# PREDICTING STOCK PRICE CHANGES OF HEALTHCARE COMPANIES BASED ON NEWS HEADLINES

By Lee Shi Min
DSI-23 | An NLP project

**01**

**OBJECTIVES**

Background
Focus on Healthcare
Problem Statement

**02**

**METHODOLOGY**

Data

**03**

**RESULTS ANALYSIS**

Model Performance
Other Findings

**04**

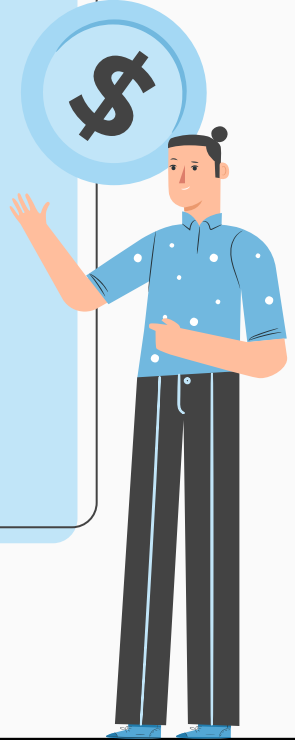**CONCLUSIONS**

Learning Points
Recommendations for
Future Project

# PREDICTING STOCK PRICES IS HARD...

- Technical versus fundamental analysis
- Evidence of post-news drift
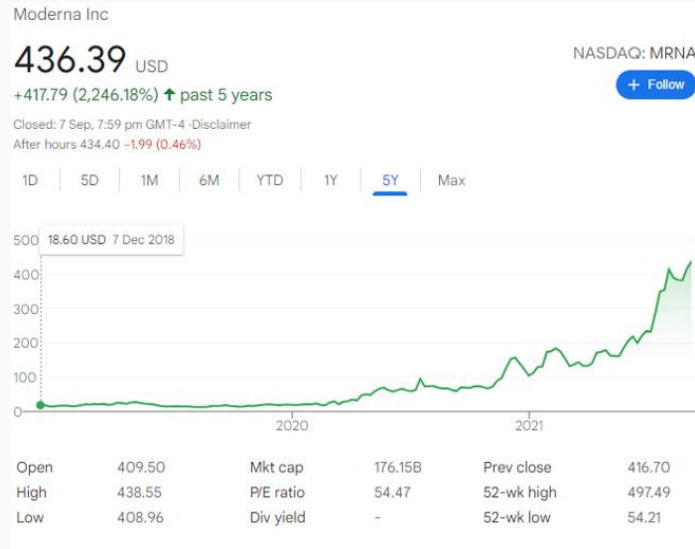- "Buy the rumor, sell the news"

Moderna Inc

436.39 USD

NASDAQ: MRNA

+417.79 (2,246.18%) ↑ past 5 years

+ Follow

Closed: 7 Sep, 7:59 pm GMT-4 -Disclaimer
After hours 434.40 −1.99 (0.46%)

| 1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max |

18.60 USD 7 Dec 2018

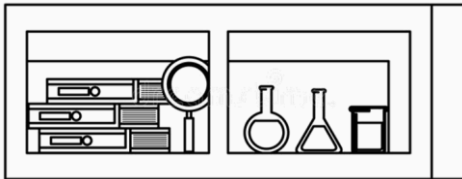| | | | | | |
|---|---|---|---|---|---|
| Open | 409.50 | Mkt cap | 176.15B | Prev close | 416.70 |
| High | 438.55 | P/E ratio | 54.47 | 52-wk high | 497.49 |
| Low | 408.96 | Div yield | - | 52-wk low | 54.21 |

# Tremendous Growth & Potential

- Development of vaccines and therapeutics highly valued by society (Moderna stock increased by > 2,000% in the last 5 years)
- Outperforming S&P500 in some instances
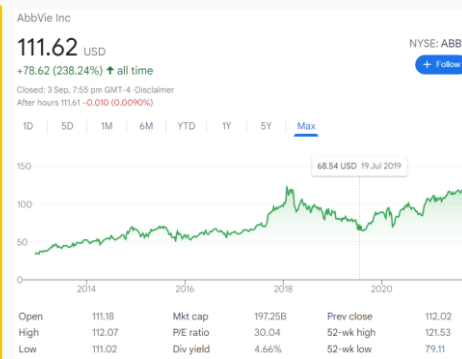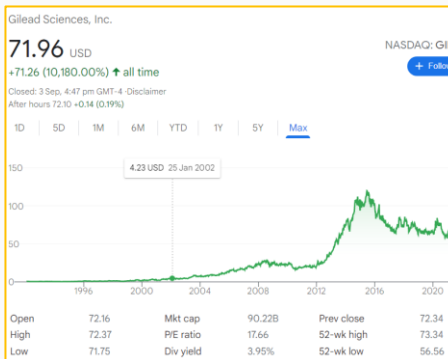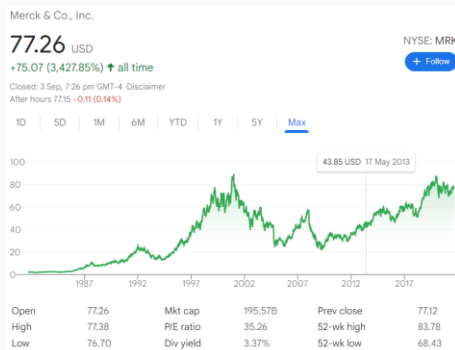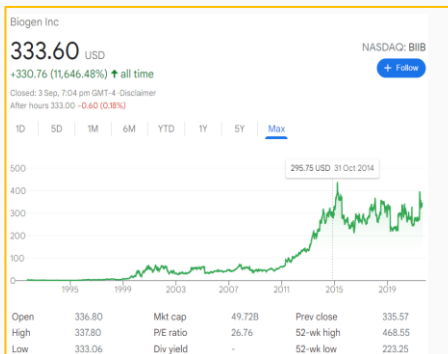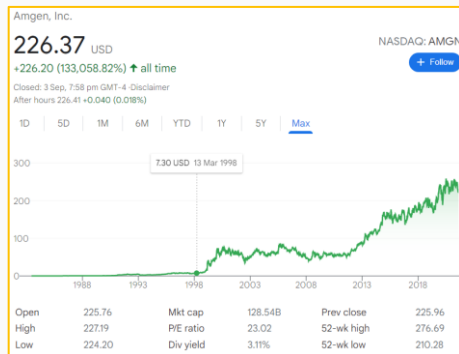
# Knowledge Gap

- Finance knowledge is not enough
- Understanding of drug development process, commercial implications of scientific technologies and regulatory frameworks needed

