

# SafeSwipe

## CC Fraud Detection

Brian Woo





# Agenda

01

Business Problem

04

Conclusion

02

Data Analysis

05

Next Steps

03

Modeling

06

Demo



# 01

## Business Problem

Financial Institutions are under pressure to combat rising fraud  
Fraud harms both finances and customer trust

Issues:

- 60% of U.S. card holders have been victimized by fraud
  - 45% have experienced fraud multiple times
- Last year: 52 million Americans faced fraudulent charges
  - Losses exceed \$5 billion



# 02

## Data Analysis

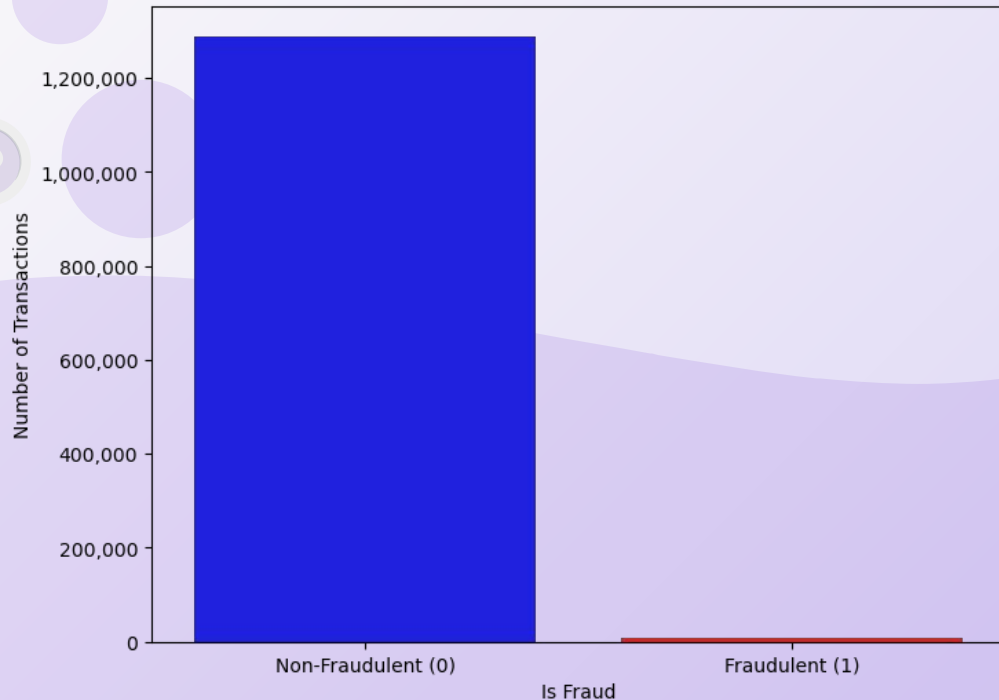
Simulated Dataset contains ~1.3 Million rows of data

- trans\_date\_trans\_time
- category
- amt
- state
- zip
- city\_pop
- lat
- long
- merch\_lat
- merch\_long
- job
- dob
- Target: is\_fraud



Source: [Credit Card Fraud Detection](#)

