# Stock Market Comparisons and Predictions

James P. StClair

Northwest Missouri State University, Maryville MO 64468, USA
S555947@nwmissouri.edu

**Abstract.** I began this project trying to see what impact COVID may
have had on several fields of work- namely, Insurance, Healthcare, and
Delivery, and attempt to predict changes from there. I reviewed a small
data set of several companies stock prices based on what they were at the
closing of per day, from early 2020 to late 2022. I downloaded the csvs for
each company, removed any missing or superfluous data, then trained and
tested the full data remaining. Based on the results I received, I would
believe that COVID still had a significant impact on each field, as the
stock prices were very clearly affected at the beginning of 2020, but as
time went on those changes either expanded or reversed. For Healthcare,
their stock price rose and then continued to rise based on the data, and
should continue to do so. Insurance on the other hand, had an initial
decrease at the beginning of 2020 but then started to stabilize again
by mid 2021, and should continue to rise from there, albeit slower than
Healthcare. Delivery as a field was the most impacted I believe, as it
initially saw a rather large rise as more people were staying home, but
as it started approaching the end of 2022, it started going back down, in
rare cases even lower than it had started, and should continue to lower
a small amount further.

**Keywords:** COVID · Stock · Delivery · Insurance · Healthcare

## 1 Introduction

The domain that I have chosen for my project is Stock Market Predictions, as
I want to look into the differences in Stock price of several fields, especially
since COVID-19. I want to try to compare and contrast differences in the fields
of Insurance, Healthcare, and Delivery services(such as Amazon), to show the
differences, what trends exist, and make predictions for these companies.

I am planning on looking into Kaggle, Google Datasets, and AWS to try to
find data-sets that would have what I am looking for- stock markets for these
fields over time, including showing any impact during and since COVID.

### 1.1 Goals of this Project

I want to analyze the difference in some stock market fields since COVID, and
show whether, as a whole, they are trending upwards since then, or downwards,
and make predictions based on the data obtained. The fields I will be looking at

are Insurance, Healthcare, and Delivery Services, to try to showcase some of the fields that would have a direct impact from COVID and its spread. I will want to show predictions for each field as a whole, and several individual companies, to show how they have declined or gone up since COVID.

## 1.2   Steps in the Project Process

First I would create a Github repo and Overleaf outline regarding this project. Next I will look over the possible data sources on hand, to try to find as clean of data as possible, and fitting what I am attempting to look for. I will look into similar projects on order to give myself some background information, and source anything that I use I will start cleaning and assessing the data I have pulled in order to present the information I am looking for, and only use the relevant data, keeping a note of every process I use to do so. I will use the finished data to run predictive modeling of the data, to try to see what the differences will look like further on I will bring everything together, as a final presentation to show what I did, what I predict, and how I think these fields will progress from here.

## 2   Key Components and Limitations

The key components of my approach and results will be the use of accurate data regarding stock prices of the companies in these fields, and using data analysis from there to show the trends for each company as well as each field as a whole. It will require accurate data, proper use of analytical skills and tools, and predictive modeling. Some limitations however is that it would not take into account many outside factors that could happen in the meantime, for the companies individually or another event such as COVID that could affect several companies at once. Another limitation I came across was that as much data as this data set had, there were few companies actually working in the fields I was looking for in this list.

## 3   Data Set

The data set I will be using, Stock Market Data(NASDAQ, NYSE, S and P500) was found in Kaggle, and was made by Paul Mooney. This data-set appears to be several csv's, 1 for each company as part of the 3 subheaders in the title. I used batch processing to go through the data, after I made sure to confirm which companies I would need data from. I am planning to use the names of the companies, the dates, and the stock price on the end of business day, in a separate chart for each company. I will be doing this for each company, in order to have viable data for each one, before working with this data both for each company, and for each field I will split them into as a whole(healthcare, insurance, and delivery). However, for any company that works in 2 or more fields(such as UNH working in both Healthcare and Insurance), I will list it under both fields. I have use the following companies for the previously named fields:

- Delivery
  - UPS
  - FDX
  - EBAY
  - AMZN
- Healthcare
  - LLY
  - UHS
  - HCA
  - JNJ
  - UNH
- Insurance
  - ALL
  - PGRE
  - BRO
  - CB
  - MET
  - UNH
  - TRV

