area under ROC curve
  - ROC: TPR-FPR curve
- F-score

$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}.$$

- Recall (TPR) =  $TP / (TP + FN)$
- Precision =  $TP / (TP + FP)$
- FPR =  $FP / (FP + TN)$

# Average Precision Score

- The precision-recall curve shows the tradeoff between precision and recall for different threshold.
- AP score is the area under the Precision-Recall Curve



$$AP = \sum_{k=1}^N p(k) \Delta r(k)$$



## Part III: Results

- Data Preprocessing
- Experiment on Algorithms



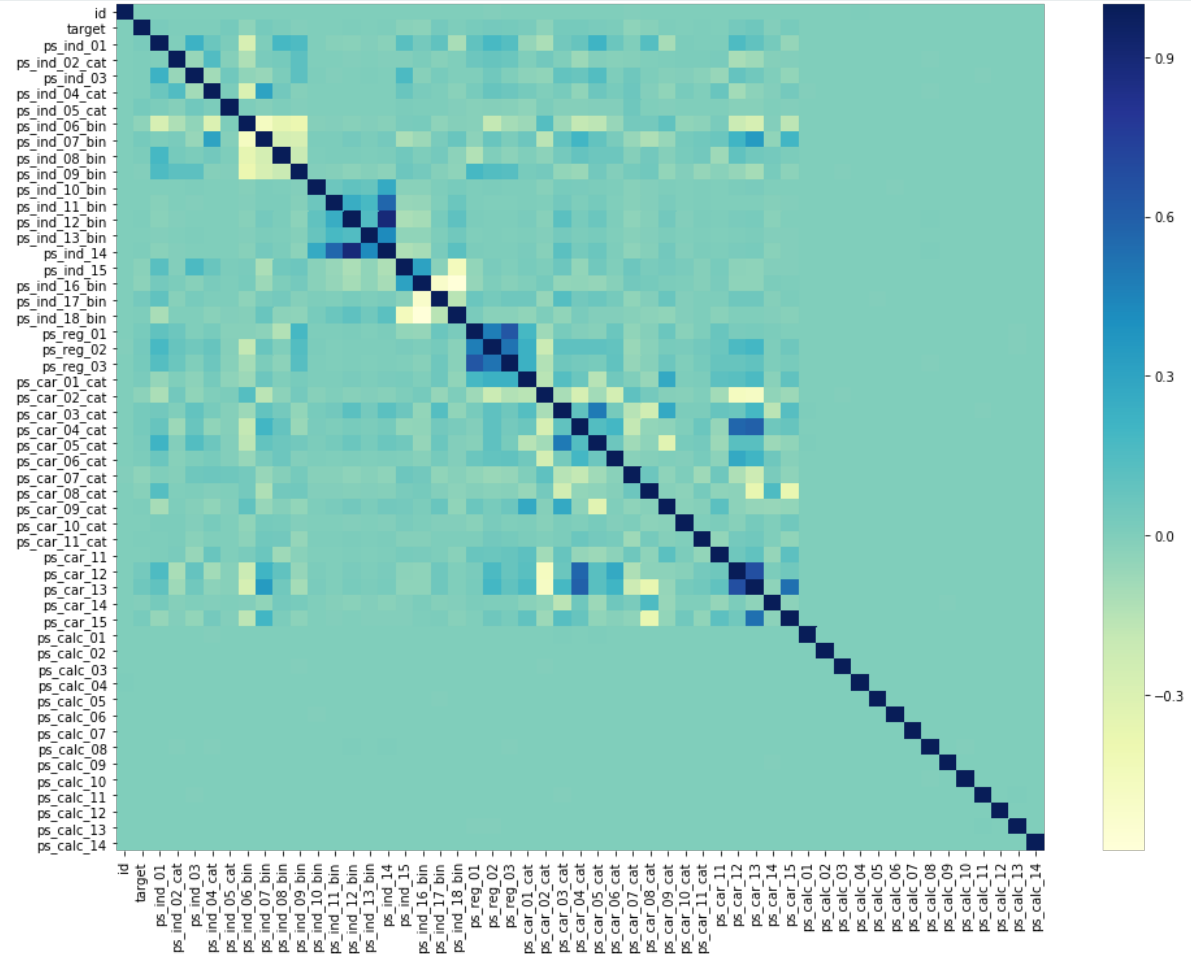


# Data Preprocessing

- Feature Selection
