

Small Business Loans

Decision Making on Loan Approval

Ed Lee
Classification, METIS, July 2022



U.S. Small Business
Administration

Background

- **SBA** audits and approves loans for small businesses
- Banks lend funds based on SBA approval



Lending partner



Small business



Background

- But some businesses end up failing to pay the loan back

→ “Charge-off”





**Which
businesses are
likely to default?**

The Dataset



Records of loans approved by SBA

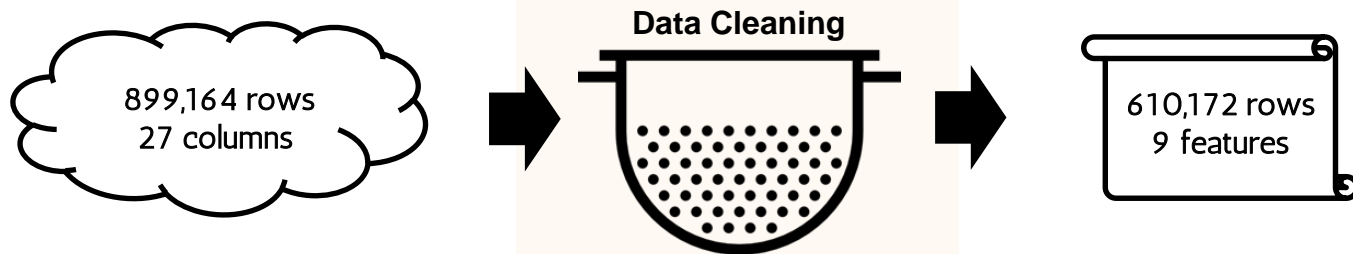
- Raw dataset contains 899,164 entries and 27 columns
- Includes traits of each business, amount approved, **loan status** (paid in full / charged off)

Tools

All the tools I have utilized were run through Python 3.8

Pandas	EDA, data cleaning, data manipulation, feature engineering
Matplotlib, seaborn	Plots, model visualizations
Scikit-learn	Building classification models (LogReg, random forest, naïve Bayes, AdaBoost), preprocessing data, evaluating model metrics
XGBoost	Building classification models (Xgboost classifier)
imblearn	Testing for class imbalance

EDA & Data Clening



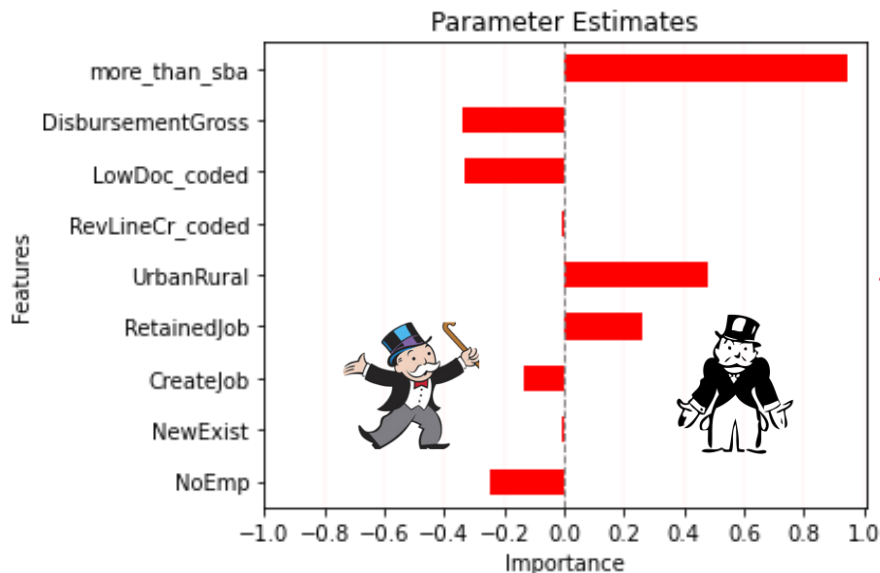
- Raw dataset

- Removing rows with invalid values
- Preprocessing
- Feature selection

- Target: **Loan Status**
- Imbalance: ~20% of all target values were positives (charged-off)

