

# Corporate Management's Bias: The Pitfalls in Financial Filings\*

Unveiling Corporate Bias in Management Discussions in Financial Filings

Nikhil Iyer

April 19, 2024

Every quarter, management must release a document outlining the performance of their company. Among this document the Management Discussion and Analysis (MD&A) section serves as an opportunity for a company's management team to discuss the inner workings of the company. However, the efficacy of these documents is often compromised by a variety of factors. Abundant filler words, complicated jargon, and an extreme positive bias make using these documents difficult. This paper conducts a statistical analysis of SEC documents from 5 of the biggest tech companies in the world. It reveals that management often has an extreme positive bias, making it difficult to identify the negative aspects of the company. This finding is important as it highlights challenges that shareholders face when assessing the true performance of a company, which could impact portfolio management strategies and investment decisions

## 1 Introduction

Big scandals make headlines, however it is often the subtle tactics that pose the greatest risk. A company must file a 10Q every quarter, and such a document is released to the public. These documents contain important information to the shareholders and the street. Oftentimes, the contents of this document along with the earnings report will dictate the company's performance in the stock market. Understanding how to use 10Qs specifically the MD&A section is essential for shareholders to make an informed investment decision.

This paper examines the Management Discussion and Analysis (MD&A) section of the Securities and Exchange Commission (SEC) filings. There is very little information on how to best use the MD&A section, however, this paper analyzes the MD&A sections of SEC filings and

---

\*Code and data are available at: <https://github.com/Niyer02/sec-market-comparison>

unveils a concerning trend: the presence of an extreme positive bias and constant downplay of negative aspects. The estimand in this paper is the extent to which management biases their discussions in the MD&A section of SEC filings while obscuring the negative aspects. By revealing the subtle strategies management uses, this study aims to empower shareholders with the knowledge and tools needed in the corporate world. This study is also presented to encourage management to engage in these practices less frequently or drop them altogether.

Firstly this paper will go over how the data (2) was retrieved, cleaned, and processed into workable features. The data retrieval and cleaning process was heavily inspired by TELLING STORIES WITH DATA (Alexander 2023). Then it will explore the model (3) intended to predict the gain or loss in a company's stock price. The results section will sum up the findings, and the paper ends with the final discussion points. The data was gathered in `Python 3` (Van Rossum and Drake 2009), then further parsed in `R` (R Core Team 2023) using `tidyverse` (Wickham et al. 2019), `dplyr` (Wickham et al. 2023), and `readr` (Wickham, Hester, and Bryan 2024). The features were plotted with `ggplot` (Wickham 2016).

## 2 Data

The data for this paper was collected in `Python 3` (Van Rossum and Drake 2009). The SEC filing data was retrieved using the SEC API (2022), and the financial data was retried using the `yfinance` package in `Python` (Van Rossum and Drake 2009). Overall the raw data set was composed of a `ticker` column, `MD&A filing date` column, and the actual MD&A text.

The financial data consisted of four main columns. The first is the difference between the stock's price 1-week after the SEC filing and 1-week before the SEC filing. This can be interpreted as the Information Given in an MD&A. A high value in this column indicates that the SEC filing had a significant impact on the stock, and the mean of this column reflects this. The remaining three columns are the 1-month, 3-month, and 6-month price changes in the stock. There were 10 tech companies chosen, and all 10Qs were retrieved from 2019 to 2023. This resulted in 15 10Qs per company, with the final data set being 150 rows.

Due to the length of the MD&A sections, they could not be used directly. Instead the **average word count**, **average positive word count**, and **average negative word count** were all computed, and z-scores were computed for each row. The result is three columns with information of how far in either direction a data point is from the mean in the three classes mentioned.

