# Using Economic Indicators and Sentiment Analysis of Economic Policies to Build a Predictive Model of S&P 500 Stock Price

Jonathan Watkins, Jimmy Zhang, Jake Jarosik
DSCI 631 Fall Quarter 2023

# Introduction and Background

- Federal Open Market Committee  (FOMC)
  - Holds 8 regular meetings a year
  - Sets federal funds interest rates and other monetary policies
  - Monitors the US economy to ensure it is working
  - They look at metrics / economic indicators like CPI and the unemployment rate
  - Previous chairman of the FOMC have shown great variance in economic philosophy
- Standard and Poor's 500 Index Fund (S&P 500)
  - A stock index fund that is a weighted representation of the 500 largest companies traded on the US stock exchange.
  - Largest sector is currently software and technology
  - Includes: {'Communication Services', 'Consumer Discretionary', 'Consumer Staples', 'Energy', 'Financials', 'Health Care', 'Industrials', 'Information Technology', 'Materials', 'Real Estate', 'Utilities'}
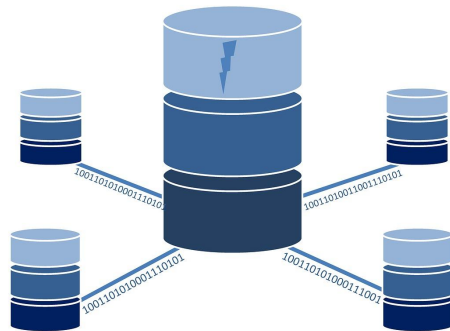
# Project Scope / Project Significance

- Build a Model to Forecast the Price of the S&P 500 Index Fund
  - Exploratory Data Analysis (EDA), specifically related to S&P 500
    - Visualization of current economic climate
  - Sentiment Analysis of FOMC textual data
    - Informs predictive model of the direction of the Federal Reserve's policies
  - Time series modeling for numeric S&P 500 stock data
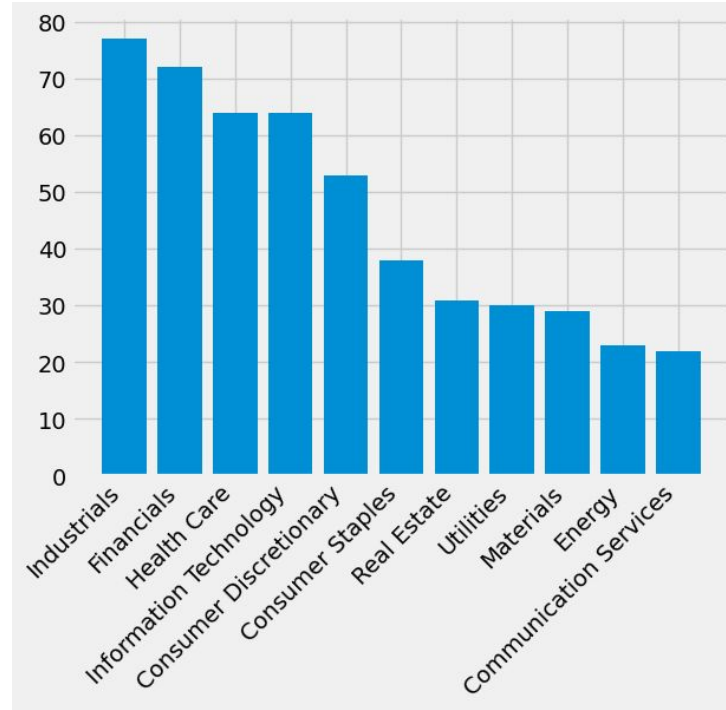    - Quantify the linearity of the model as well as predict future stock prices
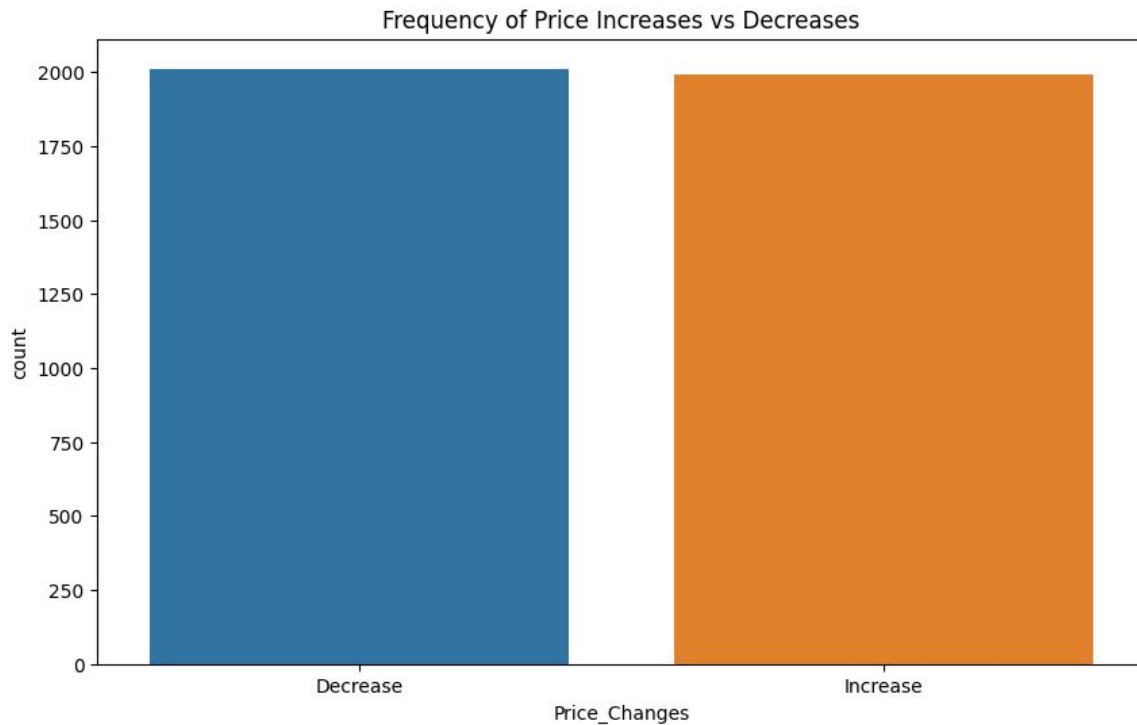
# Data Acquisition

