



# CAPSTONE 3: INCREASE CUSTOMER SATISFACTION WITH MACHINE LEARNING

JOE BOARDMAN



# PROBLEM

- Big banks will always get complaints from customers about all sorts of different products. Being able to classify complaints to help identify the root cause of the issue can be crucial in maintaining and improving customer satisfaction

# RESULT

- Effective multi-class classification model achieving a 0.94 AUC score
- Topic modeling to uncover a more accurate root cause for complaint
  - Customer Service issues, etc

# DATA

- Consumer Financial Protection Bureau (CFPB) complaint database
- Focused only on 4 large American banks
  - JP Morgan Chase
  - Bank of America
  - Wells Fargo
  - Citibank
- 600k records

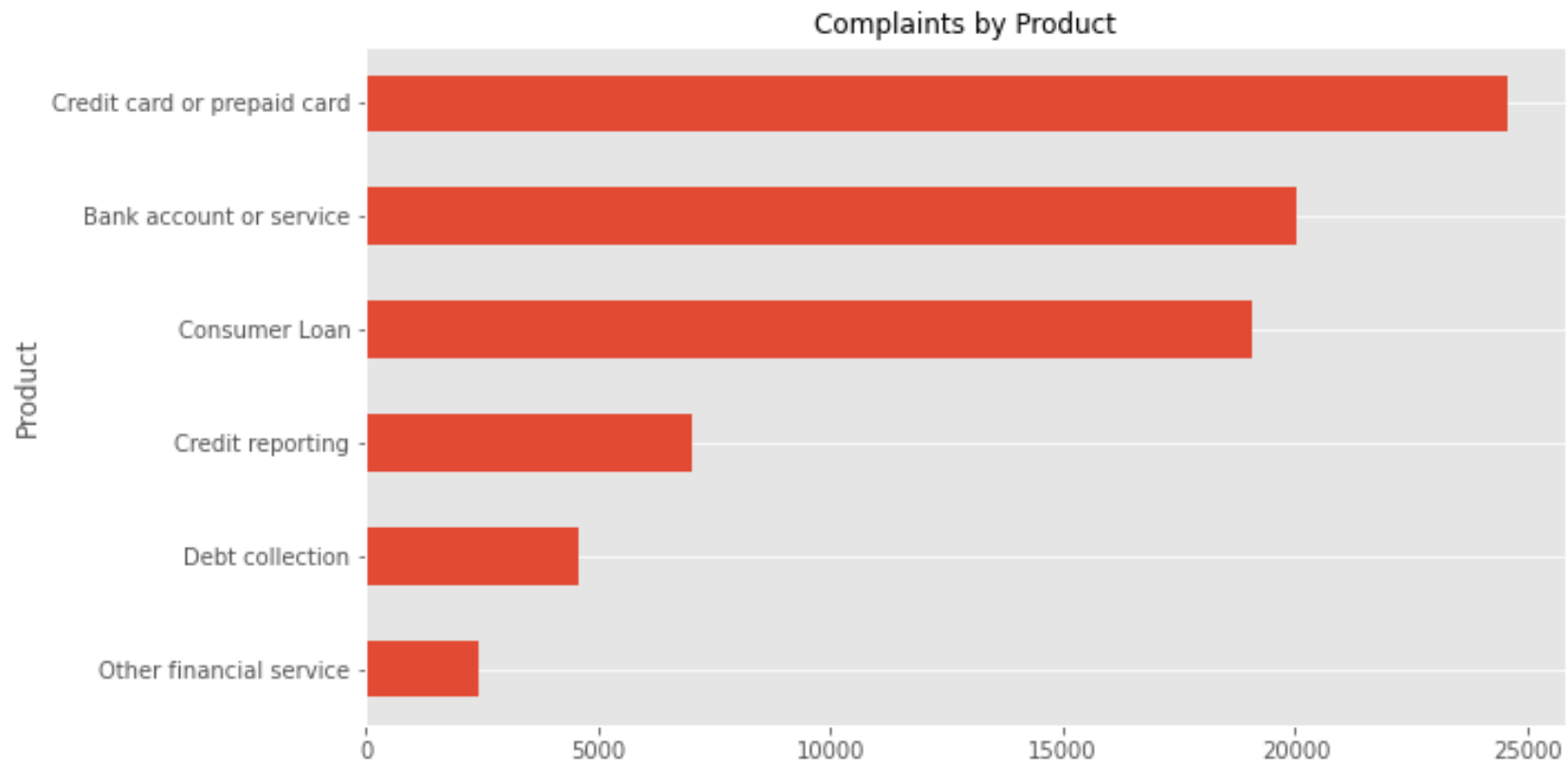
# DATA CLEANING

- Filtered records without written complaint
- Sampled from full dataset to only include 4 banks
  - Left with 77k rows
- Consolidated products to solve redundancy
- Checked for duplicates and removed if necessary