Distribution of Fraudulent Transactions



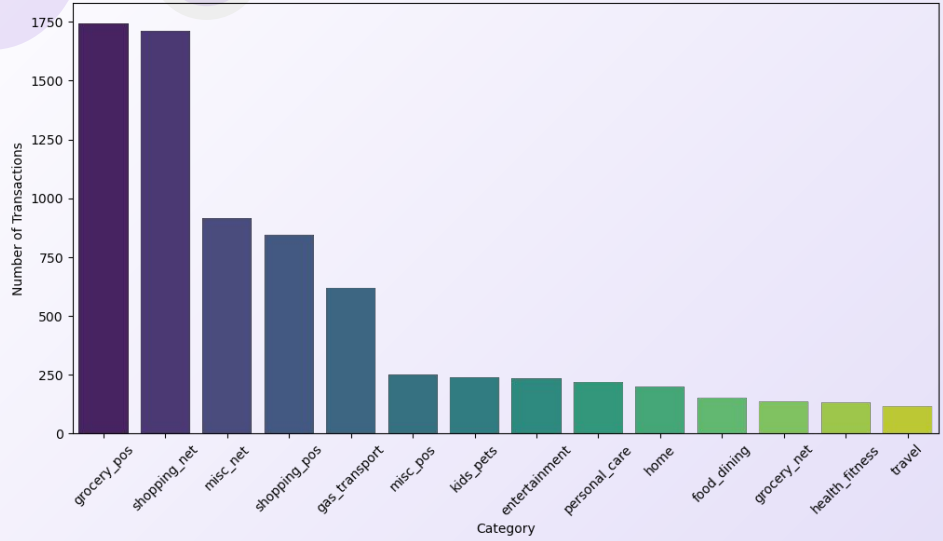
Large Class-Imbalance

7506 Fraud

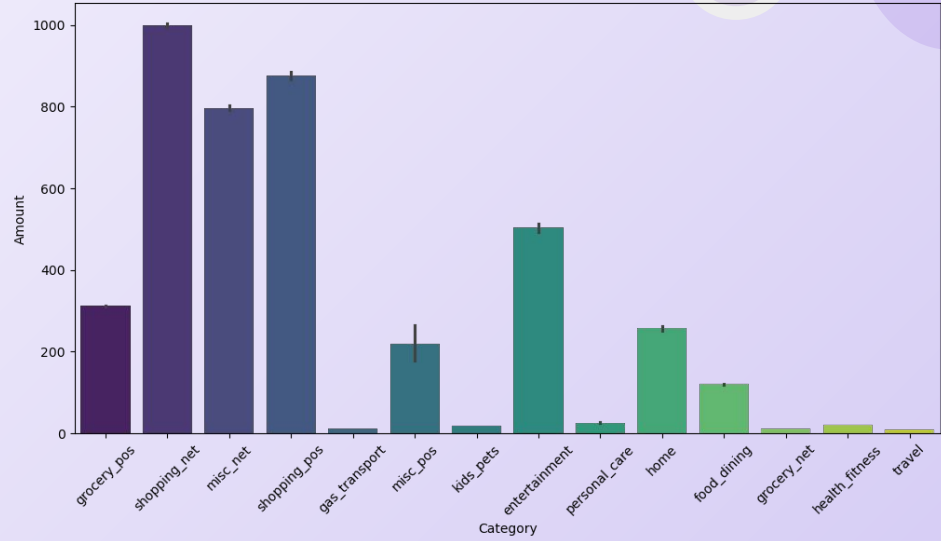
Randomly Sampled: 7506 non-fraud



Distribution of Fraudulent Transactions in Each Category

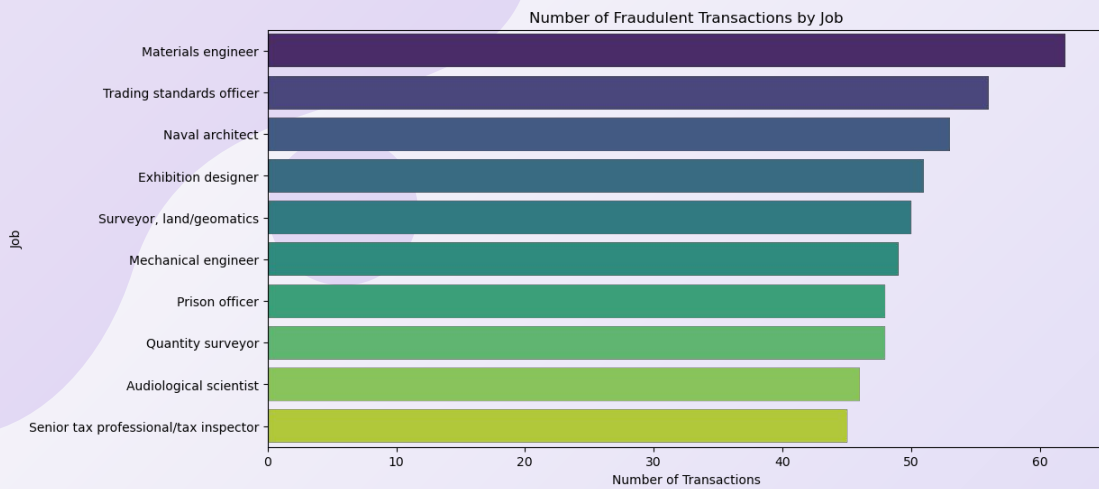


Amount of Fraudulent Transactions in Each Category

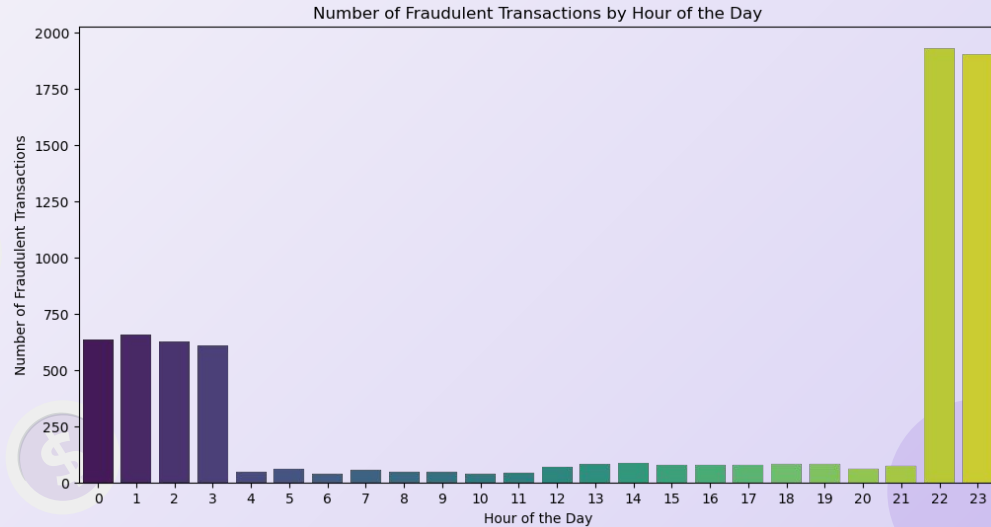


Categories have different spending amounts





Specific jobs have more fraudulent activity



Specific times during the day could indicate patterns for fraud



# 03

## Modeling

Goal: Predict whether a credit transaction is fraudulent or not