- Timeframe: From 2008 to Present
- Stock Price Data
  - Yahoo Finance API
  - Daily stock prices (Open, Close, High, Low, Volume)
  - Limited Fundamental Data (Market Cap, Moving Averages, 50 Week High and Low)
- FOMC Textual Data
  - Acquired through web scraping
  - Speeches, statements and minutes
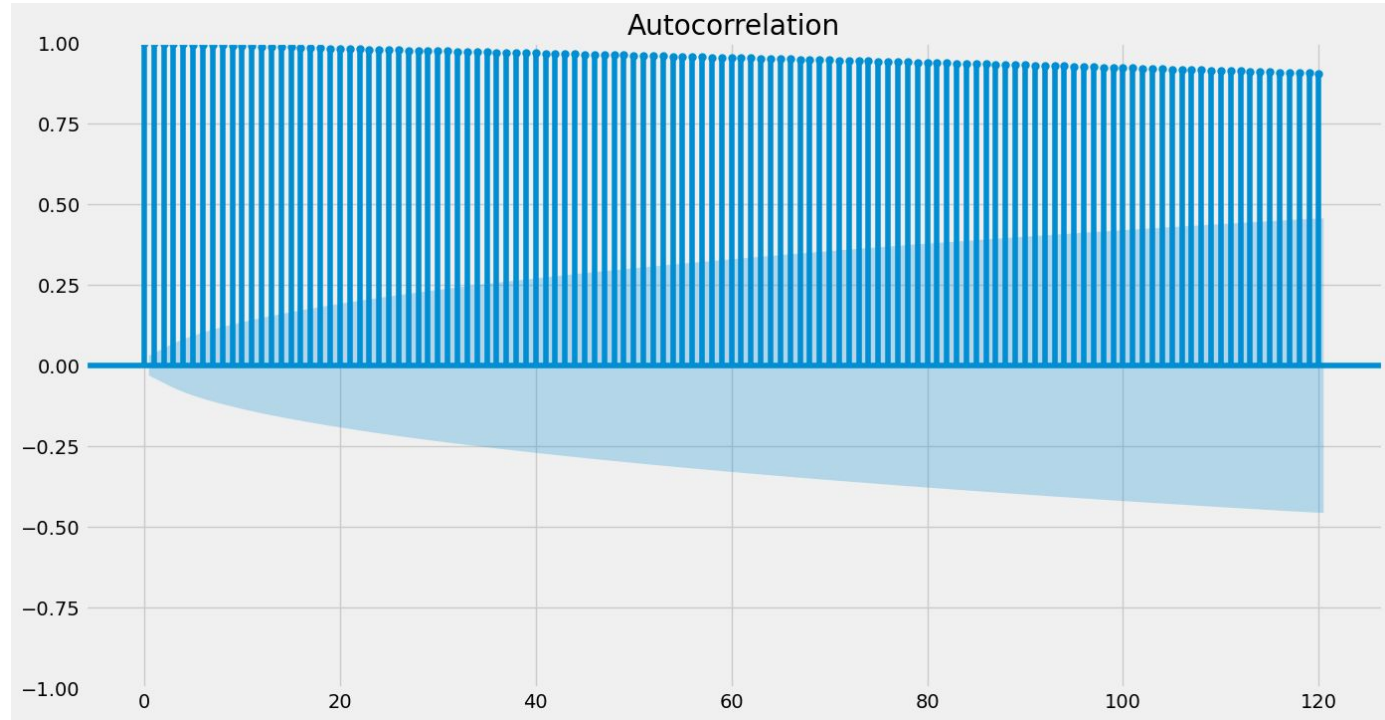- Federal Interest Data
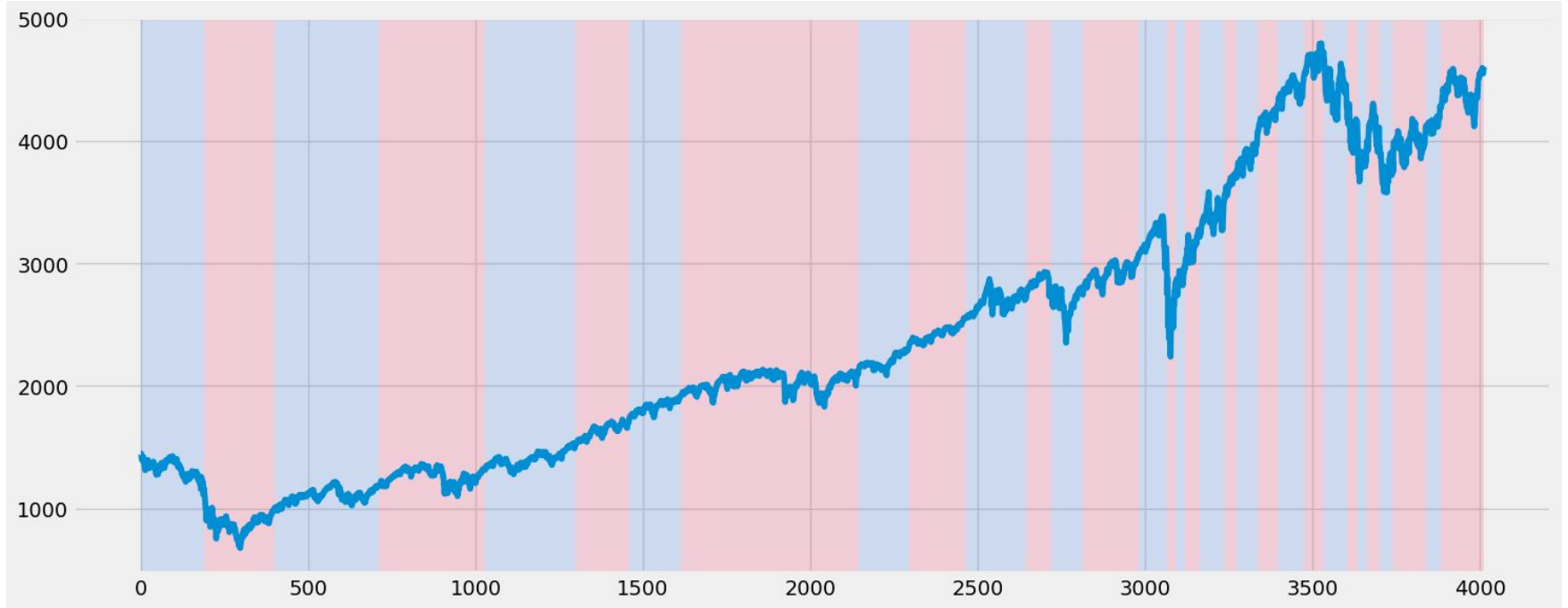- BLS Economic Data

