

Lund University School of Economics and Management
Bachelor Thesis in Finance (NEKH01)
Supervisor: Dag Rydorff
June 2022



r/wallstreetbets Influence on the Stock Market

Sentiment Analysis on r/wallstreetbets during one of the loudest and most noticeable periods of financial debate on social media

By

Gustaf Kinch
gustaf.kinch@belteberga.se

Carl Tjernberg
tjernbergcarl@gmail.com

Abstract:

This research aims to examine how much impact social media channels have on stock returns during unusual market events. We do this by measuring the sentiment and activity on the Reddit forum wallstreetbets (WSB) regarding some of the six most mentioned companies during the first four months of 2021 and connect this to their relative stock returns. The sentiment is derived through a sentiment analysis using financial dictionary analysis, capturing the sarcastic and sometimes difficult to interpret language used by the community. The sentiment is derived utilizing two financial dictionary analysis, VADER, a sentiment analysis tool designed specifically for social media, as well as our own constructed dictionary EMOJI, which is created to better analyze the emotions expressed on WSB. We collect the sentiment using Python, a high-level programming language, and empirically examine WSB effective predicting power on individual stock returns using a panel data regression. We come to the conclusion that we do not have enough significant values to draw any conclusions about how WSB influences the stock return. We, on the other hand, are capable of determining that our strategies of extracting the sentiment work differently, where EMOJI appears to be better at capturing the mood on WSB. We believe this strengthens our preconceptions about the tone on WSB and leave it for further research to proceed upon this.

Keywords:

Reddit, wallstreetbets (WSB), Stock Return, Sentiment Analysis, Dictionary Analysis

Acknowledgement

We are grateful for Dag's guidance and intelligent recommendations during this project. With his help and dedication, we were able to achieve the research goals we set before we started the project.

Contents

1. Introduction	4
2. Theoretical Framework and Previous Research	8
2.1 Theoretical Framework.....	8
2.1.1 Random Walk Theory and Efficient Market Hypothesis.....	8
2.1.2 Behavioral Finance.....	9
2.2 Previous Research.....	10
2.2.1 Sentiment Analysis Background.....	10
2.2.2 Sentiment Analysis on Social Media.....	11
2.3 Summary.....	13
3. Data	15
3.1 WSB Data.....	15
3.2 WSB Data Processing.....	16
3.3 Financial Data.....	17
4. Method	19
4.1 Research Questions.....	19
4.2 Sentiment Analysis.....	20
4.2.1 Sentiment Analysis Methodology.....	20
4.2.2 VADER.....	21
4.2.3 EMOJI.....	22
4.2.4 Daily Sentiment Score.....	23
4.3 Sentiment Time-Series.....	24
4.4 Empirical Methods.....	25
4.5 Diagnostic Checks.....	27
5. Analysis & Discussion	29
5.1 Descriptive Statistics.....	29
5.2 Diagnostic Checks.....	31
5.3 WSB Sentiment and Stock Returns.....	33
5.4 Discussion.....	35
5.5 Limitations.....	37
6. Conclusion	38
Bibliography	39
Appendix	

1. Introduction

With a new generation entering adulthood, born and raised with ubiquitous internet access and who values social networking highly, it has become increasingly important to comprehend this generation's communication style, which differs from previous generations'. According to these studies, *Generation Z*, born between 1995 and 2012, is the most impatient, acquisitive and self-directed generation yet (Dangmei & Singh, 2016). Today, more people have turned to the stock market, which may be explained by the numerous lockdowns that took place during the pandemic where people were stranded at home having more time to evaluate their investments. In 2010, retail accounts accounted for 10.1% of the total U.S. equity trading volume, a number that rose to 19.0% in Q3 2021, hitting peak levels at 24.0% in Q1 2021. Meanwhile, the daily volume of options trading in the U.S. has surged by 57%, hitting peak levels in Q1 2021, indicating that the speculation in stocks has increased significantly. Even though the rise of retail investing cannot fully explain this shift, it is safe to say that understanding the change in communication among a growing segment of the adult population, and their characteristics, is critical for banks and institutional investors still accounting for around 40% of the overall U.S. equity trading volume, as well as regulators trying to prevent stock manipulation online (Feary, Sharma, Franco & Thrasher, 2022).

The use of social media has expanded rapidly as technology has improved, making communication platforms online more available to the public. With a generation that communicates online exclusively, the use of messaging boards to manipulate stock prices has become easier than ever to exploit. All of these elements have contributed to the rise of the infamous Reddit-forum *wallstreetbets* (WSB), where users from all over the world use a lexicon of terminology aimed to persuade people to bet on specific stocks (Corbet, Hou, Hu & Oxley, 2022).

The controversial stock price increase in GameStop Corp. (GME), when the closing stock price of the American video game retailer rocketed from USD 17.25 to 325.00, during January 2021 is one of the most notorious market events in recent years. Theoretically, the huge price increase might be explained by the fact that GME's short interest was 140 % as of January 22, implying that 1.4 times GME's outstanding shares were sold short but not yet covered or closed off (Chohan, 2021). Due to the hedge funds' massive positions, they had to cover their positions by trading the stock themselves, resulting in dramatic up-movements in

the stock price as the market could not offer the liquidity these short sellers needed to prevent large losses. The phenomenon is called a *short squeeze* (Brunnermeier & Pedersen, 2005), and many attribute the lavish surge of GME's stock price and other stocks experiencing similar price movements to retail investors betting against the short-sellers by assembling on WSB and keeping the market illiquid by holding on to their positions in the stock, leading to large losses among short-selling institutions (Hu, Jones, Zhang & Zhang, 2021). WSB now has over ten million active subscribers and is by far the most popular finance-related social media forum on the internet. However, the distinct tone with extensive sarcasm and posts containing research reports unrelated to firm-fundamentals that members on the forum use, suggests that the WSB may be less informative than other finance-related forums (Bradley, Hanousek Jr., Jame & Xiao, 2022).

What factors influence people's investment decisions? Malkiel (2005) examines this in his study and simply recommends investors, both individual and institutional, to stick to index funds that comprise the market portfolio. He argues that because institutional investors rarely outperform the market, the market should be considered as efficient, with stock prices behaving randomly since news is unpredictable. Nofsinger (2005) offers another perspective on the topic in his research on how general optimism and pessimism are mirrored in investing decisions. He claims that the economy is a complicated system of human interactions, emphasizing the significance of social mood in order to understand investment behavior and the market. While the concept of an efficient market leaves the reader to understand market bubbles by themselves, Nofsinger (2005) associates bubbles with excessive social mood. He implies that if society's optimism rises too high, investors may overestimate their investment talents while underestimating risks, leading to corporate overinvestment and market bubbles.

With the tremendous growth of social media, capturing a wide social mood has never been easier than today. Since financial markets incorporate social mood changes more quickly than other markets (Nofsinger, 2005), a lot of research has been made on decoding investor sentiment in order to anticipate future stock movements. Textual analysis is an emerging area in finance with researchers actively examining the impact of qualitative information on equity valuations. According to studies, the language and tone used by media and corporate executives has an impact on stock prices (Loughran & McDonald, 2016), and academics have sought to understand investor mood and its impact on stock prices using several methodologies. The most frequent approaches to textual analysis in the field are distributing

texts based on machine learning or a dictionary (Kearney & Liu, 2014), and the evidence of its association with stock prices is rather convincing.

Today, a majority of research shows evidence of the predictive power that investor sentiment on social media platforms have on the stock market (Antweiler & Frank, 2004; Das & Chen, 2007; Bollen, Mao & Zeng, 2011; Corbet et al., 2022; Bradley et al., 2022; Hu et al., 2021). However, less research investigates social media's impact during rare market events. In this paper, we establish a sentiment analysis on WSB in order to look for its effect and eventual predictive power on stock prices. More specifically, we investigate if the sentiment in WSB posts, so-called *submissions*, along with user activity have a predictive power before, during and after the short squeeze of GME in late January 2021. Our research focuses on stocks that best satisfy the characteristics of a typical WSB stock by picking out a few well-mentioned stocks before, during, and after the hysteria that took place on the stock market in the beginning of 2021. We are also limiting our sentiment analysis to the titles of the submissions since the comment section has much recurrence, resulting in a lot of noise and inaccuracies in our data.

In our research, we find no evidence of a significant relationship between the sentiment or the activity on WSB and the respective stock returns during our sample period, covering January through April in 2021. Interfering with previous literature examining investor sentiment's predicting ability (Engelberg, 2008; Bollen et al., 2011; Hu et al., 2021; Corbet et al, 2022), our research can not find evidence that a positive tone translates to positive stock returns. On the contrary, these findings match other similar research that can not find such evidence with respect to individual stock returns (Antweiler & Frank, 2004; Das & Chen, 2007). By using a regression model with time series data in line with methods used in previous research (Kearney & Liu, 2014), our results support previous findings made by Bradley et al. (2022), suggesting that WSB has lost its informativeness due to the notorious market events in the first half of 2021 (Hu et al., 2021; Corbet et al, 2022).

The remainder of the paper is as follows. The primary empirical findings in earlier related work in the fields of sentiment analysis and its links to asset prices along with theory are presented in Chapter 2. The collection and processing of WSB data and financial data are discussed in Chapter 3. The major methodology employed in this research is described in Chapter 4, together with diagnostic checks before running the regression. The results of our

research, as well as a discussion of these results, potential improvements and further research are presented in Chapter 5. The final chapter 6, discusses our conclusions and summarizes our findings, as well as recommendations on further research.

2. Theoretical Framework and Previous Research

This chapter provides the reader with a theoretical viewpoint, beginning with a review of the theoretical framework upon which this thesis is based, and then moving on to theories about behavioral finance and sentiment analysis on social media. Finally, we provide previous findings and conclude the theory.

2.1 Theoretical Framework

2.1.1 Random Walk Theory and Efficient Market Hypothesis

In a study comparing the performance of actively managed funds versus passively managed funds in the United States, evidence shows that only 23 percent of actively managed funds outperformed the average return of passive funds from 2009 to 2019 (Johnson, 2019). The study backs up the *Random Walk Theory* (RWT), saying that financial assets behave randomly and that it is impossible to earn a higher return on the stock market consistently without taking on more risk. This means that neither analysis of companies' financial information, known as fundamental analysis, nor analysis of historical prices, known as technical analysis, are valid tools for an investor to use in order to create an asset portfolio that outperforms a portfolio composed of a random selection of assets (Malkiel, 2003).

