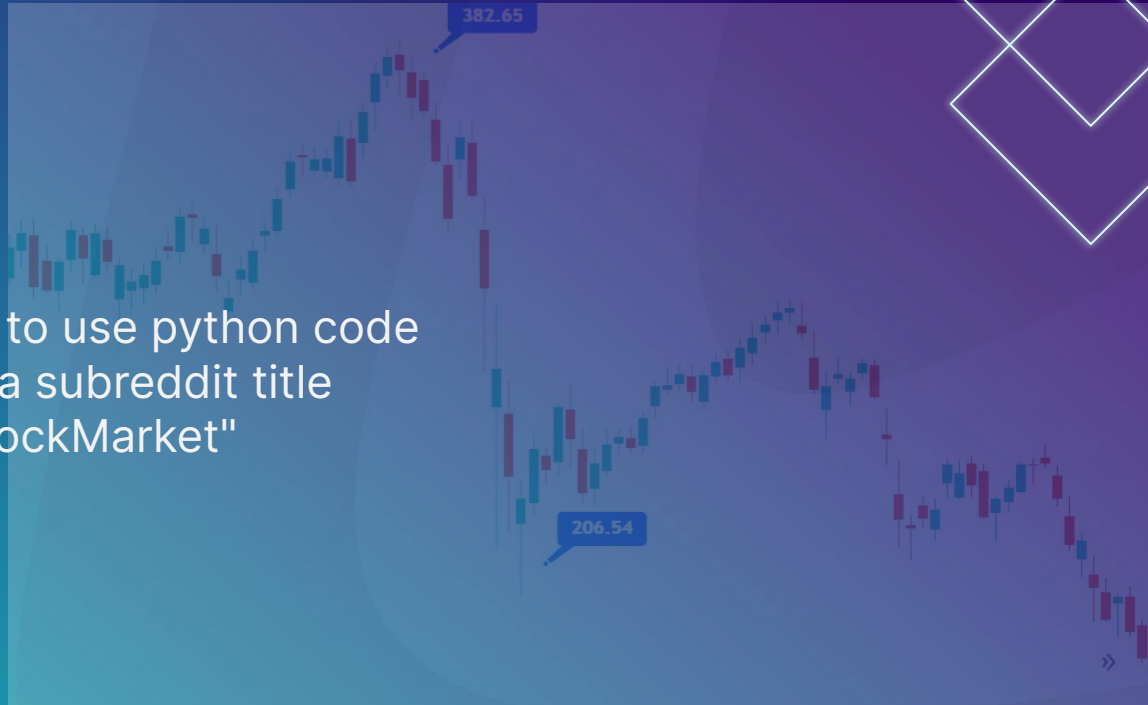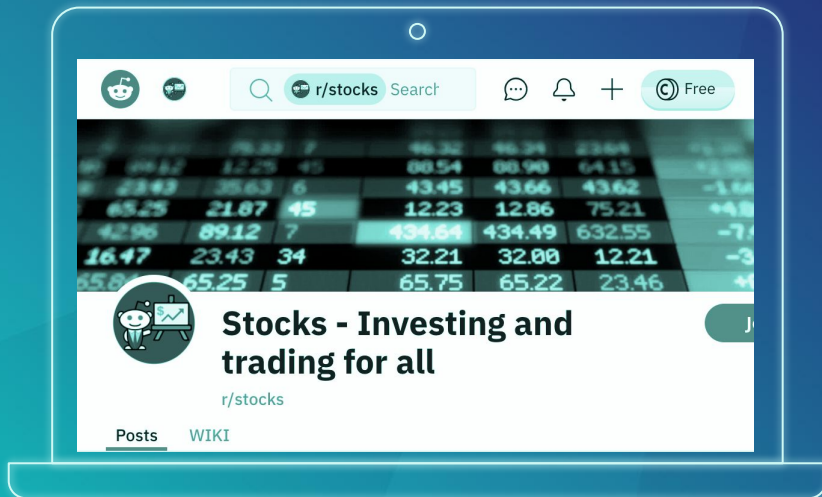# 'STOCKS' VS 'STOCKMARKET'

# PROBLEM STATEMENT

The Goal of this project is to use python code classifiers to determine if a subreddit title belongs to "stocks" or "StockMarket"
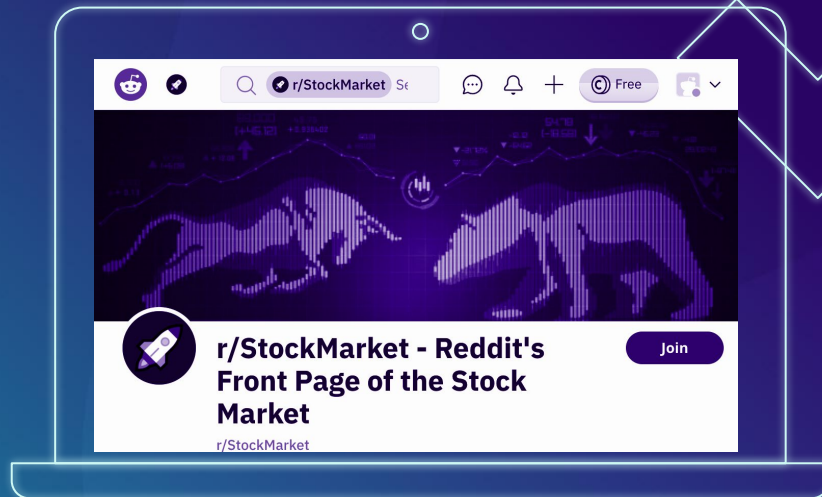
# SAMPLE OF PHRASES



stocks

'Will Beijing Supersede Hong Kong?'