# EDA: Composition of S&P 500

# EDA: Stock Price Daily Changes



Frequency of Price Increases vs Decreases
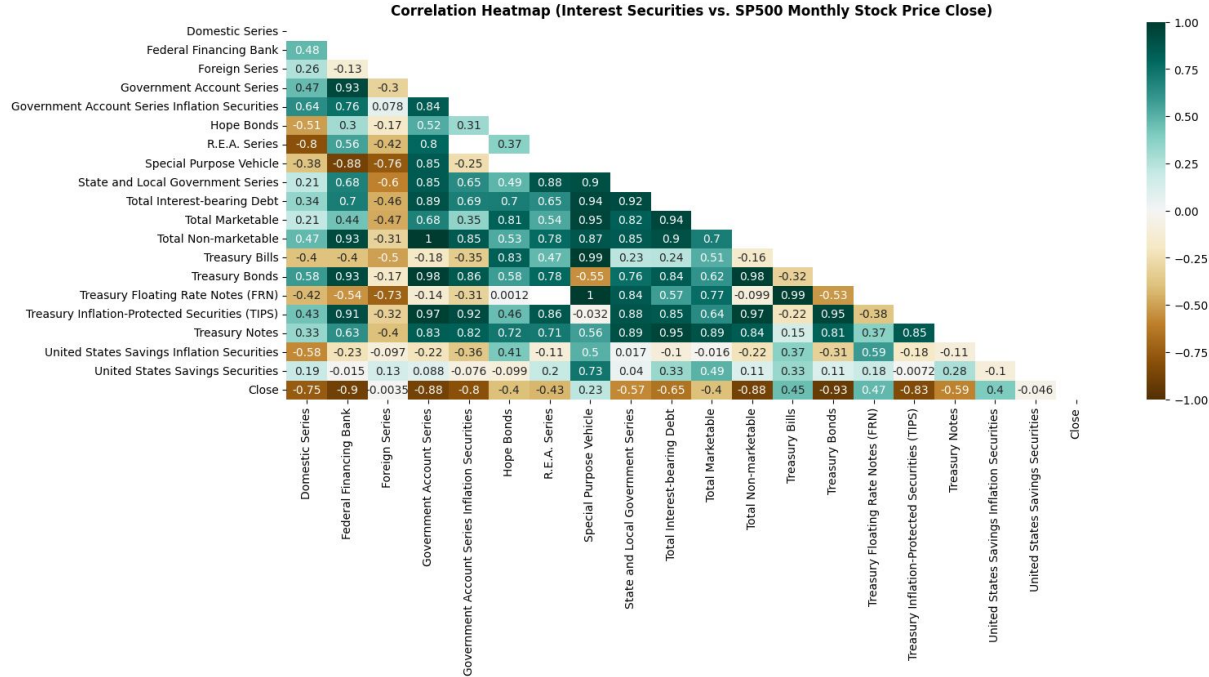
# EDA: Stock Price Autocorrelation