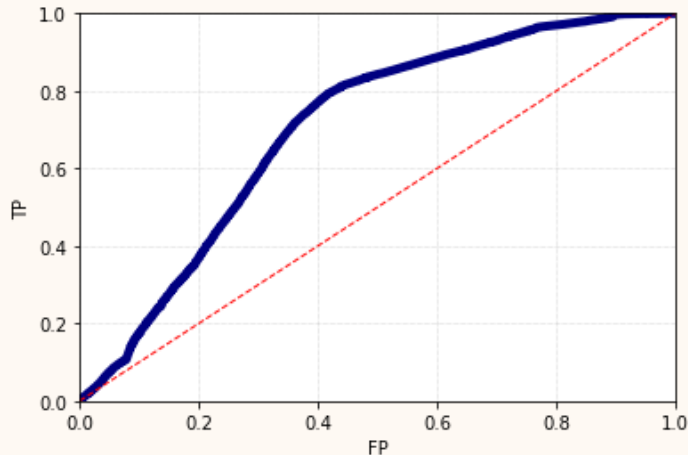
Feature Importance

Before jumping into model comparisons, let's get a rough idea of each feature's impact



Provides clients
insights on which
factors to
scrutinize upon

'Starting Point' - ROC (LogReg)



Accuracy: 0.820

Precision: 0.979

Recall: 0.002

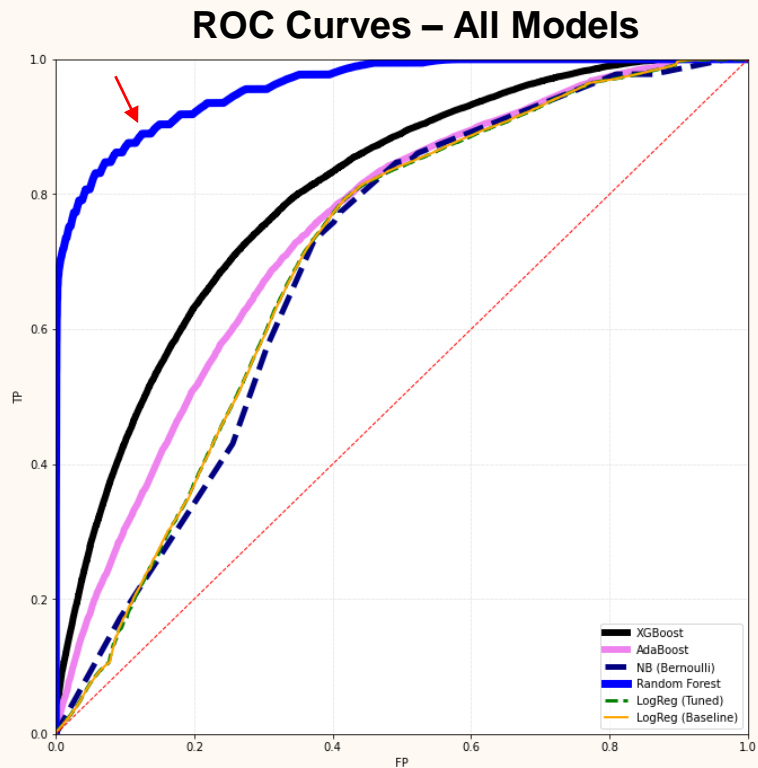
AUC: 0.703

Result - Models

Total six models were built:

1. 'Baseline' Logistic Regression (default parameter)
2. Tuned LogReg (with GridSearchCV)
3. Random Forest (tuned with GridSearchCV)
4. Naïve Bayes
5. AdaBoost
6. XGBoost

