# Reddit Financial Channels

Adam Duke
Whitney Schreiber
**{Data Mining Autumn 2021}**

# Agenda

# Executive Summary

*"At the start of the year, traders added more than $150 billion to the market cap of GameStop, AMC, and 48 other businesses" –Bloomberg*

In early 2021, select NASDAQ and NYSE listed securities that had mainly underperformed during the COVID-19 pandemic experienced suddenly radical and unstable price movements. This phenomenon was largely driven by simultaneous short position liquidation of institutional investors and mass opposing, collective retail trading activity. Popular speculation suggests observed market behavior was a consequence of observable posts, primarily those on specific Reddit channels. **We examine whether latent market volatility may be associated with text and other metadata of Reddit posts for known meme stocks, in particular GME.**

# Reddit Financial Channel Project Decision

## GOAL

Identify signals of abnormally high volatility in future stock price of a given "meme" stock using Reddit posts.

## VALUE

Advance indication of volatility can afford reduced value at risk or provide opportunities for advantageous options trading strategies.
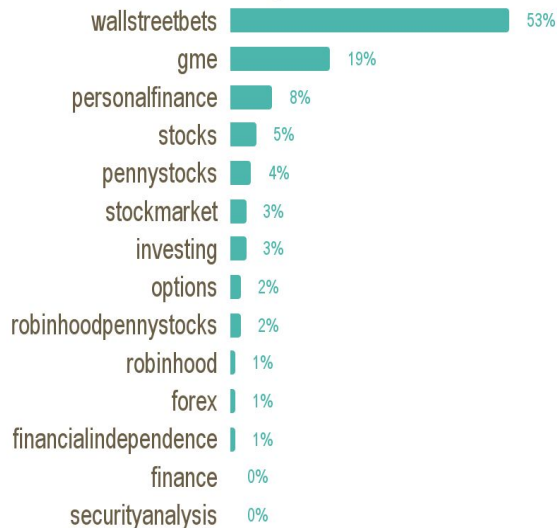
Reddit Financial Channel

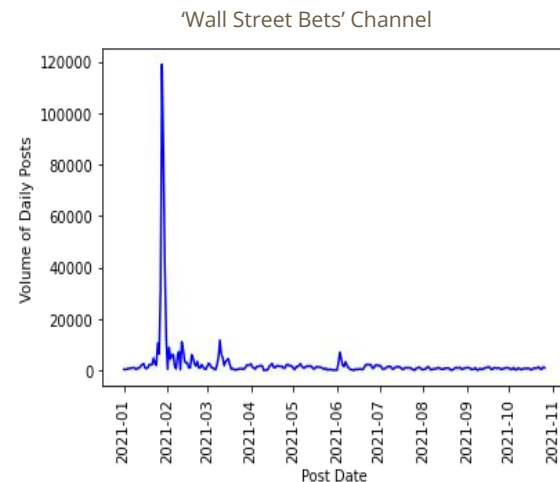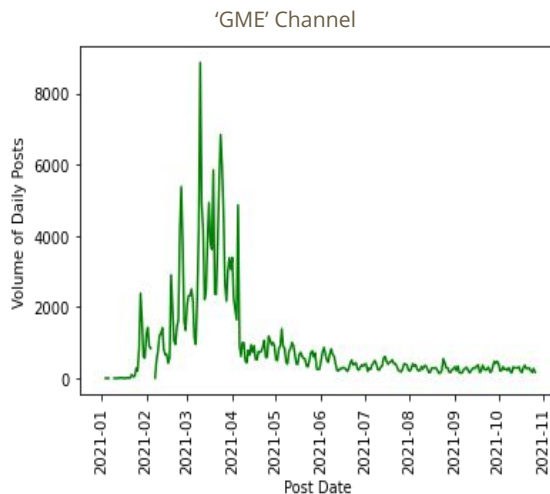# Data Exploration and Feature Engineering