In Fama's (1965) paper on the behavior of stock prices, he finds statistical evidence that stock prices can not be predicted, which supports the theory that changes in stock price follow a random walk. These findings have been a cornerstone in the development of the *Efficient Market Hypothesis* (EMH), in which asset pricing is divided into three subcategories based on how effective prices reflect available asset information. The purest form of EMH is considered as *Strong Efficiency*, which implies that all information, public and private, is reflected in the price of an asset. It is followed by *Semi-Strong Efficiency*, which implies that all public information, new and old, is reflected in asset pricing. The last form of efficiency in EMH is *Weak Efficiency*, which follows directly from RWT by stating that investors are unable to predict future price changes using historical data (Fama, 1970).

Because of the impact that EMH and RMT have had on modern finance theory and research, the two concepts have been thoroughly researched and criticized. Critics of EMH originate in

the fact that people's interpretation of available information varies a lot, which might result in prices wandering about their true value (Fielitz, 1971). The criticism regarding RWT is similar to EMH and focuses on the stock market's long-term memory, which allows investors to predict future returns through historical price changes and patterns (Lo & MacKinlay, 2002). Objections against the traditional theories are also supported by the fact that investors relying on fundamental- and/or technical analysis have been able to attain great returns historically (Greig, 1992; Griffioen, 2003).

2.1.2 Behavioral Finance

Profiting from mispricing has paved the way for researchers looking for new methods to analyze market efficiency and people's interpretation of information. One of these theories is *Behavioral Finance* (BF), which questions the traditional assumption in EMH that all investors are rational and interpret available information impeccably. The theory tries to give an explanation to why the market at times is unable to set correct prices (Barberis & Thaler, 2003), and consist of the two cornerstones *Cognitive Psychology* and *Limits of Arbitrage*, referring to how people think and the market's inefficiency, keeping prices in a non-equilibrium condition for long periods of time. Too much belief in one's own capacity and leaning too much on recent experience are common arguments for market distortion made by retail investors, however, some market misvaluation derive from institutional supply and demand imbalances. One example of this was when Yahoo was added to the S&P 500 index in 1999, making index fund managers forced to buy the stock, driving the price up by more than 50% in a week (Ritter, 2003). However, even though market inefficiency can be explained by institutional trading, a couple of BF biases are worth presenting before attempting to clarify whether social media can in fact be used by retail investors to create above market-level returns or not (Barberis & Thaler, 2003).

Overconfidence. The tendency for a person to overestimate their talents and believe that they are a better-than-average investor is known as overconfidence bias. The bias is made up of two components: *Self-Attribution* and *Hindsight*, which refers to people's tendency to ascribe success to their own abilities while blaming failure on poor luck rather than their own incompetence (Barberis & Thaler, 2003).

Representativeness. People tend to employ the representativeness heuristic when determining the probability of an event, according to extensive studies on people's

interpretation of information. While representativeness is often a useful tool, it can also lead to serious prejudices and people frequently make the error of believing that two similar events are more correlated than they are. This can lead to another bias known as *Sample Size Neglect*, which means that by believing in the correlation of two occurrences, a person might also be overly reliant on small data, missing out on the fact that small data results are more likely to be explained by high levels of variance (Barberis & Thaler, 2003).

Conservatism. The mental process in which people cling to their previous beliefs rather than acknowledging new information is called conservatism. It leads to people being slow to respond to new information and as a result drive prices excessively high, causing the market to fall into states of inefficiency (Barberis & Thaler, 2003).

Anchoring. When people make decisions, they tend to rely too much on pre-existing information or the first information they come across. This is known as anchoring bias, and in behavioral finance it means that the reference point we have at hand, or so-called '*anchor*', has a lot of impact on our decisions, and will in many cases lead to us investing irrationally (Barberis & Thaler, 2003).

2.2 Previous Research

2.2.1 Sentiment Analysis Background

Sentiment analysis is defined as the extraction of people's opinions, attitudes, and emotions about specific entities through textual analysis (Hassan, Korashy & Medhat, 2014).

Historically, campaign managers have used sentiment analysis during elections to track voters' thoughts on various issues and reactions to speeches and debates. Another approach in which sentiment analysis has been used frequently is regarding consumer product- and service reviews (Feldman, 2013). However, during the last decade researchers have also used sentiment analysis to get a better understanding of the impact that investor sentiment has on returns. This has been accomplished through a variety of methods, and analysis of investor sentiment has in the last decades been done on news articles, company reports, company press releases, analyst reports, and internet sources such as social media (Kearney & Liu, 2014).

Engelberg (2008) investigates the relationship between earnings announcements and market returns and is one researcher who refutes EMH's argument that all available information is sustained in asset pricing. By using sentiment analysis on qualitative earnings data, Engelberg (2008) discovers evidence for the predictability of asset prices. Macskassy, Saar-Tsechansky and Tetlock (2008) find similar evidence of the importance of public information in predicting returns in their paper. They show that news articles effectively capture aspects of firms' fundamentals and thus can be quickly incorporated into stock prices by fundamental investors. Kothari, Li and Short (2009), and Huang, Teoh and Zhang (2014) are other researchers who disagree with EMH and show evidence for the possibility of achieving above-market returns by doing sentiment analysis. Kothari et al. (2009) discover that favorable reports about a firm by reports made by management, analysts, and reporters correspond to a declining risk in the stock. The paper by Huang et al. (2014) investigates similar effects and finds that a positive tone in earnings press releases affects stock returns positively, but with a delayed negative reaction in the following two quarters.

2.2.2 Sentiment Analysis on Social Media

In recent decades, the flow of information about corporations on the internet has intensified. For example, the number of messages about Amazon Inc. on Yahoo's message board increased from 70,000 to 900,000 in 1998-2005 (Das & Chen, 2007). With a 45 times increase in data flows between 2005 and 2014 (Bughin, Dhingra, Lund, Manyika, Stamenov & Woetzel, 2016), estimates indicate that the global social media market size was \$159.7 billion in 2021 and will continue to expand at a compound annual growth rate of 39% by 2026 (The Business Research Company, 2022). With the number of daily active users on Twitter doubling to 217 million in 2018-21, and around 500 million tweets being sent every day (Aslam, 2022), it is safe to conclude that the prospects are limitless if able to extract the sentiment from such a platform.

Antweiler and Frank's (2004) early work on the ability to predict stock returns by extracting sentiment from internet messages posted on Yahoo! Finance and Raging Bull is a forerunner in this field of research. The study discovers evidence that positive-toned posts can lead to negative returns, and that message's impact on stock returns should be seen as statistically but economically minor. While Antweiler and Frank (2004) examine 1.5 million text messages in 2000, Das and Chen (2007) have 145,100 messages regarding 24 tech-sector stocks present in

the Morgan Stanley High-Tech Index on Yahoo during two months in 2001, in their research sample. In the research, they find no evidence for a strong relationship between sentiment and stock prices on average for individual stocks, but find a statistical relation from sentiment to stock prices for the aggregated index. Bollen et al. (2011) research on the predicting potential of Twitter sentiment is another study that looks into the veracity of EMH and RWT, and if public sentiment has a correlation with market returns. The article examines if public mood, as represented in daily Twitter messages, can forecast movements in the *Dow Jones Industrial Average index* (DJIA) over an 11-month period in 2008. By tracking Twitter posts on a daily basis using two separate approaches with two and six dimensions respectively, they find evidence that it is possible to predict daily up and down changes in the closing values of DJIA by defining public mood.

With 430 million active users a month, the social media platform Reddit has become a significantly large stage for people with different interests to communicate with like-minded people. As one of Reddit's largest forums, *wallstreetbets* (WSB), focusing on finance-related topics has more than ten million subscribers, an intensified amount of research on the subreddit's ability to predict stock price movements have been made. Corbet et al. (2022) find in their study that the selection of stocks discussed on WSB and similar messaging boards tend to be about companies within the tech-sector. That is because these types of stocks have ambiguous and hidden internal mechanics, making it easy for rumors and disinformation to prevail, especially when it comes to potential product and technological development and advancements, with little evidence available in the public forum to contradict. The research shows that with the increase of options trading, the possibility for small groups of traders to act on misinformation distributed over the internet has increased, exposing illiquid equities to coordinated acts. Similar to Corbet et al. (2022), Bradley et al. (2022) find that posts regarding research on WSB emphasize risky investments, with high volatility and short interest. In the research, Bradley et al. (2022) find that before the short squeeze of GME in 2021, the sentiment in WSB investment research reports, so-called *Due Diligence* (DD) reports, could be used to forecast returns one month ahead. With DD reports being able to forecast media sentiment, earning surprises, and earnings forecast revisions, WSB has been able to provide useful information about the future in the past.

Other research by Hu et al. (2021) backs up the claim that WSB can be used to forecast stock returns. In their study on the effect of WSB on retail investors and short-sellers' role in price

movements, they discover that higher traffic, more positive tone in posts and comments, and higher connectedness lead to greater returns, higher retail order flow, and lower shorting flows in the future. However, while WSB activity appears to encourage retail buying behavior and discourage shorting, Hu et al. (2021) discover a contradiction in this relationship. They find that even though high WSB traffic discourages shorting, shorting flows become more informational and can predict stock returns even better during times of high WSB activity.

2.3 Summary

To summarize, research through sentiment analysis has been conducted in various fields in order to explore the relationship between the retail investors' mood and stock returns. Literature shows that the methods differ, and that different platforms' popularity during the sample period, as well as their ability to capture retail investor sentiment, play a large role in the selection. The most common analysis methods for enclosing the provided sentiment are the dictionary-based approach and machine learning. While these two methods being the most commonly used when attempting to extract investor sentiment, the linear regression model on time series data is without doubt the most commonly used model for testing the relationship (Kearney & Liu, 2014). Additionally, the evidence for sentiment analysis's predictive power in previous literature is overwhelming, but not conclusive.

Bollen et. al. (2011) find similar evidence to Antweiler and Frank's (2004) statistical findings that Yahoo!Finance during the early 2000s had an impact on stock returns. They discover that defining the sentiment can help predict up and down movements in the public sentiment on company posts on Twitter. Das and Chen (2007), on the other hand, contradict these findings in their study of internet messages on Yahoo's effect on the stock market in 2001. While the previous papers find statistical proof of a relationship, Das and Chen's (2007) study shows that a statistical relationship can only be found between sentiment and an index, DIJA in this case, and not between sentiment and individual stocks on average.

Along with the majority of the literature on online forums' effect on the stock market, several studies analyzing the sentiment on the considerably more speculative Reddit forum WSB confirms that retail investors' attitude can predict stock movement. Corbet et. al. (2022) find

evidence that an increasing volume of options being traded among a stock occurring on the forum can be explained by the stocks' level of liquidity and short interest. According to Corbet et. al (2022), WSB posts generally contain information regarding stocks in the tech industry, frequently with high short interest and limited liquidity. Publications by Hu et al. (2021) and Bradley et al. (2022) also find evidence for WSB forecasting capability, with Bradley et al. (2022) demonstrating that DD reports prior to 2021 were able to predict returns one month in advance, and Hu et al. (2021) showing that intense WSB traffic increases retail buying activity.