# EDA: Stock Price Change Point Detection

# EDA: Interest Rate Correlation Matrix



Correlation Heatmap (Interest Securities vs. SP500 Monthly Stock Price Close)
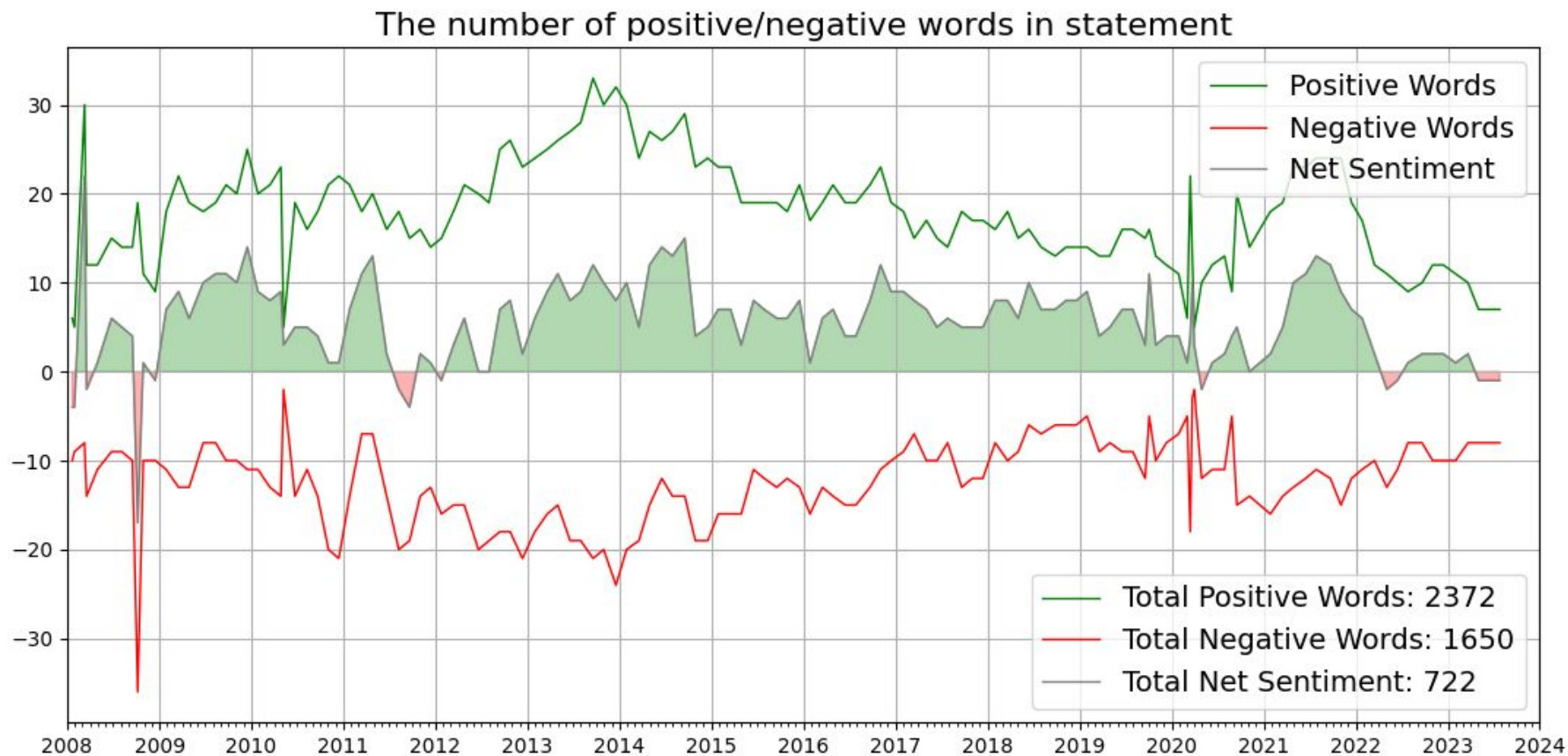
# EDA: FOMC Statement Word Cloud

# Sentiment Analysis: Loughran and McDonald Word List and Usage

```
lmdict = {'Negative': ['abandon', 'abandoned', 'abandoning', 'abandonment', 'abandonments', 'abandons', 'abdicated',
                       'abdicates', 'abdicating', 'abdication', 'abdications', 'aberrant', 'aberration', 'aberrational',
                       'aberrations', 'abetting', 'abnormal', 'abnormalities', 'abnormality', 'abnormally', 'abolish',
                       'abolished', 'abolishes', 'abolishing', 'abrogate', 'abrogated', 'abrogates', 'abrogating',
                       'abrogation', 'abrogations', 'abrupt', 'abruptly', 'abruptness', 'absence', 'absences',
                       'absenteeism', 'abuse', 'abused', 'abuses', 'abusing', 'abusive', 'abusively', 'abusiveness',
'Positive': ['able', 'abundance', 'abundant', 'acclaimed', 'accomplish', 'accomplished', 'accomplishes',
             'accomplishing', 'accomplishment', 'accomplishments', 'achieve', 'achieved', 'achievement',
             'achievements', 'achieves', 'achieving', 'adequately', 'advancement', 'advancements', 'advances',
             'advancing', 'advantage', 'advantaged', 'advantageous', 'advantageously', 'advantages',
negate = ["aint", "arent", "cannot", "cant", "couldnt", "darent", "didnt", "doesnt", "ain't", "aren't", "can't",
          "couldn't", "daren't", "didn't", "doesn't", "dont", "hadnt", "hasnt", "havent", "isnt", "mightnt", "mustnt",
          "neither", "don't", "hadn't", "hasn't", "haven't", "isn't", "mightn't", "mustn't", "neednt", "needn't",
```