# Reddit Financial Channel Raw Data

- **14** Channels
- **1,377,932** records (posts)

- Date range:
  **1/1/2021 - 11/21/2021**

### Percent of Total Posts by Channel

| Channel | Percent |
|---|---|
| wallstreetbets | 53% |
| gme | 19% |
| personalfinance | 8% |
| stocks | 5% |
| pennystocks | 4% |
| stockmarket | 3% |
| investing | 3% |
| options | 2% |
| robinhoodpennystocks | 2% |
| robinhood | 1% |
| forex | 1% |
| financialindependence | 1% |
| finance | 0% |
| securityanalysis | 0% |

**Post Volume**

'GME' Channel

'Wall Street Bets' Channel

# Data Profile

**Number of gilded awards**
`gilded`
(integer)

**Channel**
`channel`
(string)

**Author Name**
`author`
(string)

**Post Creation Date**
`created`
(datetime)

**Total awards**
`total_awards_received`
(integer)

**Score**
*(number of upvotes)*
`score`
(numeric)

**Link Flairs**
`link_flair_text`
(string)

44.2k

🔘 r/GME · Posted by u/kbme  Simple Lurking Ape  8 months ago
⚡🐵🚀🐵4 💖🍔🐂🐢😱💫23 🐢24 🔥🐵Ⓢ21 💕🐻33 🍵🐵2🎁2🐵⚠6 🐸

**Post Title**
`title`
(string)

If you lurk r/GME, and don't post anything ever, but you own GME, it is very possible that you are literally reading this title.

Discussion

I don't know. Please don't comment. The more of you the better.

Ken has his market manipulation tactics; we have invisible ape technologies. (Please see attached images.)

**Number of Comments**
`num_comments`
(integer)

💬 2.4k Comments      ↗ Share      🔖 Save      👁 Hide      🏳 Report

88% Upvoted

**Upvote Ratio**
`upvote_ratio`
(numeric)

**Post was removed**
`removed`
(binary)

**Archived**
`archived`
(binary)

📥 **This thread is archived**
New comments cannot be posted and votes cannot be cast

**Post was deleted**
`deleted`
(binary)

**Post was pinned**
`pinned`
(binary)

**Body text**
`self_text`
(string)

**Post was locked**
`locked`
(binary)

**Number of Crossposts**
`num_crossposts`
(integer)

**Image post thumbnail**
`thumbnail`
(string)

**Post's short url**
`shortlink`
(string)

Post is a text
`is_self`
(binary)

Post is a video
`is_video`
(binary)

Post set as original content
`is_original_content`
(binary)