StockMarket

'How to Value a Company with Multiples'

STOCKS

STOCKMARKET

# TRAINING
# MODEL

# PARAMETERS

## DATA

15130 data points
Evenly from both
subreddits

## TOKENS

CountVectorizer,
TfidfVectorizer
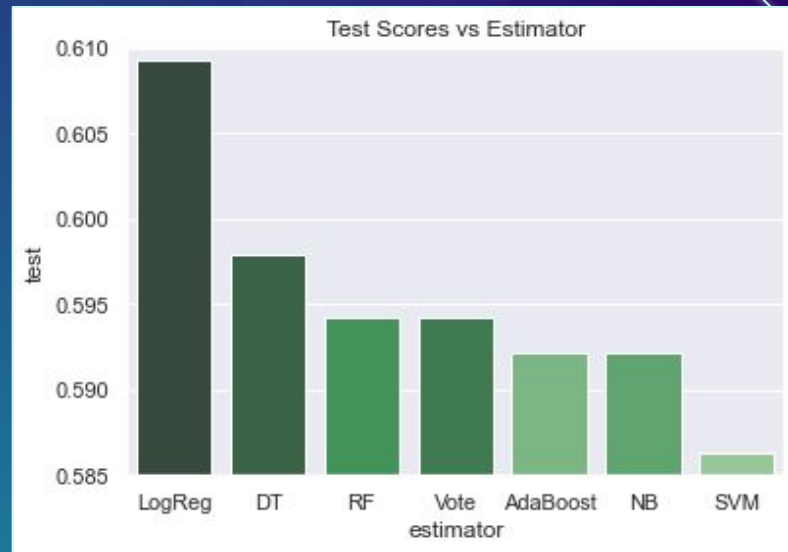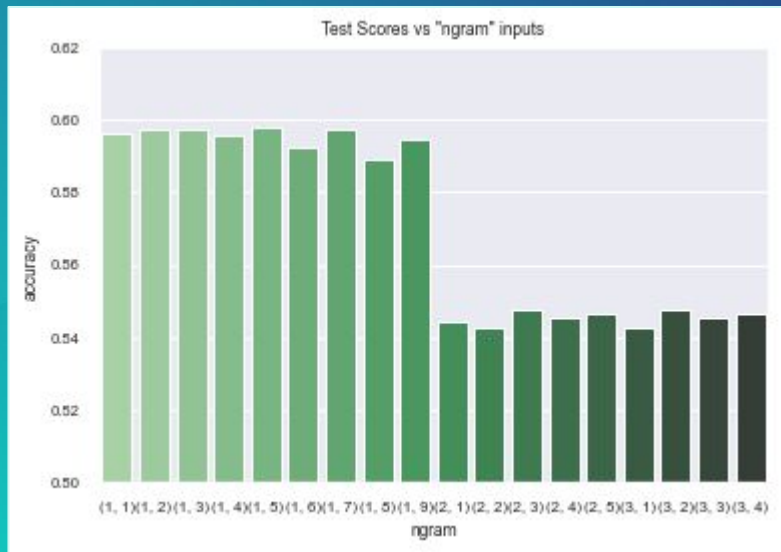
## ESTIMATOR

NB, RF Adaboost,
VotingClassifier, SVM,
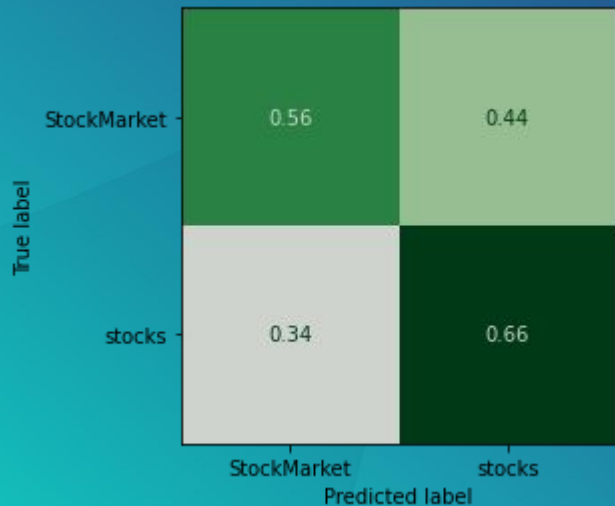LogReg

# NGRAM AND ESTIMATORS



Test Scores vs "ngram" inputs



Test Scores vs Estimator

GRADING MODEL

# LOGISTIC REGRESSION

Confusion matrix, normalized "True label"



## SCORES
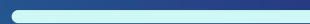
Accuracy **61%**

Sensitivity **66%**

Specificity **56%**

# LOGISTIC REGRESSION - COEFFICIENTS

| | TOP 3 COEFFICIENTS | |
|---|---|---|
| **FEATURE** | **COEFFICIENT** | **ODDS** |
| **STOCKS** | 0.57 | 1.76 |
| **QUESTION** | 0.56 | **1.75** |
| **ADVICE** | 0.53 | 1.70 |

# LOGISTIC REGRESSION - COEFFICIENTS

## Counts of titles with word 'question'

■ stocks  ■ StockMarket

Based on data there is 2.15 times as many `stocks` titles with the word "question"

Every word "question" in a subreddit title is 1.75 times as likely to be considered from `stocks` subreddit

# CONCLUSION

## MODEL

Logistic Regression is the best trialled model

## ACCURACY

Model can be useful to up to 61%

## RANKING

model coefficients understandable

## SMART

Is the model better than human classification?

# THANKS!

Do you have any questions?

changjulian17@gmail.com
github.com/changjulian17

/julian-chang/