

# Project 3: Web APIs & NLP

By: Jayson Villena  
01/18/2023

# Problem

A coworker was tasked with gathering data of individual investors and current thoughts about the market. He gathered the data from two subreddits, WallStreetBets and StockMarket. The senior data scientist likes the content but he doesn't want to utilize data from wsb, however the worker didn't label where the subreddits came from.

## Problem Statement:

Collect data from the two subreddits and create a model that can predict where the data came from and match it with the original data.



# What do we know?

- Original data
- Equal Amount

# How?

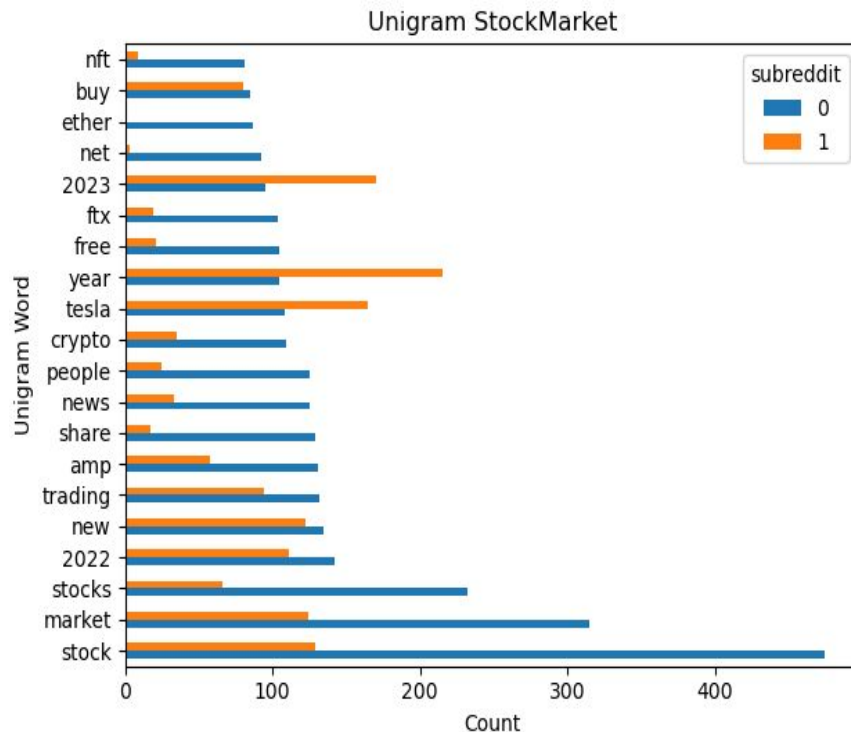
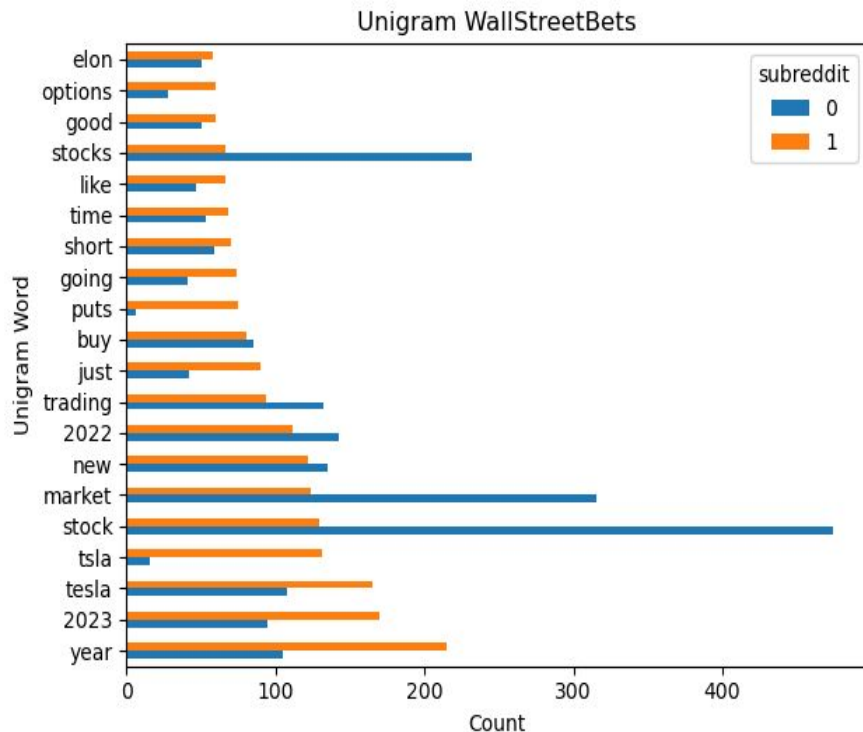
1. API
2. Preprocess and EDA
3. Model Optimization