- Objective: Account for simple negation only for positive words
  - Counting Positive and Negative Words with Negation Check
  - Simple Negation:
    - Occurs when negate words appear within three words before a positive word
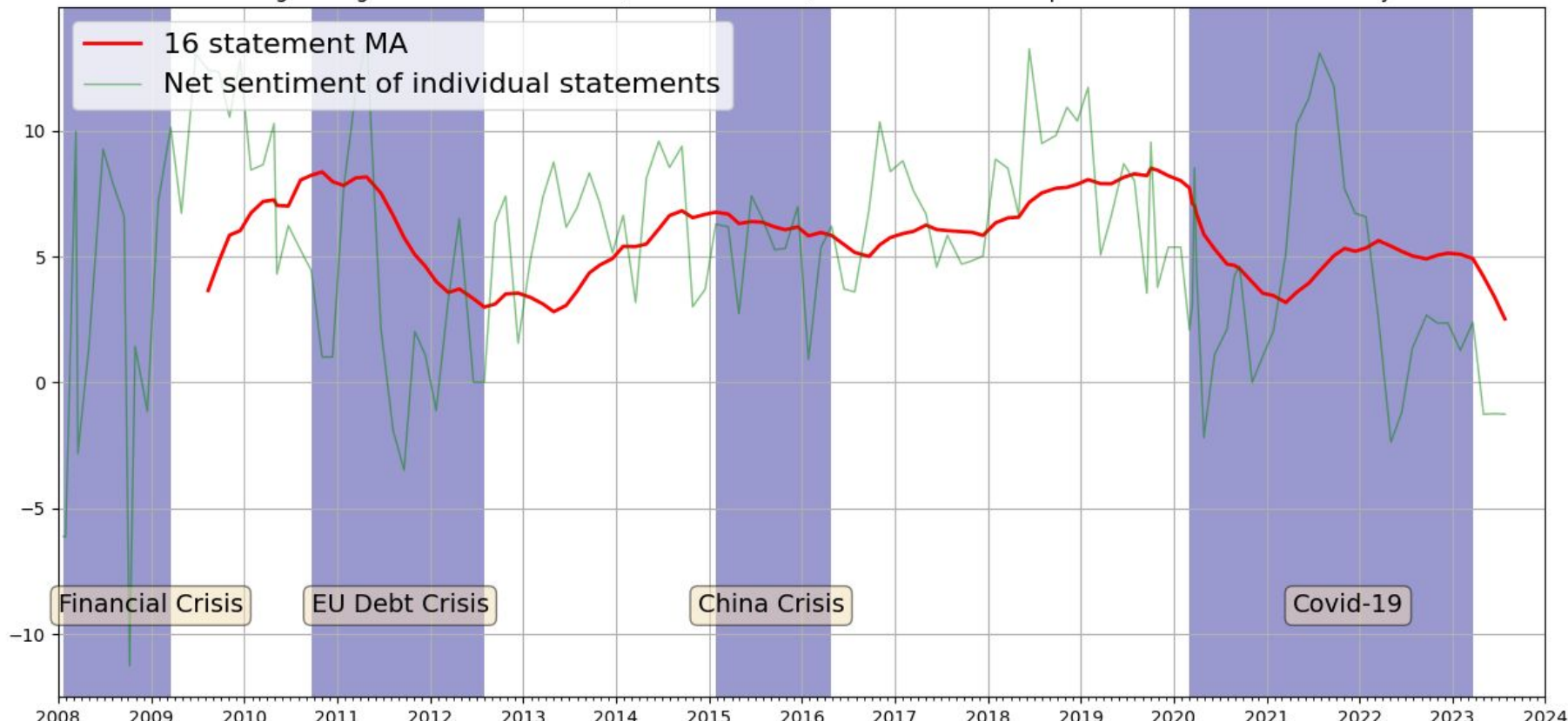
# Sentiment Analysis Plots: Tracking Positive and Negative Word Counts