**FOCUS ON HEALTHCARE COMPANIES**

6 out of 8 companies in list outperformed S&P 500 (~3,616% all time increase)

# PROJECT FOCUS

Using NLP/Data Science to help layperson or amateur investors learn about healthcare investing

| 1. NLP Project to Predict Stock Price Jumps given News Headlines | |
|---|---|
| 1. Use past data for prediction | **ROC AUC** |
| 2. Accurate identification of price jumps on news | **Accuracy** |
| 3. Precise prediction of price jumps | **Precision** |

| 2. Deep dive into the headlines that are indicative of price jumps |
|---|
| Study feature importances |

# DATA



- **Date**
- News headlines

- **Date**
- Adjusted Closing Price (Dividends, Stock Splits, S&P500)
- Trading Volume

| Feature | Description |
|---------|-------------|
| Date | Date during which the stock and news information was pulled. |
| adjusted_abs | **Absolute percentage change** for the adjusted closing price of the stock with respect to the previous trading day. |
| adj_direction | Denote the directionality of the stock price changes. |
| cleaned | Daily news headlines from Reuters pertaining to the particular ticker concatenated into one single string. |
| headline_word | Word count of news headlines |
| dict_score | Score to denote the news sentiment for healthcare domain. Negative values are assigned to bad news (e.g. failed trial) and positive values are assigned to good news (e.g. fda approval). |
| target_var | Classification label: **1** denotes when stock prices changed by more than 1.31% (P75), **0** for otherwise. |

# WORKFLOW

**Data Cleaning**

- Removal of outliers
- Text cleaning – define new stopwords, removal of stopwords, lemmatization

**Feature Engg & EDA**

- Word Count for News Headlines
- Polarity Scores
- Vectorization – Count & Tf-idf

**Choice of Model**

- Using PyCaret to compare models
- Choose model based on precision, accuracy & AUC ROC scores

**Model Tuning & Evaluation**

- GridSearchCV
- Confusion Matrix, AUC ROC Curve

# TEXT CLEANING - AMENDING THE LIST OF STOPWORDS

```python
#add words that aren't in the NLTK stopwords list
new_stopwords = ['reuters', 'country', 'population', 'government', 'united', 'states', 'company', 'economy',
                 "said", 'say', 'inc', 'data', 'business', 'one', 'two', 'three', 'four', 'five',
                 'january', 'february', 'march', 'april', 'may', 'june', 'july', 'august', 'nyse',
                 'september', 'october', 'november', 'december', 'weekly', 'monthly', 'quarter', 'year',
                 'abbvie', 'pfizer', 'gilead', 'merck', 'eli', 'lilly', 'amgen', 'biogen', 'regeneron',
                 'bristol', 'myers', 'squibb', 'nasdaq', 'new', 'york', 'exchange', 'earning', 'price',
                 'medicine', 'healthcare', 'dow', 'jones', 'index', 'dji', 'investor', 'percent', 'market',
                 'drug', 'share', 'health', 'biotech', 'per', 'cent', 'co', 'sp']

new_stopwords_list = stop_words.union(new_stopwords)

#remove words that are in NLTK stopwords list
not_stopwords = {'no', 'not', 'up', 'further', 'above', 'down', 'under', 'over', 'through', 'off', 'below'}
final_stop_words = set([word for word in new_stopwords_list if word not in not_stopwords])
```
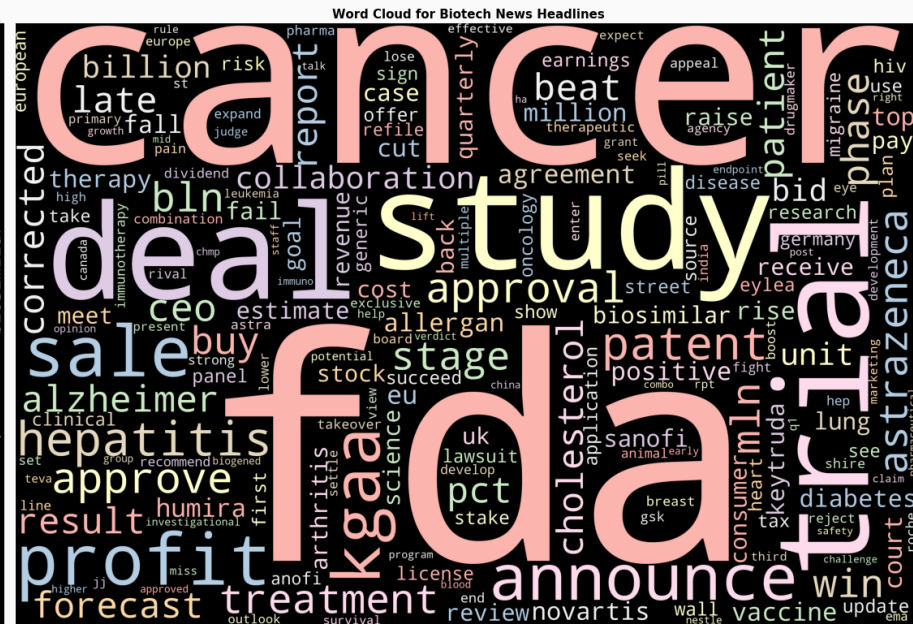
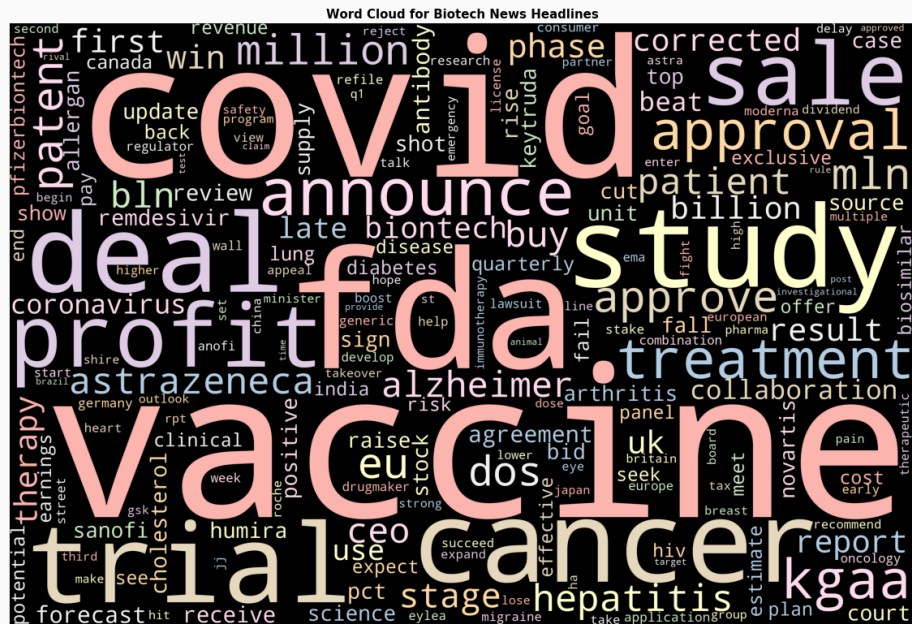# REMOVAL OF OUTLIERS BY YEAR (COVID ADJUSTMENT)

