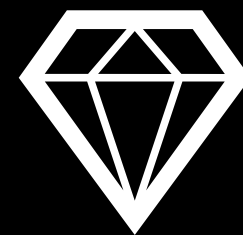




# An Analysis of r/Wallstreetbets's effect on stock prices

Capstone Project - Ryan



# Business Problem Statement



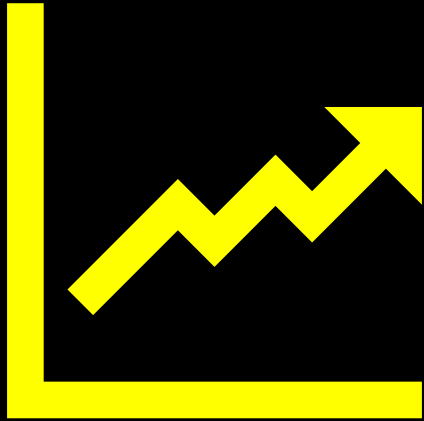
▶ Johnny   is a MD of a foreign bank



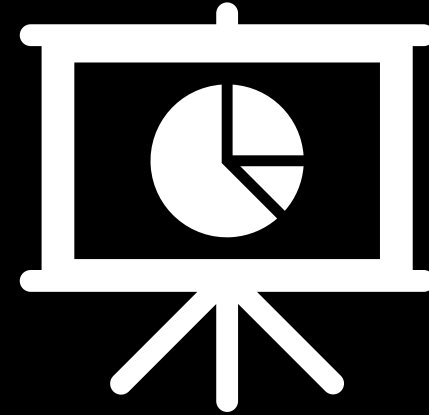
- ▶ Seen WSB influence on gme
- ▶ Wants to use WSB's activity to give the bank a trading edge



# Data Problem Statements

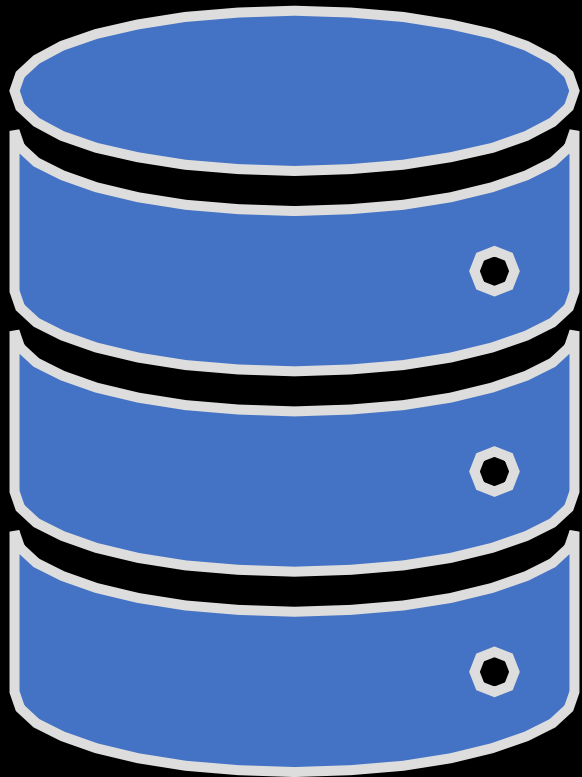


Is there a relationship between the sentiment or number of comments on WSB and stock prices?



Can we use stock data and WSB comments data to predict the next day's price movement?