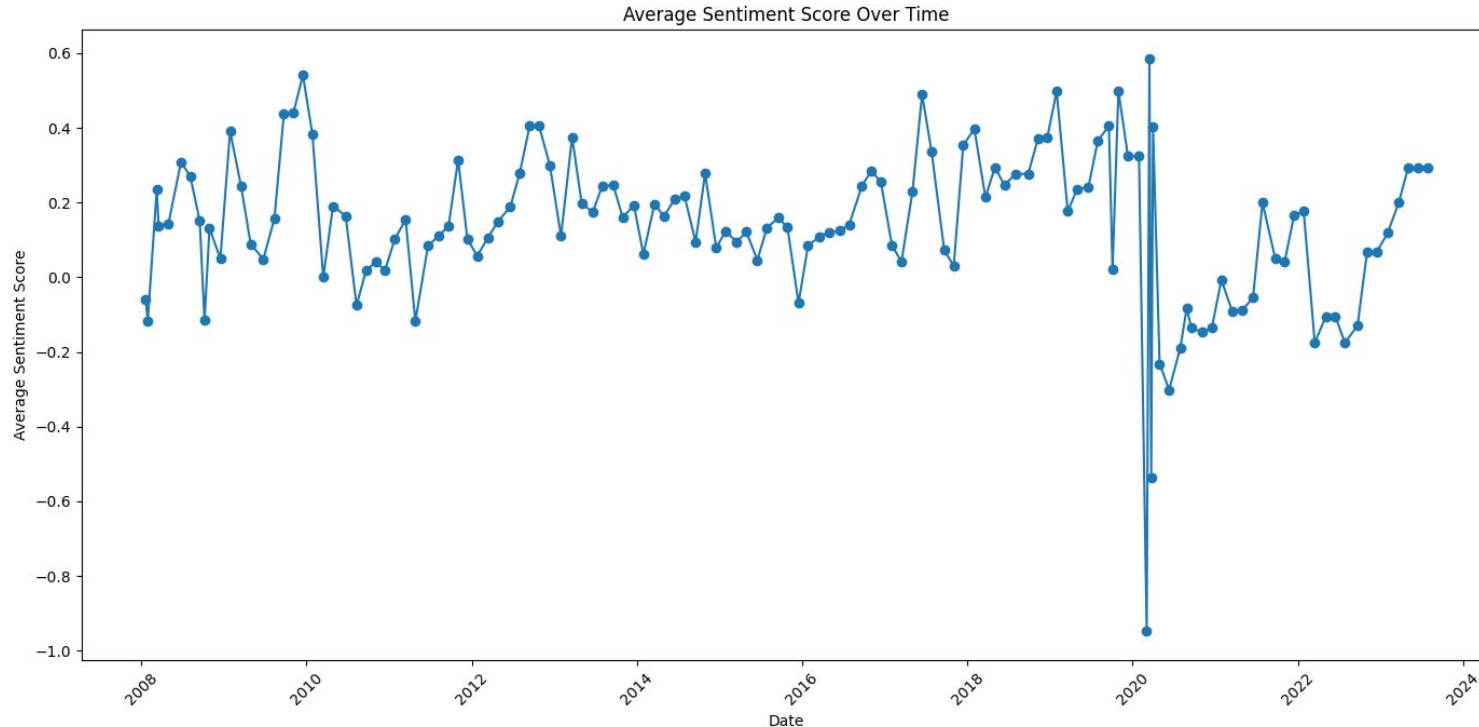The number of positive/negative words in statement

# Moving Averages



Moving average of last 16 statements (~2 Year Window) seems to match with periods of economic uncertainty
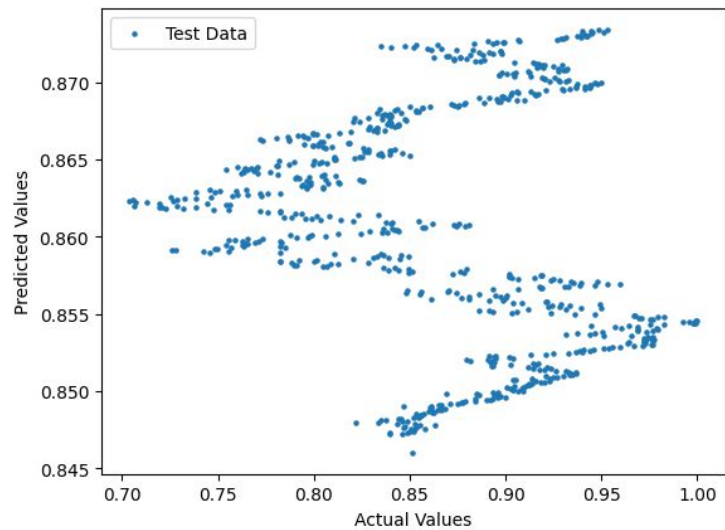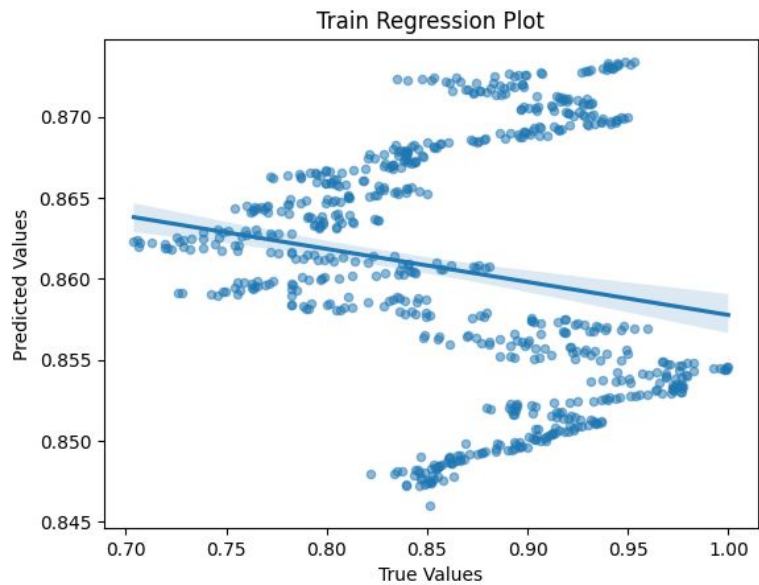
# Improving upon Sentiment Analysis with FiGAS approach and with Senticnet