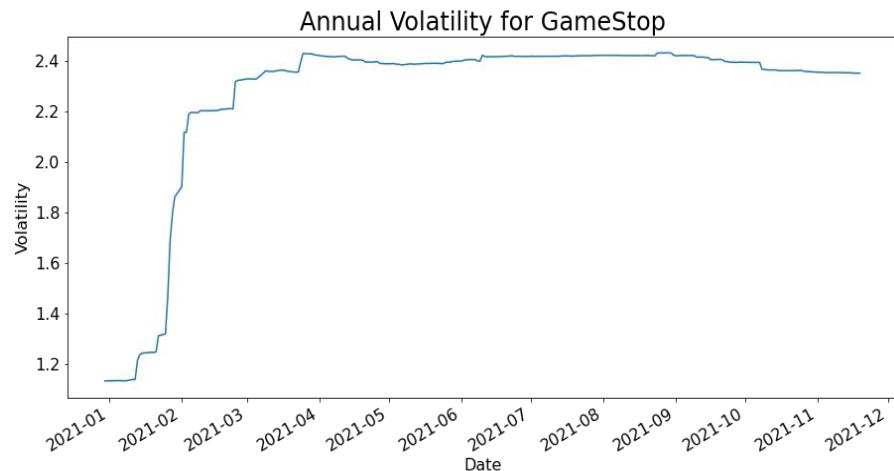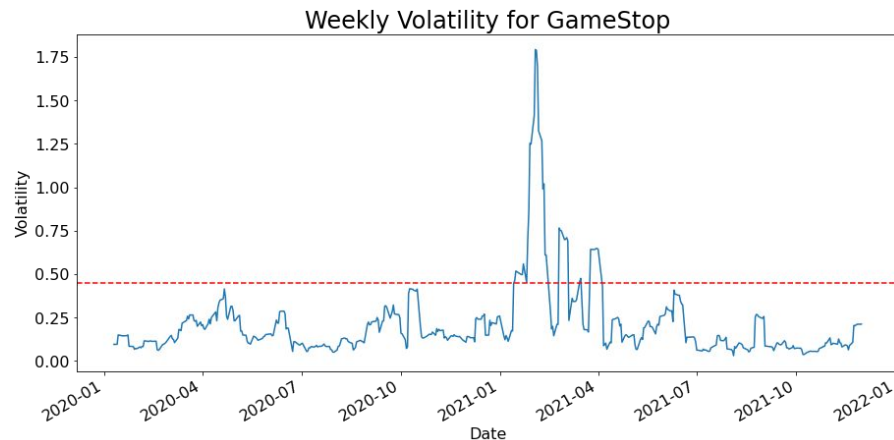# Data Profile

Retrieve stock market data using **yfinance**

- Ticker Symbols: **GME**
- Date range: **1/1/2021 - 11/21/2021**
- Metrics: Open, High, Low, Close, Adj Close, Volume
- Calculate **volatility** using Adjusted Closing price

```
return = log(closing_price_tomorrow / closing_price_today)
```

```
volatility = 7_day_rolling_std(return) * sqrt(trading days)
```

**Target Variable**

- Define **Mania**, a binary variable describing volatility
  - **True:** volatility >= threshold (0.45)
  - **False:** volatility < threshold (0.45)



Weekly Volatility for GameStop



Annual Volatility for GameStop

# Methodology



**Cleaning**

Rake and Combine
`df.alltext`

**Feature Engineering (Text Mining)**

TF-IDF · CountVectorizer

Truncated SVD · Truncated SVD

K-means · Hierarchical · DBSCAN · K-means · Hierarchical · DBSCAN

**Nonlinear Modelling**

DT · RF · DT · RF · DT · RF · DT · RF · DT · RF · DT · RF

`link_flair_text`

# Feature Engineering

One-hot encode the top 10 most common link flairs

Posted by u/SpaceMillionaire 🚀🚀Buckle up🚀🚀
12 hours ago
🟢 😎 Ⓢ 2

r/GME Megathread for Monday -
December 06, 2021 🚀 Megathread 🚀

*Link Flair*

'GME' Reddit Channel
## Percent of Posts with Link Flair

| 💎🙌 | 14% |
| Discussion | 10% |
| 💎🙌 | 7% |
| Shitpost | 6% |
| 🐵 Discussion 💬 | 6% |
| 😂 Memes 🐕 | 5% |
| ☁ Fluff 💪 | 5% |
| Fluff | 5% |
| Hedge Fund Tears | 4% |
| Memes ⚡ | 4% |

| | link_flair_text | | | link_flair_text_DD | link_flair_text_ | link_flair_text_Discussion | link_flair_text_Shitpost | link_flair_text_Memes | ... |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Discussion | | 0 | False | False | True | False | False | ... |
| 1 | Shitpost | | 1 | False | False | False | True | False | ... |
| 2 | Discussion | | 2 | False | False | True | False | False | ... |
| 3 | Fluff | | 3 | False | False | False | False | False | ... |
| 4 | God Tier DD | | 4 | False | False | False | False | False | ... |