**BEFORE: 2012 to 2021**
About a decade's worth of news heavily biased due to high frequency reporting on recent COVID-19 developments

**BEFORE: 2012 to 2019**
Remove news headlines dated to 2020 and 2021 to avoid having too much noise regarding COVID/vaccination situation



Word Cloud for Biotech News Headlines



Word Cloud for Biotech News Headlines

# REMOVAL OF OUTLIERS - NUMBER OF NEWS > P99*



Number of Daily News Report(s) across Time

* Days with number of news report > P99 (i.e. 7) for a particular company are also removed

# DISTRIBUTION OF WORD COUNT AND PRICE CHANGE VARIABLES



IQR of the absolute % price change from previous day;
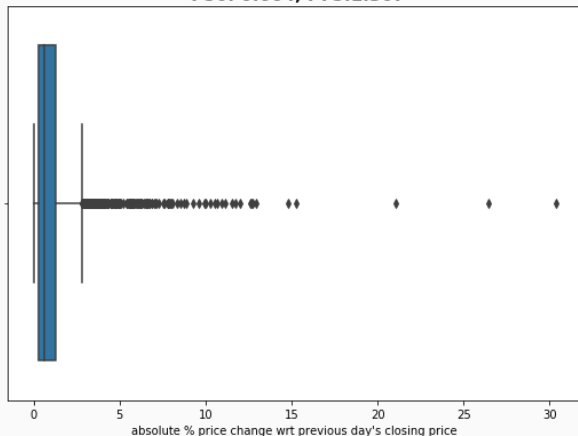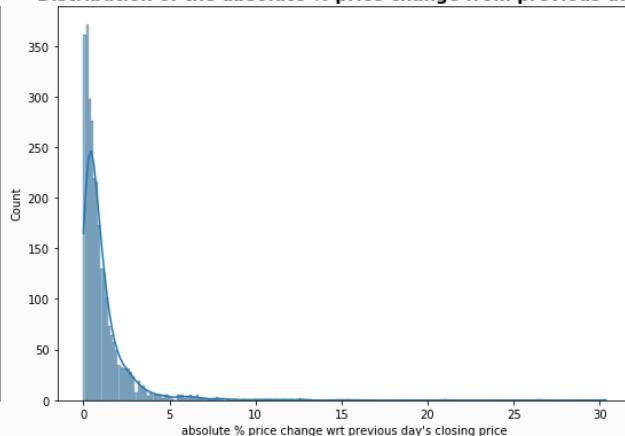P50: 0.664, P75:1.307

absolute % price change wrt previous day's closing price



Distribution of the absolute % price change from previous days

absolute % price change wrt previous day's closing price

## % Price Change

- Right skewed distribution with a long tail

- **P75** to be referenced as **threshold** of 'price jump' for target variable

## Word Count

- Right skewed distribution with a long tail

- Mean = 11, median = 8 words; question of sufficiency of text info for NLP analysis



Distribution of Word Count in News Headlines

mean = 11
median = 8

No of Words

# MORE NEWS REPORT SEEM MILDLY ASSOCIATED WITH HIGHER TRADING VOLUME AND STOCK PRICE CHANGES



Plots showing how News Count correlates with Trading Volume and Stock Price Changes

