

# CREDIT CARD FRAUD DETECTION

Noah Monastersky  
| Data Scientist

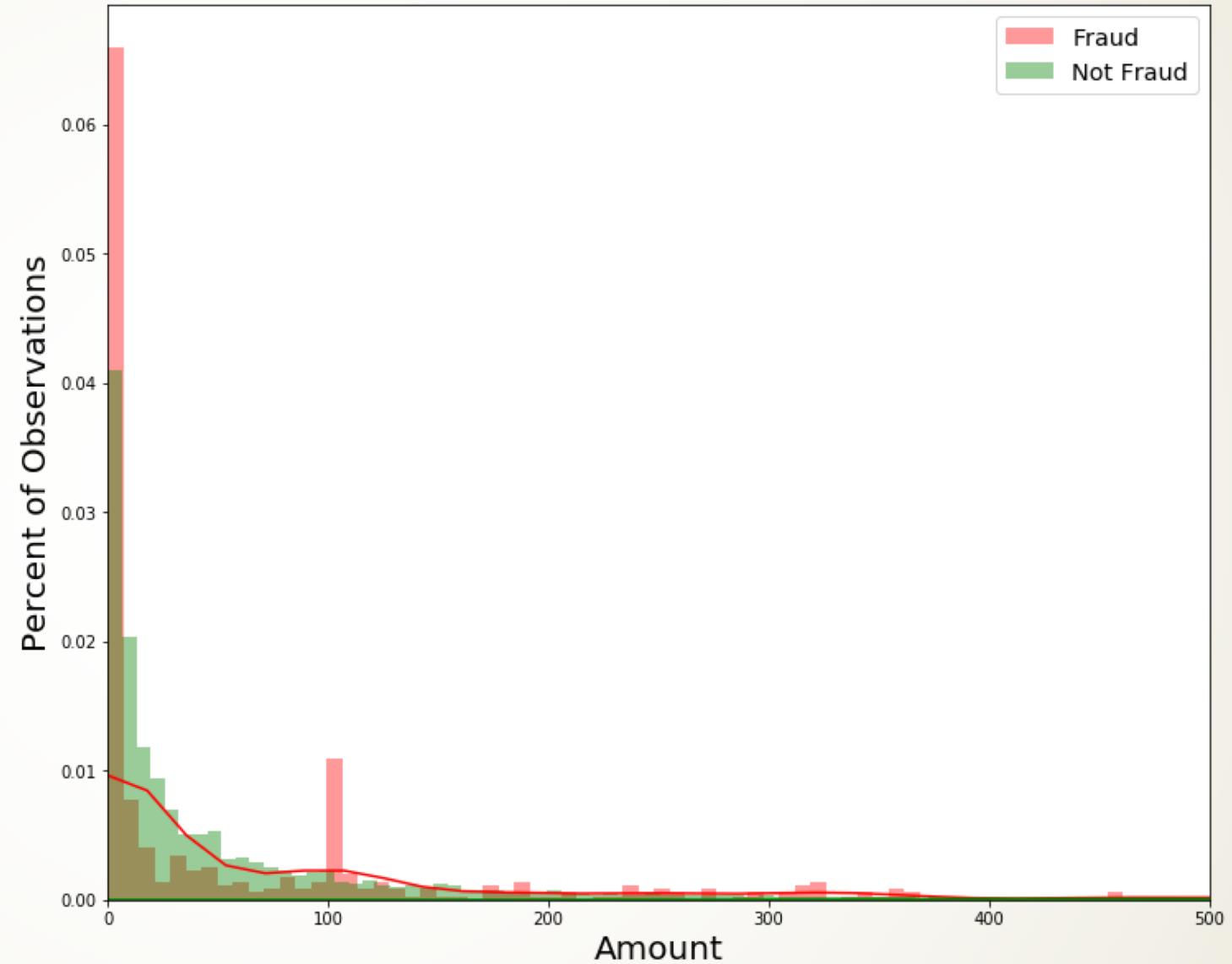
# Problem Statement

- Can I build a model that predicts whether a credit card charge is fraudulent or not?
- This has a clear use cases because financial institutions are liable when fraudulent charges occur.