Traditional theories such as RWT and EMH claim that stock prices cannot be forecasted, however modern theories, such as BF, argue that this is possible because investors are unable to act rationally on available information due to multiple biases. The theories' disagreement on the price mechanism's function are surely the backdrop to studies in this discipline, and how the price of financial assets are determined will always be scrutinized.

The growing interest and activity in online stock market discussions in recent years, with the number of users on WSB increasing by 60% hitting approximately 12 million users in March 2022 (Corbet et al., 2022), are motivations behind this research. With our research we hope to contribute complementary perspectives to existing research investigating the relationship between social media and market returns.

3. Data

In this chapter, we explain how we collect submissions from wallstreetbets (WSB), process the data, and identify important factors to investigate if the submissions on WSB affect the company's stock value over the period we choose.

3.1 WSB Data (1)

We collect submissions from Reddit and its subreddit WSB using the high-level computer language *Python* and an application programming interface (API) called *Pushshift.io*. When constructing the code to download submissions from WSB, we mainly use a Python library called *Pandas*. We gather the submissions during four months, from January 4, 2021, to April 30, 2021, receiving the following information, among others, in our dataset for each submission: author, *Coordinated Universal Time* (UTC), submission-id, number of comments, score, title etc, collecting a total of ~ 650,000 submissions. Table 3.1 provides a selection of how the data is provided at first glance when downloaded.

Table 3.1

Submissions

Title	Author	Author_fullname	created_utc	ID	Num_comments	Score
\$GME 🐛	moon_buzz	t2_95dn3zt	1609819805	kqq7em	29	1
Selling \$AAPL	Medical_LSD	t2_4tdbhbaw	1609819778	kqq73p	0	1
Holy shit, buy \$AMC 🐛🐛	xhdt	t2_8j88j1ni	1609819749	kqq6q7	11	1
What do you all think about \$NOK?	jinpiss	t2_plaov	1609819678	kqq5zp	13	1
To you all dorks that dont belive in \$GME, just bought for 100k	Petty-officer4	t2_z3m1c	1609819562	kqq4s2	0	1

Description: Examples of how submissions look at first sight when downloaded from WSB via the programming interface Pushshift.io.

When the submissions have been downloaded, we continue our filtering process by choosing six companies with various market capitalizations as they have the common denominator of being among the most frequently discussed companies on WSB during our time frame¹. Our choice of only having six companies in the sample is motivated by the fact that other mentioned companies on WSB would not contribute enough submissions on a daily basis to the dataset and thus add uncertainty to our research. The firms we choose are GameStop,

¹ Datasource: <https://swaggystocks.com/>

AMC Entertainment Holdings, Tesla, Nokia, Apple, and Palantir Technologies, who together received a total of ~ 430,000 submissions during the chosen period.

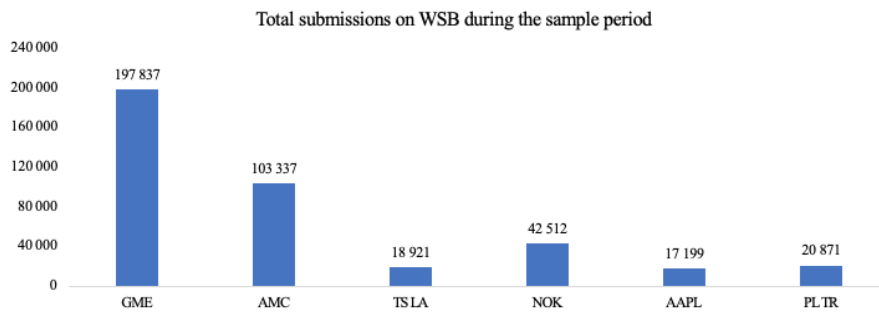
3.2 WSB Data Processing (2)

When the data is collected, we filter the submissions based on our selected company's ticker symbol, an arrangement of letters in english representing a specific asset or security listed on a stock exchange or listed publicly (e.g. GameStops ticker symbol is GME), the submission's UTC, and its title. We filter on the company's ticker in order to extract submissions regarding our selected companies as this is a commonly used method when discussing different stocks on social media and because it is in line with previous research (Mao, Liu, Wang & Wei, 2012; Challa, Majhi, Pagolu & Panda, 2016). Meanwhile, we convert UTC to *Greenwich Mean Time* (GMT) +02:00 in order to distribute the submissions later on as we will have to sort them after the market's opening hours.

Our choice of only saving the title in the dataset is due to the fact that a lot of the submissions' other qualities have been removed, with the title being the most retained part. While filtering the data, we notice that March 18th to 28th are missing, which could be due to a scraping error or to reasons relating to submissions being deleted. We double-checked the data for this problem by manually searching for submissions on WSB during these days, but did not find the missing data.

After the data thereby is filtered, the final stage before we can perform a sentiment analysis on the data is preprocessing and data cleansing. Before this final stage we look for duplicates and reduce the words to their root form (e.g. Training to Train) with the key advantage of filtering being that it reduces the amount of the data while maintaining the content that is valuable for the research. We end up with a total of 400,677 submissions after filtering for potential spam and excessive usage of cashtags, seeing that GME dominates the other selected stocks in terms of total submissions during the whole sample period, seen in Figure 3.2.

Figure 3.2



Description: WSB submissions between the 4th of January and the 30th of April 2021 sorted after ticker: 1) GME (i.e. GameStop) 2) AMC (i.e. AMC Entertainment Holdings) 3) TSLA (i.e. Tesla) 4) NOK (i.e. Nokia) 5) AAPL (i.e. Apple) 6) PLTR (i.e. Palantir Technologies).

3.3 Financial Data (3)

We use Yahoo!Finance² to get the daily closing prices for the six firms and the S&P 500 index between the 6th of January 2021 and the 30th of April 2021. After we have downloaded this data, we calculate the logarithmic returns for our selected stocks, as this is the most commonly used method in previous sentiment analysis studies (Mao et al., 2012; Bergdorf & Wolf, 2019). This is to give a consistent scale for comparing our stock market indicators and forecasters, and since it is beneficial from a numerical integration property perspective, both for time series and cross section viewpoints (Hudson & Gregoriou, 2015). The logarithmic equation is presented in Equation 3.1.

Equation 3.1

$$\text{Logarithmic stock return: } LN r_{i,t} = LN\left(\frac{\text{Close price}_{i,t}}{\text{Close price}_{i,t-1}}\right)$$

Where $LN r_{i,t}$ is the logarithmic return for company i 's closing stock price between trading day t and $t-1$.

Using logarithmic returns can, on the other hand, be misleading since returns during our time frame reach unusually high values owing to the short squeeze occurrences in some of

² Datasource: <https://finance.yahoo.com/>

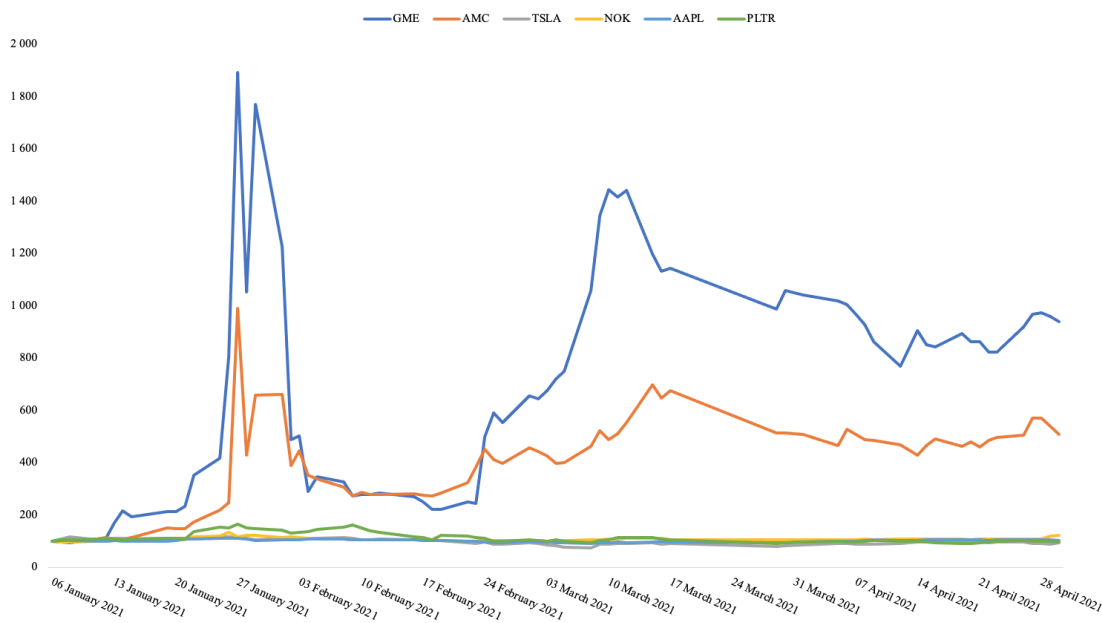
the selected companies, as shown in Figure 3.3. That is because the logarithmic formula is based on a mathematical framework called *Taylor Series*, which tells us that the logarithmic formula works best when the difference in absolute numbers between day t and $t-1$ is small. In other words, this means that the formula gets less precise as the difference between the values grows larger. Due to us presuming errors in the logarithmic formula, we choose to form another complementary equation of measuring the returns as well. As our sample period is also quite small, the two equations hint that our results may differ depending on which formula we use (Hudson & Gregoriou, 2015). The complementary equation is presented in Equation 3.2.

Equation 3.2

$$\text{Simple return: } r_{i,t} = \frac{\text{Close price}_{i,t}}{\text{Close price}_{i,t-1}} - 1$$

Where $r_{i,t}$ is the simple return for the company i 's closing stock price between trading day t and $t-1$.

Figure 3.3



Description: Stock prices between the 6th of January 6 2021 and the 30th of April 2021 for the companies, where the closing prices are indexed and set to 100 on the 6th of January. The companies are the following: 1) GME (i.e. GameStop) 2) AMC (i.e. AMC Entertainment Holdings) 3) TSLA (i.e. Tesla) 4) NOK (i.e. Nokia) 5) AAPL (i.e. Apple) 6) PLTR (i.e. Palantir Technologies).