## What do we need to do?

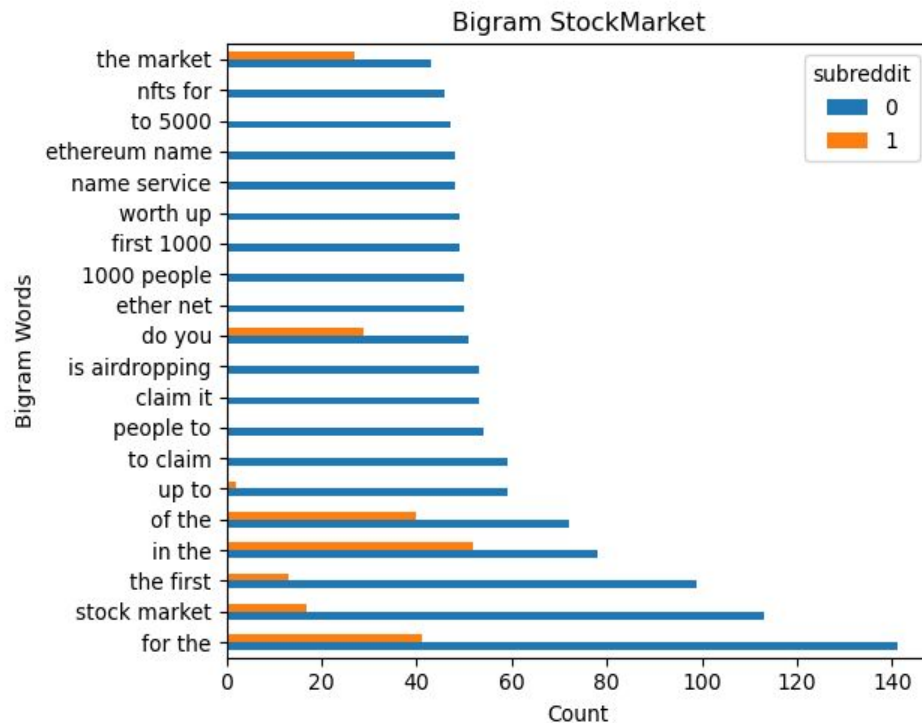
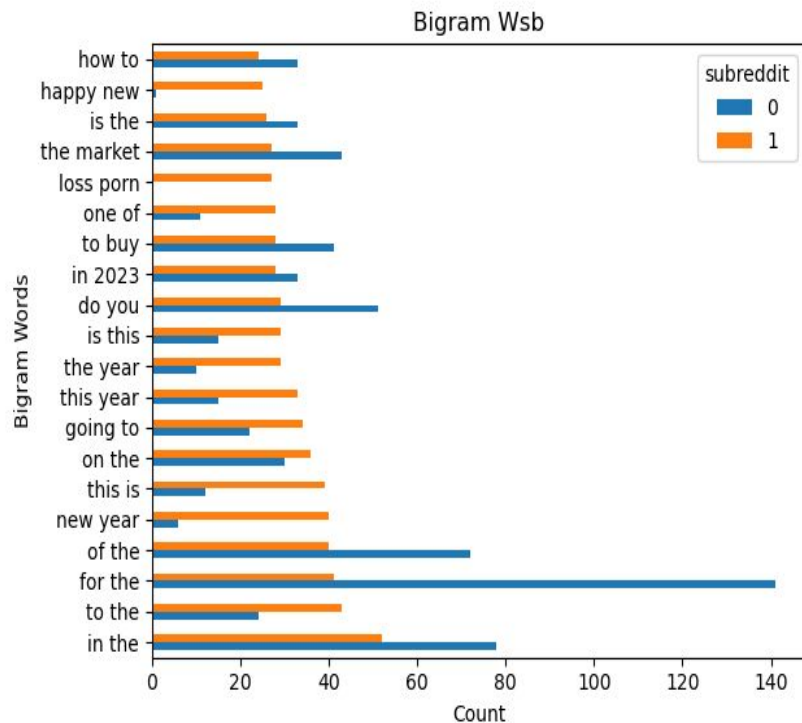
- Gather new data
- Create a predictive model
- Use model on original data