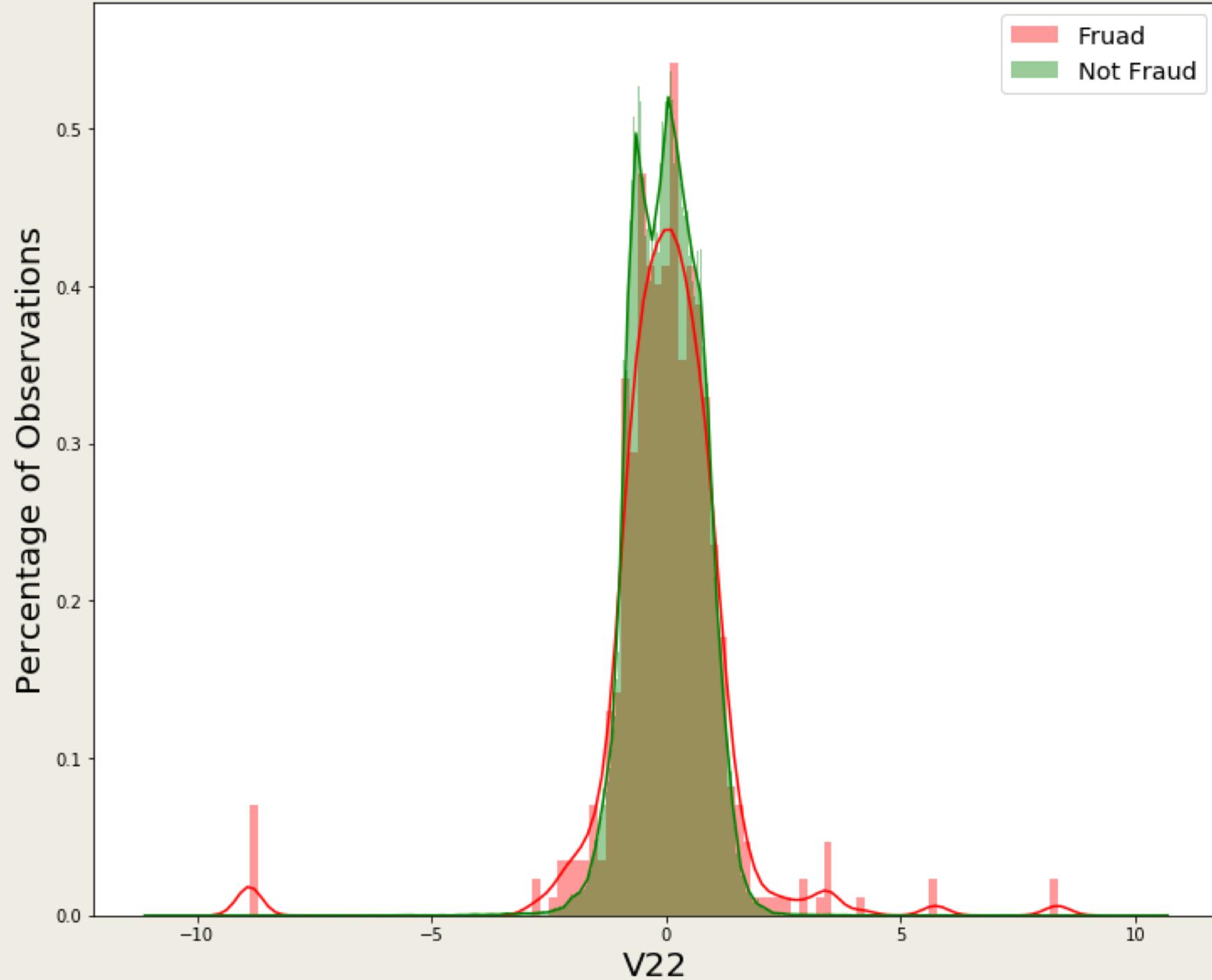
# The Data

- 284,807 observations from European cardholders in 2013
- 492 Fraud | 284,315 non-fraud
- Amount and Time column
- V1 – V28 PCA'd columns.