4. Method

In this section we present more data processing, sentiment analysis, sentiment time-series building, and our empirical testing methods on our wallstreetbets (WSB) dataset. To make it easier for readers to follow, we provide an outline of the steps as illustrated in the Figure [4].

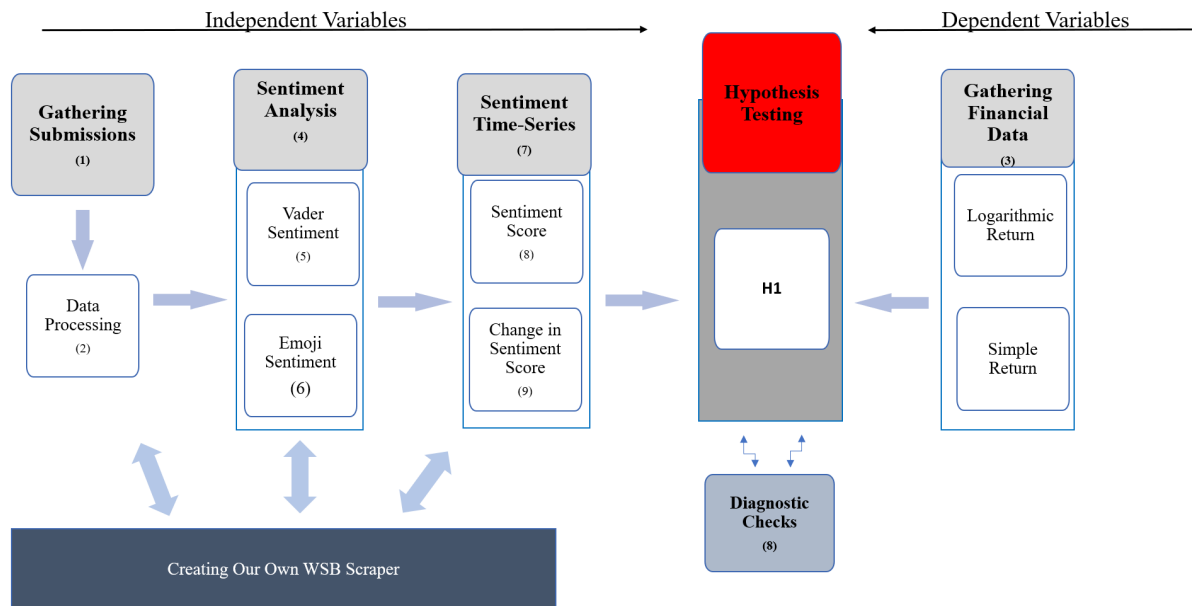
4.1 Research Questions

H1: Wallstreetbets Has a Predictive Power on Company Level Return.

The major goal of this paper is to examine the determining factors of the financial research provided on WSB during the first four months of 2021. Given the significantly growing interest of WSB in recent years, we contribute to current research on individual stock prediction using a sentiment and activity analysis. Furthermore, we want to contribute to the method of extracting sentiment because there appears to be no consensus in previous literature on which method is superior. Thus, we provide the reader with a comparison of the two different sentiment analysis methods VADER and our own built EMOJI, and a comparison in performance between the methods logarithmic and simple return.

The null hypothesis in our research is that neither the sentiment nor the activity regarding popular companies on WSB between the 4th of January and the 30th of April 2021 have predictive power on stock returns. The null hypothesis is therefore two-sided, thus we will only confirm significant coefficients at a five percentage level or lower. This is because we are not completely sure that a more positive sentiment or more activity on WSB will translate into greater stock returns after the growth in number of subscribers on the forum.

Figure 4.1: Steps of collecting and processing data



Description: The figure describes the processing of WSB and financial data before hypothesis testing. In Chapter 3.1-3.3, we describe the initial data collection process, including the gathering and processing of WSB submissions (steps 1-3) and the treatment of financial data (step 9). In the following sections, we first describe our sentiment analysis in Chapter 4.2 (steps 4-6), which is then followed by Chapter 4.3 where we construct our time-series (steps 7-9). Finally, we present the regression model we use in our hypothesis testing in Chapter 4.4 (step 10) and the diagnostic checks of the model in Chapter 4.5 (step 11).

4.2 Sentiment Analysis

4.2.1 Sentiment Analysis Methodology (4)

Sentiment analysis is a term in the field of natural language processing that refers to the process of extracting and identifying sentiment in a text using computational linguistics and textual analytics (Jacobsen & Pedersen 2021). In order to study WSB investors' thoughts on the stock market and its predictive power on popular stocks discussed on the forum, we implement a sentiment analysis on the titles we obtain from the downloaded data and filtered submissions. As mentioned in previous research, dictionary-based analysis, and machine learning are the two most common methodologies for this kind of textual sentiment

classification (Kearney & Liu, 2014). When using the former, the sentiment is classified by using a predefined dictionary, which is built on prior dynamics. The quality of the dictionary, as well as how the words are weighted and constructed, will thereafter determine the outcome. Underlying the choice of using a dictionary based method, prior studies indicate that applying machine learning algorithms over the simpler dictionary-based strategy for data classification, particularly for social media sentiment classification, offers no substantial advantage (Gilbert & Hutto, 2014). As there is little or limited data accessible, it is also difficult to train machine learning on a pre-set collection of Reddit-data. It is, however, worth mentioning that there is no consensus regarding which dictionary-based sentiment analysis technique that will perform best. Instead, what matters is what the sentiment classification's purpose is (Gilbert & Hutto, 2014).

Another commonly used dictionary is *LM* which was created by Lougrahan and McDonald (2016) and typically used in finance literature. The dictionary is an enlarged list from the Harvard/GI word list. However, the problem of implementing LM is that it was designed to analyze larger texts and hence is not ideal for WSB titles that are not very long. Another downside of using LM is that it has a hard time understanding sarcasm, and does often count positive phrases as negative and the other way around (e.g. "It is not very bad" will be counted as negative). Therefore, we do not use LM as the dictionary does not fit our assumptions about the content on WSB.

4.2.2 VADER (5)

The first method that we use to extract the sentiment from a submission's title on WSB is a dictionary called *Valence Aware Dictionary and Sentiment Reasoner* (VADER). The method has surpassed numerous well-known dictionary-based approaches as well as machine learning techniques from its very beginnings (Gilbert & Hutto, 2014). It has a number of benefits over other models and uses a mix of a sentiment lexicon and a list of lexical properties that are commonly categorized as positive or negative according to their semantic orientation on social media platforms such as Reddit. VADER will also consider capitalization writing, which alters the sentiment's power and strengthens adverbs (e.g. Incredibly good) (Bergdorf & Wolf, 2019).

Because of the numerous positive elements of VADER, we integrate this method into our code. When running VADER through our submission titles, the method analyzes the sentences as a whole, resulting in us receiving a sentiment score for each submission between -1 and +1, as seen in the examples in Table 4.1. The interpretation of the score is that a negative value indicates a negative sentiment for that submission, while a positive score means that the submission title is positive (Gilbert & Hutto, 2014).

Table 4.1

Submissions			
Title	GMT	\$Ticker	VADER sentiment score
Buy \$AMC LFG go to the moon	2021-01-28 14:50:40	\$GME	0.77
SELL \$AAPL Dorks	2021-01-28 14:50:39	\$AAPL	-0.52
\$GME 🚀🚀🚀	2021-01-28 20:58:38	\$AMC	0.00

Description: Examples of how the submissions look like when filtered and given a sentiment score by VADER. The closer a value gets to 1, the more bullish is the title, while a value closer to -1 indicates bearish content in the title.

4.2.3 EMOJI (6)

Since the tone and language on WSB is persuasive and expressive where users typically are sarcastic and use emojis frequently to convey a state of mind (e.g. 🚀 is a bullish expression for going to the moon, implying that the stock price will rise) we choose to create a complementary dictionary to characterize the sentiment of the submissions. We choose to call this dictionary EMOJI as it only takes commonly used emojis (e.g. 🚀 is assigned with +1 to the submission's total sentiment score) and expressions (e.g. GUH, meaning that the user has lost a lot of money, is assigned with -1 to the submission's total sentiment score) to extract the emotion in a submission. By doing this, EMOJI is built to understand the sarcastic environment on WSB and the intention of EMOJI is to aid the research with a better interpretation of the mood on WSB than VADER or any other existing dictionaries can.

By the mentioned procedure, EMOJI analyzes a submission's title word by word instead of as a whole, giving a bullish word +1 to the sentiment and a bearish word -1. This results in EMOJI receiving either a positive (bullish content), neutral (both bullish and bearish content) or negative total score (bearish content), as seen in Table 4.2.

Table 4.2

Submissions				
Title	GMT	\$Ticker	EMOJI sentiment score	
Buy \$AMC LFG go to the moon	2021-01-28 14:50:40	\$GME	1	
SELL \$AAPL Dorks	2021-01-28 14:50:39	\$AAPL	-1	
\$GME 🚀🚀🚀	2021-01-28 20:58:38	\$AMC	3	

Description: Examples of how the submissions look when filtered and given a sentiment score by EMOJI. The more positive the sentiment score is, the more bullish is the title, while a more negative value indicates more bearish content in the title (e.g Three rocket-emoji gives a score of 3)

4.2.4 Daily Sentiment Score

Since we want to use VADER and EMOJI separately, and still investigate their relative performance, we need to categorize them in equal amounts of dimensions. In previous research, a common approach is to narrow the sentiment into different categories according to moods (Kearney & Liu, 2014). We choose to categorize the submissions as bearish (i.e. negative content), neutral, and bullish (i.e. positive content) in line with previous research (Gilbert & Hutto, 2014). To do this, we give submissions with a VADER score less than -0.2 a value of -1 indicating bearish content, a score above 0.2 is given a score of 1, and a score in between is given a score of 0. The result of this is seen in Table 4.3. To categorize EMOJI similarly, a positive total score for a submission is translated into a sentiment score of 1, a neutral total score becomes 0, and a negative total score becomes -1, as seen in Table 4.3. Since VADER is able to evaluate broader financial content, and EMOJI is able to recognize emotions and the sarcastic environments in the submissions more effectively, we assume that the two methods will provide information that will be interesting to analyze.

Table 4.3

Submissions				
Title	GMT	\$Ticker	VADER sentiment score	EMOJI sentiment score
\$GME 🚀	2021-02-02 21:58:44	\$GME	0	1
Selling \$AAPL	2021-02-03 21:58:43	\$AAPL	-1	-1
Holy shit, buy \$AMC 🚀🚀	2021-02-04 21:58:43	\$AMC	1	1
What do you all think about \$NOK?	2021-02-05 21:58:40	\$NOK	0	0
To you all dorks that dont belive in \$GME, just bought for 100k	2021-02-06 21:58:38	\$GME	1	0

Description: Final view of the submissions when classified as bullish (+1), bearish (-1) or neutral (0) by EMOJI and VADER.