The individual data, as well as each process for the individual companies and the fields as a whole, can be found here, on the Github Repository: Capstone-Project. https://github.com/jpstclair/CapstoneProject.git

### 3.1   Cleaning the Data

The data-set I had was already pretty well organized by company into many separate csvs, but it did also contain a large amount of data that for this project I did not need. The first thing I needed to do was remove the columns for the attributes I did not need, and then add into each one a new column simply to state which of the 3 fields I was looking into was. After that, I had 1 main attribute remaining for each field- which was Date. Each field also included several different individual companies that had data in this data-set, each labelled by their stock symbol(such as ALL for Allstate), as well as one last field labelled Average, which was the average of each of these previous fields. Any dates that had missing values were removed, as they would not be able to affect the data and would simply be clutter where it wasn't necessary. These were primarily removed by utilizing Excels own programs to remove rows that had missing values. This left me with roughly 299 dates for each company, with which I will attempt to plot both a line plot of the individual companies and the averages(shown in 3 figures below, labelled as each fields stock price from 2020 to 2022) to see if the stock prices appear to generally be going up or down, as well as applying that data to attempt to see if I can predict what will occur with that stock after the end of that period.

### 3.2   Analyzing the Data

After looking over the data for some time, I was able to find a few things. As a whole, most of these companies dipped down in early 2020, though Johnson and Johnson(Healthcare), Eli Lily and Co(Healthcare), Progressive(Insurance), Regeneron(Healthcare), and United Health Group Inc(Healthcare) are all exceptions to this, as they have all been generally rising up gradually since 2020. This is, of course accounting for the general oscillations that stock prices tend to have, as while there were a lot of ups and downs, the trend lines were absolutely positive for these companies throughout. On the whole, it appears that COVID may have affected some of these companies quite a bit. For example, delivery companies saw large decreases in early 2020, and massive rises in stock in late 2020, but as of late 2022 saw that stock come back down for almost every single one of them, though Fedex and UPS are slight outliers here in that they did not dip back down. Insurance companies and healthcare companies, on the other hand have been a bit more varied. They all generally had a small dip in 2020, but some went further down by 2022, some went back to roughly the same price by 2022, and some have kept increasing to well over where they started, showing that while COVID did have some affect on some of these companies, it didn't seem to affect them as much as delivery companies on the whole were affected. Specifically, I used line plots and to get these basic ideas in place, as well as looking at the line plots to get a sense of when they started going downwards or upwards. The source for this data can be looked at here:

https://www.kaggle.com/datasets/paultimothymooney/stock-market-data

## 4   Results and Predictions

The first thing that I did was split the data I had into training data and test data. After selecting the data involved, I ran a few different machine learning algorithms across the average results from each field as a whole, to better show the impact COVID had on those fields. I had immediately ruled out linear regression, as while it could have been good to use as a trend-line, with how often stocks oscillate even on the same week, I decided that it likely wouldn't be as useful for predicting the stocks in the future. As shown in the following figures, I ran each average through Polynomial Regression and Elastic Net, in order to try to find approximate estimates, or at least a solid footing for finding those estimates. Based on the results that were shown, I would believe that these trainings were at least mostly successful with showing how the average for each field will progress. Specifically, for the Insurance data, the Elasticnet generally fit better, with the R2 being .9309; for the Healthcare data, the Elasticnet also fit better, with the R2 being .9455; for the Delivery dataset, it was less accurate though still had an R2 of .7997. As such, I would believe any predictions made with these results to be mostly accurate, though with some possible variance based on events and how the individual companies may be perceived. I find it interesting that Healthcare started going down for a small amount of time at the beginning of COVID before generally rising up, while Insurance on the whole
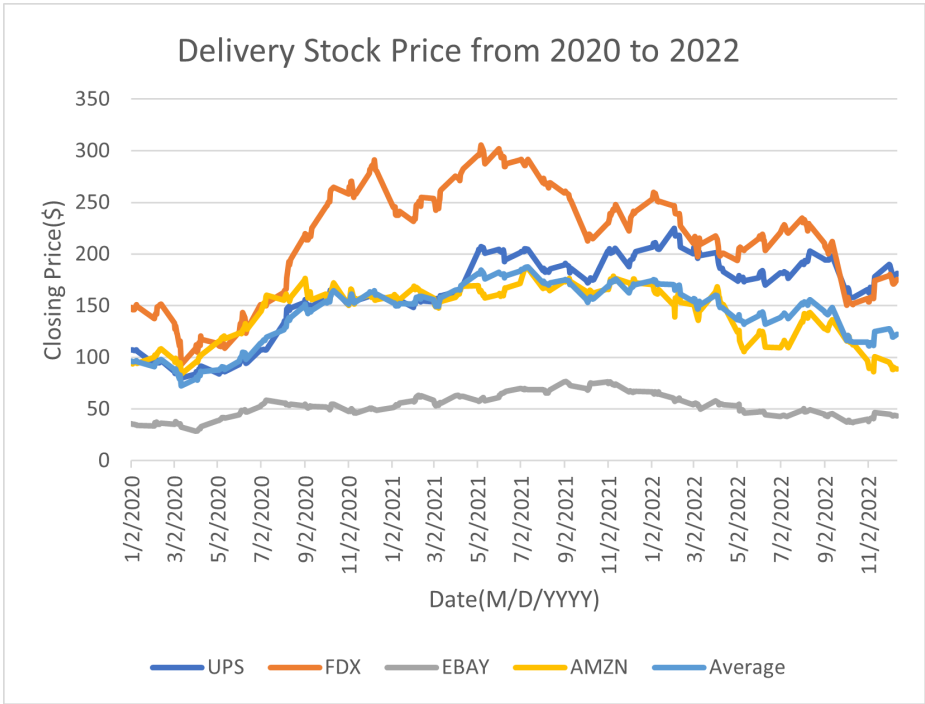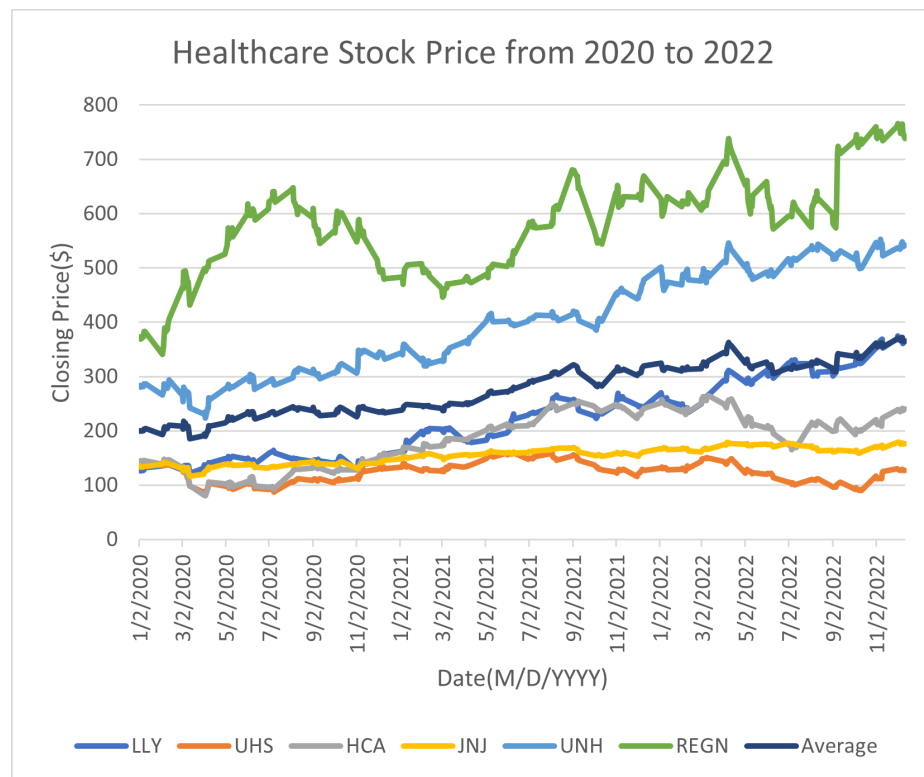
**Fig. 1.** Delivery Stock Prices