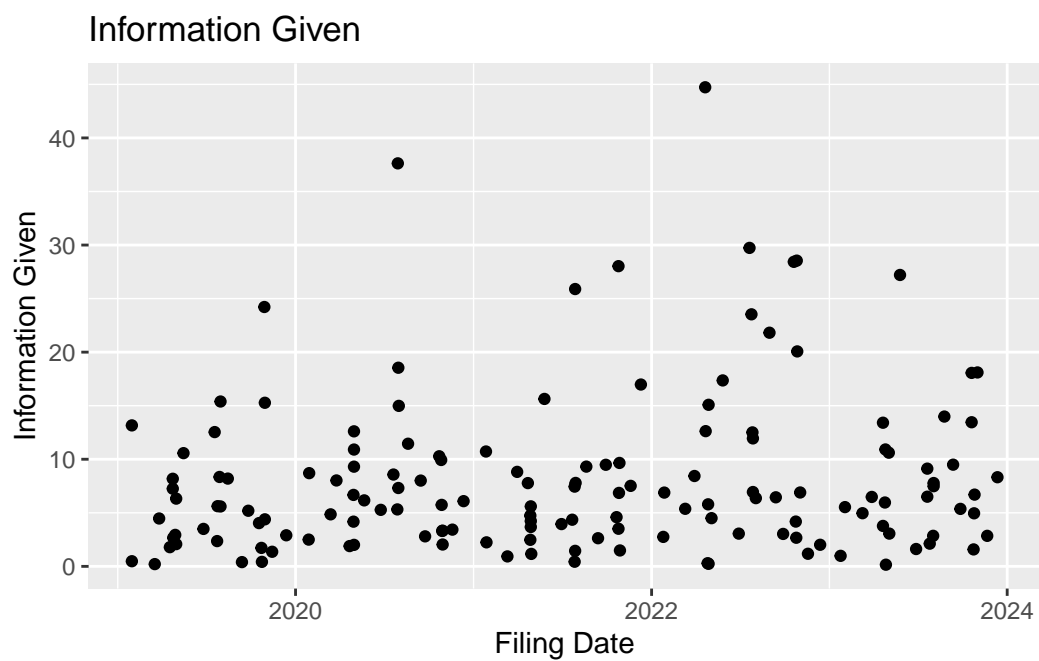


Figure 1: Information Given Distribution

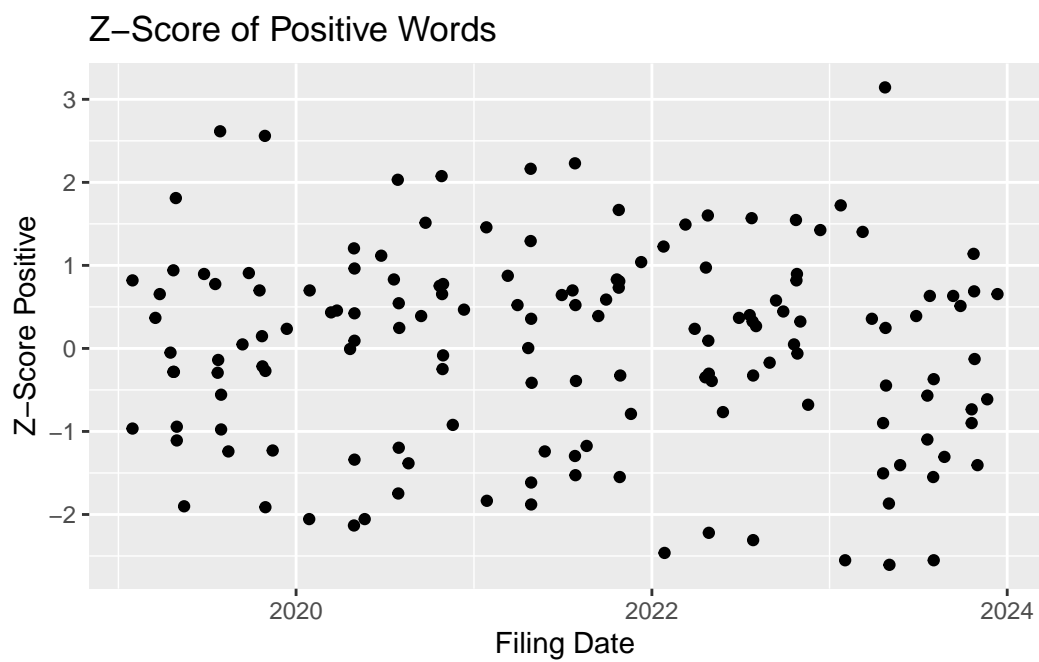


Figure 2: Positive Word Count Distribution

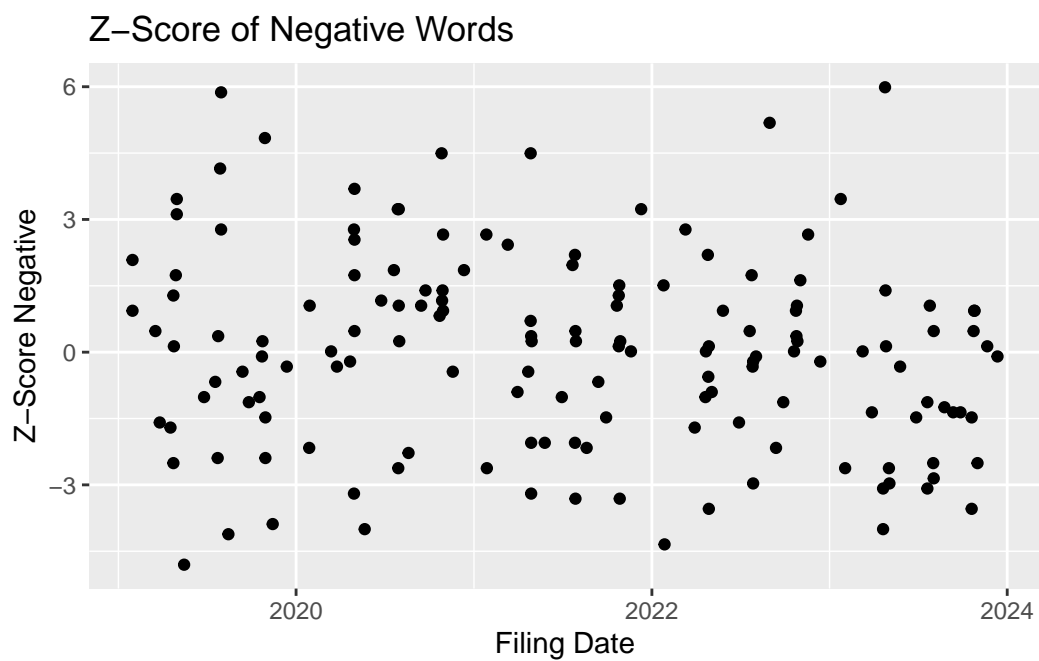


Figure 3: Negative Word Count Distribution

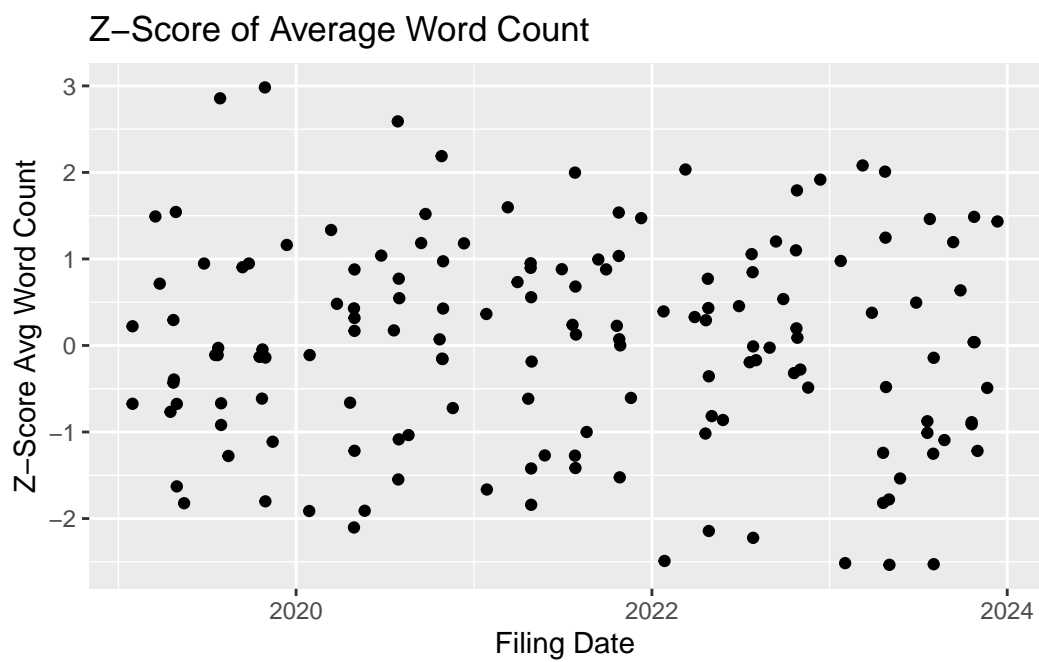


Figure 4: Average Word Count Distribution

Figure 1, Figure 2, Figure 3, and Figure 4 shows the scatter plots of 4 features. Figure 1 shows the one-week price change in stock during the one-week window of the SEC filing. As seen in the@fig-data1, an SEC filing has a large impact on the stock price, as in such a short window within the SEC filing release the stock prices jumped in almost all cases. All of the graphs above show very little, if any correlation. This is a surprising result. The stock price increased 10%-20% in most cases, however, there was no change in the number of uses of positive or negative words as explained by Figure 5.