Metrics:

- Precision
  - Correctly predicts fraud out of all predicted fraud
  - Reduces false alarms
- Recall
  - Identifies actual fraud cases
  - Catches most fraudulent transactions
- F1-Score
  - Balance between Precision and Recall







## Dataset for Training Models

- category
- amt
- state zip
- is\_fraud



### (Feature Engineered)

- trans\_day
- trans\_month
- trans\_hour
- age\_at\_transaction
- city\_pop\_group
- job\_category
- time\_since\_last\_trans
- avg\_transaction\_amount
- category\_transaction\_count
- unique\_transactions\_day
- user\_merchant\_distance



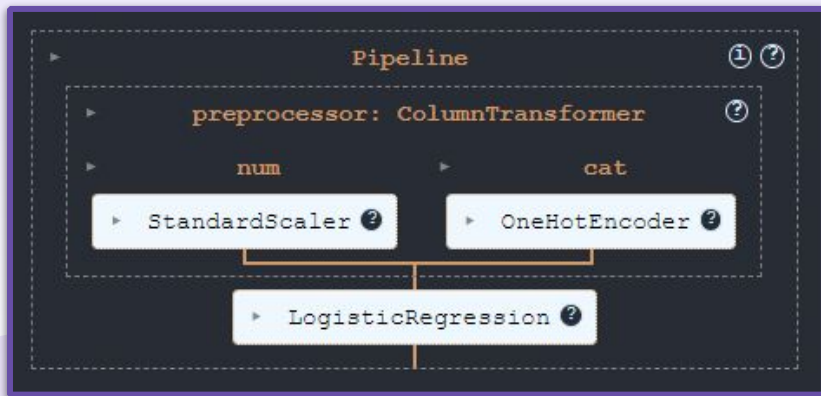
## Models:

### Baseline:

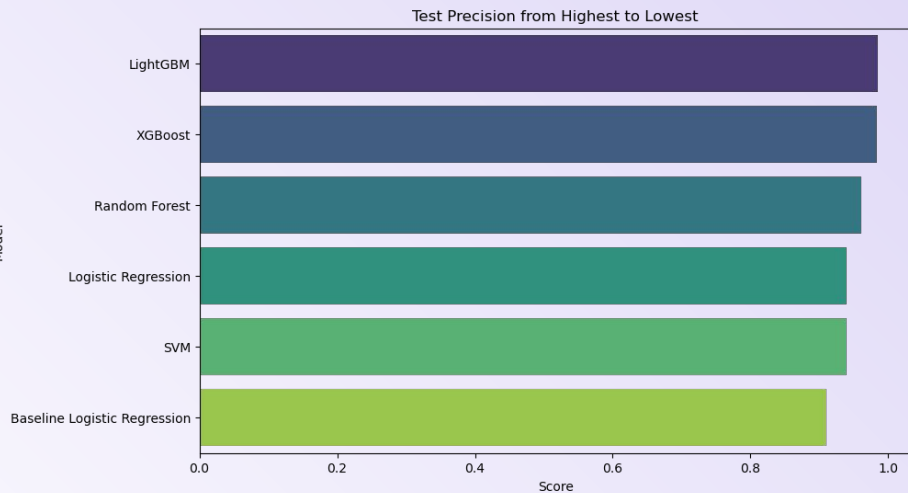
1. Logistic Regression

### Tuned:

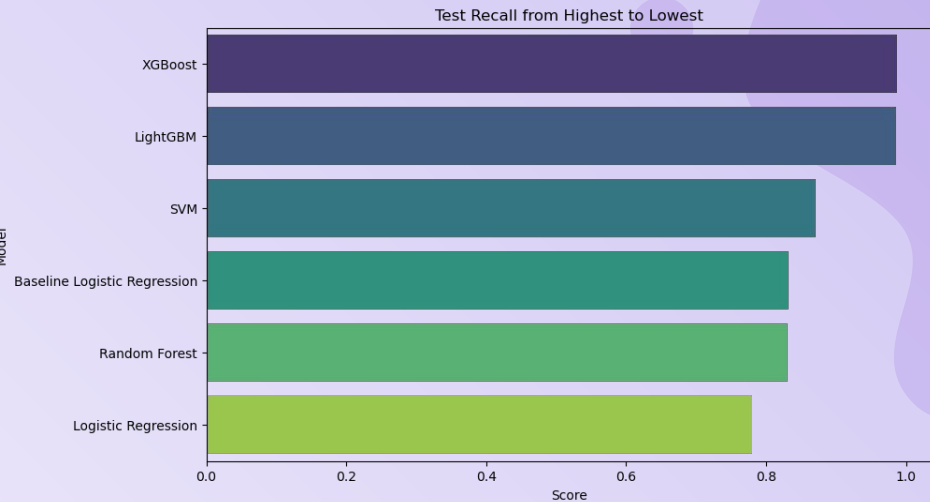
1. Logistic Regression
2. LightGBM
3. XGBoost
4. Random Forest
5. Support Vector Machines