- One hot encoding
- Missing Values

# Pearson Correlation

$$\rho_{X,Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$





## What is “One Hot Encoding”?

CompanyName	Categoricalvalue	Price
VW	1	20000
Acura	2	10011
Honda	3	50000
Honda	3	10000

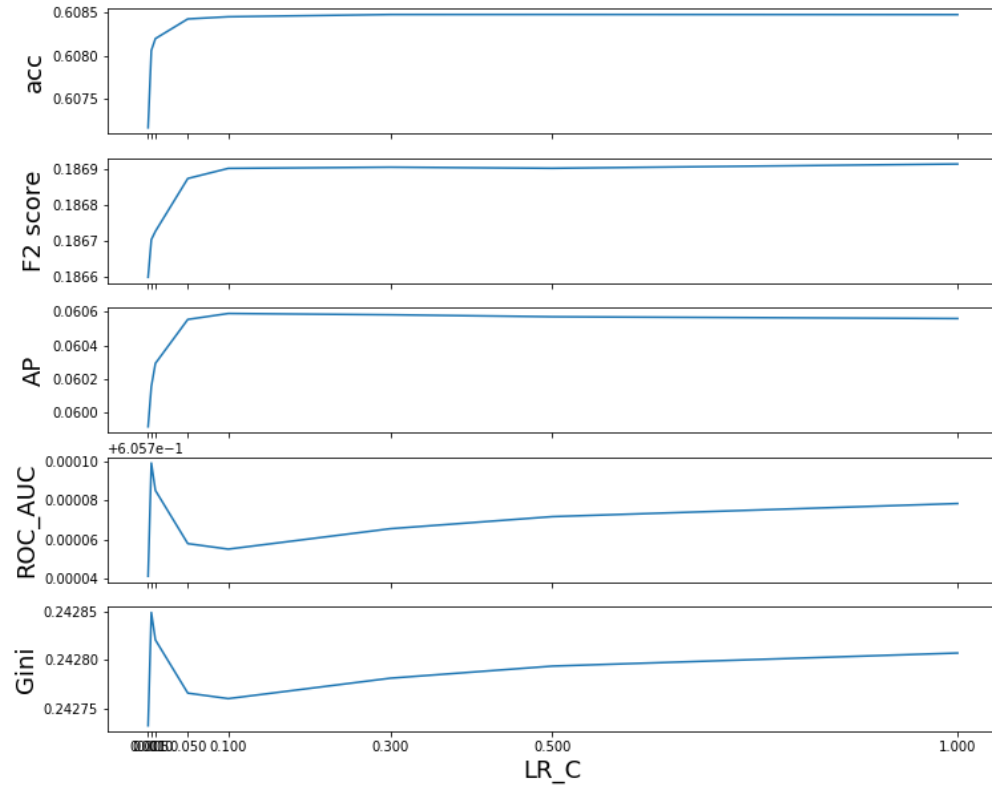
VW	Acura	Honda	Price
1	0	0	20000
0	1	0	10011
0	0	1	50000
0	0	1	10000



# Logistic Regression

- Package: `sklearn.linear_model.LogisticRegression`
- Missing Value
  - Fill the missing values with mean for continuous data
- Data Standardization (0 mean and 1 std)
- Parameter: Inverse of regularization strength  $C = 0.1$