Table 1: Correlation Matrix between IG and z-scores

	IG	z_score_positive	z_score_negative	z_score_avg_word_count
IG	1.0000000	-0.0798141	-0.0003675	-0.0933501
z_score_positive	-0.0798141	1.0000000	0.6677007	0.9187923
z_score_negative	-0.0003675	0.6677007	1.0000000	0.6070120
z_score_avg_word_count	-0.0933501	0.9187923	0.6070120	1.0000000

We confirm this hypothesis by examining the correlation matrix in Table 1. With relation to Information Given, the z-score columns do not correlate at all, however, we do see a strong positive correlation between the positive word count and average word count. Such correlation is the basis of this study. From the figures in Figure 5, we observe that the distribution of the growth variables is very similar. These graphs tell an interesting story when compared to the previous features. Although stock price is increasing, there was no increase in the count of positive words, and no decrease in the count of negative words, even among large price increases.

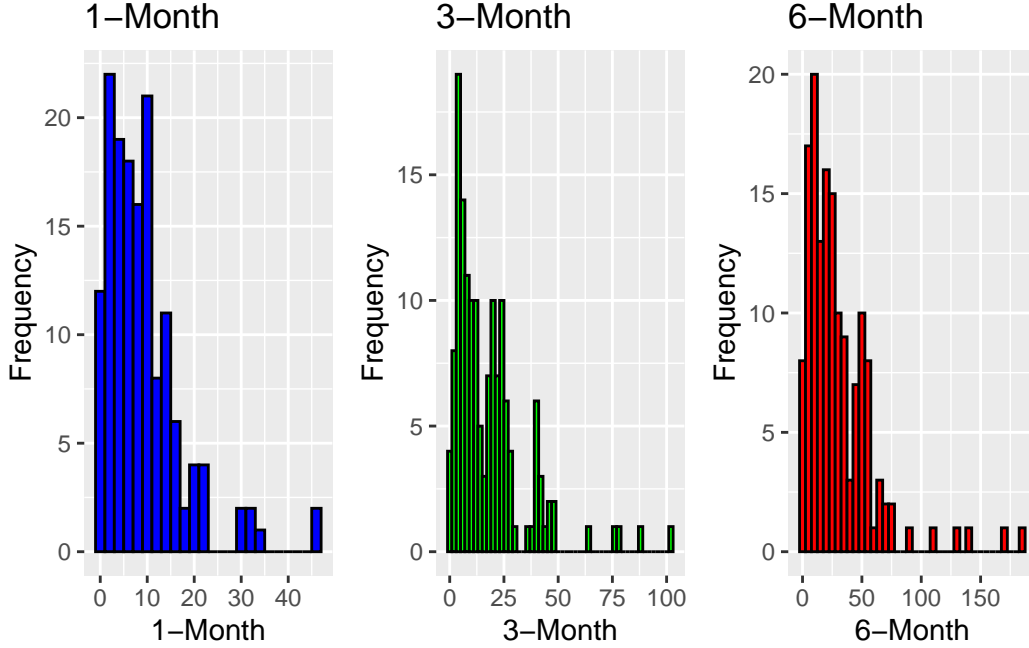


Figure 5: Distribution of Stock Price Growth over Time Periods

From examining the data in [?@fig-data](#) and Figure 5, we can see that the lack of correlation is worrying. An increase in stock price should lead to more positive comments and a decrease in negative comments. However this is not the case, and this is further reinforced in the model (3) chosen.

### 3 Model

The goal of our modeling strategy is to attempt to find a correlation between Z-Score of Positive Word Count, Z-Score of Negative Word Count, and Z-Score of Average Word Count, and to then build a model to predict the percentage a stock will increase or decrease in a given time frame. For this, a Random Forest Model was settled on to predict the amount of change in a stock price, however, due to its simplistic nature, it was not able to accurately model the data. Linear models as well as Gaussian models were also trained (Goodrich et al. 2024), however, they performed worse overall than the Random Forest. The model set-up section (3.1) will go over the setup of a Random Forest, however the data itself can be better modeled by a linear regression model, where we can see the true scale of the lack of correlation between the features.

For the purpose of this model, the 3-month and 6-month features were omitted in an attempt to obtain the highest accuracy model possible. Thus, the Random Forest model was

trained to predict 1-Month given the predictors: `z_score_positive`, `z_score_negative`, and `z_score_avg_word_count`.