4.3 Sentiment time-series (7)

After giving a sentiment score to each submission for both methods, we must create a sentiment time-series of daily observations in order to sort the submissions according to the market's opening hours. To do this, we must define the time-thresholds before we can create the time series. When it comes to defining time-thresholds, there has been no consensus in past literature (Kearney & Liu, 2014). As a result, we will define it as the stock market's closing hours, which means that all submissions received after the New York stock exchange closes (22:00:00 GMT +01.00) will be rolled into the next trading day. When it comes to weekends and national holidays affecting the opening hours, we simply take the entire time period into consideration and roll it into when the stock exchange opens (e.g. submissions after Friday 22:00 will be counted for on Monday GMT +01:00). The main reason for this is that any submission made outside of trading hours will be reflected in the next trading day and will have no effect on the past. It is also worth mentioning that we take the summertime into consideration since it changes differently in the US and will affect the stock market's closing hour between the 14th and the 28th of March. During this time, we adjust the closing time (21:00:00 GMT +01.00).

To connect the submissions sentiment score to its corresponding trading day, we use an aggregation method to decide the daily sentiment score for VADER and EMOJI separately. By using this method, we incorporate neutral mood into our sentiment analysis, as this may also contribute useful information to the sentiment (Smailović, Grčar, Lavrac, Znidarsic, 2013). The equation for the daily sentiment score is presented in Equation 4.1:

Equation 4.1

$$\text{Average Daily Sentiment Score: } Score_{i,t} = \frac{1}{n} \sum_{j=1}^n s_{i,t,j}$$

Where $s_{i,t,j}$ represents the score for a submission j with ticker i during trading day t , and n represents the total number of submissions for ticker i during that day.

With Equation 4.1 producing the average daily sentiment score for each one of our companies, we want to examine the changing factor in WSB sentiment and its correlation with stock returns. Therefore, we need to define such an equation which we can use in our model to test the relationship. Due to researchers investigating the sentiment on diverse areas,

having different size data and time series, there is no distinct equation for this kind of variable (Kearney & Liu, 2014). Because of this, we choose to measure the daily change in sentiment in Equation 4.2, where we use the difference for the daily change in sentiment between day t and $t-1$ since the relative change will attain too big numbers as the sentiment sometimes approaches 0. As we also want to measure the increase and decline in the number of posts' correlation with stock returns later on, we define an equation for this in Equation 4.3, which we can use as an independent variable in our testing model.

Equation 4.2

Daily Change in Sentiment: $Sent_{i,t} = Score_{i,t} - Score_{i,t-1}$

Where $Sent_{i,t}$ measures the change in sentiment score for ticker i between day t and $t-1$.

Equation 4.3

Daily Change in Number of Submissions: $Subm_{i,t} = No. Subm_{i,t} - No. Subm_{i,t-1}$

Where $Subm_{i,t}$ measures the change in number of submissions for ticker i between day t and $t-1$.

4.4 Empirical Methods (11)

When we have set up the daily sentiment and compiled the amount of submissions for the companies on each trading day, we must test their relationship with stock returns. The range of methods to model this relationship is wide, and with sentiment analysis and its predictive power on stock returns being a relatively new area of research, there is no shortage of innovative approaches. However, the most common method of investigating this relationship is by using a linear regression model on time series data, while also accounting for some general market indicator (Kearney & Liu, 2014).

The regression model that is usually used in previous literature is a linear autoregressive distributed lag model with panel data, which has the advantage of incorporating delays from the sentiment into the model (Li, 2006, Das & Cen, 2007; Macskassy et al., 2008; Davis, Piger, Sedor, 2011; Doran, Peterson, McKay Price, 2012; Hu et al., 2021; Corbet et. al., 2021;

Bradley et. al., 2022). Therefore, we choose to form the regression Equation 4.4, where the S&P 500 index is included as an approximation of the market portfolio, and used as an independent variable with one time lag.

Equation 4.4

$$r_{i,t} = \alpha_i + \beta_i^r r_{t-1,i} + \sum_{j=0}^1 \beta_{i,t-j}^{Sent} Sent_{i,t-j} + \sum_{j=0}^1 \beta_{i,t-j}^{Subm} Subm_{i,t-j} + \sum_{j=0}^1 \beta_{t-j}^{rm} rm_{t-j} + \varepsilon_{i,t}$$

where $r_{t,i}$ is the return of the stock belonging to company i between day t and $t-1$, $Sent_{i,t}$ is the daily change in sentiment score between day t and $t-1$, $Subm_{i,t}$ is the daily change in number of submissions between day t and $t-1$, rm_t is the market return from the S&P 500 index between day t and $t-1$, and $\varepsilon_{i,t}$ the error term in the regression model for ticker i during day t .

We choose to regress the model with only one lag because WSB, like other frequently used forums, is an internet forum with a high frequency of posts, with popular submissions being exchanged on a daily basis, falling quickly out of readers' eyes. Therefore, we do not assume that the daily sentiment from more than one day prior to the daily return will have an impact on daily returns. In the results, we expect one of the coefficients $Sent_{t,i}$, $Subm_{t,i}$ to be significant if mood and/or activity on WSB will have any significant relationship with stock returns. The regressions we will perform in order to evaluate our hypothesis are presented in Table 4.4.

Table 4.4

<i>Regression</i>	1)	2)	3)	4)
<i>Dependent variable</i>	Stock return, t		LN Stock return, t	
<i>Independent variable</i>				
LN Stock return, t-1			x	x
Stock return, t-1	x	x		
VADER, t	x		x	
VADER, t-1	x		x	
EMOJI, t		x		x
EMOJI, t-1		x		x
# of submissions, t	x	x	x	x
# of submissions, t-1	x	x	x	x
LN S&P Return, t			x	x
LN S&P Return, t-1			x	x
S&P Return, t	x	x		
S&P Return, t-1	x	x		

Description: The four different regressions we are going to investigate where Columns 1) and 2) examine the predictive power of our sentiment analysis methods on stock returns when stock returns are measured using a simple return method, whereas Columns 3) and 4) apply the logarithmic return approach. The x:s mark if the variable is included in the regression equation.

4.5 Diagnostic Checks (10)

Since our research uses a linear autoregressive distributed lag model with multiple explanatory variables, it is necessary to test if the model obtains robustness. To test if the model may cause misleading results, we test for the following biases using the econometric program called *Gretl*.

- (1) Heteroscedasticity**, a form of inconsequence in the variance of an explanatory variable. It occurs when for example bigger values in the variable have greater variance in average than smaller values of the variable. In order to check for

heteroscedasticity and inefficient estimates of beta-values in the regression, we apply a test called *White's test* on the model (Waldman, 1983).

- (2) **Multicollinearity**, when there is a significant correlation between two or multiple explanatory variables. This happens when two data points after one another are strongly correlated with each other, leading to unreliable coefficient values with high levels of variances and standard errors (Mansfield & Helms, 1982).
- (3) **Autocorrelation**, the degree of similarity between a given time series and the lagged version of itself during a period. Since we are dealing with a distributed lag model in our regression, we need to test this bias, and we do this by using the so-called *Durbin Watson test* (Savin & White, 1977).
- (4) **Stationarity**, when a data set does not have any trend during the sample period. This is a common assumption when using time series data sets and it means that a stationary process' mean, variance and autocorrelation structure do not change over time. To test if the data series are stationary we perform a panel unit root test called the *Levin Lin Chu test* instead of the more commonly used unit root test called *Augmented Dickey Fuller test*. This is because our model is dealing with panel data and therefore needs a panel root test to examine this bias (Barbieri, 2005).

Looking at the data from the period and the research questions, the possibility that the selected stocks suffer from the mentioned distortions at times during the sample period and that we need to adjust the regression is rather high. This is due to the fact that some of the stocks experience massive changes in stock returns, as well as sentiment score and volume of posts, which motivates the tests

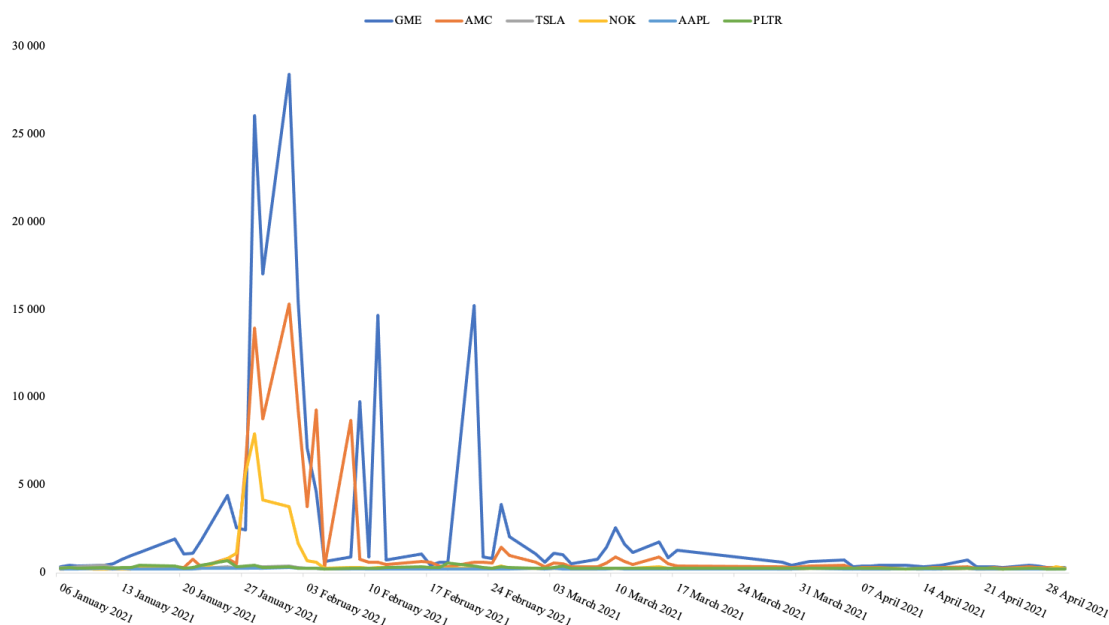
5. Analysis & Discussion

The following chapter analyzes and discusses the outcome of this paper. We present our regression analysis and discuss the results presented in this chapter and end with possible limitations.