**ALL HEADLINES**

**Most Common Words in News Reported**

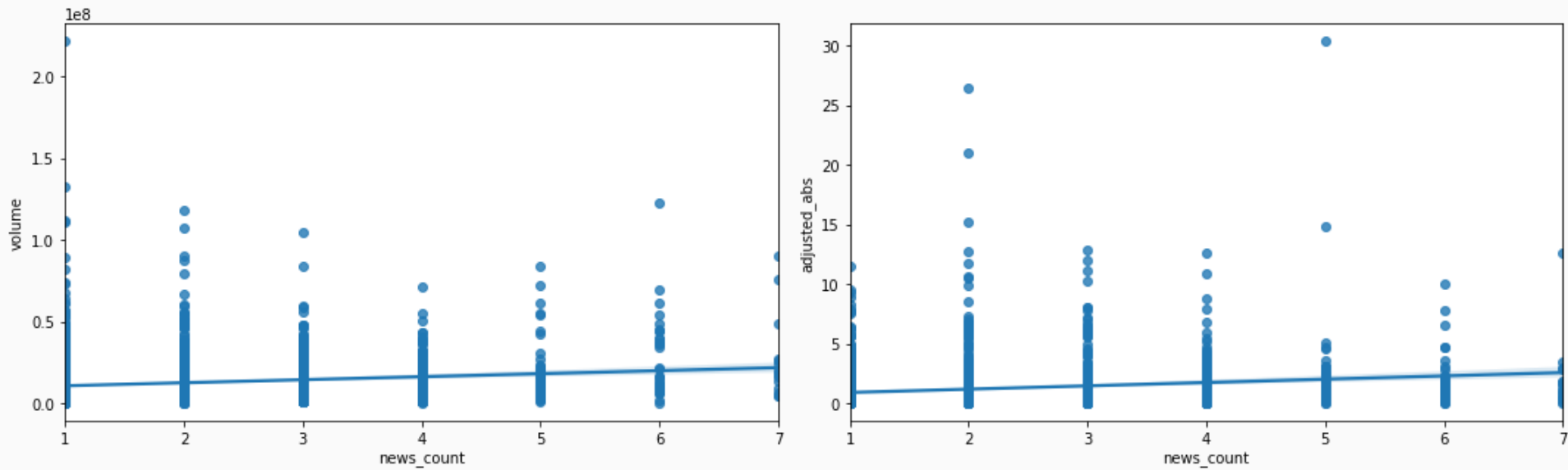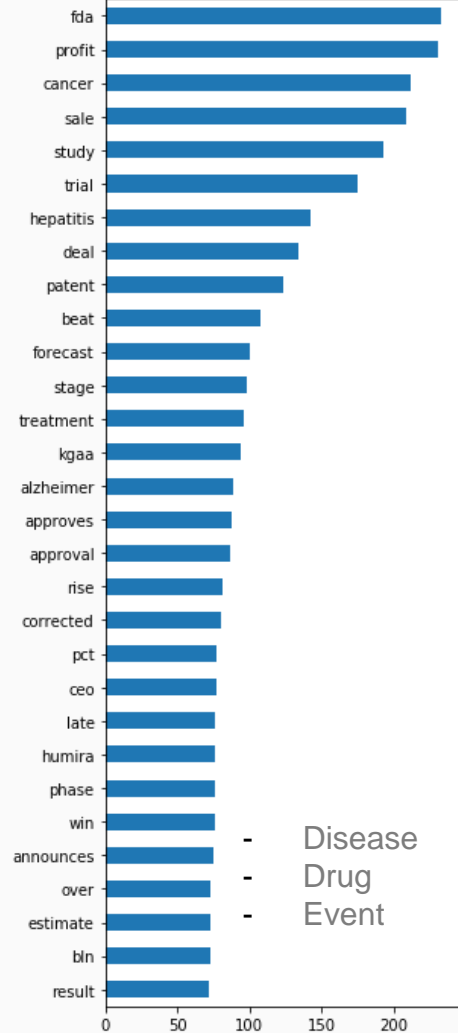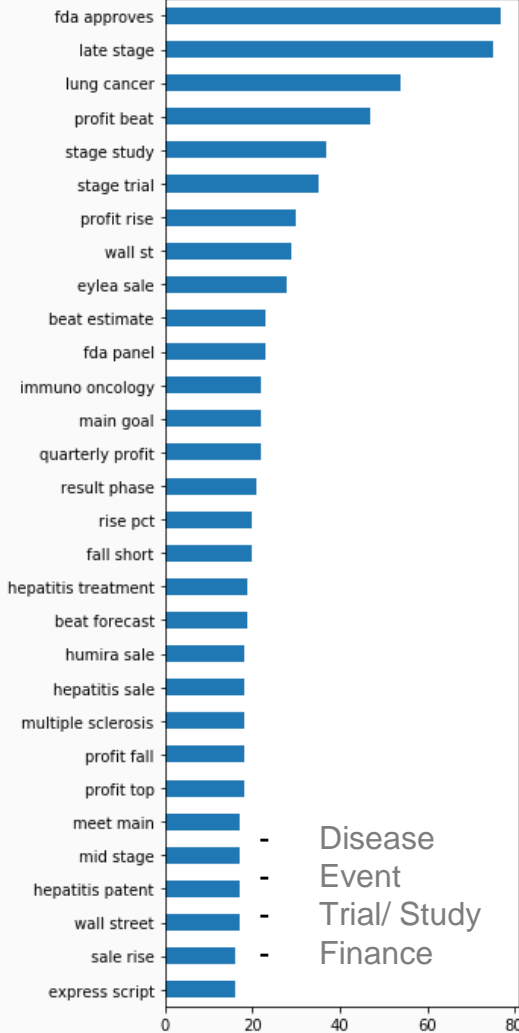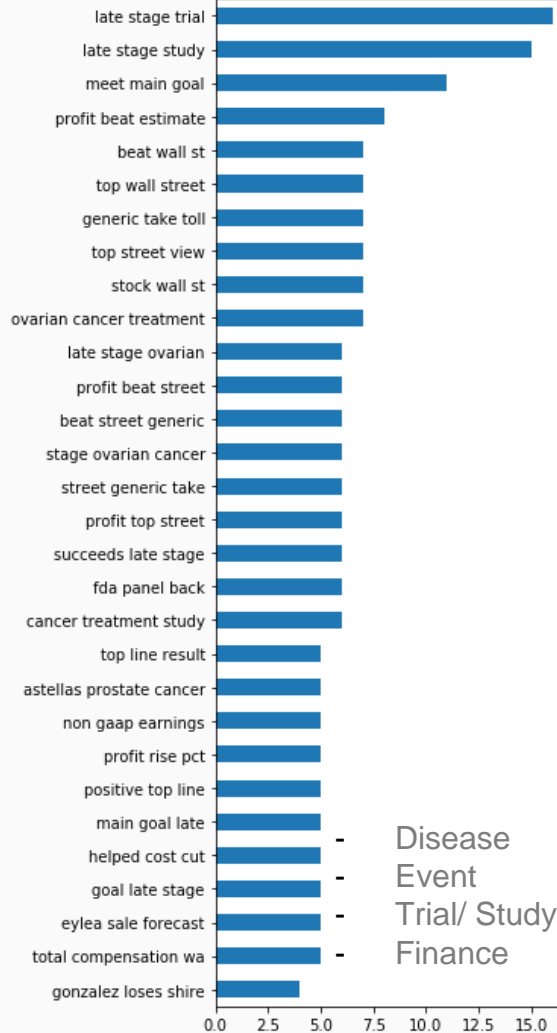| Word | Value |
|---|---|
| fda | ~230 |
| profit | ~225 |
| cancer | ~210 |
| sale | ~205 |
| study | ~190 |
| trial | ~175 |
| hepatitis | ~140 |
| deal | ~130 |
| patent | ~120 |
| beat | ~105 |
| forecast | ~100 |
| stage | ~98 |
| treatment | ~95 |
| kgaa | ~93 |
| alzheimer | ~88 |
| approves | ~86 |
| approval | ~85 |
| rise | ~82 |
| corrected | ~80 |
| pct | ~78 |
| ceo | ~77 |
| late | ~76 |
| humira | ~75 |
| phase | ~75 |
| win | ~75 |
| announces | ~73 |
| over | ~72 |
| estimate | ~72 |
| bln | ~71 |
| result | ~70 |

Legend:
- Disease
- Drug
- Event

**Most Common Words in News Reported**

| Word | Value |
|---|---|
| fda approves | ~77 |
| late stage | ~75 |
| lung cancer | ~54 |
| profit beat | ~47 |
| stage study | ~37 |
| stage trial | ~35 |
| profit rise | ~30 |
| wall st | ~29 |
| eylea sale | ~28 |
| beat estimate | ~23 |
| fda panel | ~23 |
| immuno oncology | ~22 |
| main goal | ~22 |
| quarterly profit | ~22 |
| result phase | ~21 |
| rise pct | ~20 |
| fall short | ~20 |
| hepatitis treatment | ~19 |
| beat forecast | ~19 |
| humira sale | ~18 |
| hepatitis sale | ~18 |
| multiple sclerosis | ~18 |
| profit fall | ~18 |
| profit top | ~18 |
| meet main | ~17 |
| mid stage | ~17 |
| hepatitis patent | ~17 |
| wall street | ~17 |
| sale rise | ~16 |
| express script | ~15 |