❖ Class imbalance was ultimately addressed via class weights



Result – Models

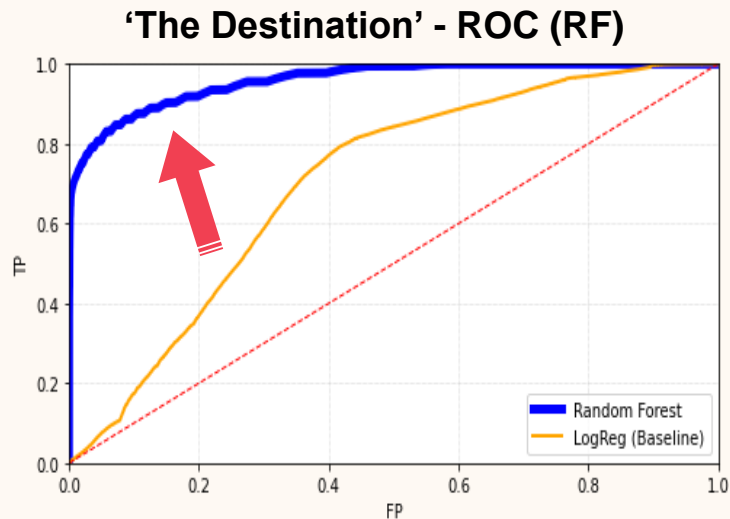
Choosing the best model:

- Model performances were compared via plotting ROC curve and evaluating its area (ROC AUC)
- **Random Forest** model yielded the best model

The Champion:



Random Forest Model



Metrics improvement over the baseline:

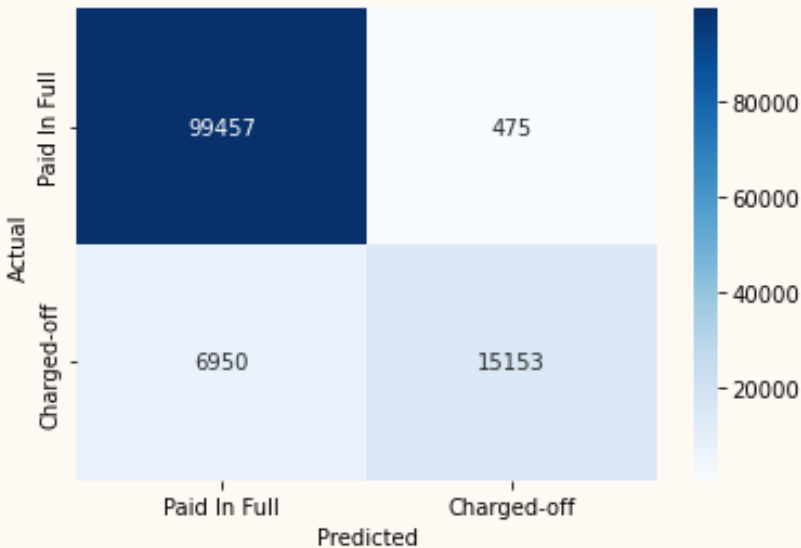
Accuracy: 0.820 → **0.938**

Precision: 0.979 → 0.970

Recall: 0.002 → **0.681**

AUC: 0.703 → **0.961**

Confusion Matrix



The Champion:



Random Forest Model



Metrics improvement over the baseline:

Accuracy: 0.820 → **0.938**

Precision: 0.979 → 0.970

Recall: 0.002 → **0.681**

AUC: 0.703 → **0.961**



<https://www.bamboohr.com/small-business/>

Further Works

Further optimizations and visualizations

- Employing more features from the dataset
ex) Business type, franchise, location...
- Visualizing predicted values, possibly via geocoding with tableau
- More thorough feature engineering
- Further fine-tuning (hyperparameters, etc)



U.S. Small Business
Administration

Thank you!

Any questions?

CREDITS: This presentation template was created by
Slidesgo, including icons by **Flaticon**, infographics & images
by **Freepik**

Please, keep this slide for the attribution