# Feature Engineering
## Generate Basket of Key Words

Parallel **rake** 'title' and 'selftext' fields and combine the keywords

| Post title `title` | Post body text content `selftext` | Basket of words generated by rake `alltext` |
|---|---|---|
| GME is FINALLY going to the moon, this technical analysis looks very nice 🚀🚀🚀 | After some downwards movement, I think everybody needs some good news, and here they are. We are seeing some very strong technical indications showing that we are in fact on the verge of breaking up once again, the target is $26! :) [https://youtu.be/JlwXg5-H7cg] (https://youtu.be/JlwXg5-H7cg) | gme finally going moon technical analysis looks nice 🚀🚀🚀 downwards movement think everybody needs good news seeing strong technical indications showing fact verge breaking target 26 :) https :// youtu jlwxg5 h7cg ]( |

# Feature Engineering

## Transform Unstructured Text

**alltext**

gme moon
🚀🚀

need see
gme 🚀🚀
🚀🚀🚀
watching
took
position
rig...

short
squeeze
incoming
🚀🚀🚀
🚀🚀

convinced
💰 gme
extreme
pump
coming guy
explai...

already
know must
brothers
sisters
submit
comp...

## TF-IDF

Transform text to an array by comparing the word frequency in a doc and the number of docs with the word

**Output:** sparse array
```
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
```
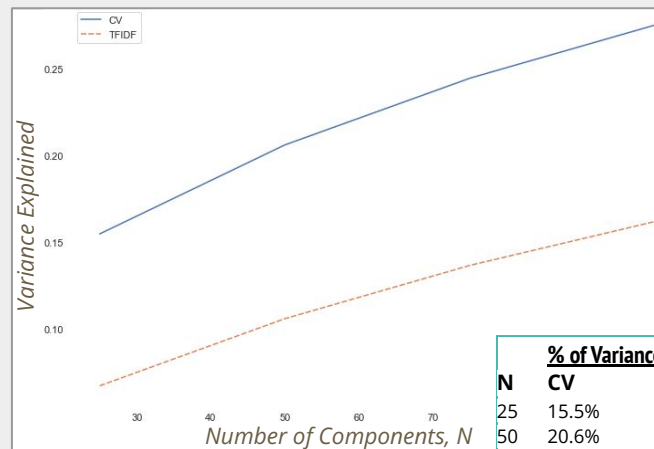
## CountVectorizer

Produce array of the frequency of each unique word that occurs in the entire text

**Output:** sparse array
```
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
```

## Truncated SVD

Linear dimensionality reduction of sparse matrix



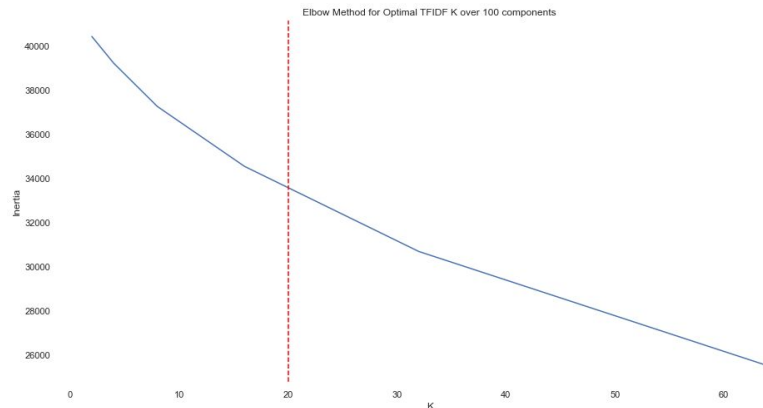| | % of Variance Explained | |
|---|---|---|
| **N** | **CV** | **TFIDF** |
| 25 | 15.5% | 6.7% |
| 50 | 20.6% | 10.6% |
| 75 | 24.5% | 13.7% |
| 100 | 27.5% | 16.2% |