Legend:
- Disease
- Event
- Trial/ Study
- Finance

**Most Common Words in News Reported**

| Word | Value |
|---|---|
| late stage trial | ~16 |
| late stage study | ~15 |
| meet main goal | ~11 |
| profit beat estimate | ~8 |
| beat wall st | ~7 |
| top wall street | ~7 |
| generic take toll | ~7 |
| top street view | ~7 |
| stock wall st | ~7 |
| ovarian cancer treatment | ~7 |
| late stage ovarian | ~6 |
| profit beat street | ~6 |
| beat street generic | ~6 |
| stage ovarian cancer | ~6 |
| street generic take | ~6 |
| profit top street | ~6 |
| succeeds late stage | ~6 |
| fda panel back | ~6 |
| cancer treatment study | ~6 |
| top line result | ~5 |
| astellas prostate cancer | ~5 |
| non gaap earnings | ~5 |
| profit rise pct | ~5 |
| positive top line | ~5 |
| main goal late | ~5 |
| helped cost cut | ~5 |
| goal late stage | ~5 |
| eylea sale forecast | ~5 |
| total compensation wa | ~5 |
| gonzalez loses shire | ~4 |

Legend:
- Disease
- Event
- Trial/ Study
- Finance
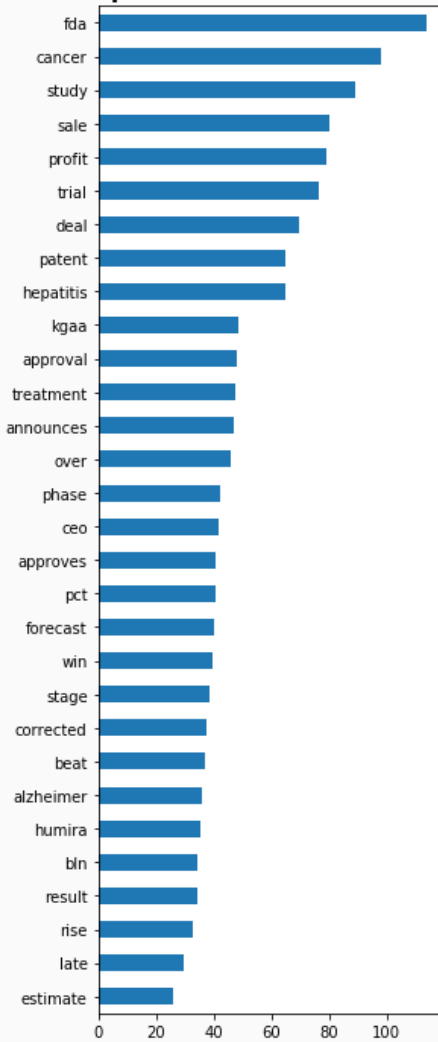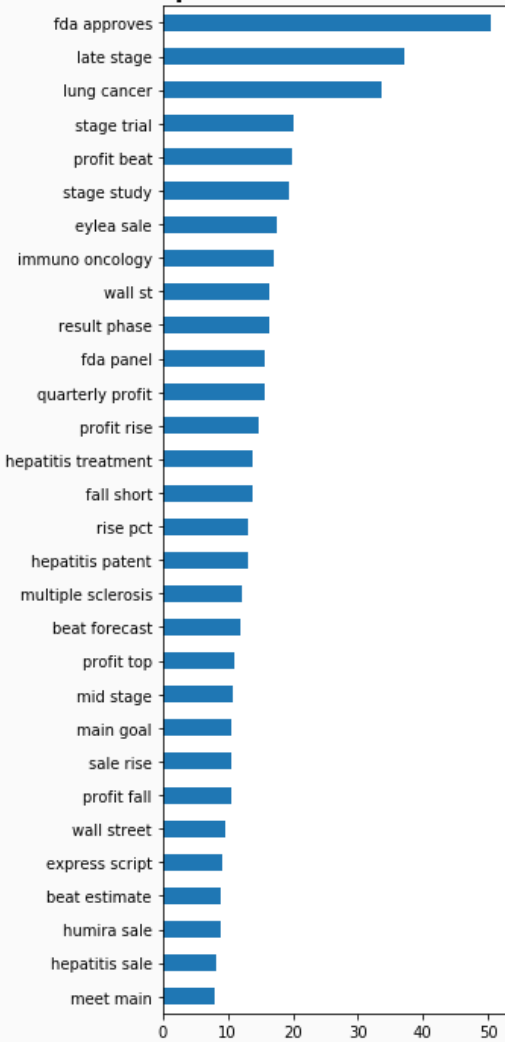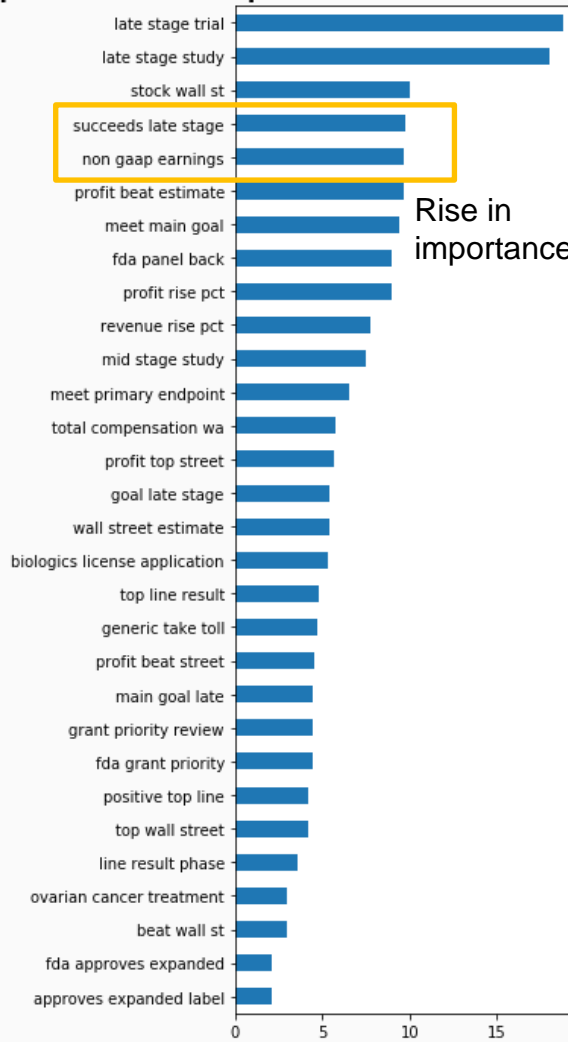
Most Important Words in News Reported