5.1 Descriptive statistics

We can see a couple of interesting things by looking at the statistics on our sentiment analysis and the measured number of submissions across the firms. The first thing we notice about the dataset is that the number of submissions every day is particularly centered around specific dates, as shown in Figure 5.1. We can also notice the difference in the volume of conversations sparked by the various companies, with GME attracting the most interest, with up to 30,000 submissions on a few occasions, while on other days, it drops below 5,000. When compared to individual firm's stock returns, a correlation between the number of submissions and stock returns during the time frame seems to appear.

Figure 5.1



Description: In this figure, the number of collected submissions about each company between the 6th of January 2021 and the 30th of April 2021 is shown. The companies are; 1) GME (i.e. GameStop) 2) AMC (i.e. AMC Entertainment Holdings) 3) TSLA (i.e. Tesla) 4) NOK (i.e. Nokia) 5) AAPL (i.e. Apple) 6) PLTR (i.e. Palantir Technologies).

Another interesting feature of the dataset seen in Table 5.1 is that the sentiment analysis differs based on the approach used. When compared to the EMOJI method, the method of applying VADER in the analysis has a lower average sentiment across all firms (e.g 0.15 against 0.21 for GME). Despite the lower values that VADER exhibits, the gap in daily average sentiment score between VADER and EMOJI appears to remain consistent, with the average VADER sentiment changing in lockstep with EMOJI across the firms. Other than this observation, we can not observe any significant difference in standard deviation, skewness, or kurtosis. Because of the datasets' incomprehensible behavior this may not be worth studying further, but nevertheless important to note and perhaps consider when interpreting the following regressions.

Table 5.1

<i>Ticker</i>	GME	AMC	TSLA	NOK	AAPL	PLTR
# of submissions						
Average	2 786.43	1 455.45	266.49	598.76	242.23	293.95
Median	775.74	369.74	253.74	250.74	238.74	267.74
Standard deviation	5 577.39	3 098.06	38.21	1 259.23	10.48	74.28
Kurtosis	10.36	9.55	5.74	20.49	11.69	11.15
Skewness	3.21	3.16	2.32	4.42	3.30	2.95
VADER sentiment score						
Average	0.15	0.16	0.16	0.16	0.16	0.19
Median	0.15	0.17	0.15	0.12	0.18	0.18
Standard deviation	0.06	0.11	0.23	0.25	0.30	0.17
Kurtosis	-0.35	1.56	3.47	7.33	3.26	1.83
Skewness	-0.14	-0.85	0.78	-0.40	-0.60	-0.37
EMOJI sentiment score						
Average	0.21	0.22	0.23	0.24	0.23	0.28
Median	0.22	0.23	0.25	0.24	0.19	0.29
Standard deviation	0.05	0.16	0.17	0.14	0.24	0.11
Kurtosis	2.99	12.38	1.16	-0.34	0.37	0.07
Skewness	-0.94	2.51	-0.26	-0.04	0.74	-0.65

Description: Some statistics about a few independent variables in the regression model for each company in the research, which are; 1) GME (i.e. GameStop) 2) AMC (i.e. AMC Entertainment Holdings) 3) TSLA (i.e. Tesla) 4) NOK (i.e. Nokia) 5) AAPL (i.e. Apple) 6) PLTR (i.e. Palantir Technologies)..

5.2 Diagnostic Checks (10)

Before we run our regressions and look at the results, we need to clarify whether the regressions suffer from biases or not. First out is the Levin Lin Chu test, which we perform to check if the time-series satisfies the condition of stationarity. When performing this test, we have the null hypothesis that the panels contain unit roots, meaning that the data is non-stationary. In Table 5.2, we see that the t-statistic for each variable in the test obtains a rather high value, meaning that we can reject the null hypothesis and continue with our regression.

Table 5.2

Levin Lin Chu test	
Variable in regression	t-statistic
Stock return, t	-29.541
LN Stock return, t	-23.292
VADER, t	-35.542
EMOJI, t	-34.581
# of submissions, t	-28.372
S&P 500 return, t	-23.735
LN S&P 500 return, t	-23.719

Description: The t-statistic for each independent variable that is used in the regressions when runned in the Levin Lin Chu test. A high t-statistic in absolute values indicates a low p-value, meaning that we can reject the null hypothesis and confirm that our variables are stationary.

Since we are working with a time series with lag, we continue by testing the dataset for multicollinearity and autocorrelation. Firstly, we perform the multicollinearity test in Gretl, where we observe that our regression models do not suffer from the bias, shown in Table 5.3. Thereafter, we test for autocorrelation using the Durbin Watson test. In this test, a test statistic value of 2 indicates that the model does not have autocorrelation, with values less than two indicating positive autocorrelation and values above two indicating negative autocorrelation. Our test statistics show values slightly below two, seen in Table 5.4, indicating that the variables in the regressions almost have no autocorrelation at all.

Table 5.3

Collinearity	
Regression	Yes/No
1)	No
2)	No
3)	No
4)	No

Description: A test in Gretl for collinearity in the regressions, giving us information about the equation that we can interpret as 'Yes' if there is multicollinearity in the equation, and 'No' if there is not.

Table 5.4

Durbin Watson test	
Regression	test statistic
1)	1.918
2)	1.915
3)	1.929
4)	1.927

Description: The test statistic for each independent variable that is used in the regressions when runned in the Durbin Watson test. A test statistic value close to 2 translates to no autocorrelation, while a value between 0 and 2 indicate positive autocorrelation, and a value between 2 and 4 indicate negative autocorrelation.

Finally, we test our explanatory variables for heteroscedasticity, seen in 5.5, so that we can interpret the regressions correctly. The method of doing this is using White's heteroscedasticity test, in which we have the null hypothesis that the regression is homoscedastic. Since all regressions obtain a p-value above our significance level of 0.05, we can accept our null hypothesis and conclude that the dataset is homoscedastic. This indicates that the variances of the error terms in the regression equations are constant.

Table 5.5

White's test	
Regression	p-value
1)	0.076
2)	0.085
3)	0.110
4)	0.121

Description: The p-value for each independent variable that is used in the regressions when runned in White's test. Here, a p-value above 0.05 indicates that we can accept the null hypothesis at a 5 percent significance level and confirm that the model is homoscedastic.

5.3 WSB sentiment and Stock Returns

H1: WSB activity and sentiment and stock returns (9)

In Table 5.6, we see our four different regressions. The coefficients of the independent variables are displayed in the table, with a parameter value of -0.04082 indicating that for every one unit increase in that variable, the dependent variable decreases by 0.04082 units. Looking at our obtained p-values, we can see that neither VADER, EMOJI, nor the amount of submissions appear to be significant, meaning that we can conclude that our sentiment analysis does not have predictive power over the stock returns of the companies we use in our sample. However, in Table 5.6 we see that our self-built EMOJI sentiment analysis seems to be more significant than VADER.

In terms of our investigation regarding whether we get different outcomes when using simple returns instead of logarithmic returns or not, we find evidence that the logarithmic method performs slightly better. Looking at the p-values in Table 6, we see that the stock returns on day $t-1$ are significant for the returns on day t , apart from also showing evidence of the S&P 500 return on day t 's significance similar to the simple return method.

Table 5.6

Regression	1)	2)	3)	4)
Dependent variable	Stock return, t		LN Stock return, t	
Independent variable				
LN Stock return, t-1			-0.1114* (0.0502)	-0.1148** (0.0436)
Stock return, t-1	-0.04082 (0.4642)	-0.04290 (0.4418)		
VADER, t	-0.002176 (0.9573)		0.004395 (0.8751)	
VADER, t-1	0.01410 (0.7283)		0.009458 (0.7348)	
EMOJI, t		0.04658 (0.4002)		0.04850 (0.2028)
EMOJI, t-1		0.02915 (0.5922)		0.02335 (0.5329)
# of submissions, t	-1.254e-07 (0.9815)	-1.297e-07 (0.9809)	-2.626e-06 0.4923	-2.622e-06 (0.4922)
# of submissions, t-1	-2.3103e-06 (0.6331)	-2.402e-06 (0.6193)	-4.656e-06 (0.1592)	-4.743e-06 (0.1509)
LN S&P Return, t			-2.260*** (0.0031)	-2.286*** (0.0027)
LN S&P Return, t-1			-0.06201 (0.9363)	-0.1065 (0.8907)
S&P Return, t	-3.553*** (0.0014)	-3.572*** (0.0013)		
S&P Return, t-1	-0.1990 (0.8602)	-0.2505 (0.8246)		

Description: In the table, the values in ordinary style are the coefficient values for the regressions, while the values standing underneath in parenthesis and with cursive style are the p-value for the coefficient value above. The table presents the four different regressions we can see in Table 4.4, and the number of '*' represents the level of significance for the coefficient value.

5.4 Discussion

Our results indicate that daily sentiment and daily sentiment with one lag have no significant impact on predicting ability, however an increase in average daily sentiment score on days $t-1$ and t has a positive but not significant effect on the stock return on day t for the selected stocks. However, these findings back up previous research that claims investors concentrate too much on coordinated trading tactics during the hectic period in the beginning of 2021, presumably at the expense of studying corporate fundamentals (Bradley et al. 2022)

Another noteworthy takeaway from the findings is that our strategies for determining sentiment seem to work differently. Meanwhile EMOJI does not show any significant level of predicting power across the selected firms, it is possible that the method is better at capturing the investor mood on WSB since it provides a more significant p-value than VADER, see Table 5.6. Since EMOJI only pays attention to expressions and emotions rather than financial information and is having these results compared to VADER, we may confirm a few hypotheses about WSB. With EMOJI merely accounting for popular finance-related emojis and abbreviations and still showing a tendency of higher significance, we believe this strengthens the preconceptions about the expressive tone on WSB, and that this is an approach to consider onwards for further research on the platform.

Part of the findings in this work are consistent with some previous studies, implying that WSB users' informativeness declines with time as a result of the platform's massive growth, which has resulted in changes in the substance of reports, lowering the value of relevance (Bradley et al. 2022). However, our findings disagree with other previous research, which suggests that the volume of messages on social media has a significant impact on stock returns and that the positive influence of social media is related with better future returns (Antweiler & Frank, 2004). The disparity in results could be explained by the fact that latter research looks at the relationship across time frames that are not affected by unexpected market events like short squeezes and major movements which affect the outcome in informativeness and thus change the outcome of the result.

With previous articles concluding that a change in sentiment can be used to predict stock returns, this validates that the random walk theory (RWT) and the efficient market hypothesis (EMH) suffer from inaccuracies. While RWT may be stable over longer periods of time

(Johnson, 2019), research suggests that stock returns can be predicted by investor mood during shorter periods of time (Engelberg, 2008; Macskassy et al., 2008; Kothari et al., 2009; Huang et al., 2014; Bollen et al., 2011; Corbet et al., 2022; Hu et al., 2021), contradicting both RWT and EMH by showing evidence that public and private information regarding stocks can be extracted from investor sentiment and used to achieve profits. On the other hand, our research follows other research showing that individual stock returns can not be predicted by the investor sentiment (Antweiler & Frank, 2004; Das & Chen, 2007).