# Feature Engineering

## Cluster Transformed Text

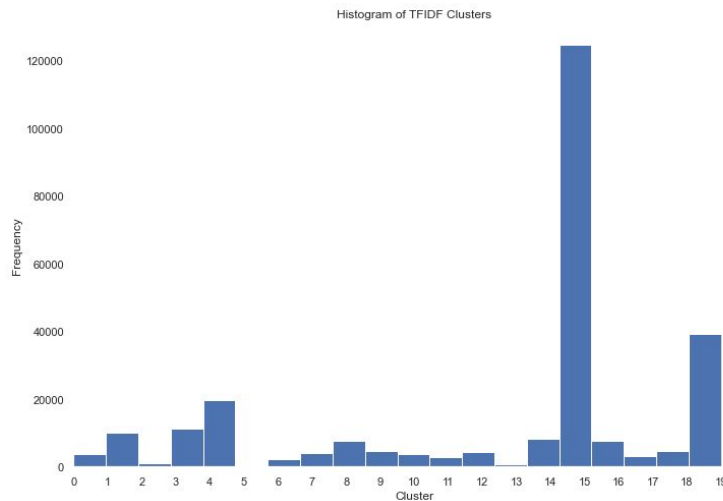### k-Means of count vectorization

Use elbow method to determine the appropriate number of clusters, k



**Use k = 20**

Resulting cluster sizes *(k = 20)*



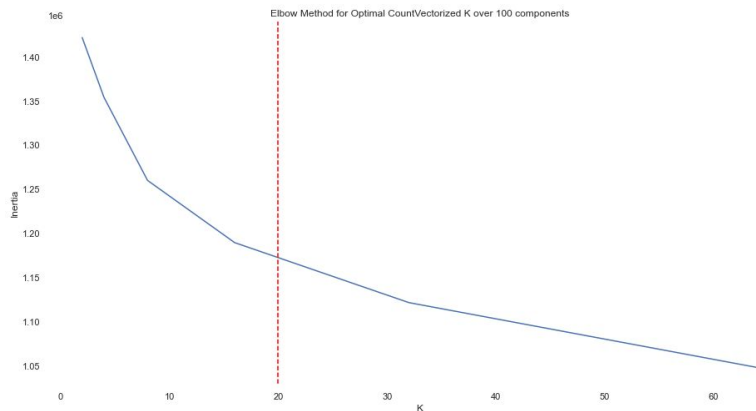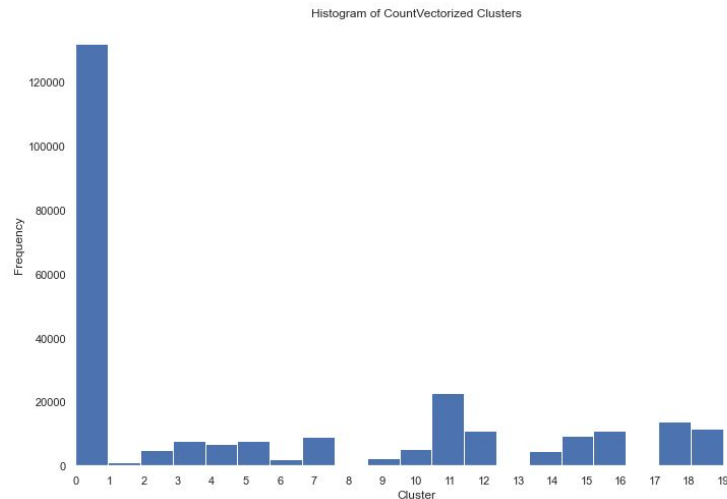| Rake & Combine Text | CountVectorizer | Truncated SVD | K-Means |