Most Important Words in News Reported

Most Important Words in News Reported

ALL HEADLINES

Rise in importance
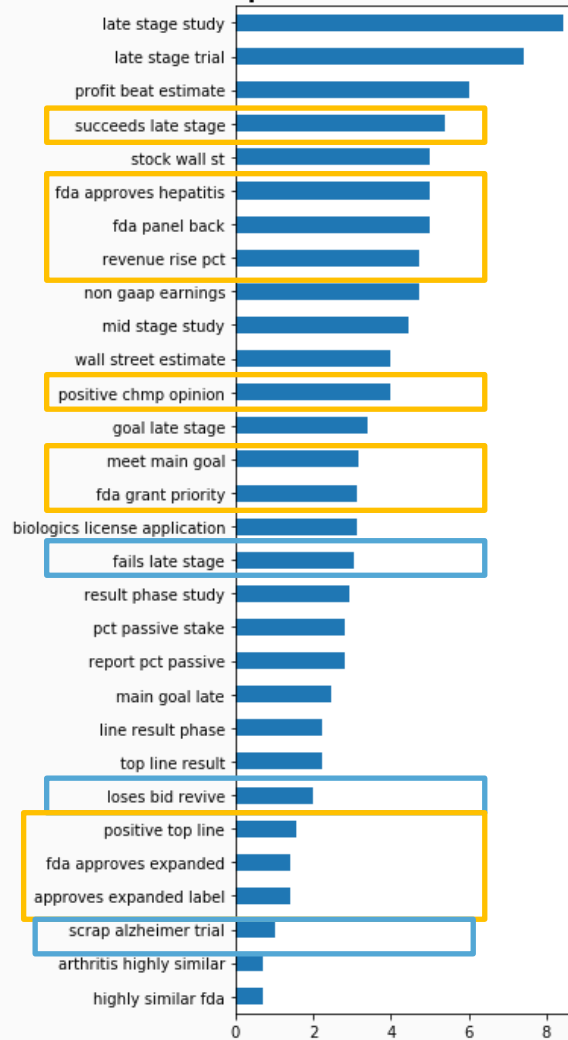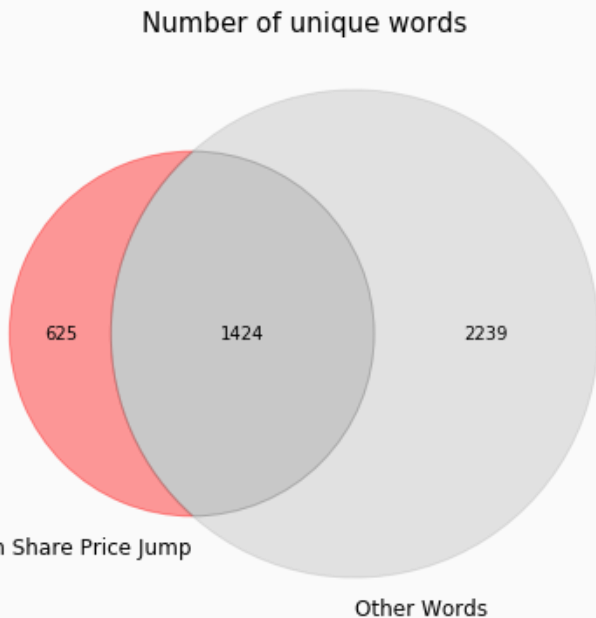
**Most Important Words in News Reported** (left chart)

HEADLINES ASSOCIATED WITH PRICE INCREASE

late stage trial
late stage study
meet main goal
stock wall st
profit rise pct
top wall street
non gaap earnings
succeeds late stage
top street view
eylea sale forecast
fda panel back
profit beat estimate
profit top street
total compensation wa
profit beat street
beat wall st
positive top line
top line result
generic take toll
gonzalez loses shire
ovarian cancer treatment
main goal late
goal late stage
helped cost cut
astellas prostate cancer
beat street generic
street generic take
late stage ovarian
stage ovarian cancer
cancer treatment study

**Most Important Words in News Reported** (right chart)

HEADLINES ASSOCIATED WITH PRICE DECREASE

late stage study
late stage trial
profit beat estimate
succeeds late stage
stock wall st
fda approves hepatitis
fda panel back
revenue rise pct
non gaap earnings
mid stage study
wall street estimate
positive chmp opinion
goal late stage
meet main goal
fda grant priority
biologics license application
fails late stage
result phase study
pct passive stake
report pct passive
main goal late
line result phase
top line result
loses bid revive
positive top line
fda approves expanded
approves expanded label
scrap alzheimer trial
arthritis highly similar
highly similar fda

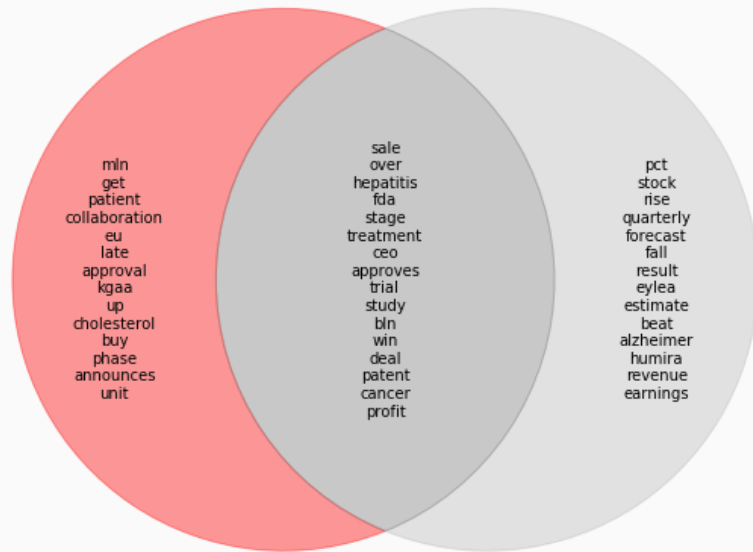Headlines associated with stock price decrease are a **mixed bag** of positive and negative events