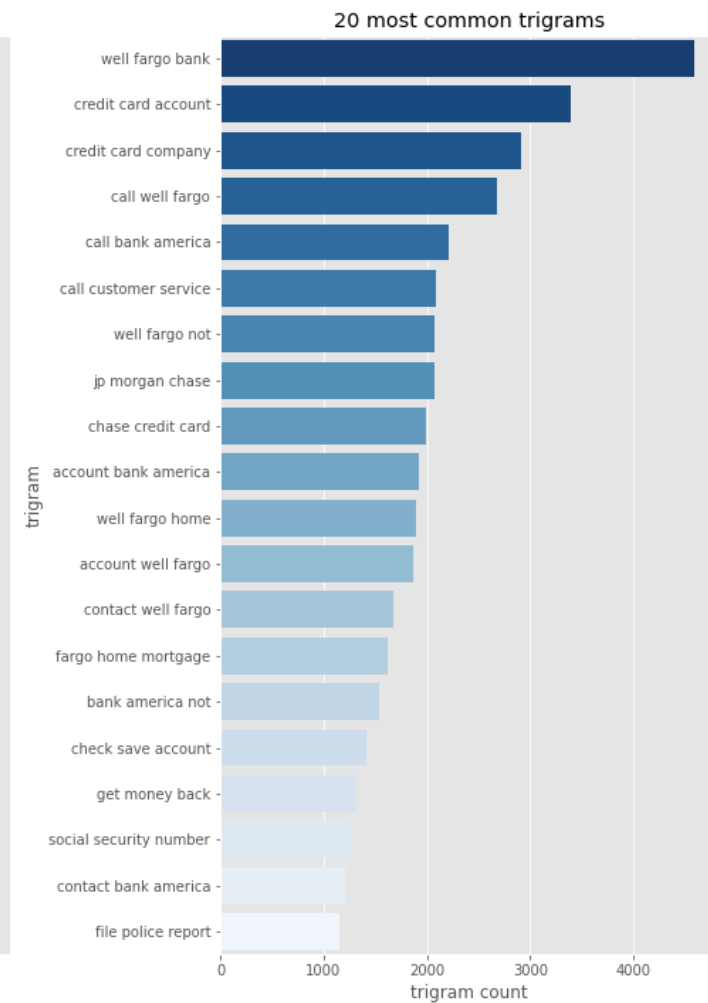
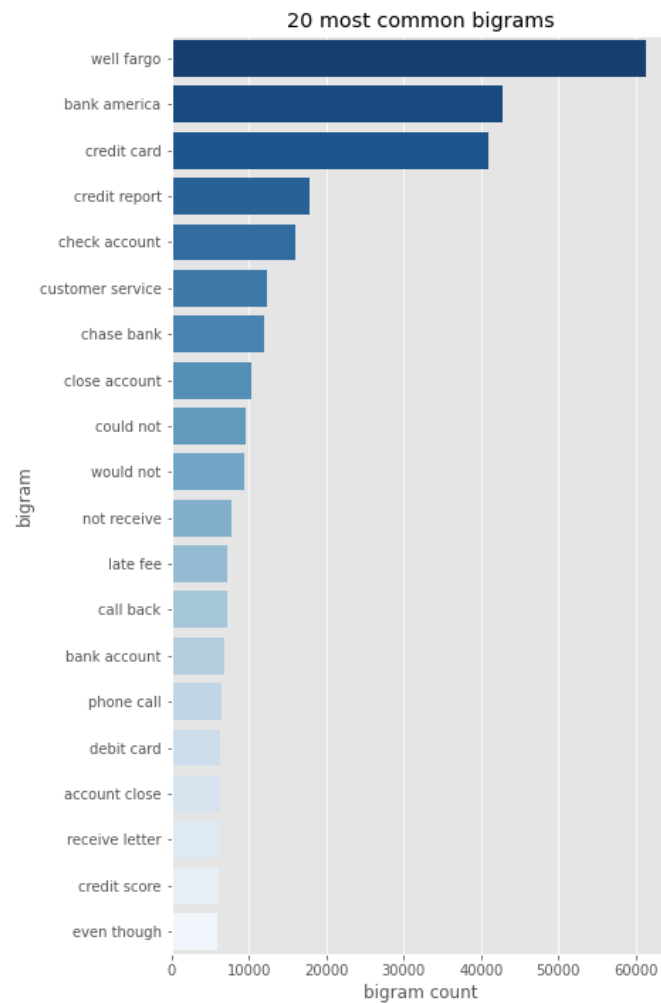
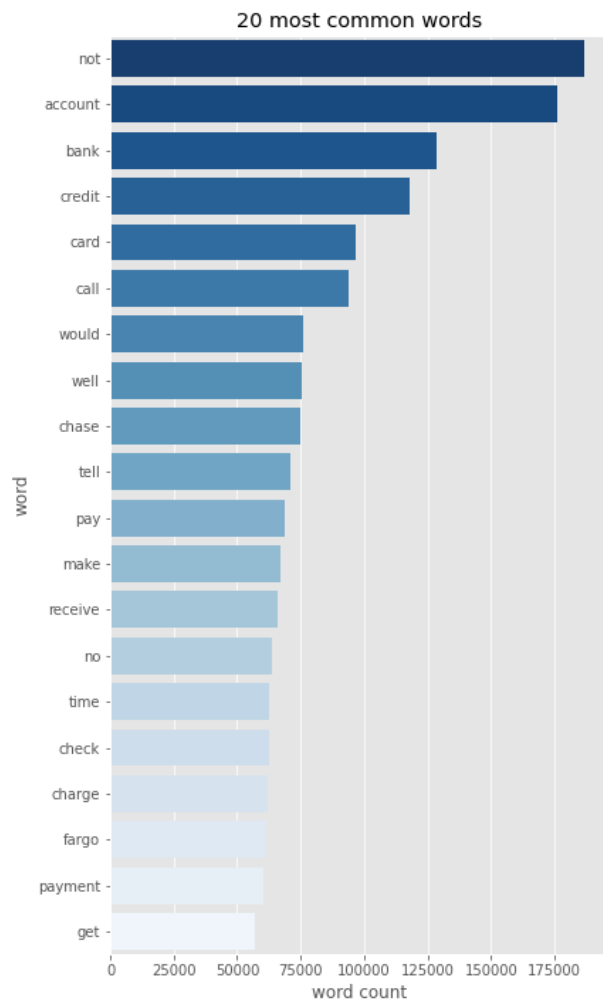
# TEXT PREPROCESSING

- Word tokenization
- Lemmatization
- Removal of stop words
- Convert back to a string