LR\_C

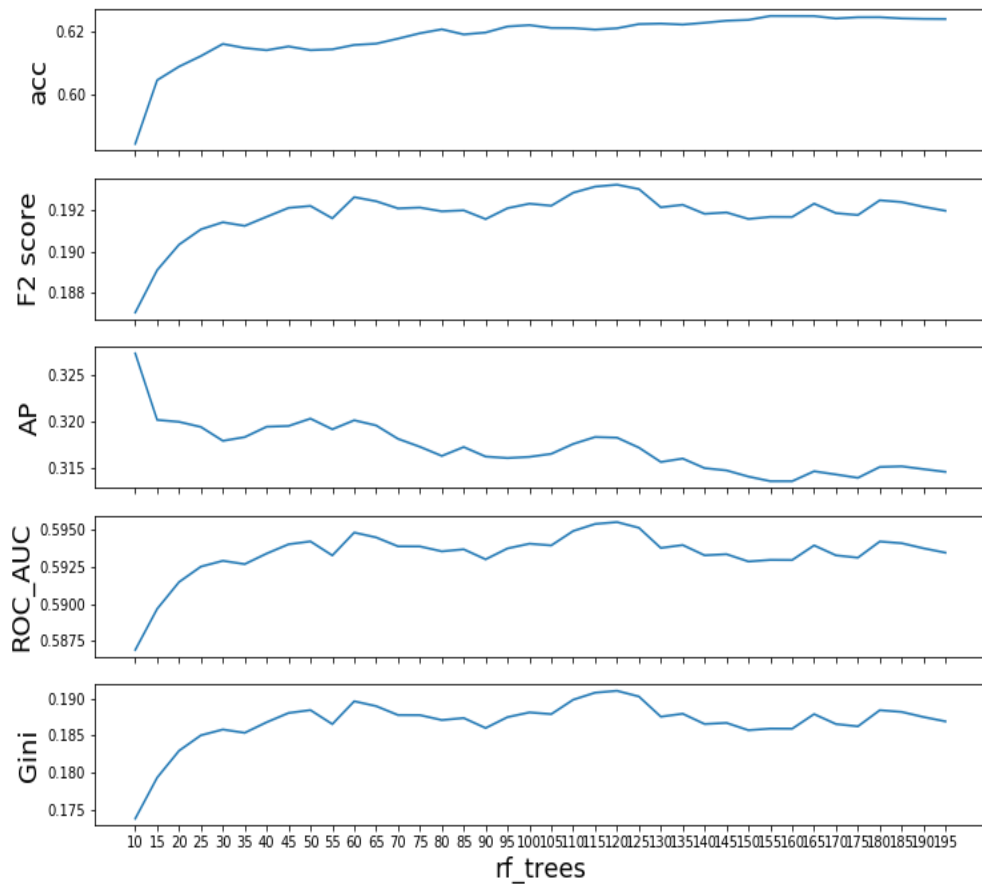




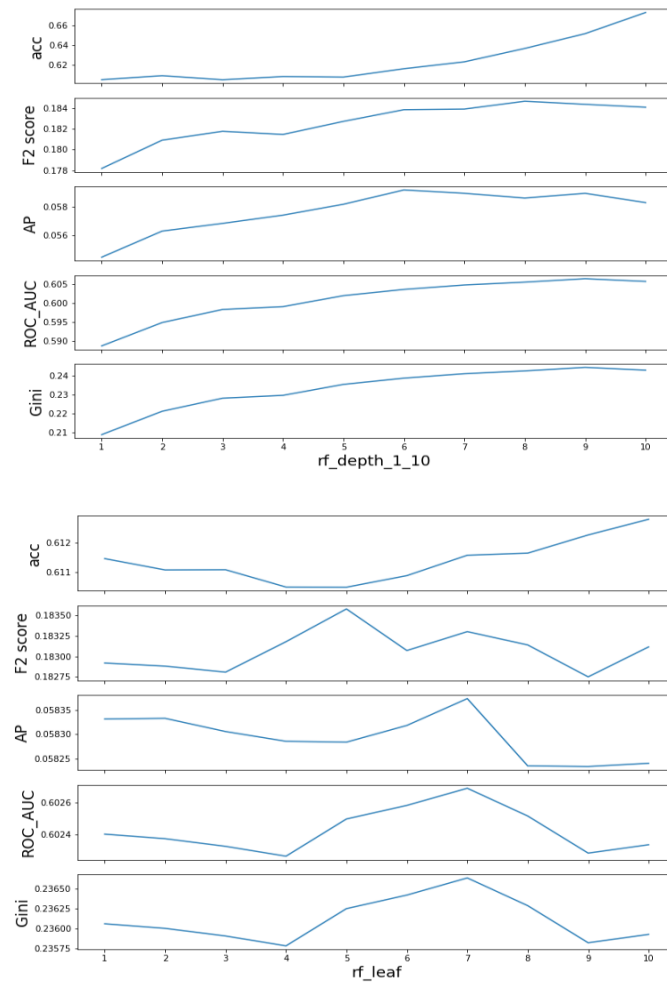
# Random Forest

- Package: `sklearn.ensemble.RandomForestClassifier`
- Optimal Parameter
  - Number of trees: 120