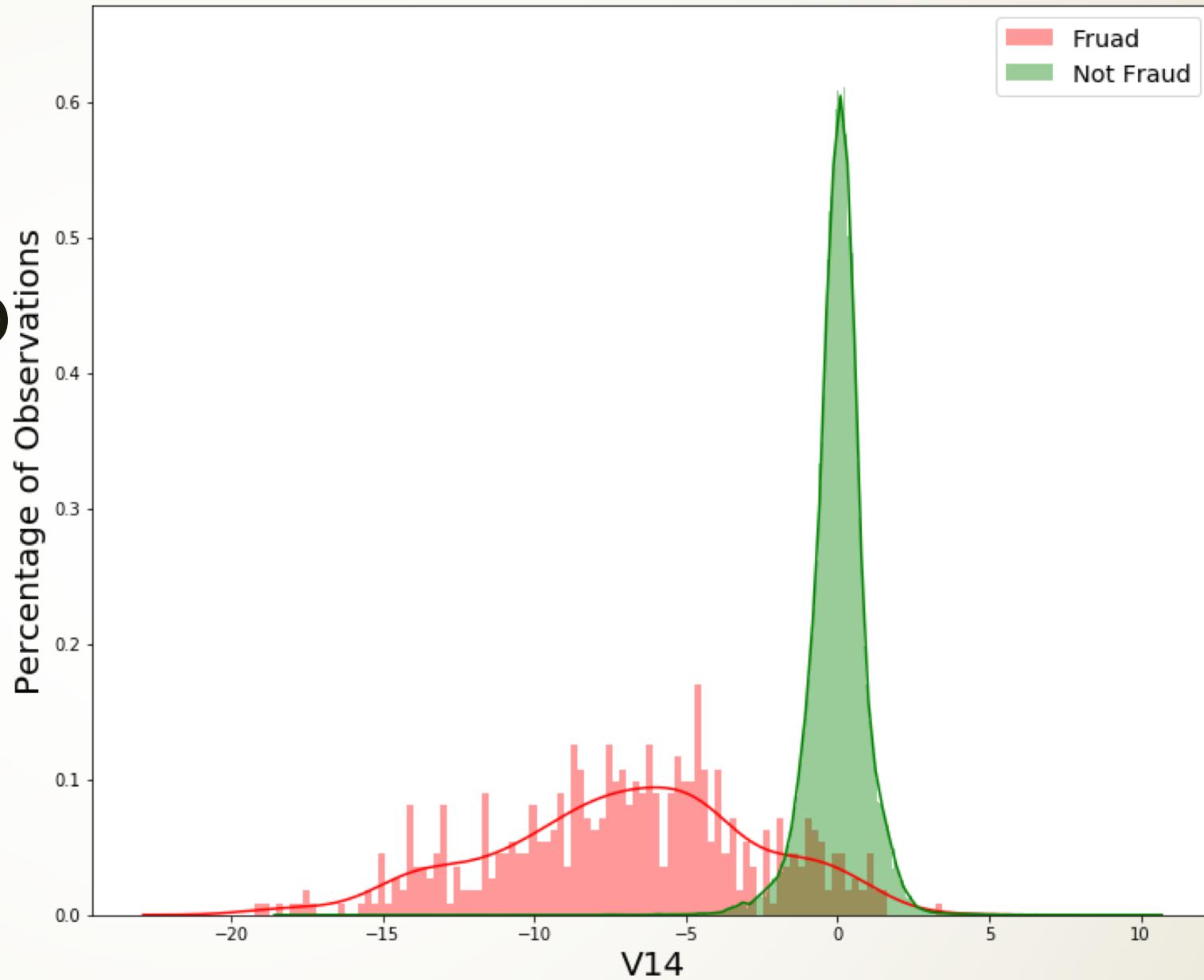
Most Charges  
are low dollar  
amount.



