

This project attempts to predict the performance of SPY ETF using financial news headlines on Reuters webpage on the previous trading day.

Data

News headlines between October 2018 and September 2020 are scraped from the Reuters financial news homepage¹. Past versions of the webpage are retrieved from the Internet Archive². The website snapshot closest to noon each day is used. Due to the limited scale of the project, data before October 2018 are not used because there are differences in the format of the HTML files.

The website is separated into 20 parts (figure 1). The text in each part is extracted and save into CSV files. Data validation is done to ensure the timestamp of the webpage is within the required date. Note that some fields may be empty since Reuters may have less than 6 stories for some days.

SPY ETF is chosen because it tracks the S&P500 index and can act as a proxy of the US stock market. The adjusted closing price and trade volume is retrieved from Yahoo Finance³.

Returns for each trading day is calculated. Two Boolean variables is also created: whether the price moves up, and whether the trade volume is greater than USD 75,000,000 (median of recent 2 years).

Analysis

Sentiment analysis is performed on the text data using TextBlob and NLTK Vader. A polarity score, showing if the statement is positive or negative, is generated for each of the 20 parts of the website.

A random forest model is fitted to predict the direction of the stock price and the relative size of the daily trade volume. Random forest model is chosen because it is non-parametric and it can automatically perform feature selection. Due to the limited scale of the project, other models are not considered. Cross validation is performed to tune the complexity of the model to avoid overfitting.

The data is split into a training and a test set, which contains 17% and 83% of the data respectively. ETF price increases in 65.1% of the test set, and trading volume is greater than USD 75M in 44.2% of the test set.

6 models are fitted in an attempt to find predictive power as shown below.

Model	Features	Response Variable	Test Accuracy
Model 1	TextBlob polarity and subjectivity	ETF price direction	56.98%
Model 2	NLTK polarity	ETF price direction	62.79%
Model 3	NLTK polarity for news contents only	ETF price direction	46.51%
Model 4	TextBlob polarity and subjectivity	Size of trade volume	54.65%
Model 5	NLTK polarity	Size of trade volume	53.49%
Model 6	NLTK polarity for news contents only	Size of trade volume	47.67%

Results

All models do not show significant predictive power on both price direction and volume of trade as their accuracy is lower than a null model (predicting using the most occurring result).

Interestingly, for model 2, the content of the 3rd story has much more predictive power than other features (figure 2). This may be caused by pure luck; however, it can be a starting point for further analysis.

Due to the limited scale, this project simplifies multiple processes. More in-depth analysis can be done by extracting more data, generate more features from the text data, and use a more iterative approach for model building.

¹ <https://www.reuters.com/finance>

² <https://archive.org>

³ <https://finance.yahoo.com/quote/SPY/history/>

Appendix

Figure 1: Reuters financial news homepage with annotations

Main Title points to the CME Group banner.

Main Content points to the main article titled "Exclusive: Biden campaign tells miners it supports domestic production of EV metals".

Featured Title 1-3 points to three articles: "Huawei third-quarter revenue rises 3.7%, ending double-digit growth streak", "Stocks hold tight ranges as U.S. election caution sets in", and "Exclusive: Wells Fargo explores sale of asset management business - sources".

Video Title 1-3 points to three video thumbnails: "Gilead's remdesivir gets U.S. FDA approval", "Stocks rally with eyes on stimulus talks", and "Goldman Sachs to pay \$3 billion over 'IMDB scandal'".

Stories Title 1-6 points to six story thumbnails: "HK brokers ready war chest for mom-and-pop bidding frenzy in Ant's mega IPO", "Intel's margins tumble as customers shift to cheaper chips, shares slide 10%", "Walmart sues federal government over opioid case", "Top BI and Data Trends of 2021 - See What's Coming Next", "California drivers sue Uber over in-app messages asking to support ballot measure", and "Major airline groups push for end to coronavirus quarantines, travel bans".