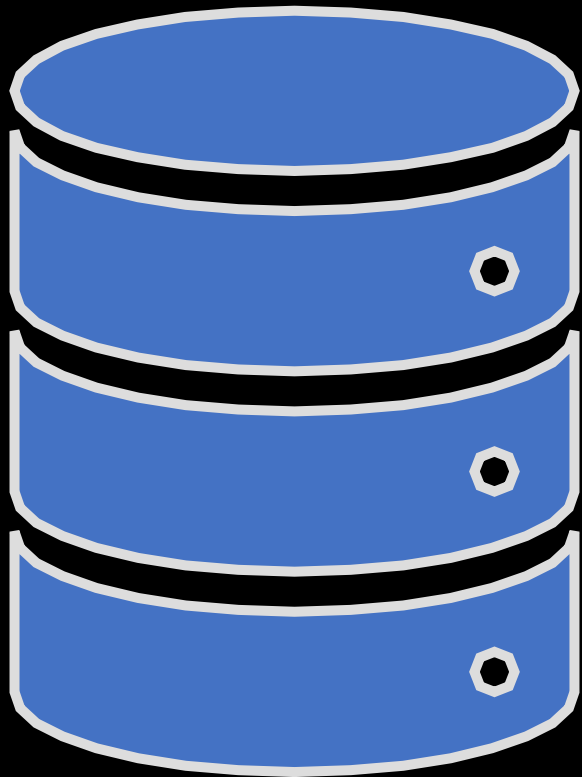
# Data sources



Yahoo Finance: GME, TSLA, PLTR, NOK, BB, SP500, AMC

Kaggle: WSB comments (Feb-2018 to Feb 2021)

# Data sources



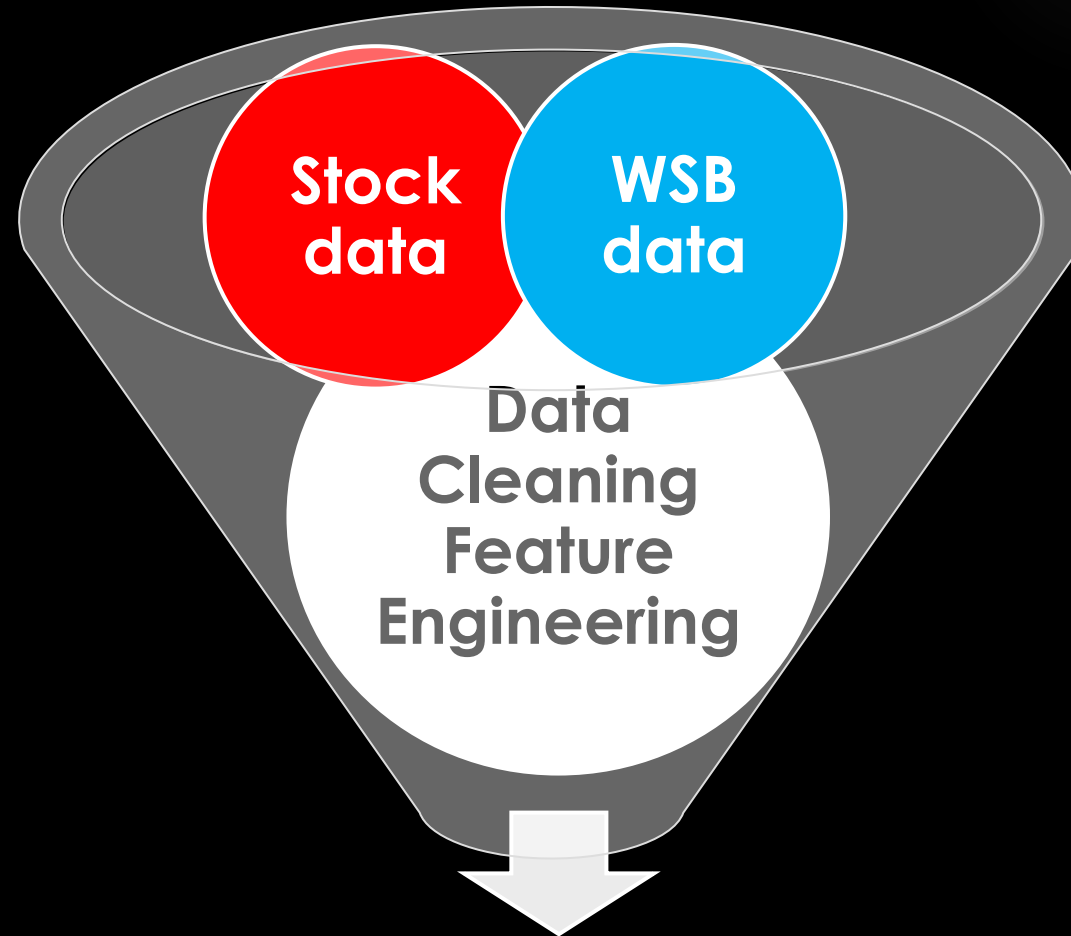
- ▶ Stock data columns:

- ▶ `Open, High, Low, Close, Volume`

- ▶ WSB comments columns:

- ▶ `created_utc, text`

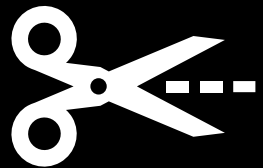
# Data preparation



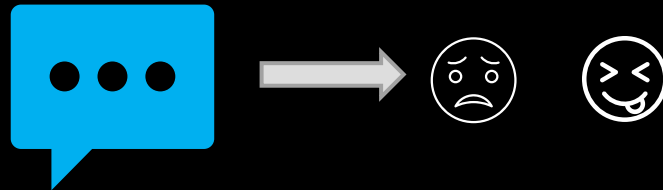
Final Dataframes

# Data Cleaning / Feature Engineering

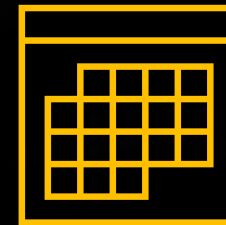
(WSB comments)



Remove NA  
values



Apply spacy's sentiment  
and subjectivity models.



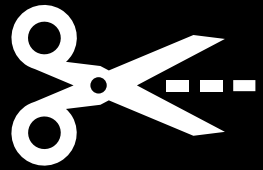
- Convert unix timestamp (**created\_utc**) to tz-aware datetime object
- Get **day\_of\_week** and month label from **created\_utc**

## Available Columns

- ▶ **created\_utc**
- ▶ **text**
- ▶ **sentiment**
- ▶ **Subjectivity**
- ▶ **month**
- ▶ **day\_of\_week**

# Data Cleaning / Feature Engineering

(Stock data)



Remove all data  
prior to 16<sup>th</sup> Feb  
2018



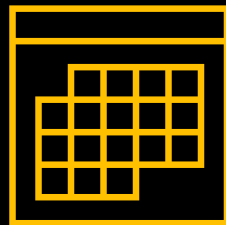
Compared **Open**  
and **Close** values  
to get **pct\_change**  
column



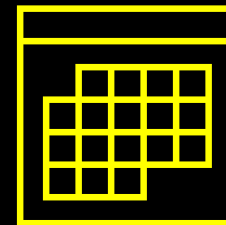
Compared Open and  
Close values to get  
bool column **up\_today**  
and **up\_tomorrow**