Average Sentiment Score Over Time

# Introduction to ARIMA

- Autoregressive Integrated Moving Average
- The most simplest of modeling approach of the three.
  - Does not take additional regressors. Only fits historical values.
  - Does not account for Seasonality
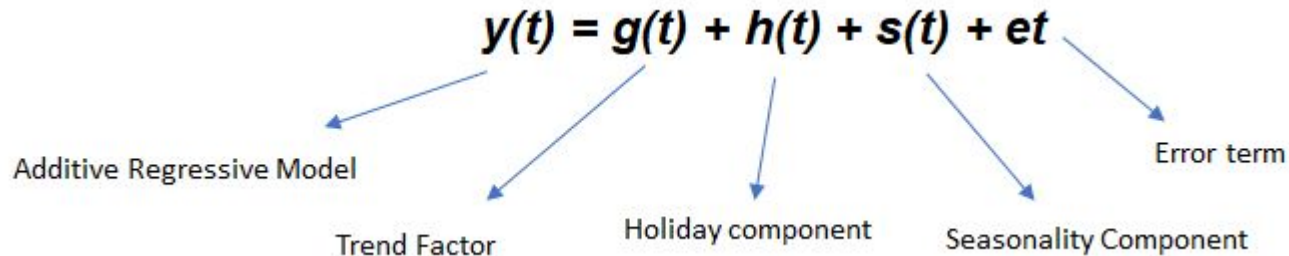  - Designed for Regular Time Spacing / Intervals
  - Assumes Stationarity

# ARIMA Results

# Introduction to Prophet

- Open Source Package By Meta / Facebook
- Prophet is a procedure for forecasting time series data based on an additive model where non-linear trends are fit with yearly, weekly, and daily seasonality, plus holiday effects.
- Simplifies Time Series Forecasting

$$y(t) = g(t) + h(t) + s(t) + et$$

Additive Regressive Model

Trend Factor

Holiday component

Seasonality Component

Error term

# Prophet Results - Autoregression

# Prophet Results - Comparing Different Model Settings (n_splits=4)

# Introduction to LSTM

- Long short-term memory network is a recurrent neural network, aimed to deal with the vanishing gradient problem present in traditional RNNs.
- Excels at capturing long-term dependencies, making it ideal for sequence prediction tasks

# LSTM Results

# LSTM Results - Time Series Split (n_splits = 4)

# LSTM Results - Final Model



Learning Curve

# LSTM Results

Train Score: 0.02 RMSE
Test Score: 0.16 RMSE

Mean Absolute Error (Train): 0.0162
Mean Absolute Error (Test): 0.1364



Comparison of values

# Key Project Findings and Insights

- EDA on S&P 500 Index Composition
  - Provides a diverse representation of different sectors in the US stock market with a focus on software and tech
- EDA on S&P 500 Index Fund Price Data
  - High autocorrelation in daily price, low autocorrelation in daily changes in price
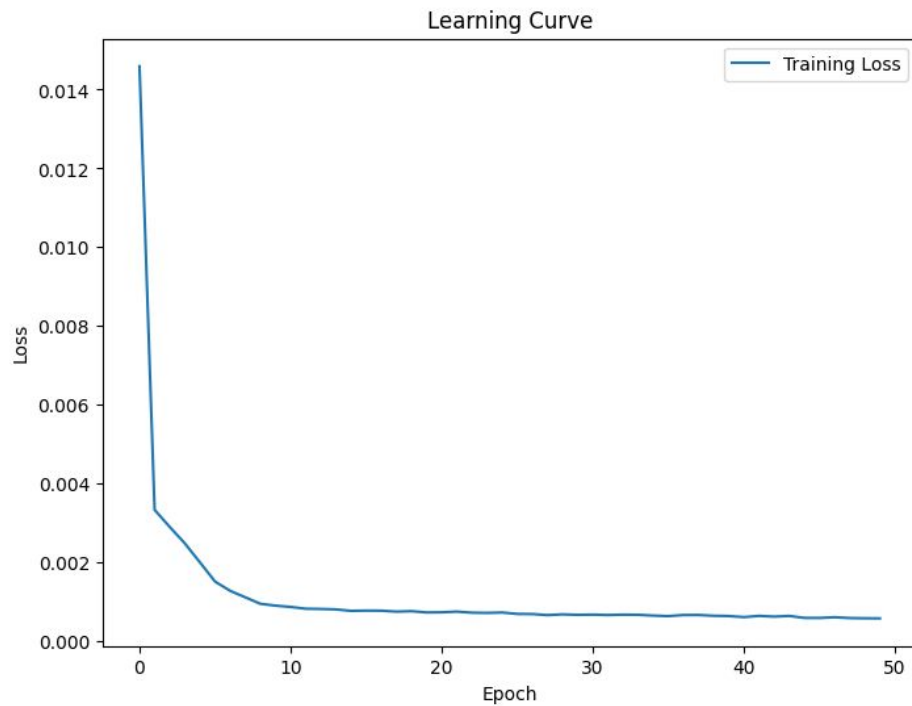- Sentiment Analysis of FOMC Statements
  - Leveraged Loughran-McDonald Master Dictionary and FiGAS approach for Sentiment Analysis
- Machine Learning Models
  - Times Series Modeling requires specialized machine learning algorithms
  - While we did not see the algorithms in lecture, modeling a time series uses the same concepts

