

Fraud Detection Strategy

Mitchell Seiter



Background

The Client:

- Any institution that handles credit card transactions (or any kind of financial transactions) and that would have some liability if fraudulent transactions happened with their service.

The Goal:

- Develop a fraud detection algorithm that ideally catches all cases of fraud without stopping any legitimate transactions from being blocked.

Data

- Data comes from a simulated data generator that mimics real life credit card transactions. For privacy purposes it was not possible to obtain “real “ credit card transaction data.
- Tools and models used were sklearn to generate and test various supervised machine learning models such as KNN, Logistic Regression, Decisions Trees, and Random Forests.
- Primary model we will look at is random forest as this provided the best metrics

Preliminary Look at the Data

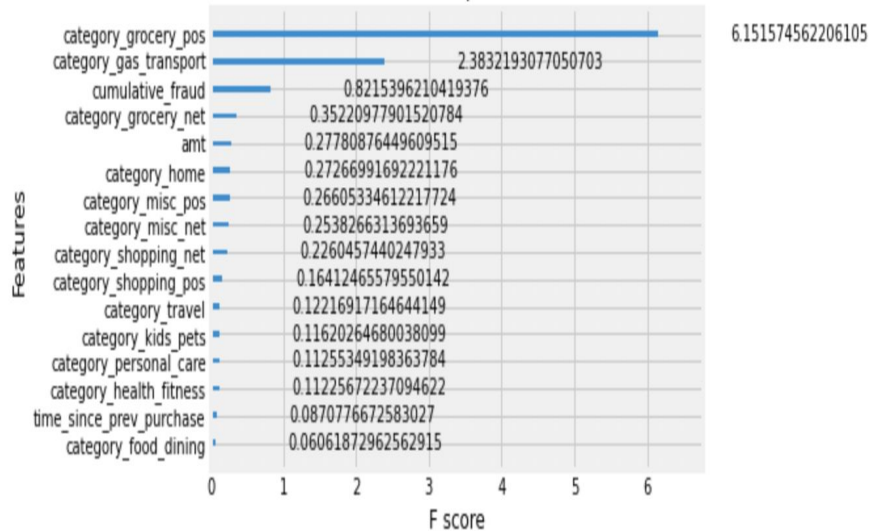
- **555,719** total transaction out of which **2,145** were fraud so **1/3 of a percent** total fraud
- **218 out of 917** people responsible for all cases of fraud
- Fraudulent transactions have a mean value of **\$529 compared to \$69** for all transactions

Some Key Features

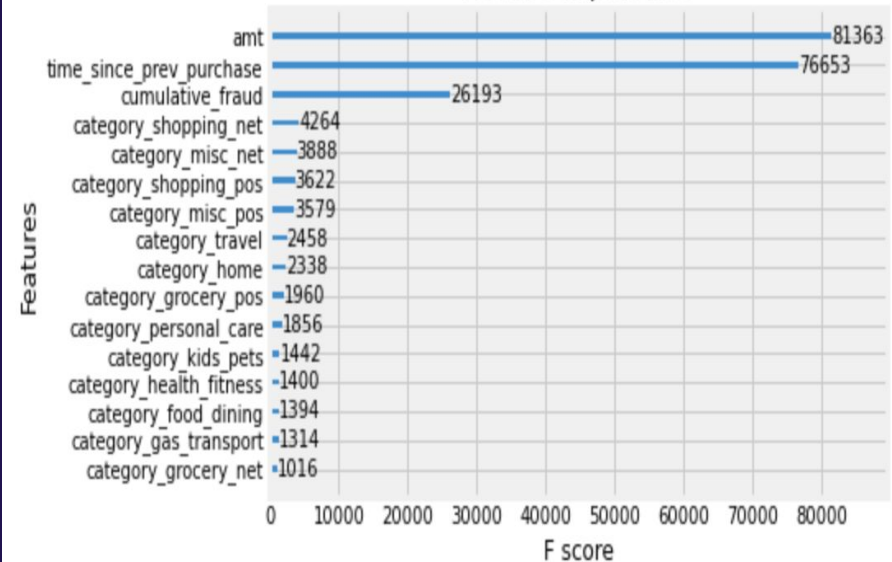
- Amount of transaction
- Category of merchant for the purchase such as: groceries, gas, entertainment, etc.
- Time since last transaction (added feature)
- Does this person have a history of fraud (added feature)