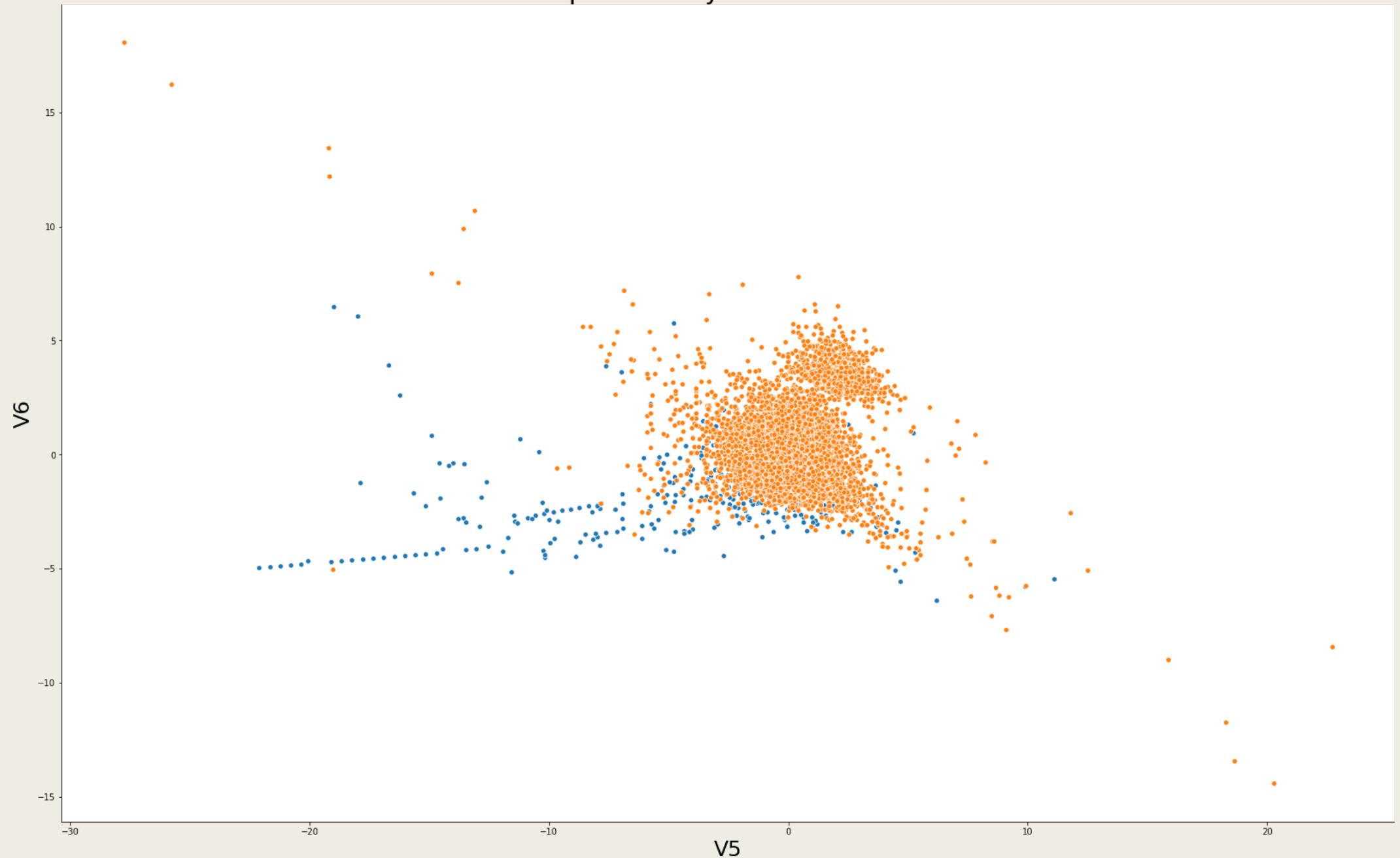
Variables with  
low correlation  
to class like  
V22 follow  
roughly the  
same shape.