# FEWER UNIQUE WORDS ARE ASSOCIATED WITH SHARE PRICE JUMP, BUT THERE ARE NO DISCERNABLE PATTERNS OR WORDS THAT SEEM INDICATIVE OF SHARE PRICE JUMP
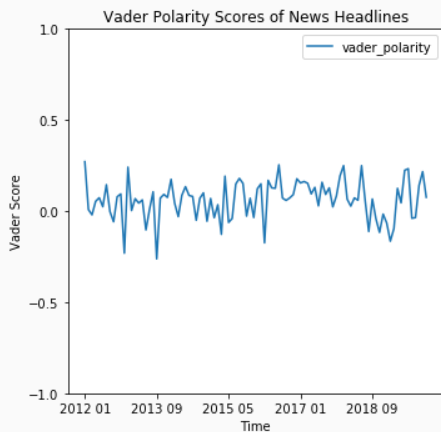


Number of unique words

625    1424    2239

Words Associated with Share Price Jump

Other Words

Top 30 words

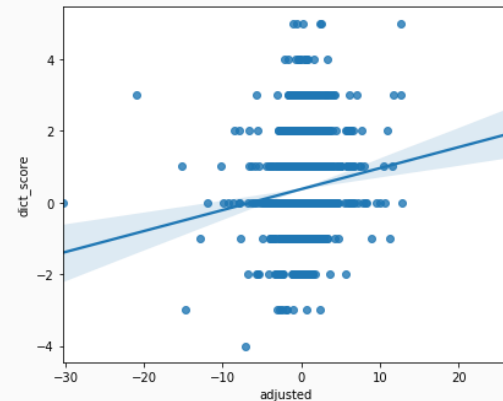| | | |
|---|---|---|
| mln | sale | pct |
| get | over | stock |
| patient | hepatitis | rise |
| collaboration | fda | quarterly |
| eu | stage | forecast |
| late | treatment | fall |
| approval | ceo | result |
| kgaa | approves | eylea |
| up | trial | estimate |
| cholesterol | study | beat |
| buy | bln | alzheimer |
| phase | win | humira |
| announces | deal | revenue |
| unit | patent | earnings |
| | cancer | |
| | profit | |

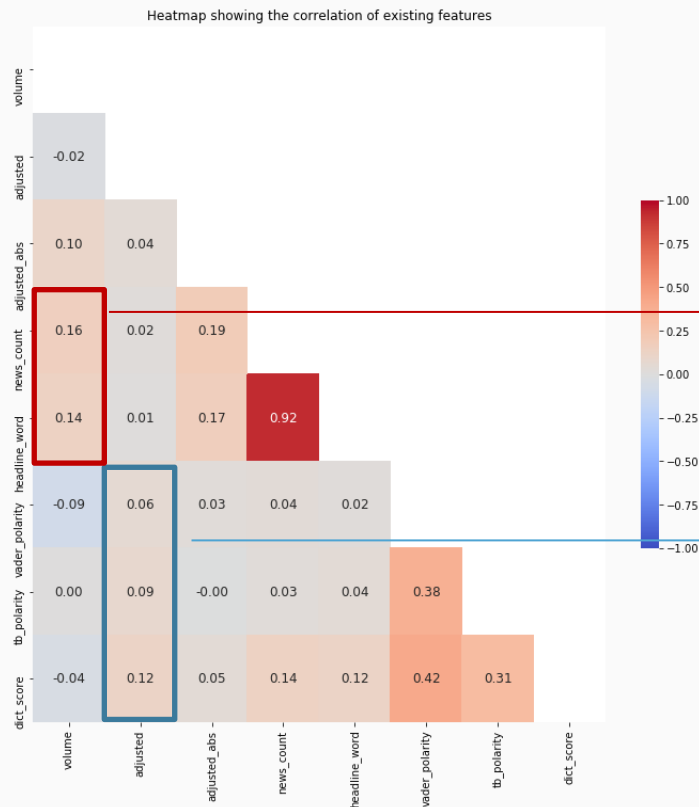Words Associated with Share Price Shock

Other Words

# SCATTER PLOT OF SENTIMENT SCORES SUGGEST EXISTING LIBRARIES/FUNCTIONS (E.G. VADER & TEXTBLOB) ARE NOT EFFECTIVE AT DIFFERENTIATING NEWS HEADLINES

# WEAK POSITIVE CORRELATION BETWEEN NEWS AND TRADING VOLUME, AND DICTIONARY SCORE AND ADJUSTED STOCK PRICE CHANGES IN %



Heatmap showing the correlation of existing features

Higher news count weakly & positively correlated with higher trading volumes

The manually determined dictionary score shows a higher correlation with the percentage change in stock prices than the vader and TextBlob polarity scores

| S/N | Model | Vectorizer | N-Gram | Accuracy | AUC | Recall | Precision | F1 |
|---|---|---|---|---|---|---|---|---|
| 1 | Light Gradient Boosting Machine | Count | 1-gram | 0.764 | 0.625 | 0.248 | 0.549 | 0.339 |
| 2 | Logistic Regression | Count | 1-gram | 0.764 | 0.651 | 0.307 | 0.529 | 0.386 |
| 3 | Gradient Boosting Classifier | Count | 2-gram | 0.747 | 0.560 | 0.068 | 0.663 | 0.120 |
| 4 | Logistic Regression | Count | 2-gram | 0.745 | 0.594 | 0.165 | 0.516 | 0.247 |
| 5 | Gradient Boosting Classifier | Count | 3-gram | 0.749 | 0.522 | 0.046 | 0.496 | 0.084 |
| 6 | Logistic Regression | Count | 3-gram | 0.747 | 0.538 | 0.060 | 0.493 | 0.106 |
| 7 | **Logistic Regression** | **Tf-idf** | **1-gram** | **0.758** | **0.654** | **0.117** | **0.672** | **0.199** |
| 8 | K Neighbors Classifier | Tf-idf | 1-gram | 0.748 | 0.568 | 0.114 | 0.577 | 0.187 |
| 9 | Logistic Regression | Tf-idf | 2-gram | 0.766 | 0.600 | 0.043 | 0.636 | 0.079 |
| 10 | Gradient Boosting Classifier | Tf-idf | 2-gram | 0.764 | 0.600 | 0.080 | 0.524 | 0.135 |
| 11 | Logistic Regression | Tf-idf | 3-gram | 0.748 | 0.570 | 0.006 | 0.150 | 0.011 |
| 12 | Ada Boosting Classifier | Tf-idf | 3-gram | 0.748 | 0.533 | 0.066 | 0.511 | 0.116 |