Compared **High**  
and **Low** values to  
get  
**pct\_volatility**  
column



Converted  
**Date** column to  
tz-aware  
datetime object



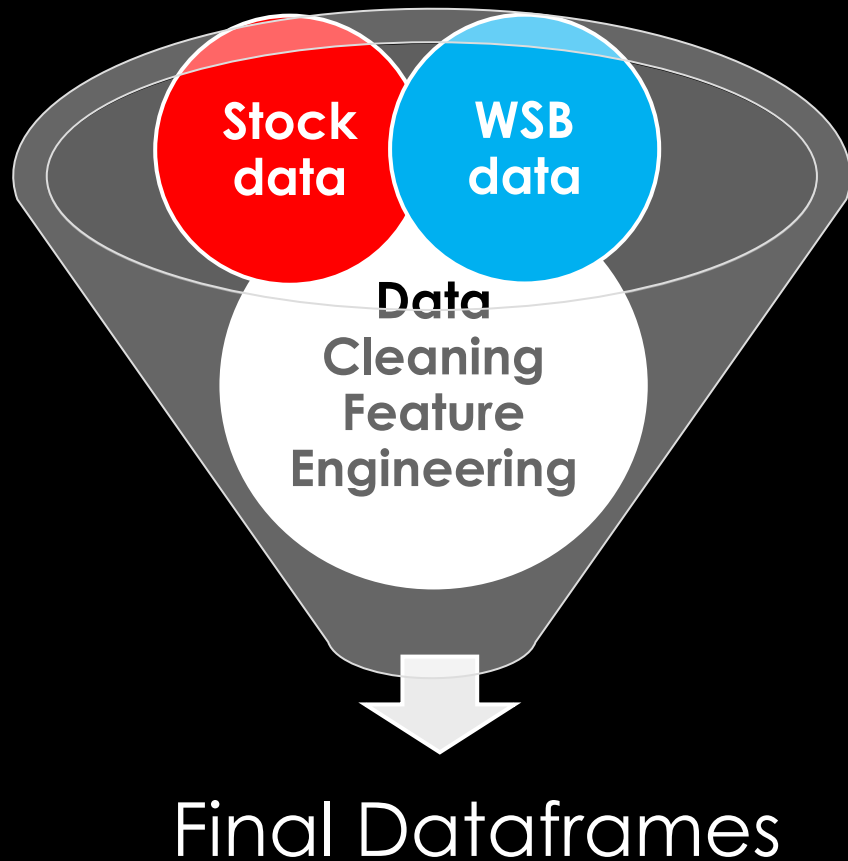
Get  
**day\_of\_week**  
and **month**  
label from **Date**

## Available Columns

- ▶ **Date**
- ▶ **Open**
- ▶ **High**
- ▶ **Close**
- ▶ **Low**
- ▶ **Adj Close**
- ▶ **Volume**
- ▶ **pct\_change**
- ▶ **up\_today**
- ▶ **up\_tomorrow**
- ▶ **pct\_volatility**
- ▶ **day\_of\_week**
- ▶ **month**