Importance of Features

Feature importance

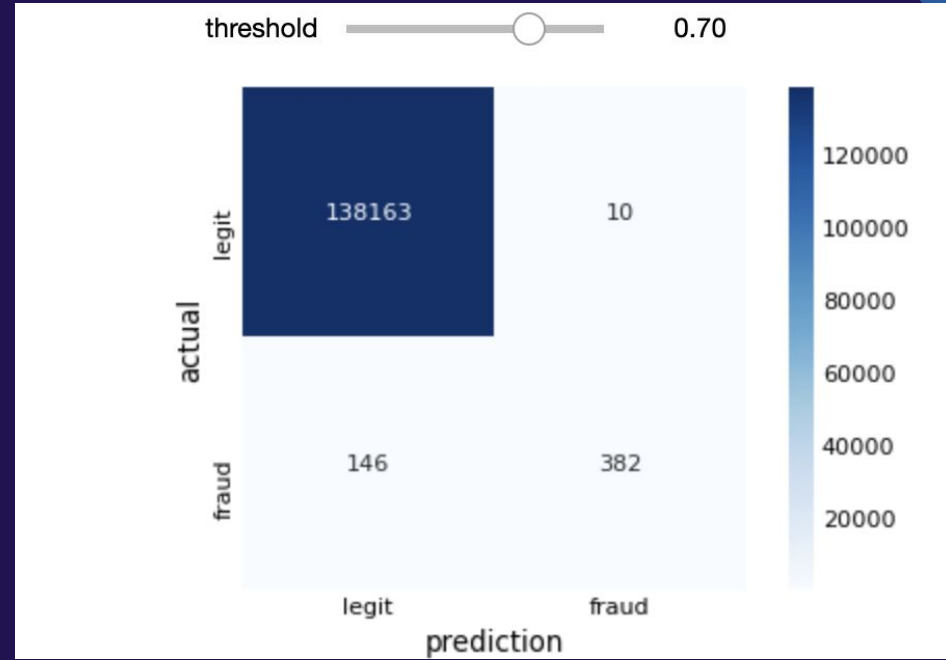


Feature importance



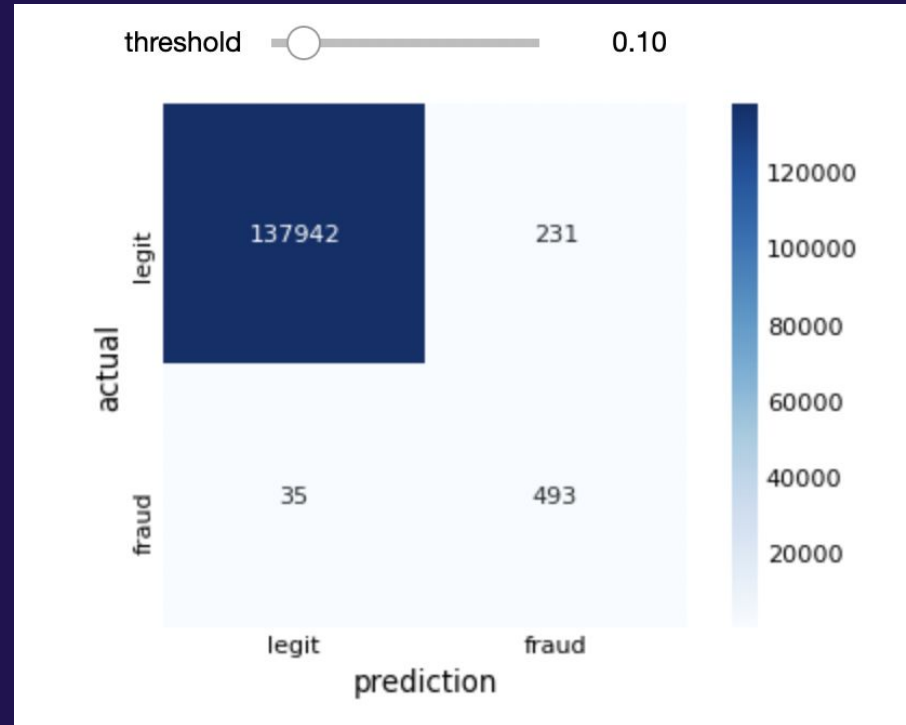
Focus on not alerting customers needlessly

- Set a threshold of 70%
- Catches 72% of all fraud
- Virtually no customers will receive false alerts