# Challenges & Limitations

- Data Acquisition and Feature Selection:
  - Reliability of Free Stock Price Data and Data Retrieved via Web Scrape
- Sentiment Analysis:
  - Loughran and McDonald's Sentiment Word List still quite "basic" and misses context, intensity, offering only a limited view of sentiment through our raw scores/net sentiment
  - Flipping the sentiment for negation may not capture true causative factors behind the FED rate changes
- ARIMA
  - Model assumptions and requirements are too stringent for our Dataset
- Prophet
  - Lack of detailed documentation. While less restrictive than ARIMA, it lacks more advanced machine learning algorithms for additional regressors. And, the Holiday term is limited to model for a narrow range of events.
- LSTM
  - The most complex model, and most challenging to understand - made it difficult to use more advanced training validation methodologies (such as rolling window cross validation)
  - Shows a lot of promise, with sufficient tuning

# References

Board of Governors of the Federal Reserve System. (n.d.). Federal Open Market Committee (FOMC). Federal Reserve. Retrieved [July 22, 2023], from https://www.federalreserve.gov/monetarypolicy/fomc.htm

Consoli, S.; Barbaglia, L.; and Manzan, S. (2020). Fine-grained, aspect-based semantic sentiment analysis within the economic and financial domains. In Proceedings - 2020 IEEE 2nd International Conference on Cognitive Machine Intelligence, CogMI 2020, 52 – 61.

Consoli, S.; Barbaglia, L.; and Manzan, S. (2021). Explaining sentiment from Lexicon. In CEUR Workshop Proceedings, volume 2918, 87 – 95.

Consoli, S.; Barbaglia, L.; and Manzan, S. (2022). Fine-grained, aspect-based sentiment analysis on economic and financial lexicon. Knowledge-Based Systems, 247: 108781

Conti-Brown, P. (2016). The Power and Independence of the Federal Reserve. Princeton: Princeton University Press. https://doi.org/10.1515/9781400873500

Esri. (2023). How change point detection works. *ArcGIS Pro Documentation.* Retrieved from https://pro.arcgis.com/en/pro-app/latest/tool-reference/space-time-pattern-mining/how-change-point-detection-works.html

Facebook. (2023). Quick Start. In *Prophet Documentation*. Retrieved from https://facebook.github.io/prophet/docs/quick_start.html

Hansen, K. B. (2021). Model Talk: Calculative Cultures in Quantitative Finance. Science, Technology, & Human Values, 46(3), 600–627. https://doi.org/10.1177/0162243920944225

Jabbri Maleki, F. (2023). Metrics. In *Time Series Exploration with Python: A Journey from Traditional to Advanced Forecasting Models.* Retrieved from https://faridjb.github.io/Time-Series-Exploration/chapters/chapter_09.html

# References Continued

Kutzkov, K. (2023, July 31). ARIMA vs Prophet vs LSTM for Time Series Prediction. *Neptune.ai*. Retrieved from https://neptune.ai/blog/arima-vs-prophet-vs-lstm

Krieger, M. (2021, February 19). Time series analysis with Facebook Prophet: How it works and how to use it. *Towards Data Science.* Retrieved from https://towardsdatascience.com/time-series-analysis-with-facebook-prophet-how-it-works-and-how-to-use-it-f15ecf2c0e3a

Ma, E. (2018, November 3). Understand how to transfer your paragraph to vector by doc2vec. Towards Data Science. https://towardsdatascience.com/understand-how-to-transfer-your-paragraph-to-vector-by-doc2vec-1e225ccf102

Mehdian, S., Rezvanian, R., & Stoica O. (2019). "The Effect Of The 2008 Global Financial Crisis On The Efficiency Of Large U.S. Commercial Banks," Review of Economic and Business Studies, Alexandru Ioan Cuza University, Faculty of Economics and Business Administration, issue 24, pages 11-27, December.

Mikolov, T., Chen, K., Corrado, G. & Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. *CoRR*, abs/1301.3781.

Mukhiya, S. K., & Ahmed, U. (2020). *Hands-on exploratory data analysis with Python: Perform EDA techniques to understand, summarize, and investigate your data*. Packt Publishing.

Takahashi, Y. (2020, November 13). FedSpeak — How to build a NLP pipeline to predict central bank policy changes: A guide to analyse policymakers' conversations by deep neural network. Towards Data Science. https://towardsdatascience.com/fedspeak-how-to-build-a-nlp-pipeline-to-predict-central-bank-policy-changes-a2f157ca0434

Tensorflow. (2023). Time series forecasting. *TensorFlow Core.* Retrieved from https://www.tensorflow.org/tutorials/structured_data/time_series

Stocker, M., Lakatos, C., Ohnsorge, F., & Kose, M. A. (2017, February 27). *Understanding the global role of the US economy*. CEPR. https://cepr.org/voxeu/columns/understanding-global-role-us-economy

# Team and Contributions

- Jake Jarosik
  - ARIMA
  - LSTM

- Jonathan Watkins
  - Sentiment Analysis
  - EDA

- Jimmy Zhang
  - Data Acquisition and EDA
  - Prophet