# EXPLORATORY DATA ANALYSIS



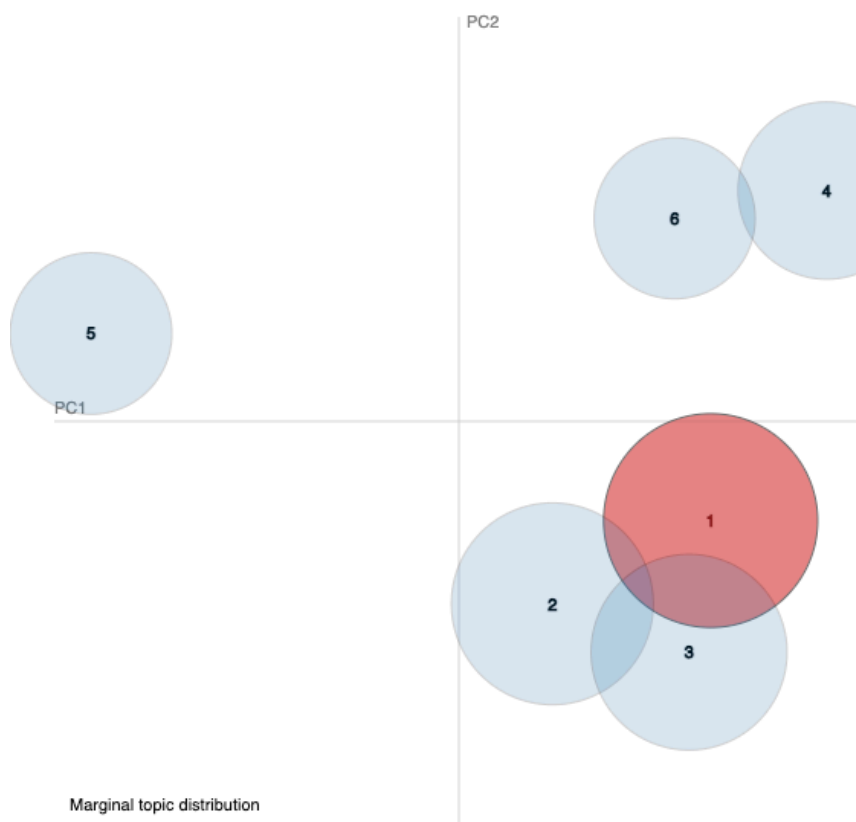
# EXPLORATORY DATA ANALYSIS





# TOPIC MODELING

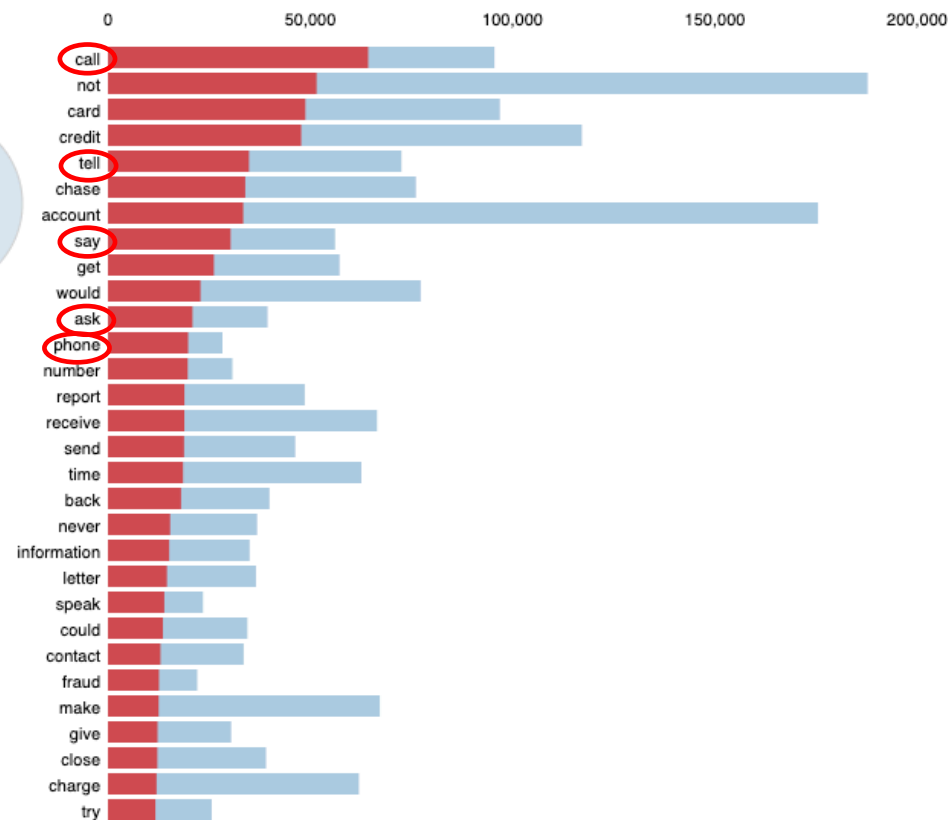
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 1 (22% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