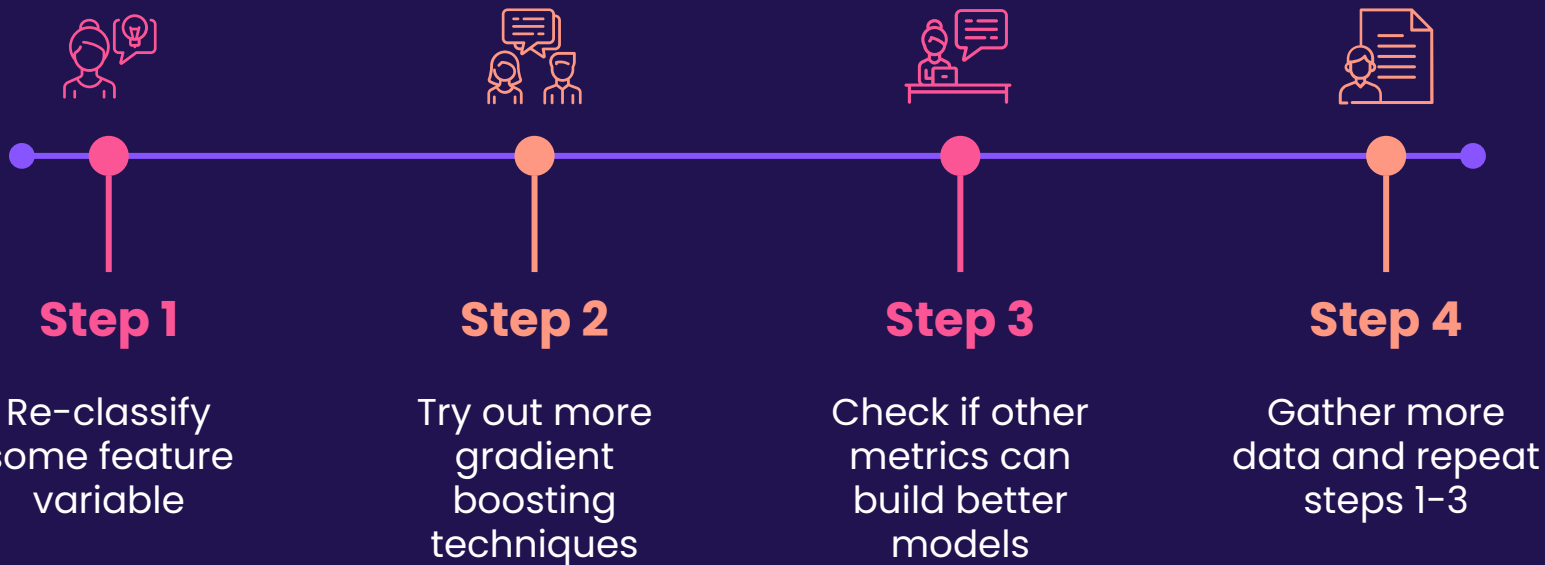


Focus on catching as much fraud as possible

- Set a threshold of 10%
- Catches 93% of all fraud
- 23x more customers impacted but still less than 1%