Variables  
with high  
correlation to  
class like  
V14 follow  
different  
shapes for  
fraud and  
non-fraud



## Spatial analysis: V5 and V6



# What does my model want to optimize

- False negative: Fraud occurred and was not caught
- False positive: Fraud did not occur but we acted as it did.
- False negatives are worst case scenario so we want a model that optimizes our recall score.

# Training and Testing my Data

- I created a subset of my data by randomly selecting 50,000 rows and training my data and then testing my models on the entire dataset.
- I also used SMOTE to underrepresent the majority class and overrepresent the minority class.
- Scaled my data.

# Modeling

---

Logistic Regression

---

Random Forest

---

Gradient Boosting

---

XG Boost

---

K Nearest Neighbors

---

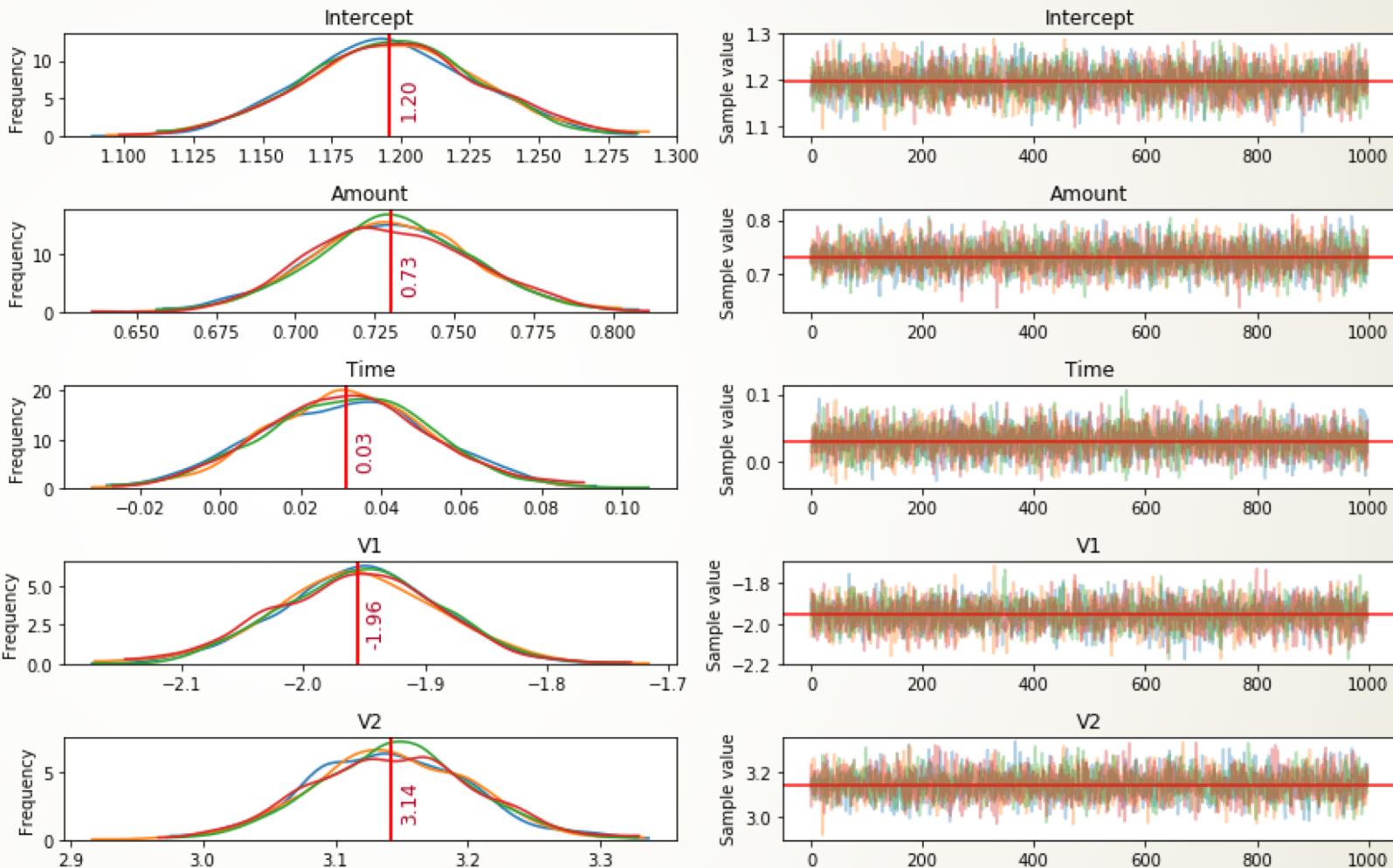
Support Vector Machine

---

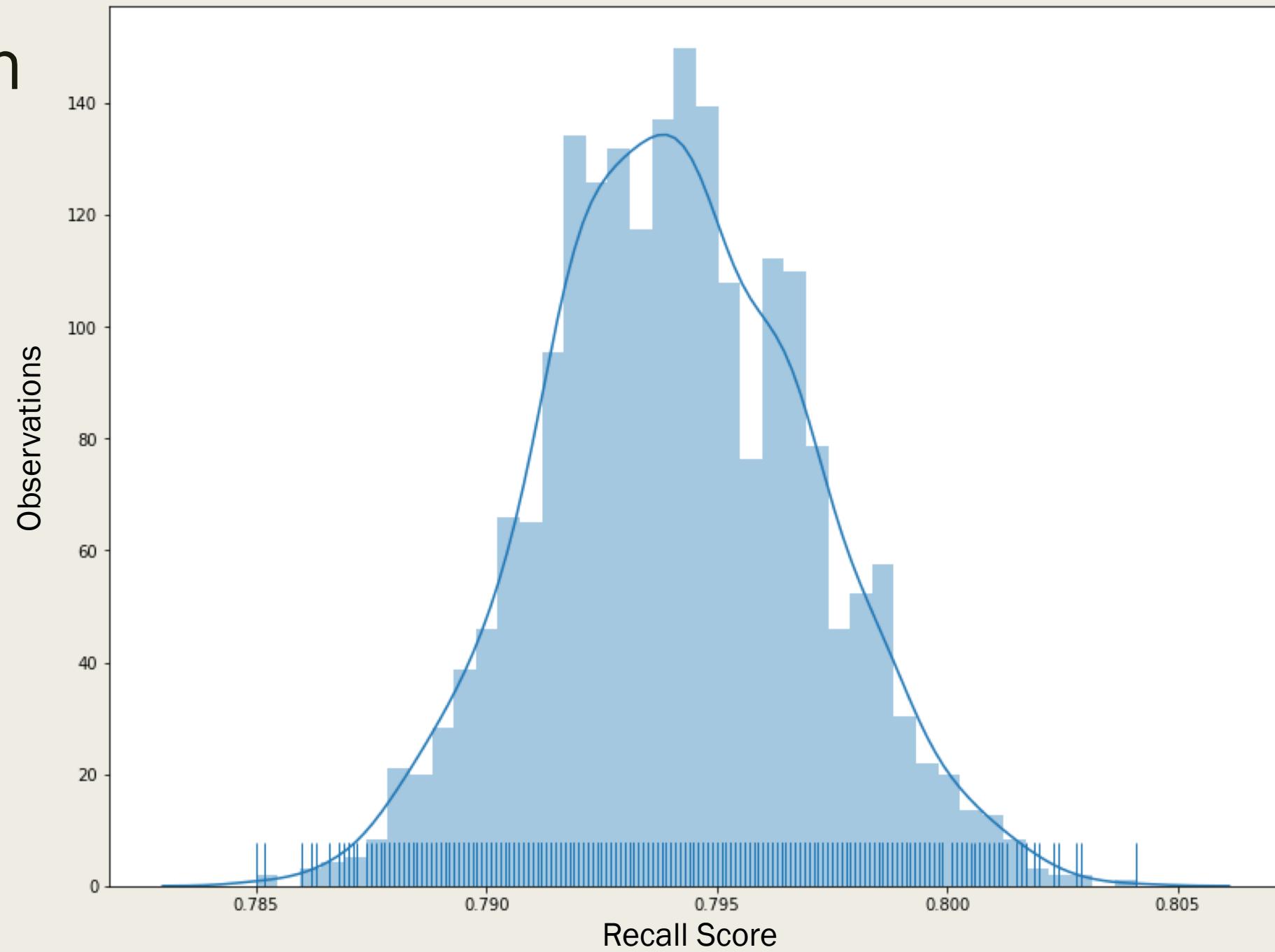