# MODELS USED

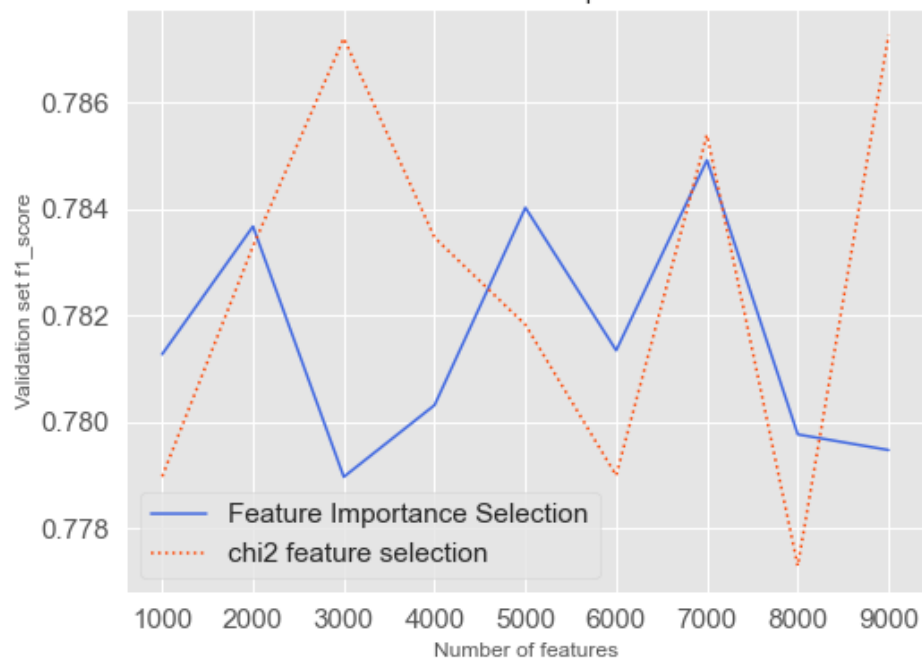
- Random Forest
- Multinomial Naïve Bayes
- Logistic Regression