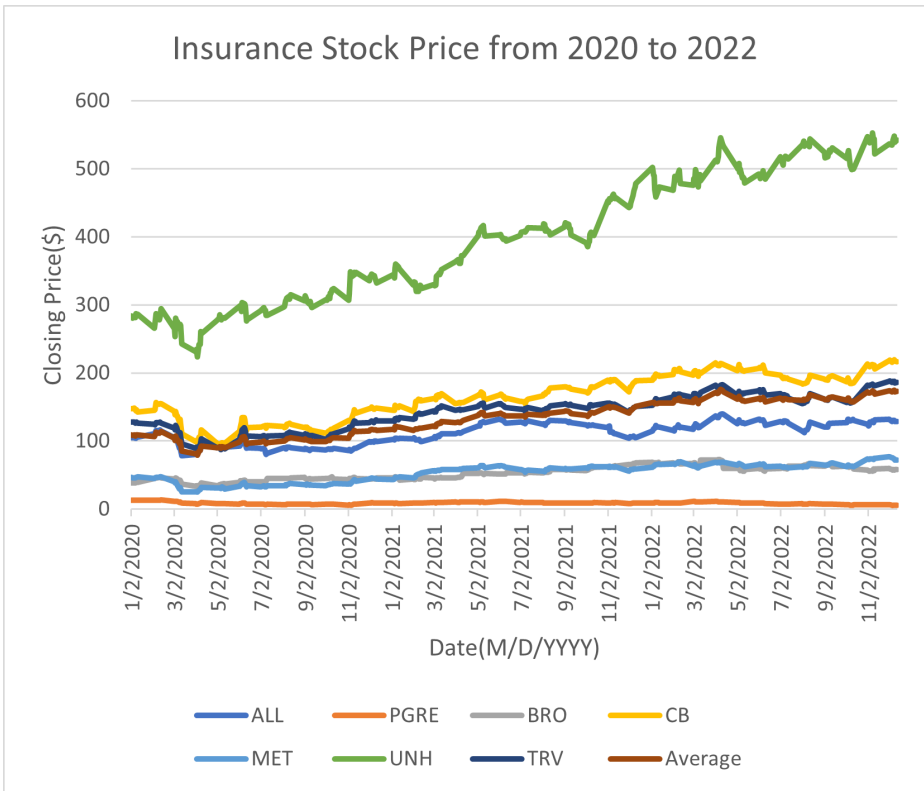**Fig. 2.** Healthcare Stock Prices

**Fig. 3.** Insurance Stock Prices

generally stayed around the same, or just went down a little bit across the entire 2020 to 2022 timeframe. At the same time, Delivery stock went up for the beginning of the worst of the COVID timeframe, but then as people started being more active and going out more often again it seemed to fall back down.

These trainings, as well as line plots for each individual company and the average for each field can be found here:

https://github.com/jpstclair/CapstoneProject.git

## 5    Conclusions and Future Work

I would believe that based on the results I received overall, that each of these fields were impacted in some way by COVID, as several companies either rose or fell at the beginning of COVID, and several companies that rose started going back down once COVID started to go down. It is interesting to me how these fields were affected, such as Healthcare stock prices rising as more attention was given to it and more was needed from it, but almost more interesting was the way that Delivery stock prices went up initially, but then did start to noticeably decay during mid 2022. For future work, there are 2 different things I would like to do. first, I would want to find more dataset in this same timeframe for more companies in these fields, to get a better idea for the whole of the impact on these fields. Second, I would also want to start looking into other fields as well, to see what impact may have been had on some of those fields as well, as there could be some fields that were impacted that, at first glance, we would think wouldn't be as affected by COVID.