K-Means

# Bayesian Logistic Regression

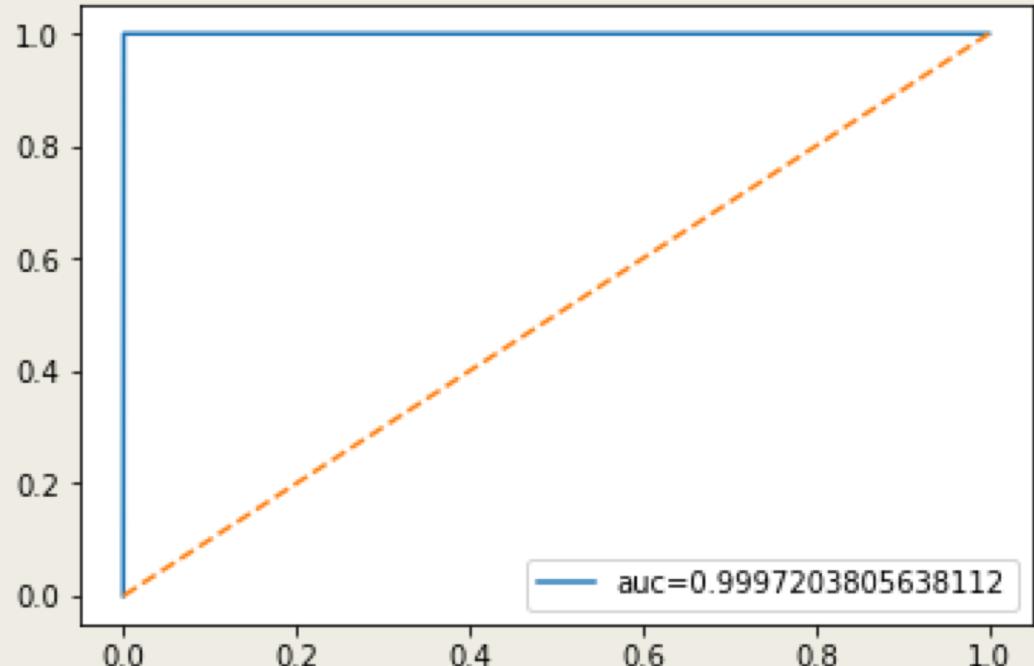
- For each Beta there is a distribution for possible values it could observe.



# Distribution of Models



# Conclusions



- Use Random Forest to model as the accuracy score is .9994 and I am able to get 0 false negatives.
- Bayesian models did not outperform regular logistic regression.
- In the future, I would love to get the anonymized data and explore those variables.