# FEATURE SELECTION

## ■ Chi-squared vs Feature Importances

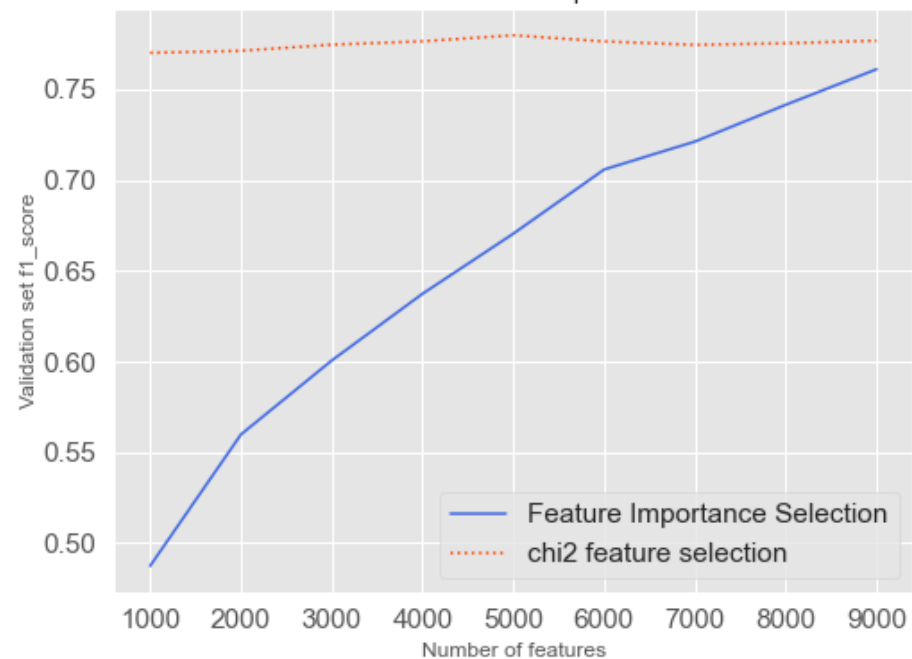
### Random Forest

Feture Selection: Feature Importances vs Chi2



### Naïve Bayes

Feture Selection: Feature Importances vs Chi2

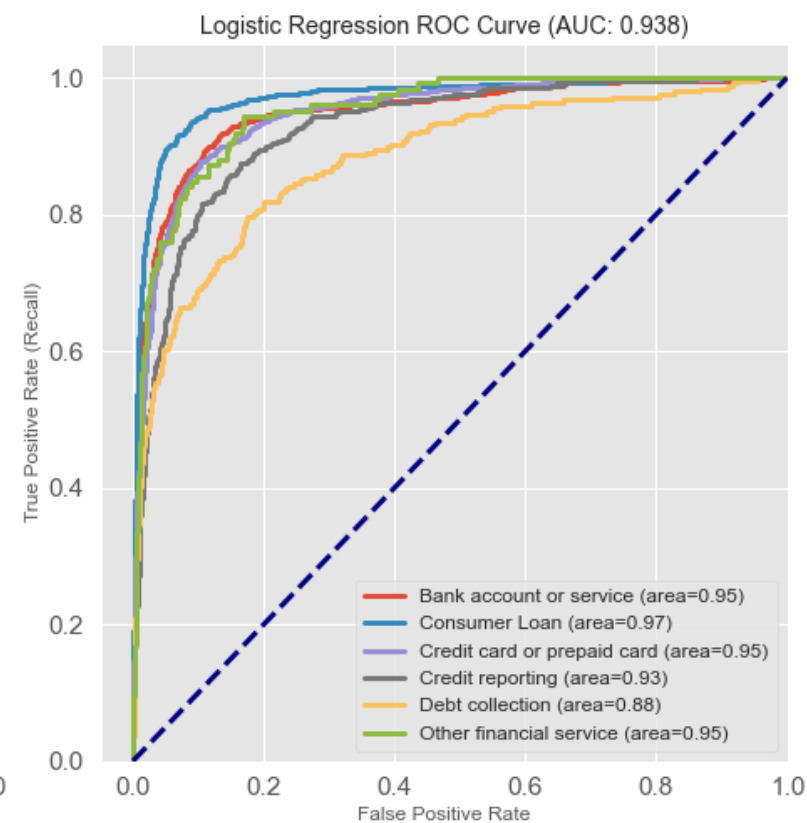
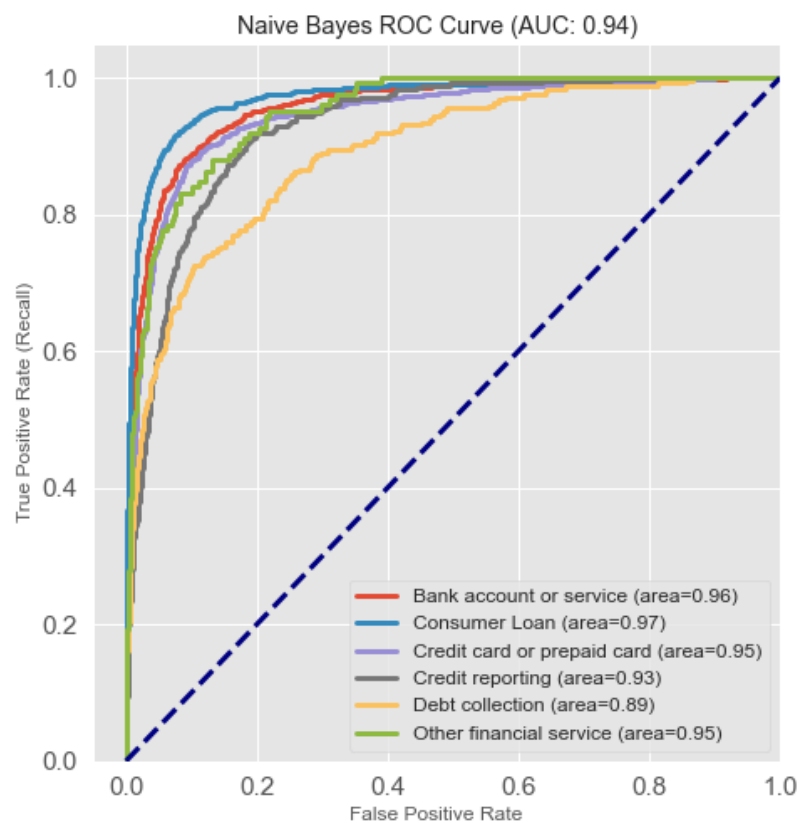
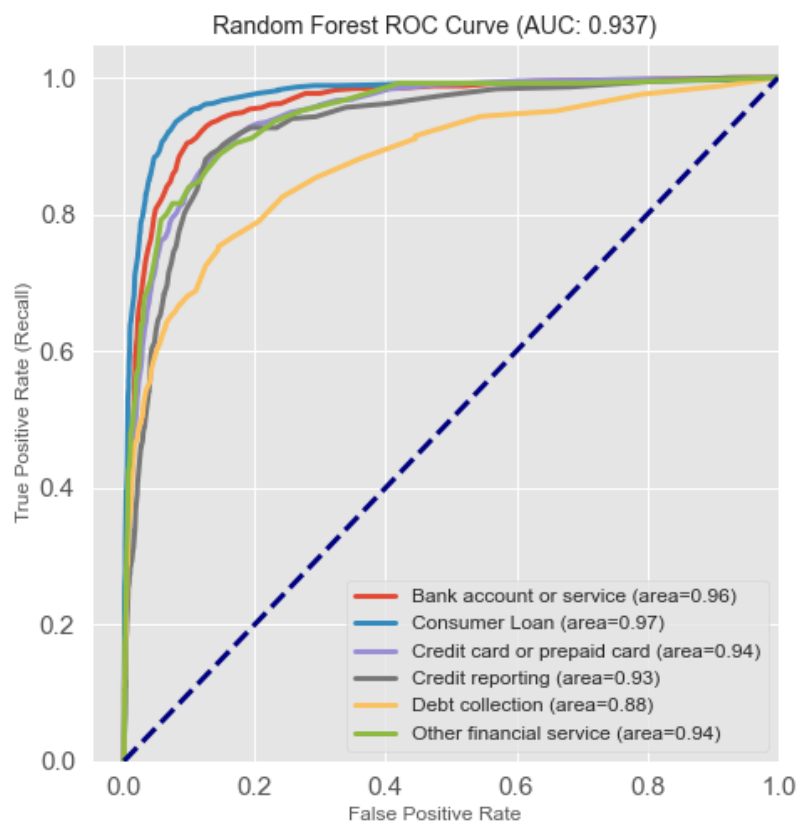


# EVALUATION

- My recommendation would be to use Naïve Bayes
  - Similar performance to Random Forest and Logistic Regression
  - Very computationally efficient

	Weighted f1	Runtime (mins)	AUC Scores
Model			
Random Forest Classification	0.787	1.0	0.937
Naive Bayes Classification	0.786	0.0	0.940
Logistic Regression	0.786	22.9	0.938

# EVALUATION



# VALUE TO CUSTOMER

- Effectively identify product alignment of complaint
- Compare with the assigned topic to understand a possible other driver

## EXPLORE FURTHER

- Neural Network
- Try more techniques to involve Sentiment Analysis



QUESTIONS?