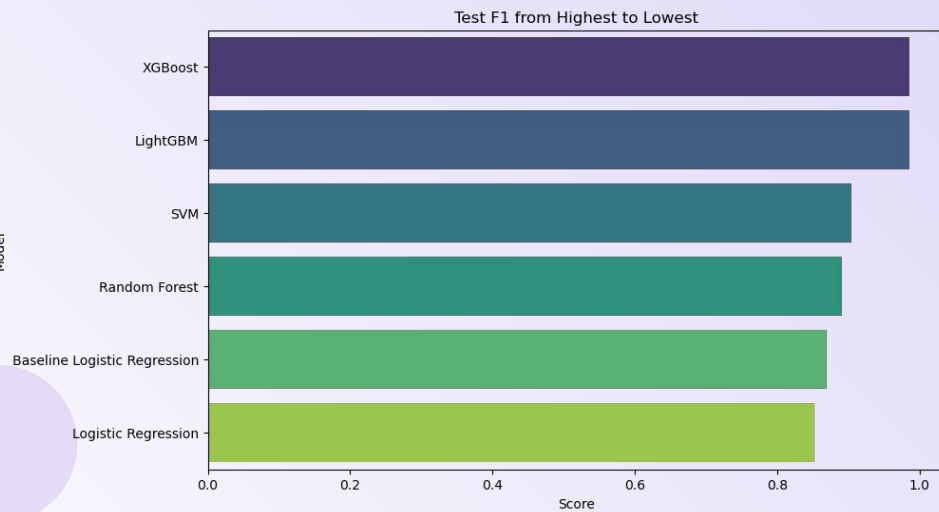
Model



Model



Model



# 04

## Conclusion



	Model	Train Precision	Test Precision	Train Recall	Test Recall	Train F1	Test F1
3	LightGBM	0.998475	0.984472	0.997145	0.985346	0.997810	0.984909
2	XGBoost	0.993143	0.983621	0.992387	0.986679	0.992765	0.985147
4	Random Forest	0.967514	0.960493	0.833270	0.831261	0.895388	0.891216
1	Logistic Regression	0.942632	0.940011	0.775600	0.779307	0.850997	0.852149
5	SVM	0.932290	0.939597	0.872668	0.870337	0.901494	0.903642
0	Baseline Logistic Regression	0.907161	0.910593	0.829463	0.832149	0.866574	0.869606

Best Model: XGBoost

- 2nd highest Precision score  
BUT
- Highest Recall and F1-Scores
- Highest Overall Scores



# 05

## Next Steps



### Financial Institutions

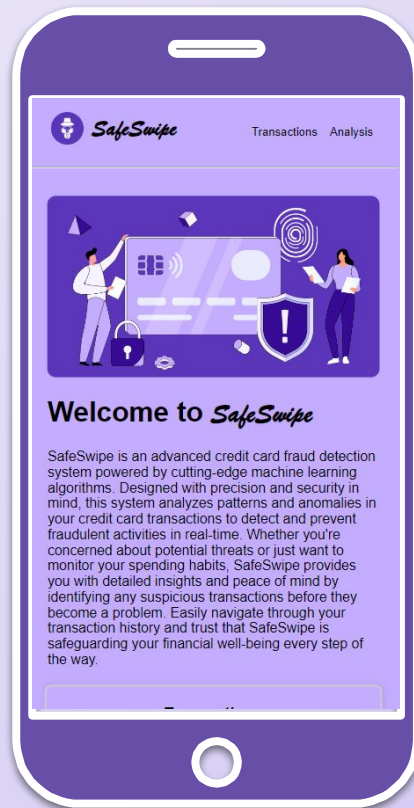
- Build Account Database: managing user accounts
- Flexible Input: allow optional fields for easier use
- Implement Real-Time Alerts

# 06

Demo

## Check out our app!

<https://safeswipe-e7d39aac3b48.herokuapp.com/>



# THANKS!

Do you have any questions?  
brianhwwoo@gmail.com  
<https://github.com/Haoweee>



Credits: This presentation template was created by  
**Slidesgo**, and includes icons by **Flaticon**, and  
infographics & images by **Freepik**

Please keep this slide for attribution