  - Tree depth: 9
  - Leaf number: 7

rf\_trees



rf\_depth\_1\_10





## XGBoosting

$$\text{obj}(\theta) = \sum_i^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

$$\text{obj}^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2}(H_j + \lambda)w_j^2] + \gamma T$$

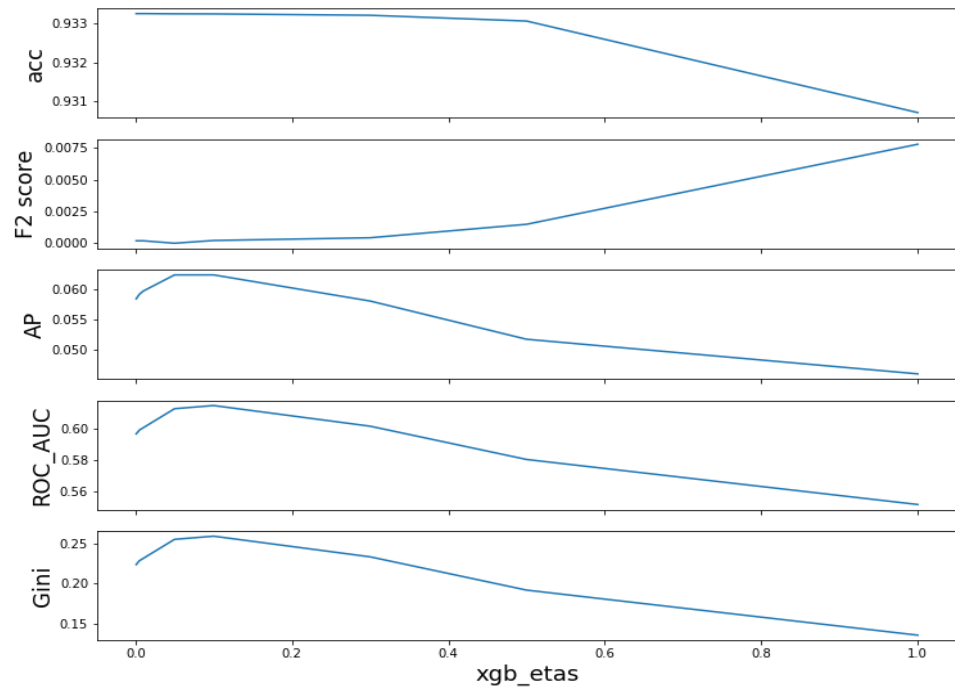
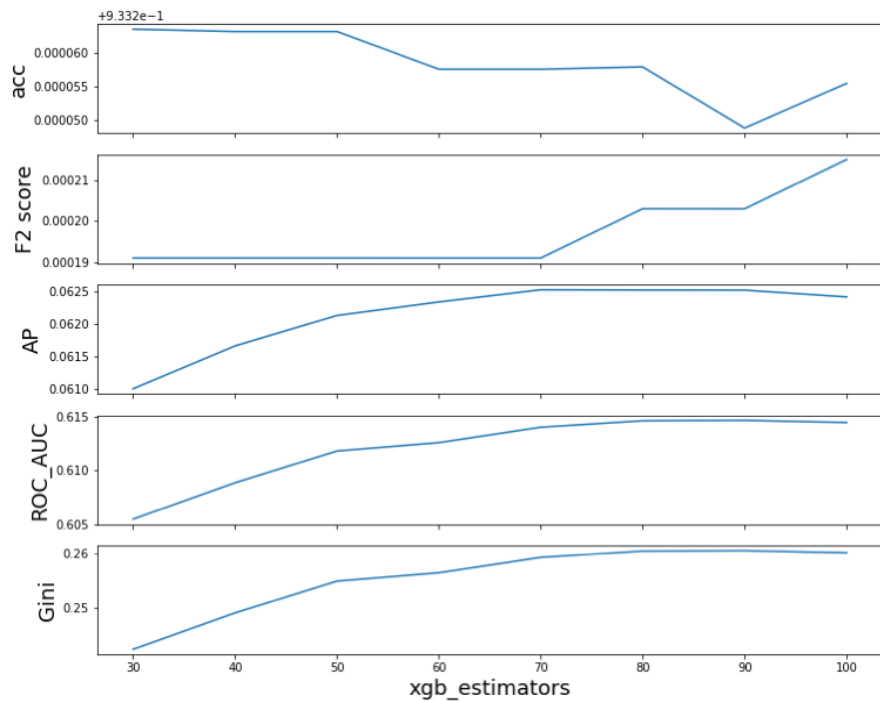
$$\text{Gain} = \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$$