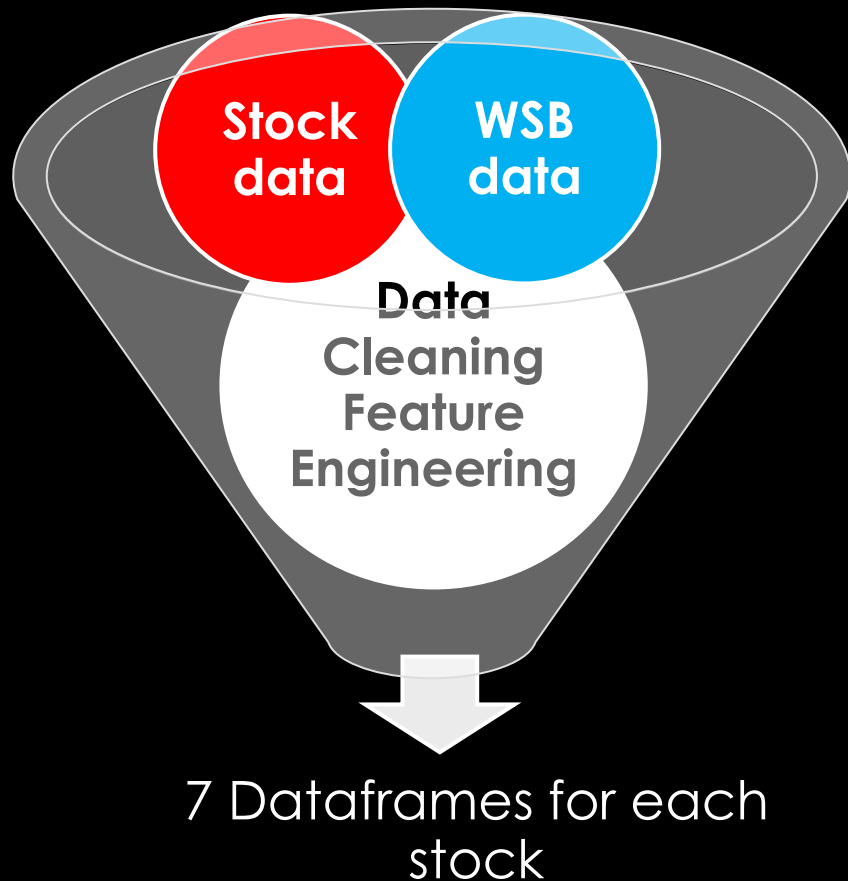


# Data preparation



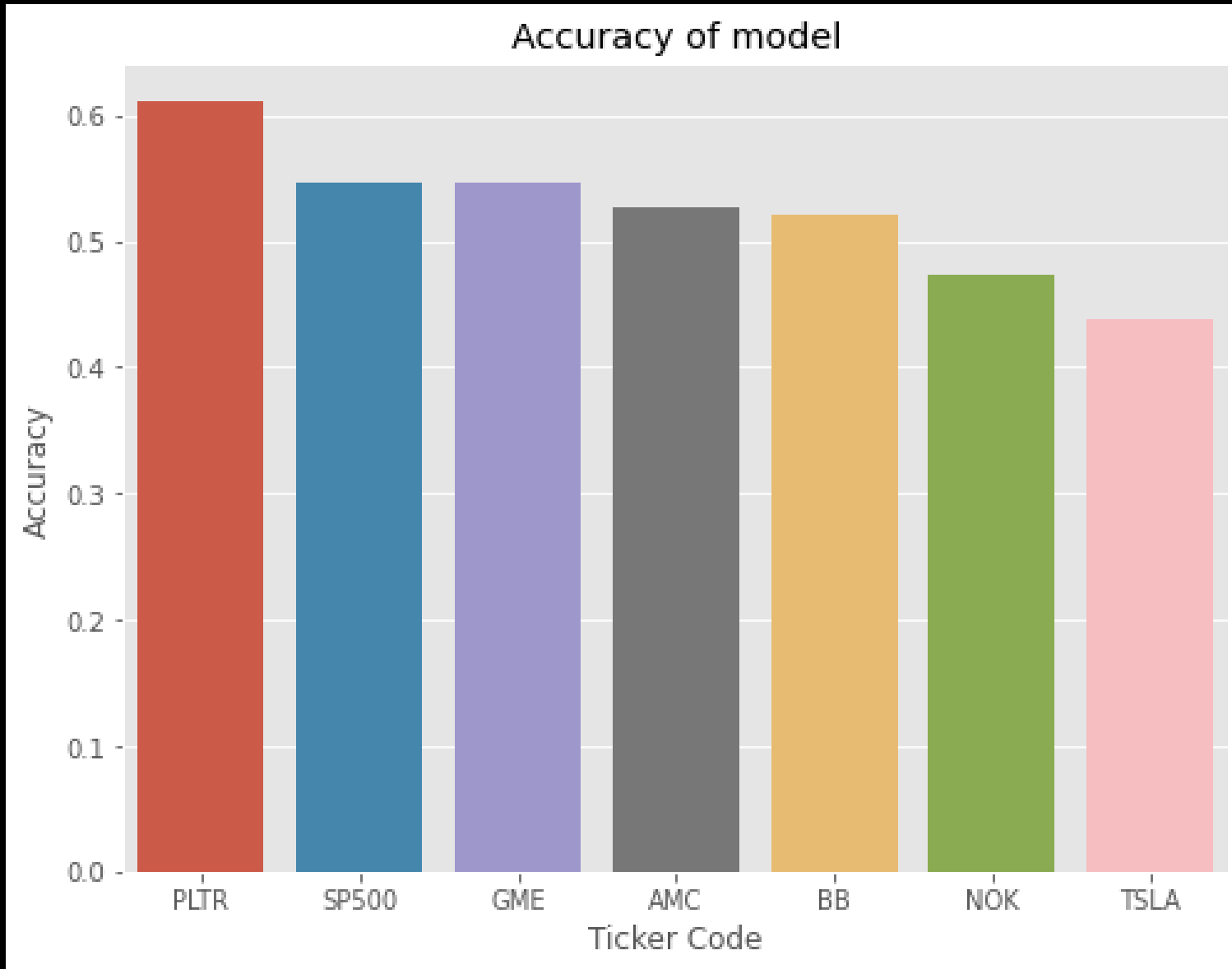
- ▶ Group wsb data by date and obtain the following columns:
  - ▶ Mean subjectivity
  - ▶ Mean sentiment
  - ▶ Count of comments
- ▶ Left join wsb data on stock data to obtain 7 different dataframes
- ▶ Target variable is **up\_tomorrow**