Our findings show that a rational investor will not be able to reach greater returns by taking advice from a positive sentiment in WSB submissions, which contain difficult-to-interpret knowledge regarding available information and future occurrences. This is in line with previous research on the subject giving evidence to WSB being less informative during the first half of 2021 than before (Bradley et al., 2022). The fact that retail investors should not consider WSB to be anything more than an uninformative discussion forum can be read as supporting evidence for the rise of overconfidence, representativeness, conservatism, and anchoring among investors on WSB. However, because this study does not provide information on how investors react to WSB content, it cannot be considered as complete evidence for these Behavioral Finance (BF) biases.

While a majority of the selected stocks in this research is meeting the criterias of having high volatility (Figure 3.3), and being well-mentioned on WSB throughout our time frame (Figure 5.1), there will still be a chance for investors who understand how to evaluate and incorporate WSB's content into their investment decisions to make significant profits from WSB. However, with this study concluding that neither an increase in the number of posts, nor a more positive tone, lead to higher stock returns, one could argue that investors should not be taking advice from WSB. With WSB developing into a world wide discussion hub for investors to share information about certain companies in order to get large stock returns, the platform can be considered to have been transformed into a more publicly available forum, leading to more noise and less informativeness. Thus, one could claim that people taking advice from WSB suffer from BF biases.

During the last years, online communication tools have increased rapidly in terms of usage and technology, where information flows are much faster and more efficient today than they were before. With WSB reaching over ten million users and becoming a widely known forum

during 2021, one may claim that the information in submissions on the platform could be seen as public information and that it has already been incorporated in stock prices, in contrast to research on the forum before 2021 (Corbet et al., 2022; Hu et al., 2021). With our study showing evidence that you can not increase returns by investing based on the sentiment nor the activity on WSB in line with Bradley et al. (2022), one could also argue that public information is already incorporated in stock prices regarding stocks being mentioned on platforms with a lot of users. Moreover, this indicates that the stock market regarding these stocks shows a tendency of Semi-Strong Efficiency and confirming the Random Walk Theory (RWT). On the other hand, previous findings oppose this statement by saying that the sentiment on WSB could predict stock returns during times when it had less subscribers. This leaves us with the outcome that social media platforms focusing on investments may work as an informative forum while they are still uncharted publicly, but that they tend to become more driven by opinions and emotions as it attracts more users.

5.5 Limitations

To end this section, we want to discuss the potential limitations of this research. Firstly, we are only investigating a short number of companies who are among the top mentioned companies on the platform during the period, which might influence our results as the dataset becomes smaller.

Another limitation is the short time frame, in which we look at some of the months when WSB was most energetic and influenced by the attention of millions of people. Another constraint is the absence of some days in March, which may have an impact on our findings.

Lastly, 2021 was a historically good year for the stock market, with the S&P index rising 26.89 percent (Miao & Macheel, 2021). This could have altered our results and their ability to be extrapolated to other time periods, and it would therefore be interesting to look into the results if the market was more negative, such as in early 2022 when the index was down -8.80% in just April (Silverblatt, 2022).

6. Conclusion

The purpose of this study is to see if *wallstreetbets* (WSB) sentiment may accurately forecast company-level returns. Since WSB is a relatively new platform, often neglecting the fundamentals of investing, there is little previous sentiment analysis research upon this. However, in line with existing research on other social media platforms, we present two different sentiment analysis approaches VADER and EMOJI and can indicate that both methods fail to produce significant results. However, our findings endorse previous research's results on WSB with the extremely strong increase in users in the beginning of 2021 leading to noise and coordination of trading tactics (Bradely et al. 2022)

The primary takeaway from our findings is that sentiment over the chosen period gives useful information about retail investors behavior due to the number of subscribers, and that this can be built upon in further research. From our results, we can also conclude that it is possible that EMOJI captures the investor mood on WSB better than VADER since it provides more significant p-values in the regression models.

In recent years, WSB has seen a huge increase in users, and it continues to develop in tandem with the growth and simplification of buying and selling securities online. As the increase in users may lead to a more efficient market where the information will be spread faster, this may also lead to more noise, raising the difficulty of interpreting WSB content as more irrelevant and rubbish posts will be present. We, however, believe and hope that our research will be useful for developing future investment techniques that incorporate social media sentiment, particularly on WSB. For further research, would it be interesting to perform research over a longer timespan to see whether WSB has predictive power on the returns on company level during events of less volatility. To be able to detect the sentiment on WSB in a better way, there is always room for improvement in the dictionary analysis, where future researchers may be able to provide a better dictionary to scrape the sentiment resulting in a significant result.

Lastly, we find it interesting to further investigate the various associated companies' short interest. Since the users on WSB frequently seek to discuss heavily shorted stocks in order to assemble their investments and thereby reduce the profit from different hedgefunds.

Bibliography

Antweiler, W., Frank, M. Z. (2004). Is All That Talk Just Noise? The Information Content of Internet Stock Message Boards, *The Journal of Finance*, vol. 59, no. 3, pp. 1259-1294

Aslam, S., 2022, *Twitter by the Numbers: Stats, Demographics & Fun Facts*, Available online: <https://www.omnicoreagency.com> [Accessed 7 May 2022]

Barberis, N., Thaler, R. (2003), A survey of behavioral finance, *Handbook of the Economics of Finance*, vol. 1, no. 1, pp. 1053-1128

Barbieri, L. (2005). Panel Unit Root Tests: A Review, Working Paper, *Catholic University of the Sacred Heart*, January 2005

Bergdorf, O., Wolf, F. (2019). Twitter Sentiment and Stock Returns, Master's Thesis, *Lund University*, 8 August 2019

Bollen, J., Mao, H., Zeng, X. (2011). Twitter mood predicts the stock market, *Journal of Computational Science*, vol. 2, no. 1, pp. 1-8

Bradley, D., Hanousek Jr., J. Jame, R., Xiao, Z. (2022). Place your bets? The market consequences of investment research on Reddit's Wallstreetbets*, Working Paper, University of South Florida, March 2022

Brunnermeier, M. K., Pedersen, L. H. (2005). Predatory Trading, *The Journal of Finance*, vol. 60, no. 4, pp. 1825-1863

Bughin, J., Dhingra, D., Lund, S., Manyika, J., Stamenov, K., Woetzel, J., (2016). *Digital Globalization: The New Era of Global Flows*, McKinsey Global Institute

Challa, K. N. R., Majhi, B., Pagolu, V. S., Panda, G. (2016). Sentiment Analysis of Twitter Data for Predicting Stock Market Movements, Working paper, *Indian Institute of Technology*, 28 October 2016

Chohan, U. W. (2021). Counter-Hegemonic Finance: The Gamestop Short Squeeze, Working Paper, *UNSW Business School*, 4 February 2021

Corbet, S., Hou, Y., Hu, Y., Oxley, L. (2022). We Reddit in a forum: The influence of messaging boards on firm stability, *Review of Corporate Finance*, vol. 2, no. 1, pp. 151-190

Dangmei, J., Singh, Dr. A. P. (2016). Understanding the Generation Z: The Future Workforce, *South Asian Journal of Multidisciplinary Studies*, vol. 3, no. 3, pp. 1-5

Das, S. R., Chen, M. Y. (2007). Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web, *Management Science*, vol. 53, no. 9, pp. 1375-1388

Davis, A. K., Piger, J., Sedor, L. M. (2011). Beyond the Numbers: Measuring the Information Content of Earnings Press Release Language, *Contemporary Accounting Research*, vol. 29, no. 3, pp. 845-868

Doran, J. S., Peterson, D. R., McKay Price, S. (2012). Earnings Conference Call Content and Stock Price: The Case of REITs, *The Journal of Real Estate Finance and Economics*, no. 45, pp. 402-434

Engelberg, J. (2008). Costly Information Processing: Evidence from Earnings Announcements, Working Paper, *University of California*, 21 March 2008

Fama, E. F. (1965). The Behavior of Stock-Market Prices, *The Journal of Business*, vol. 38, no. 1, pp. 34-105

Fama, E. F. (1970). Efficient Capital Markets: A Review of Theory and Empirical Work, *The Journal of Finance*, vol. 25, no. 2, pp. 383-417

Feary, R., Sharma, C., Franco, M., Thrasher, C., 2022, *The Investor Landscape: four evolving themes and their implications*, Available online: <https://www.credit-suisse.com> [Accessed 5 May 2022]

Feldman, R. (2013). Techniques and applications for sentiment analysis, *Communications of the ACM*, vol. 56, no. 4, pp. 82-89

Fielitz, B. D. (1971). On the Random Walk Hypothesis, *The American Economist*, vol. 15, no. 1, pp. 105-107

Gilbert, E., Hutto, C. (2014). VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text, *Proceedings of the Eight International AAAI Conference on Weblogs and Social Media*, vol. 8, no. 1, pp. 216-225

Greig, A. C. (1992). Fundamental analysis and subsequent stock returns, *Journal of Accounting and Economics*, vol. 15, no. 2-3, pp. 413-442

Griffioen, G. A. W. (2003). Technical Analysis in Financial Markets, Working Paper, *University of Amsterdam*, 3 March 2003

Hassan, A., Korashy, H., Medhat, W. (2014). Sentiment analysis algorithms and applications: A survey, *Ain Shams Engineering Journal*, vol. 5, no. 4, pp. 1093-1113

Hu, D., Jones, C. M., Zhang, V., Zhang, X. (2021). The Rise of Reddit: How Social Media Affects Retail Investors and Short-sellers' Roles in Price Discovery, Working Paper, Northwestern University, 2 April 2021

Huang, X., Teoh, S. H., Zhang, Y. (2014). Tone Management, *The Accounting Review*, vol. 89, no. 3, pp. 1083-1113

Hudson, R. S., Gregoriou, A. (2015). Calculating and comparing security returns is harder than you think: A comparison between logarithmic and simple returns, *International Review of Financial Analysis*, vol. 38, pp. 151-162