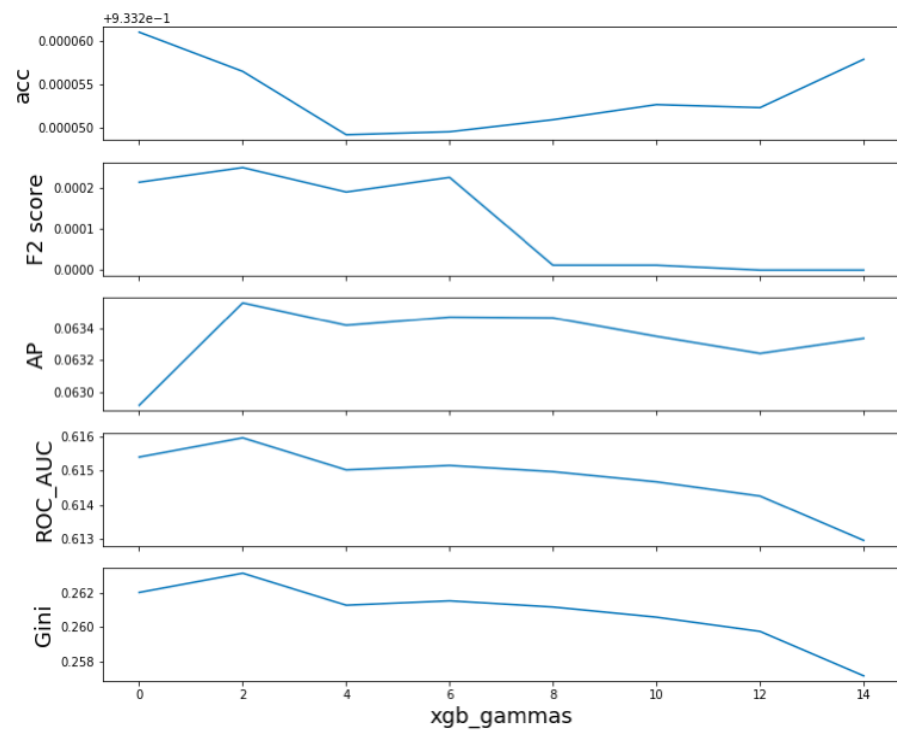
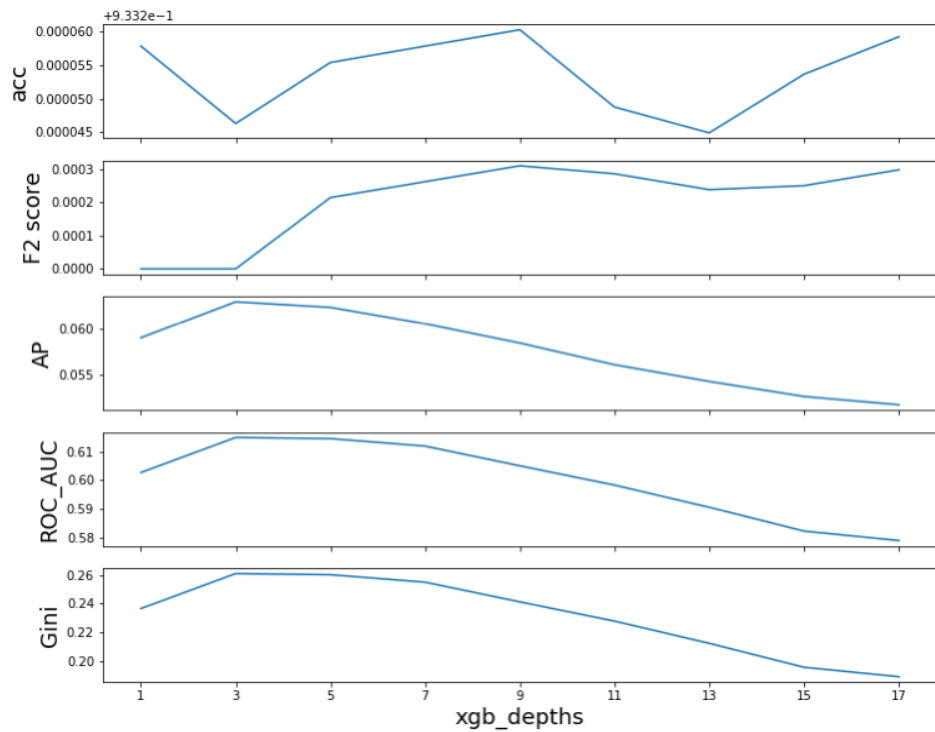