# Final Dataframe Columns



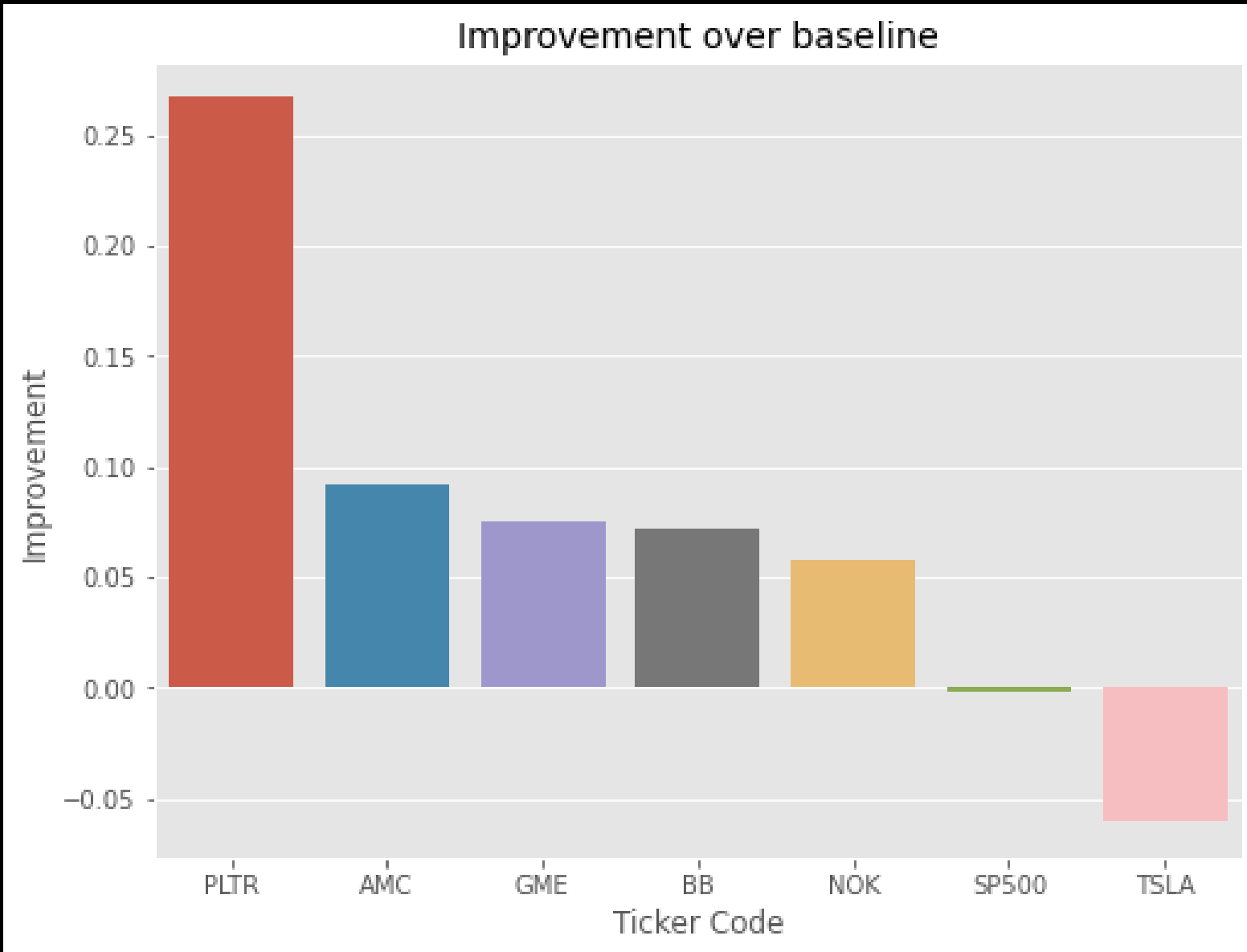
- ▶ Features:
  - ▶ Open : float64
  - ▶ High : float64
  - ▶ Low : float64
  - ▶ Close : float64
  - ▶ Volume : float64
  - ▶ up\_today : bool
  - ▶ pct\_change : float64
  - ▶ sentiment : float64
  - ▶ subjectivity : float64
  - ▶ comment\_count : int64
  - ▶ day\_of\_week : category
  - ▶ month : category
- ▶ Prediction column:
  - ▶ up\_tomorrow : bool