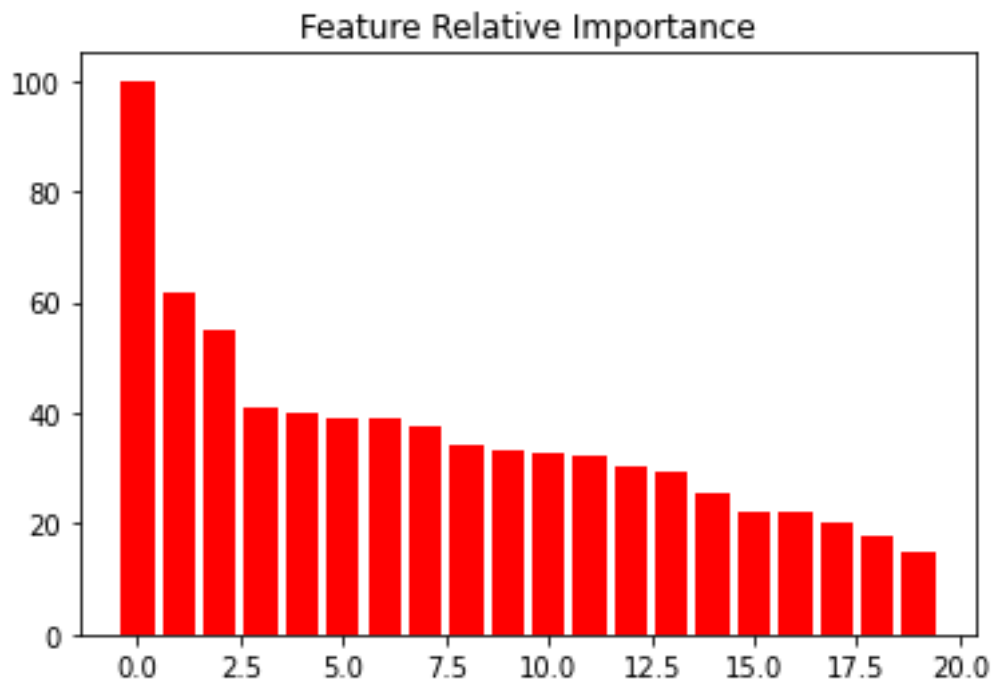
Stories Content 1-6 points to the corresponding text content for the six stories listed above.

Market Data Sidebar:

- STOCKS**
 - S&P 500: 3,493.49 (+0.52%)
 - Dow Jones: 28,363.66 (+0.54%)
 - FTSE 100: 5,785.65 (+0.16%)
 - Nikkei 225: 23,555.51 (+0.35%)
- BONDS**
 - US 10YR: +0.853 (+0.003)
 - DE 10YR: -0.567 (+0.001)
 - JP 10YR: +0.036 (+0.005)
- CURRENCIES**
 - USD / EUR: 0.8473 (+0.15%)
 - USD / GBP: 0.7655 (+0.17%)
 - USD / JPY: 104.6900 (-0.12%)

Footer: Follow Reuters: Twitter, Facebook, YouTube, LinkedIn. Subscribe: Newsletters | Podcasts | Apps. Reuters News Agency | Brand Attribution Guidelines | Advertise with Us | Careers | Reuters Editorial Leadership | Reuters Fact Check.

Figure 2: Model 2 variable importance chart and table



Features	Relative Importance	Importance
stories_content_3	100%	13.74%
stories_title_0	62%	8.46%
video_title_1	55%	7.56%
stories_title_1	41%	5.65%
stories_content_4	40%	5.49%
stories_title_4	39%	5.38%
stories_content_5	39%	5.35%
feature_title_2	37%	5.15%
video_title_2	34%	4.72%
main_title	33%	4.57%
main_content	33%	4.47%
stories_content_0	32%	4.45%
feature_title_1	30%	4.16%
stories_content_1	29%	4.02%
stories_content_2	26%	3.53%
stories_title_2	22%	3.05%
feature_title_0	22%	3.02%
video_title_0	20%	2.76%
stories_title_3	18%	2.43%
stories_title_5	15%	2.02%

[Attachment](#)

Filename	Description
Step_1_scrape_web.py	Scrape text data from web and save into CSV files
Step_2_analyse_data.py	Data pre-processing and fitting random forest models
web_data/web_data_yyyymm.csv	Text data from Reuters financial news homepage
SPY.csv	Historic ETF prices from Yahoo Finance