# XGBoosting

- Package from <https://github.com/dmlc/xgboost/>
- Optimal Parameters
  - number of estimators: 90
  - tree depth: 5
  - learning rate: 0.1
  - Gamma: 2

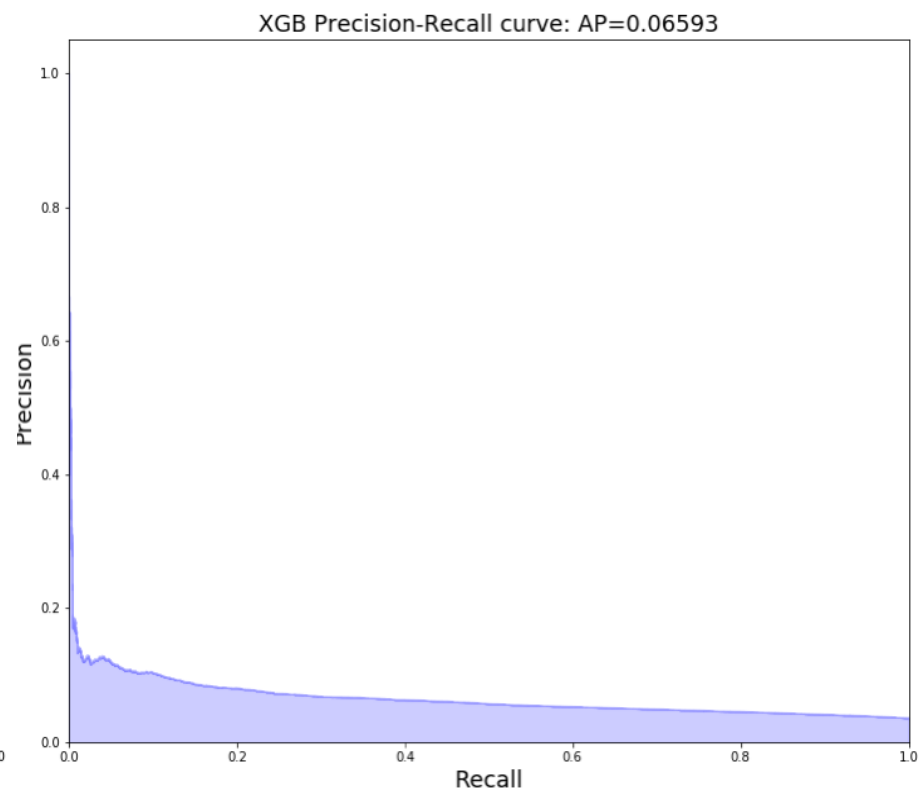
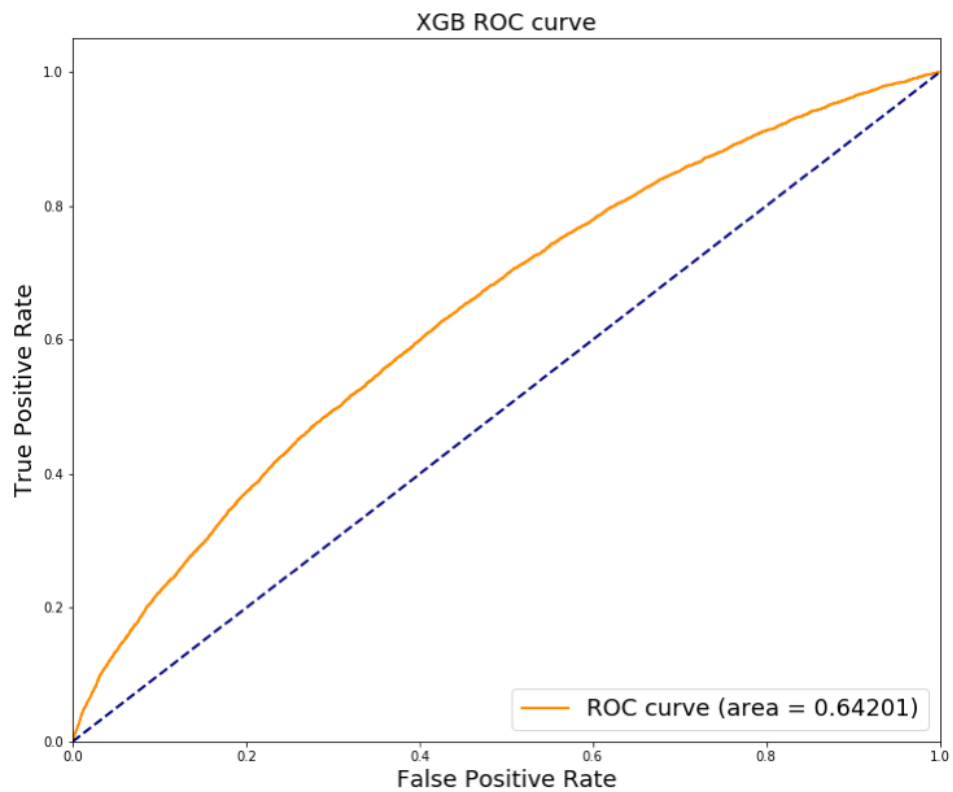






## Comparison

	<b>Accuracy</b>	<b>F2 score</b>	<b>AP</b>	<b>ROC_AUC</b>	<b>Gini</b>
<b>XGB</b>	0.964127	0.000936	0.065808	0.642011	0.284023
<b>RF</b>	0.665444	0.192977	0.061271	0.633579	0.267158
<b>LR</b>	0.627770	0.192531	0.063220	0.632558	0.265116





## Part IV: Conclusions

- Something to share:
  - AP Score as an extra evaluation criterion
- Challenges we encountered:
  - Feature selection
  - No time for grid search



**Q & A:**

Thank You!