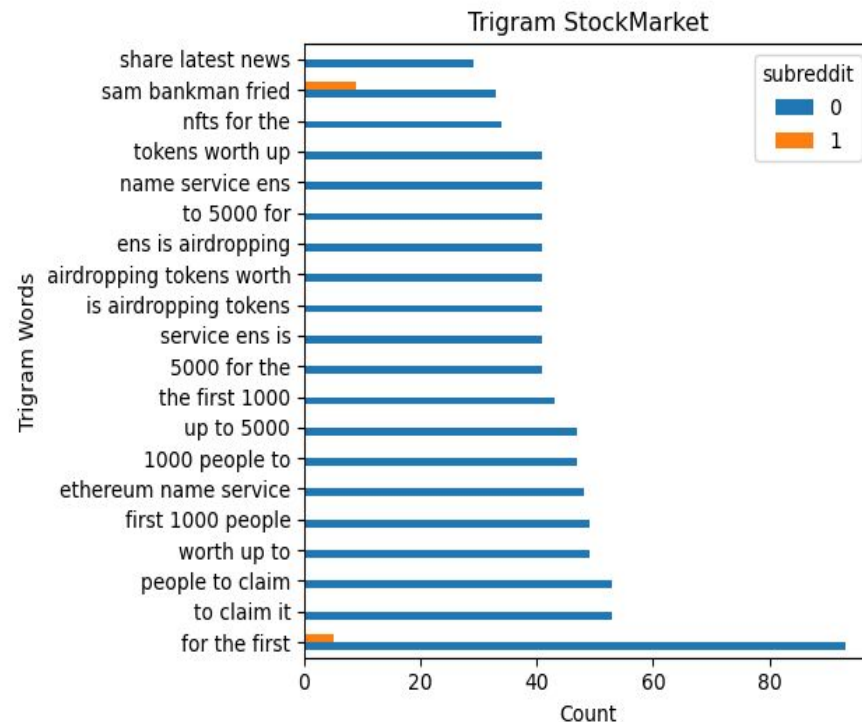
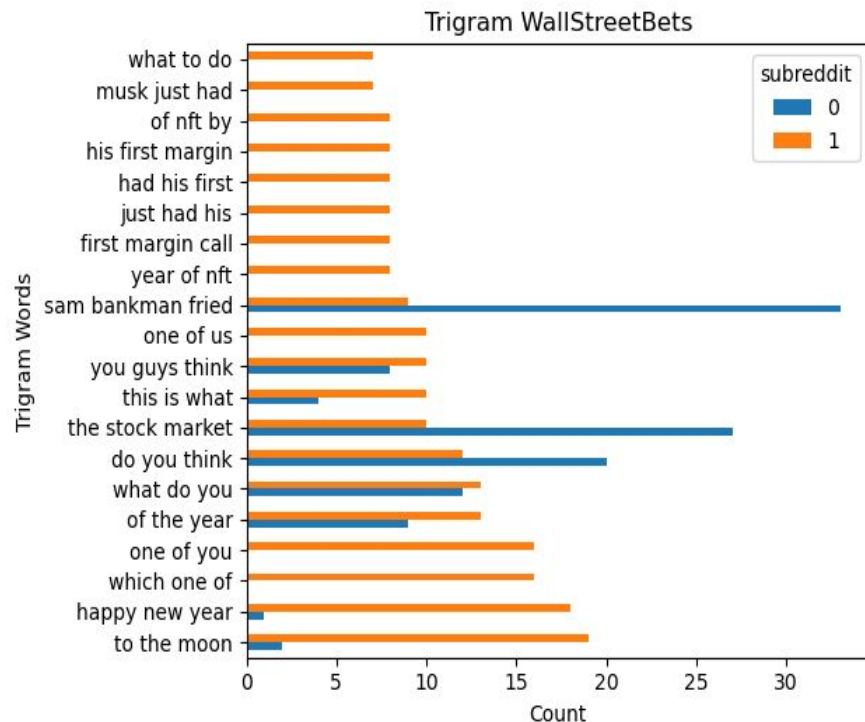
# What are we looking for?