**Fig. 4.** Splitting the data and first part of Polynomial Regression - Delivery

**Fig. 5.** Polynomial Regression and results - Delivery
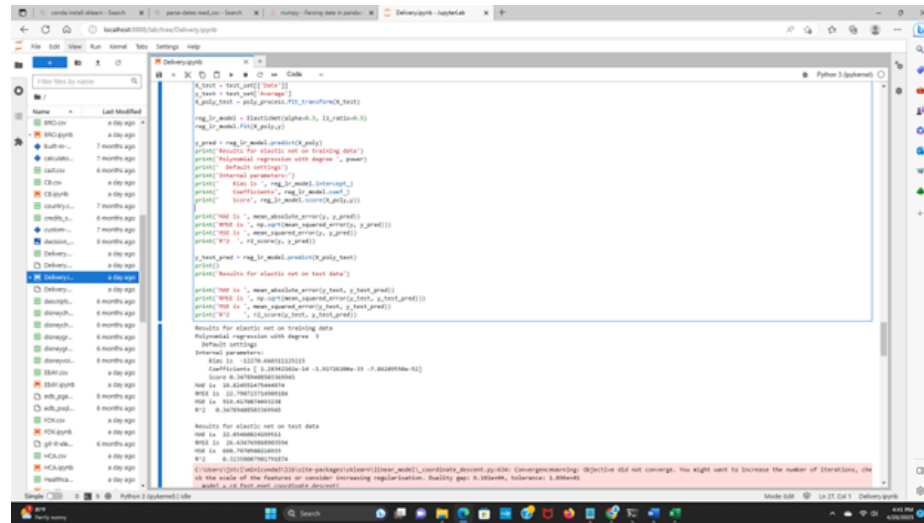


**Fig. 6.** First part of ElasticNet- Delivery
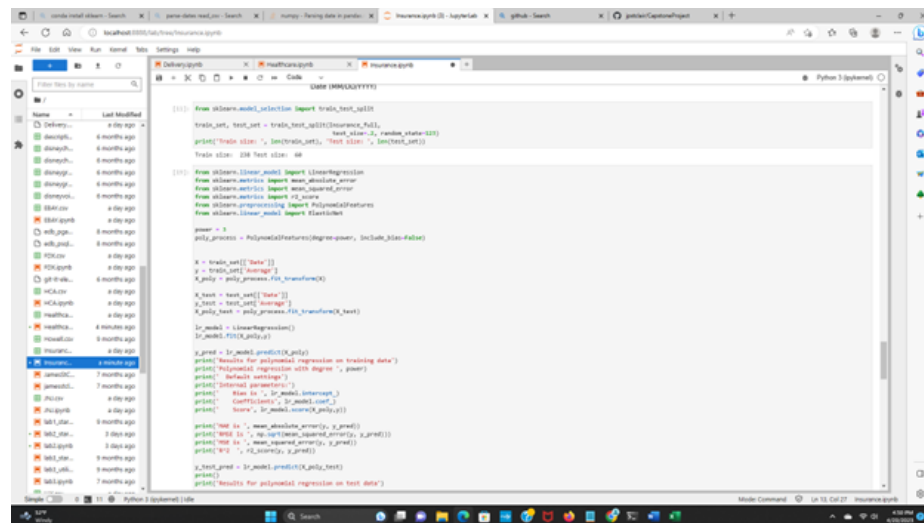
Fig. 7. Second part of ElasticNet and Results- Delivery



Fig. 8. Splitting the data and first part of Polynomial Regression - Insurance

**Fig. 9.** Second part of Polynomial Regression and Results- Insurance



**Fig. 10.** First part of ElasticNet- Insurance

**Fig. 11.** Second part of ElasticNet and results- Insurance



**Fig. 12.** Splitting the data and first part of Polynomial Regression - Healthcare

**Fig. 13.** Second part of Polynomial Regression and Results- Healthcare



**Fig. 14.** First part of ElasticNet- Healthcare

**Fig. 15.** Second part of ElasticNet and results- Healthcare

# References

1. Mehrotra, N.: The aftermath and impact of covid-19 on stock markets, https://www.forbes.com/sites/theyec/2023/02/10/the-aftermath-and-impact-of-covid-19-on-stock-markets/?sh=4bdb943ec120
2. Mooney, P.: Tstock market data(nasdaq, nyse, sp500), https://www.kaggle.com/datasets/paultimothymooney/stock-market-data
3. Wang, K.M., Lee, Y.M.: Are life insurance futures a safe haven during covid-19? Finance Innov. 2023 **9**(1), 13 (2023). https://doi.org/10.1186/s40854-022-00411-z

[1]
[3] [2]