# Feature Engineering

## Cluster Transformed Text

### k-Means of TF-IDF vectorization

Use elbow method to determine the appropriate number of clusters, k



**Use k = 20**

Resulting cluster sizes *(k = 20)*



| Rake & Combine Text | TF-IDF | Truncated SVD | K-Means |

# Modeling Data Frame

# Modelling

# Decision Tree Classifier

*from CountVectorizer*



Confusion matrix



## Feature Importance

| Feature | Importance |
| --- | --- |
| score | 25.5% |
| upvote_ratio | 21.6% |
| num_comments | 18.5% |
| link_flair_text_Memes | 6.1% |
| link_flair_text_Fluff | 5.1% |
| link_flair_text_Discussion | 3.1% |
| total_awards_received | 2.7% |
| cluster_0 | 2.7% |
| link_flair_text_DD | 1.6% |
| link_flair_text_Shitpost | 1.6% |

| | precision | recall | f1-score | support |
| --- | --- | --- | --- | --- |
| False | 0.69 | 0.97 | 0.81 | 35877 |
| True | 0.57 | 0.08 | 0.14 | 16796 |
| | | | | |
| accuracy | | | 0.69 | 52673 |
| macro avg | 0.63 | 0.53 | 0.47 | 52673 |
| weighted avg | 0.65 | 0.69 | 0.60 | 52673 |

Rake & Combine Text    CountVectorizer    Truncated SVD    K-Means    Decision Tree

# Decision Tree Classifier

*from TF-IDF*



Confusion matrix

## Feature Importance

| Feature | Importance |
|---|---|
| link_flair_text_Memes | 25.2% |
| link_flair_text_Fluff | 23.8% |
| link_flair_text_Discussion | 13.6% |
| num_comments | 12.9% |
| score | 10.5% |
| cluster_11 | 9.7% |
| upvote_ratio | 2.7% |
| cluster_15 | 1.6% |
| cluster_8 | 0.0% |
| cluster_7 | 0.0% |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| False | 0.68 | 1.00 | 0.81 | 35877 |
| True | 0.63 | 0.02 | 0.03 | 16796 |
| accuracy |  |  | 0.68 | 52673 |
| macro avg | 0.65 | 0.51 | 0.42 | 52673 |
| weighted avg | 0.67 | 0.68 | 0.56 | 52673 |

Rake & Combine Text   TF-IDF   Truncated SVD   K-Means   Decision Tree

# Random Forest Classifier

*from CountVectorizer*

## Feature Importance



| | score | 25.5% |
| upvote_ratio | 21.6% |
| num_comments | 18.5% |
| link_flair_text_Memes | 6.1% |
| link_flair_text_Fluff | 5.1% |
| link_flair_text_Discussion | 3.1% |
| total_awards_received | 2.7% |
| cluster_0 | 2.7% |
| link_flair_text_DD | 1.6% |
| link_flair_text_Shitpost | 1.6% |

### Confusion matrix



|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| False | 0.69 | 0.97 | 0.81 | 35877 |
| True | 0.57 | 0.08 | 0.14 | 16796 |
| accuracy |  |  | 0.69 | 52673 |
| macro avg | 0.63 | 0.53 | 0.47 | 52673 |
| weighted avg | 0.65 | 0.69 | 0.60 | 52673 |