# Results



- ▶ Model used:
- ▶ XGBoost Classifier:
- ▶ 80% : 20%  
**train\_test\_split**
- ▶ Best accuracy performer is PLTR, over 60% accuracy
- ▶ Model predicted the right movement 60% of the time

# Results



- ▶ Best improved model over baseline is PLTR
- ▶ Baseline: predict stock is up every day

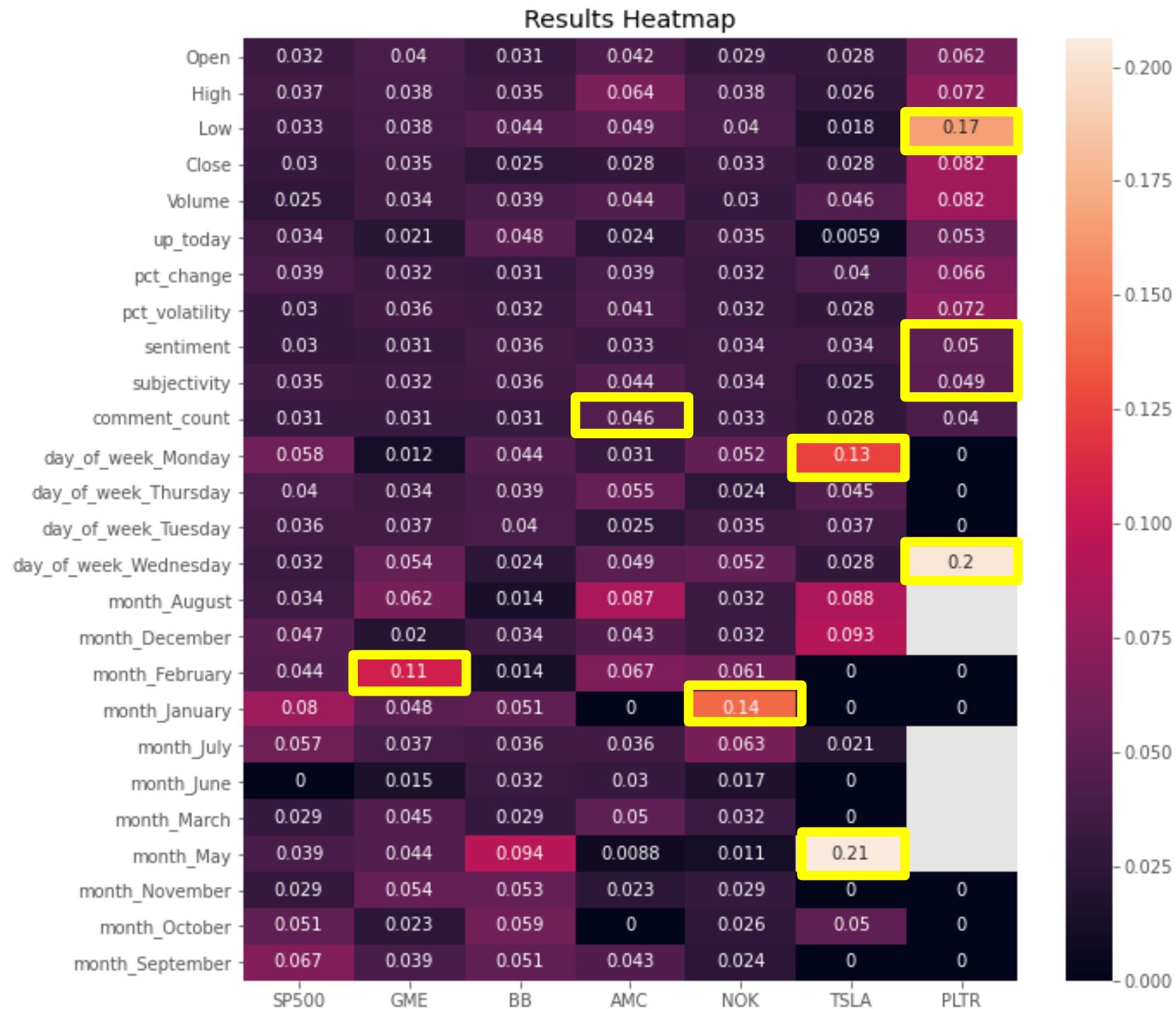
# Analysis of results

- ▶ PLTR IPO'd recently and stock data is only from 2<sup>nd</sup> Oct 2020 onwards
- ▶ For all other stocks, the analysis runs from Feb-2018 to Feb-2021
- ▶ Aug 2020 was the start of WSB's huge membership growth
- ▶ WSB's influence initially is low due to low membership count



# Feature Importances

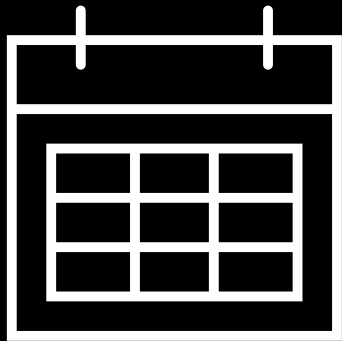
- ▶ PLTR – affected by **low**, and whether or not it is **wednesday**
- ▶ TSLA – affected by whether or not it is **monday** and **may**
- ▶ Most affected by WSB **sentiment** & **subjectivity**– PLTR
- ▶ Most affected by WSB **comment\_count** – AMC
- ▶ Most affected by the month of January – GME
- ▶ Most affected by the month of February



# Future work



- ▶ Train a BERT Model on labelled data for sentiment analysis
- ▶ (Bidirectional Encoder Representations from Transformers)



- ▶ Narrow the timeframe down to Aug-2020 onwards as this is the period of WSB's explosive growth
  - ▶ Downsides: lack of data

# Conclusion

- ▶ The results of PLTR show that WSB is starting to have influence on the stock market
  - ▶ The cult-like behavior in WSB is contagious
  - ▶ WSB encourages people to hold stocks through memes – something never done before
  - ▶ There is also the desire to get back at hedge funds for the 08 crisis
- ▶ The investing landscape has changed significantly since the Warren Buffet days
  - ▶ People were expecting a crash but no crash came
  - ▶ S&P broke 4000 during the middle of a pandemic
  - ▶ The economy and the stock market has decoupled