### 3.1 Prediction Model set-up

Define  $y$  as 1-Month. Then define  $x$ ,  $w$ ,  $t$ , as `z_score_positive`, `z_score_negative`, and `z_score_avg_word_count` respectively.

Random Forest Model:  $y = f(x, w, t)$

- $y$  represents the target variable (dependent variable).
- $f$  represents the Random Forest model function.
- $x$ ,  $w$ ,  $t$  represent the predictor variables (independent variables).

We run the model in R (R Core Team 2023) using the `randomforest` package of `randomForest` (Liaw and Wiener 2002).

#### 3.1.1 Model justification

Table 2: Random Forest Model Evaluation Metrics

Metric	Value
MAE	4.8699451
RMSE	7.3134372
R-squared	0.4361033

We can see that in Table 2, the Mean Absolute Error (MAE) and the RMSE (Root Mean Squared Error) are relatively low when looking at financial data. Additionally, we have a relatively high R-squared value, however, this is not consistent. The variance of the Random Forest model is very high, so the metrics are not consistent, and when extrapolated to larger data will fail to accurately predict. However, this was the best-performing model, without delving into Transformers.

## 4 Results

The study's result, based on evaluation metrics of a Random Forest model applied to financial data, as well as correlation inferences indicates poor predictive performance when looking at the MD&A alone. This outcome was expected, when looking at the features and plots in (2)

we saw early on that the correlation between the features did not exist, thus making it almost impossible for a model to accurately find a pattern within the data. The Random Forest parameters can be seen in Table 3.

Table 3: Random Forest Model Parameters and Attributes

Parameter	Value
Number of Trees	11
Variable Importance	
mtry	1

## 5 Discussion

### 5.1 What was learned

Despite the positive market trends, management tends to keep a positive tone in their Discussions and Filings. This study, although failing to provide an accurate model, suggests that management is making a deliberate effort to downplay negative aspects and constantly keep a positive narrative despite the actual conditions. Another finding is the lack of correlation between the financial data and predictor variables. Typically, higher stock prices can be attributed to positive events in the company, however, such events are not reflected in the MD&A. Due to the constant positive bias, positive events within the company do not get flagged as anything out of the usual. This makes finding negative events extremely difficult as well. There is also an active attempt by management to downplay the negative aspects, which is the reason training a model to predict stock price on biased text data is difficult.

### 5.2 Weaknesses of this study

Weaknesses in this study start primarily with the data gathered. The data-gathering process sampled only 10 companies, which were some of the largest technology companies in the world. This resulted in a stock price increase across the entire data set, making a Boolean classification impossible. Additionally, the complex nature of financial data is not able to be modeled by such simple models. Highly advanced Transformers tend to perform better on text data, however, they were out of the scope of this study.

The study’s choice of predictor variables also may have overlooked important features. Z-scores provide valuable insight into the relative performance of an MD&A with respect to the data set, however, such simple approaches may have led to a loss of crucial information. Such losses could have been the reason that the model was not as accurate as initially desired. The study also did not account for external factors such as interest rates, inflation, political events, etc.



Finally, the study's reliance on historical data from the specified five-year time frame may limit its applicability to future market conditions.

### **5.3 Next steps**

The next step in this study would be to look at a more complex model. Transformers and LLMs have risen in popularity and for good reason. Transformers can capture the intricacies of text data much better than a simple model like a Random Forest. Exploring models like BERT would be the next step in attempting to understand MD&A's.

## References

2022. *SEC Emblem*. <https://www.sec.gov/edgar/searchedgar/companysearch>.
- Alexander, Rohan. 2023. *TELLING STORIES WITH DATA with Applications in r*. Boca Raton: CRC Press.
- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2024. “Rstanarm: Bayesian Applied Regression Modeling via Stan.” <https://mc-stan.org/rstanarm/>.
- Liaw, Andy, and Matthew Wiener. 2002. “Classification and Regression by randomForest.” *R News* 2 (3): 18–22. <https://CRAN.R-project.org/doc/Rnews/>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Van Rossum, Guido, and Fred L. Drake. 2009. *Python 3 Reference Manual*. Scotts Valley, CA: CreateSpace.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *Dplyr: A Grammar of Data Manipulation*. <https://dplyr.tidyverse.org>.
- Wickham, Hadley, Jim Hester, and Jennifer Bryan. 2024. *Readr: Read Rectangular Text Data*. <https://readr.tidyverse.org>.