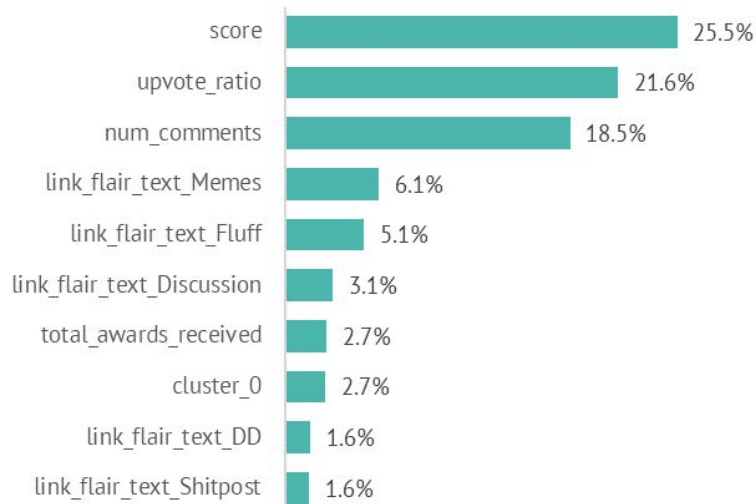Rake & Combine Text    CountVectorizer    Truncated SVD    K-Means    Random Forest
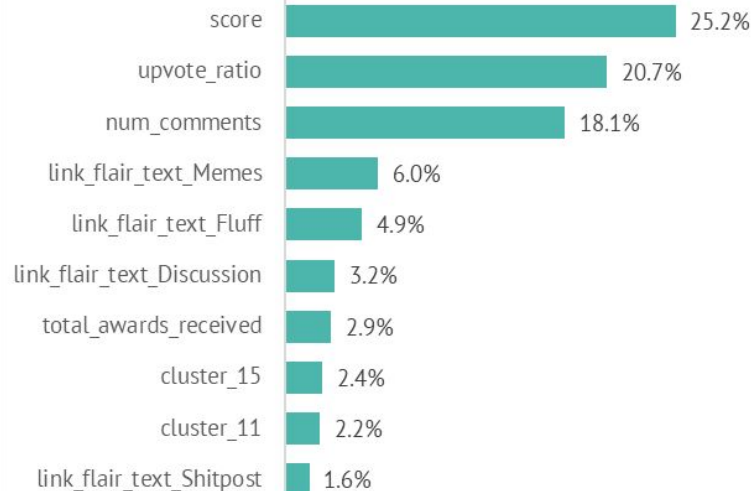
# Random Forest Classifier

*from TF-IDF*



Confusion matrix



Feature Importance

| score | 25.2% |
| upvote_ratio | 20.7% |
| num_comments | 18.1% |
| link_flair_text_Memes | 6.0% |
| link_flair_text_Fluff | 4.9% |
| link_flair_text_Discussion | 3.2% |
| total_awards_received | 2.9% |
| cluster_15 | 2.4% |
| cluster_11 | 2.2% |
| link_flair_text_Shitpost | 1.6% |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| False | 0.69 | 0.98 | 0.81 | 35877 |
| True | 0.58 | 0.07 | 0.13 | 16796 |
| accuracy |  |  | 0.69 | 52673 |
| macro avg | 0.64 | 0.52 | 0.47 | 52673 |
| weighted avg | 0.66 | 0.69 | 0.59 | 52673 |

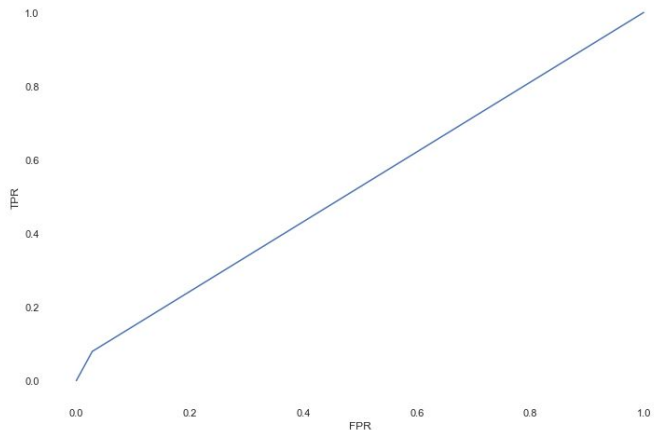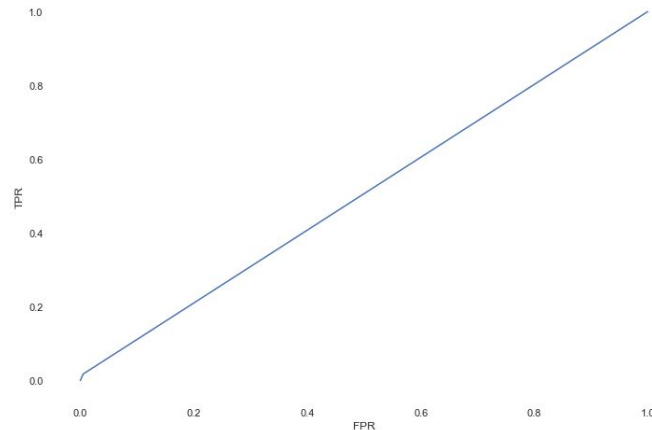Rake & Combine Text | TF-IDF | Truncated SVD | K-Means | Random Forest