- Jacobsen, T., Pedersen, T.F. (2021) *WallstreetBets on Wall Street: An Empirical Analysis of the Market Power of WallstreetBets*, Norwegian School of Economics Bergen, 9 September 2021
- Johnson, B. (2019). *Active Funds vs. Passive Funds: Which Fund Types Had Increased Success Rates?*, Available online: <https://www.morningstar.com> [Accessed 3 April 2022]
- Kearney, C., Liu, S. (2014). Textual sentiment in finance: A survey of methods and models, *International Review of Financial Analysis*, vol. 33, pp. 171-185
- Kothari, S. P., Li, X., Short, J. E. (2009). The Effect of Disclosures by Management, Analysts, and Business Press on Cost of Capital, Return Volatility, and Analyst Forecasts: A Study Using Content Analysis, *The Accounting Review*, vol. 84, no. 5, pp. 1639-1670
- Li, F. (2006). Do Stock Market Investors Understand the Risk Sentiment of Corporate Annual Reports?, Working Paper, *Shanghai Jiaotong University*, 26 April 2006
- Lo, A. W., MacKinlay, A. C. (2002). *A Non-Random Walk Down Wall Street*, Princeton: Princeton University Press
- Loughran, T., McDonald, B. (2016). Textual Analysis in Accounting and Finance: A Survey, *Journal of Accounting Research*, vol. 54, no. 4, pp. 1187-1230
- Macskassy, S., Saar-Tsechansky, M., Tetlock, P. C. (2008). More Than Words: Quantifying Language to Measure Firms' Fundamentals, *The Journal of Finance*, vol. 63, no. 3, pp. 1437-1467
- Malkiel, B. G. (2003). The Efficient Market Hypothesis and Its Critics, *The Journal of Economic Perspectives*, vol. 17, no. 1, pp. 59-82
- Malkiel, B. G. (2005). Reflections on the Efficient Market Hypothesis: 30 Years Later, *Financial Review*, vol. 40, no. 1, pp. 1-9
- Mansfield, E. R., Helms, B. P. (1982). Detecting Multicollinearity, *The American Statistician*, vol. 36, no. 3a, pp. 158-160
- Mao, Y., Liu, B., Wang, B., Wei, W. (2012). Correlating S&P 500 stocks with Twitter data, Working Paper, University of Connecticut, August 2012
- Miao, H., Macheel, T., 2021, *S&P 500 ends 2021 with a nearly 27% gain, but dips in final trading day*, Available online: <https://www.cnbc.com> [Accessed 6 May 2022]
- Nofsinger, J. R. (2005). Social Mood and Financial Economics, *Journal of Behavioral Finance*, vol. 6, no. 3, pp. 144-160
- Ritter, J. R. (2003). Behavioral finance, *Pacific-Basin Finance Journal*, vol. 11, no. 4, pp. 429-437

Savin, N. E., White, K. J. (1977). The Durbin-Watson Test for Serial Correlation with Extreme Sample Sizes or Many Regressors, *The Econometric Society*, vol. 45, no. 8, pp. 1989-1996

Silverblatt, H. (2022). *U.S. Equities Market Attributes April 2022*, S&P Dow Jones Indices, Available online: <https://www.spglobal.com> [Accessed 3 May 2022]

Smailović, J., Grear, M., Lavrac, N., Znidarsic, M. (2013). Predictive Sentiment Analysis of Tweets: A Stock Market Application, *Human-Computer Interaction and Knowledge Discovery in Complex*, vol. 7947, pp. 77-88

The Business Research Company, 2022, *Social Media Global Market Report 2022*, The Business Research Company, Available online: <https://www.thebusinessresearchcompany.com> [Accessed 5 May 2022]

Waldman, D. M. (1983). A note on algebraic equivalence of White's test and a variation of the Godfrey/Breusch-Pagan test for heteroscedasticity, *Economic Letters*, vol. 13, no. 2-3, pp. 197-200

Appendix

Regression Results

Table A.1

Model 1: Pooled OLS, using 426 observations				
Included 6 cross-sectional units				
Time-series length = 71				
Dependent variable: Returnt				
	coefficient	std. error	t-ratio	p-value
const	0.0312050	0.0105123	2.968	0.0032 ***
Returnt1	-0.0408221	0.0557241	-0.7326	0.4642
RVt	-0.00217584	0.0405922	-0.05360	0.9573
RVt1	0.0141036	0.0405669	0.3477	0.7283
RVPt	-1.25373e-07	5.41546e-06	-0.02315	0.9815
RVPt1	-2.30913e-06	4.83399e-06	-0.4777	0.6331
SPReturnt	-3.55270	1.10367	-3.219	0.0014 ***
SPReturnt1	-0.199031	1.12916	-0.1763	0.8602
Mean dependent var	0.023933	S.D. dependent var	0.207242	
Sum squared resid	17.75859	S.E. of regression	0.206118	
R-squared	0.027114	Adjusted R-squared	0.010821	
F(7, 418)	1.664207	P-value(F)	0.116044	
Log-likelihood	72.35456	Akaike criterion	-128.7091	
Schwarz criterion	-96.27360	Hannan-Quinn	-115.8964	
rho	0.010695	Durbin-Watson	1.917786	
Excluding the constant, p-value was highest for variable 23 (RVPt)				

Description: Regression results for regression equation 1) (see table 4.4). Returnt1 stands for simple return between day t-2 and t-1, RVt (RVt1) stands for daily change in VADER sentiment score between day t (t-1) and t-1 (t-2), RVPt (RVPt1) stands for daily change in number of submissions between day t (t-1) and t-1 (t-2), and SPReturnt (SPReturnt1) stands for S&P 500 index simple return between day t (t-1) and t-1 (t-2).

Table A.2

Model 2: Pooled OLS, using 426 observations				
Included 6 cross-sectional units				
Time-series length = 71				
Dependent variable: Returnt				
	coefficient	std. error	t-ratio	p-value
const	0.0314726	0.0105109	2.994	0.0029 ***
Returnt1	-0.0428919	0.0557097	-0.7699	0.4418
RVPt	-1.29743e-07	5.41163e-06	-0.02397	0.9809
RVPt1	-2.40238e-06	4.83143e-06	-0.4972	0.6193
SPReturnt	-3.57244	1.10341	-3.238	0.0013 ***
SPReturnt1	-0.250488	1.12943	-0.2218	0.8246
REt	0.0465846	0.0553136	0.8422	0.4002
REt1	0.0291546	0.0543841	0.5361	0.5922
Mean dependent var	0.023933	S.D. dependent var	0.207242	
Sum squared resid	17.73560	S.E. of regression	0.205985	
R-squared	0.028374	Adjusted R-squared	0.012102	
F(7, 418)	1.743786	P-value(F)	0.097212	
Log-likelihood	72.63054	Akaike criterion	-129.2611	
Schwarz criterion	-96.82556	Hannan-Quinn	-116.4484	
rho	0.011434	Durbin-Watson	1.915284	
Excluding the constant, p-value was highest for variable 23 (RVPt)				

Description: Regression results for regression equation 2) (see table 4.4). Returnt1 stands for simple return between day t-2 and t-1, REt (REt1) stands for daily change in EMOJI sentiment score between day t (t-1) and t-1 (t-2), RVPt (RVPt1) stands for daily change in number of submissions between day t (t-1) and t-1 (t-2), and SPReturnt (SPReturnt1) stands for S&P 500 index simple return between day t (t-1) and t-1 (t-2).

Table A.3

Model 3: Pooled OLS, using 426 observations				
Included 6 cross-sectional units				
Time-series length = 71				
Dependent variable: LNReturnt				
	coefficient	std. error	t-ratio	p-value
const	0.0143581	0.00717549	2.001	0.0460 **
RVPt	-2.62597e-06	3.82130e-06	-0.6872	0.4923
RVPt1	-4.65609e-06	3.30176e-06	-1.410	0.1592
RVt	0.00439453	0.0279297	0.1573	0.8751
RVt1	0.00945796	0.0278986	0.3390	0.7348
LNSPReturnt	-2.25991	0.758891	-2.978	0.0031 ***
LNSPReturnt1	-0.0620085	0.775712	-0.07994	0.9363
LNReturnt1	-0.111392	0.0567175	-1.964	0.0502 *
Mean dependent var	0.009496	S.D. dependent var	0.143629	
Sum squared resid	8.403076	S.E. of regression	0.141785	
R-squared	0.041556	Adjusted R-squared	0.025505	
F(7, 418)	2.589061	P-value(F)	0.012679	
Log-likelihood	231.7364	Akaike criterion	-447.4729	
Schwarz criterion	-415.0374	Hannan-Quinn	-434.6602	
rho	0.033900	Durbin-Watson	1.929123	
Excluding the constant, p-value was highest for variable 27 (LNSPReturnt1)				

Description: Regression results for regression equation 3) (see table 4.4). LNReturnt1 stands for simple return between day t-2 and t-1, RVt (RVt1) stands for daily change in VADER sentiment score between day t (t-1) and t-1 (t-2), RVPt (RVPt1) stands for daily change in number of submissions between day t (t-1) and t-1 (t-2), and LNSPReturnt (LNSPReturnt1) stands for S&P 500 index simple return between day t (t-1) and t-1 (t-2).

Table A.4

Model 4: Pooled OLS, using 426 observations				
Included 6 cross-sectional units				
Time-series length = 71				
Dependent variable: LNReturnt				
	coefficient	std. error	t-ratio	p-value
const	0.0145721	0.00716592	2.034	0.0426 **
RVPt	-2.62209e-06	3.81477e-06	-0.6874	0.4922
RVPt1	-4.74319e-06	3.29599e-06	-1.439	0.1509
LNSPReturnt	-2.28597	0.757831	-3.016	0.0027 ***
LNSPReturnt1	-0.106543	0.774986	-0.1375	0.8907
LNReturnt1	-0.114772	0.0567098	-2.024	0.0436 **
REt	0.0485046	0.0380203	1.276	0.2028
REt1	0.0233535	0.0374168	0.6241	0.5329
Mean dependent var	0.009496	S.D. dependent var	0.143629	
Sum squared resid	8.372788	S.E. of regression	0.141529	
R-squared	0.045010	Adjusted R-squared	0.029018	
F(7, 418)	2.814439	P-value(F)	0.007117	
Log-likelihood	232.5056	Akaike criterion	-449.0111	
Schwarz criterion	-416.5756	Hannan-Quinn	-436.1984	
rho	0.034956	Durbin-Watson	1.926820	
Excluding the constant, p-value was highest for variable 27 (LNSPReturnt1)				

Description: Regression results for regression equation 4) (see table 4.4). LNReturnt1 stands for simple return between day t-2 and t-1, REt (REt1) stands for daily change in EMOJI sentiment score between day t (t-1) and t-1 (t-2), RVPt (RVPt1) stands for daily change in number of submissions between day t (t-1) and t-1 (t-2), and LNSPReturnt (LNSPReturnt1) stands for S&P 500 index simple return between day t (t-1) and t-1 (t-2).