Future Work



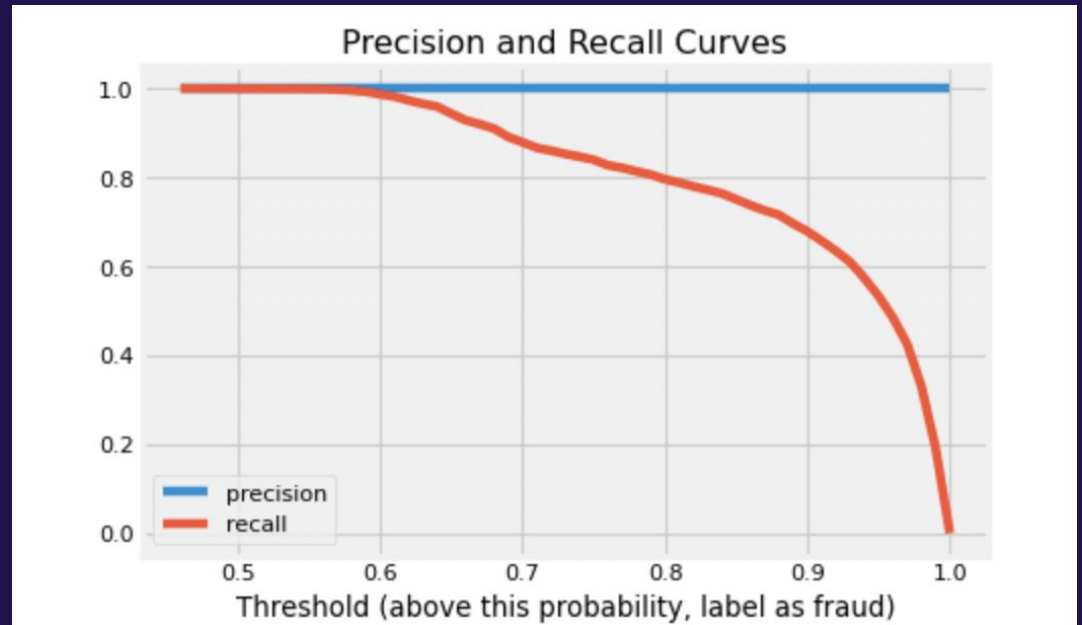
Conclusions

- Random forest provided the best model fit.
- 4 key features with purchase categories having some of the strongest influence but appeared less frequently
- Deciding on business need will help dictate which model to choose, but regardless we can catch a high amount of fraud along with having minimal unnecessary contacts

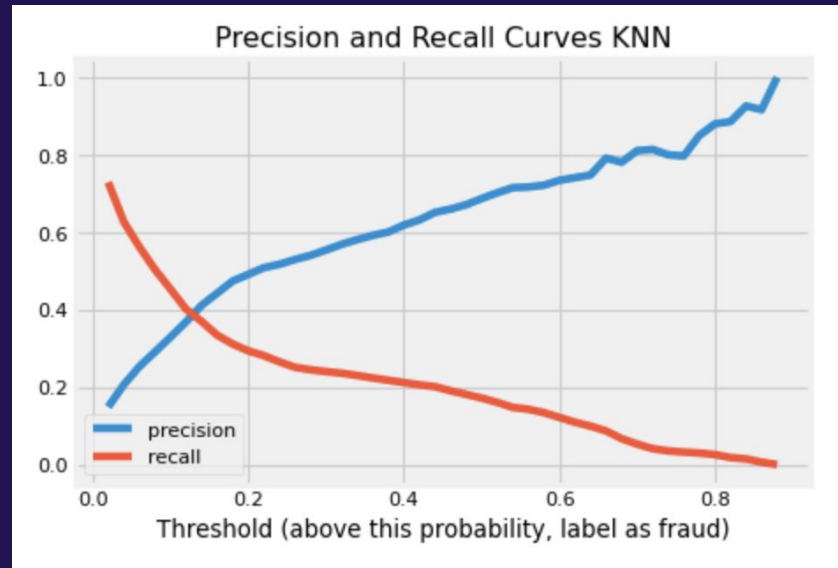
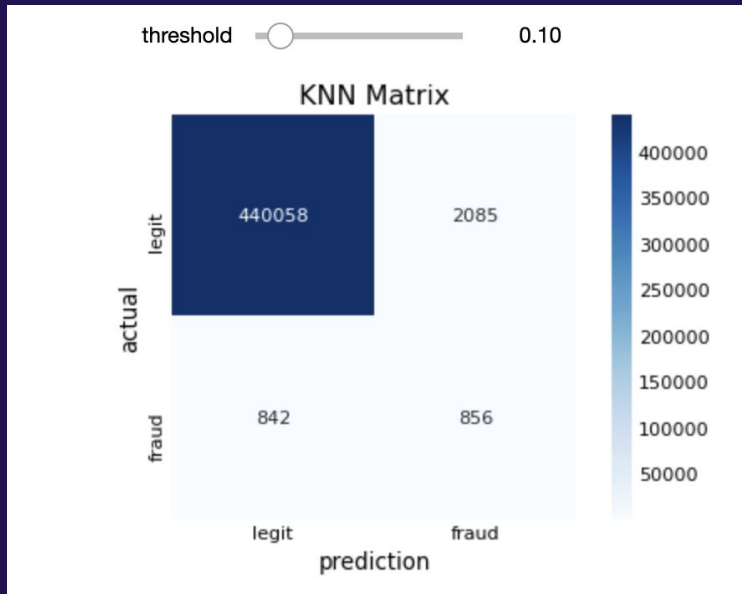