# Results

Given a nearly straight ROC curve and AUC curve near 0.5, neither model (decision tree nor random forest) trained on count or tf-idf vectorized text features is able to strongly separate true samples from false ones. Despite a relatively high accuracy, the models perform hardly better than guessing along proportion of true class of target variable.
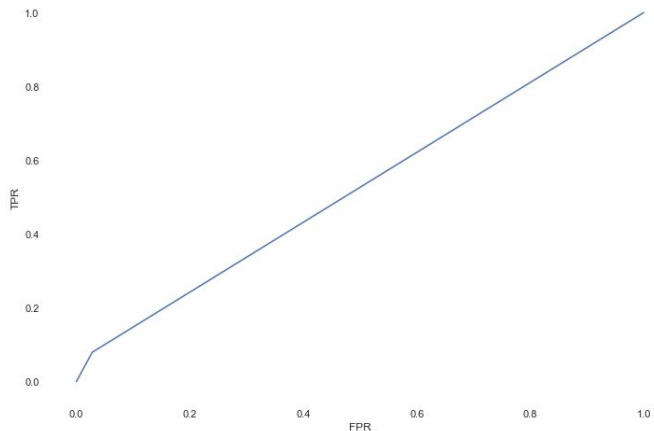


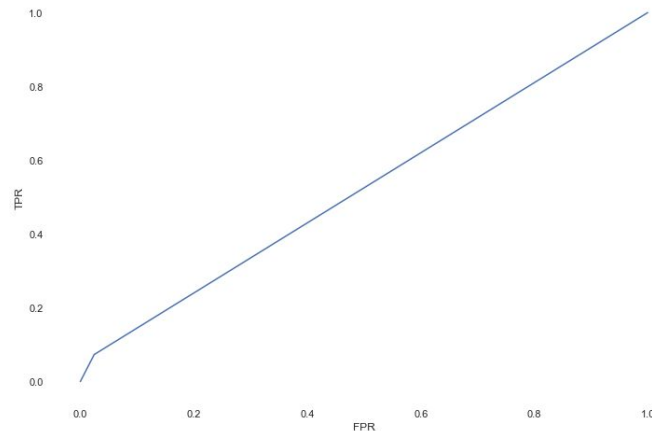CountVectorized Decision Tree ROC curve
AUC = 0.5257

TF-IDF Decision Tree ROC curve
AUC = 0.5064

CountVectorized Random Forest ROC curve
AUC = 0.5257

TF-IDF Random Forest ROC curve
AUC = 0.5244

# Takeaways

## Key Findings

- **Text mined features were not the most significant factors**
- **Models marginally improved upon guessing forward volatility**

## Challenges

- Large data
  - Over 1.3M documents
- Imbalanced data
  - Low proportion of volatile days
- Noisy data
  - Many irrelevant or spam documents
- Hardware & time constraints
  - DBSCAN and hierarchical clustering failed or crashed kernel

## Next Steps

- Enhanced data cleaning
- Various target variables
- Other securities (i.e. AMC)
- Other models (CNN, etc.)
- More granular price data
- Scrape comments text
- Classify popular posts and then pipe only that classification into a volatility prediction model

# Thank You