# Bigrams



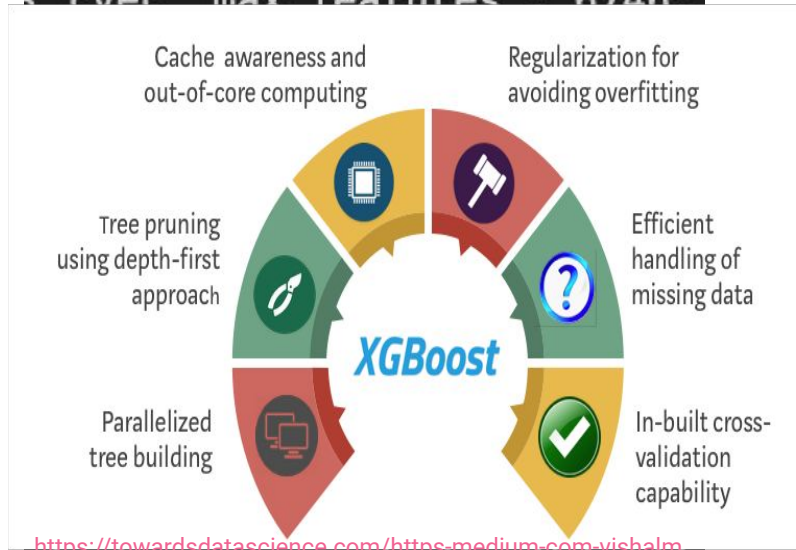
# Trigrams



# Modelling



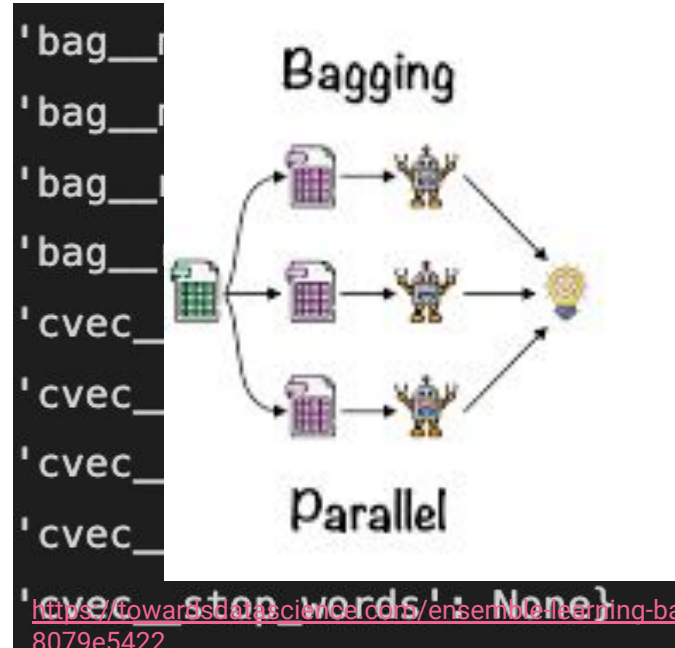
```
{'cvec_max_features': 6746
```



<https://towardsdatascience.com/https-medium-com-vishalmorde-xgboost-algorithm-long-sha-may-rein-edd9f995bc63d>

```
xgb -tree method -n 100 }
```

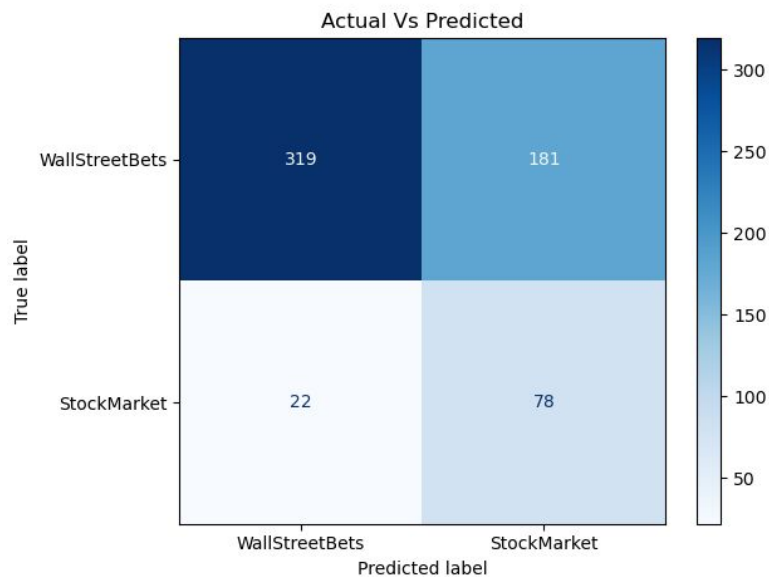
## Bagging



<https://towardsdatascience.com/ensemble-learning-bagging-boosting-8079e5422>

# Results

XgBoost: 0.71    Bagging: 0.72



# Future:

- Better understanding of the API
- Optimization



# Project 3: Web APIs & NLP

By: Jayson Villena  
01/18/2023