Appendix

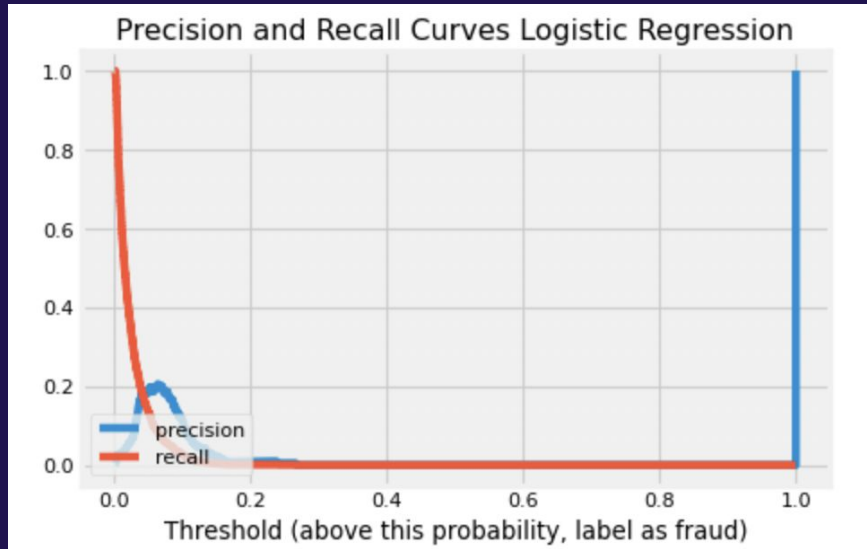
Random Forest Curves



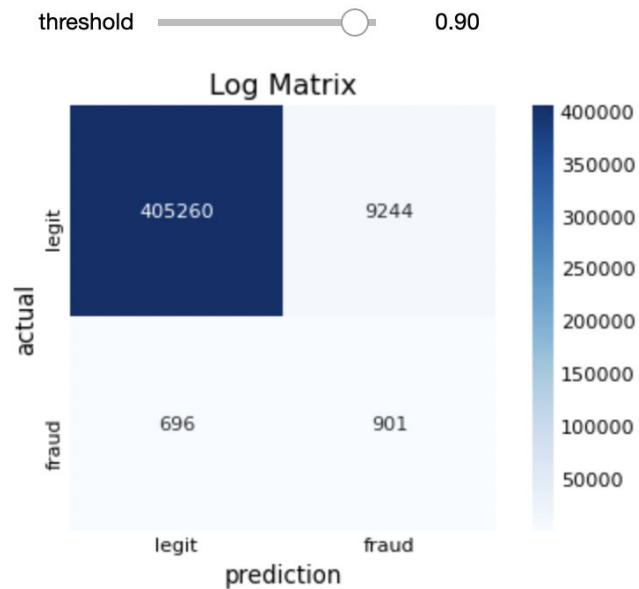
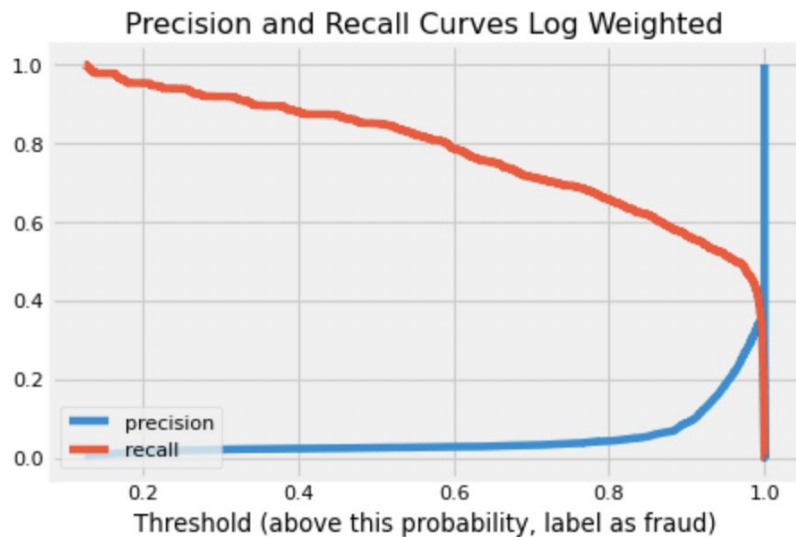
Appendix



Appendix



Appendix



Appendix