# RESULTS OF A ONE-GRAM TD-IDF LOGISTIC REGRESSION



Confusion Matrix for Logistic Regression Model (TF-IDF)
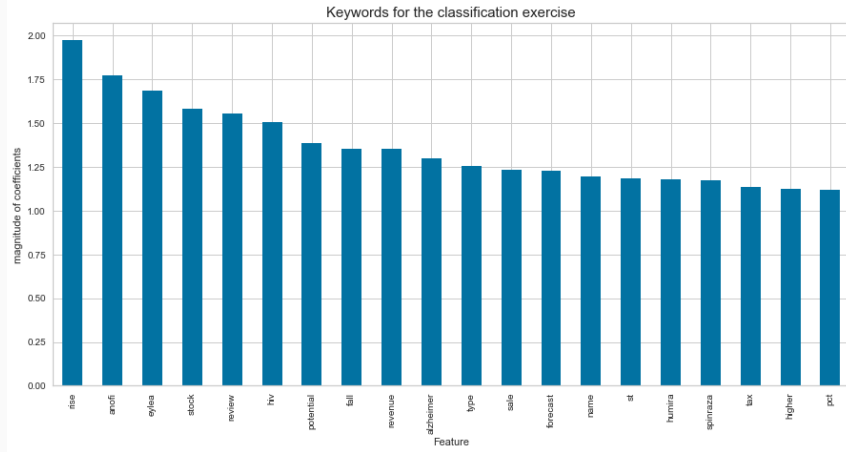Accuracy: 0.77, Precision: 0.75, Recall: 0.12



Receiver Operating Characteristic Curve of
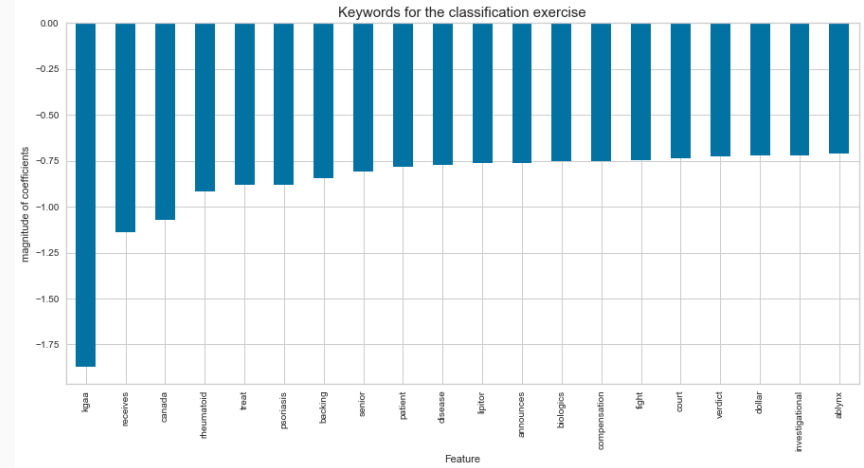Tfidf Vectorizer-Logistic Regression
ROC AUC: 0.642

The model has an accuracy score of 0.77, a ROC AUC score of 0.642 and a 0.75 precision score. While the precision score is high, it is obtained at the expense of a lower recall score of 0.12. The **model will be relatively good at predicting when a news headline would likely be a price jumper**. **However, the model may misclassify other news which are 'price jumper' as 'non-price jumper'.** This may lead to traders missing out on a significant news events that would affect stock prices significantly.

Words, e.g. '**rise**', '**stock**', '**hiv**', '**potential**', '**fall**', are predictive of stock price jumps

Words, e.g. '**kgaa**', '**receives**', '**canada**', '**rheumatoid**', are associated with constant stock prices
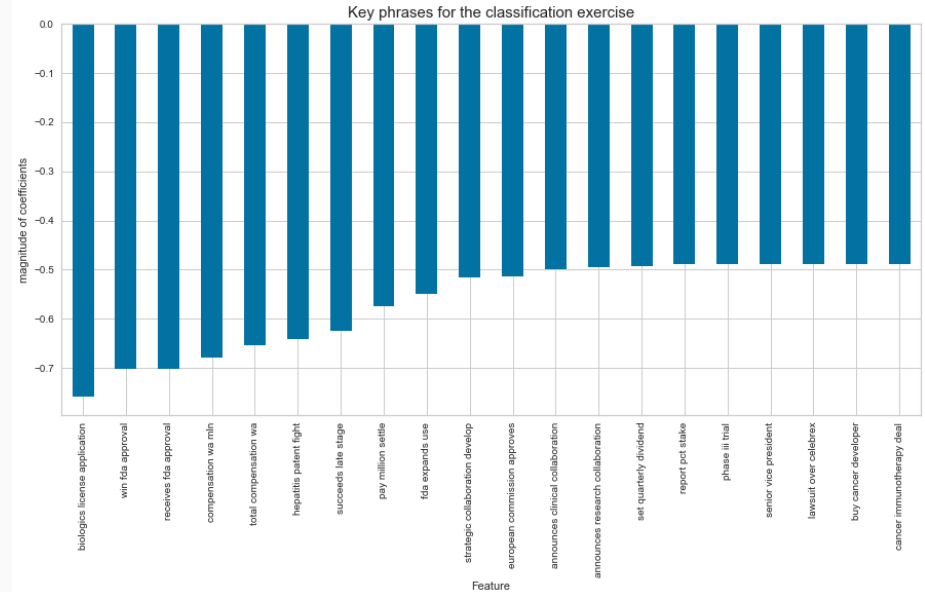


Keywords for the classification exercise



Keywords for the classification exercise

Phrases, e.g. **'revenue rise pct'**, **'profit above estimate'**, **'quarterly sale rise'**, **'mid stage study'**, are predictive of stock price jumps
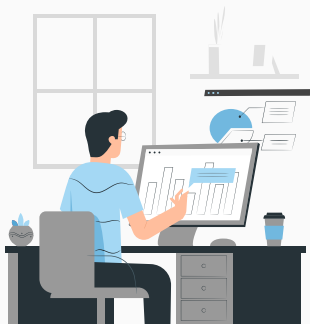
Phrases, e.g. **'win fda approval'**, **'receives fda approval'**, **'succeed late stage'**, are associated with <u>**constant**</u> stock prices



Key phrases for the classification exercise



Key phrases for the classification exercise

# CONCLUSION & RECOMMENDATIONS

1. **News headlines can be used to predict changes to stock prices.** The tf-idf logistic regression model is able to achieve > 70% accuracy and precision scores, and > 60% AUC ROC scores.

2. **The release of positive/ negative news may not necessarily lead to a jump in stock prices.** From this modelling experience, we learn that the announcement of significant milestones (such as FDA approval) may not lead to significant stock price changes.

**Considerations for Future Project(s):**
a. Aggregate news headlines from more sources (e.g. Financial Times, Seeking Alpha) for analysis
b. Study how the long-form news articles affect stock prices
c. Aggregate news headlines for more healthcare, pharmaceutical or biotech companies for analysis

# THANK YOU