

California State University, Northridge

COVID-19 Prediction

Group number 1

COMP 541 #17129 Spring 2021

Students: Alisiya Balayan, Bishoy Abdelmalik, Arian Dehghani, Ariel Kohanim, Sergio Ramirez, Natalie Weingart

Professor: Taehyung Wang

Table of Contents

| | |
|---|-----------|
| Project's objective | 2 |
| Background information | 2 |
| Data mining scope | 3 |
| Data exploration and data preprocessing | 4 |
| Modeling | 10 |
| Assessment of results | 11 |
| Summary of learning experiences | 12 |
| References | 13 |
| Appendix | 13 |

1. Project's objective

- Objective

The project's objective is to determine when will be the optimal time for the U.S. to return to normal operations and end lockdowns without putting people's lives at risk.

- Success criteria

The project's success criteria is having a road map of recovery to know when businesses and schools can return to normal operations. To obtain the roadmap of recovery, a prediction model will be built to understand COVID-19 trends.

- Goal

To predict using machine learning models and data mining techniques when the U.S. population will achieve herd immunity.

2. Background information

- Tackled Problem

Coronavirus disease (COVID-19) is an infectious disease caused by a newly discovered virus. It has affected every country and city around the world causing health, economic and psychological impacts. Therefore, predicting when the pandemic will come to an end by analyzing data about vaccinations, herd immunity, and hospitalizations is crucial for repairing our societies [1].

- Solution

The solution we attempted in this project is to use data and models to predict when the pandemic will come to an end by analyzing data about vaccinations, herd immunity, and hospitalizations. The current assumption leading scientists have is that roughly 70% of the population needs to have COVID immunity to achieve herd immunity. Either through vaccination or by natural immunity, meaning through contracting the virus and developing antibodies. As a result our model will attempt to predict how many months are remaining given the current data for the country to achieve herd immunity.

3. Data mining scope

3.1. Datasets used

- COVID-19 World Vaccination Progress
 - This dataset tells us about: which country is using what vaccine, which country's vaccination program is more advanced, and where the rate of vaccinated people per day is higher in terms of percent from the entire population.
 - <https://www.kaggle.com/gpreda/covid-world-vaccination-progress>
- Novel Coronavirus 2019 Dataset
 - This dataset has daily level information on the number of affected cases, deaths and recovery from COVID-19. It provides us with data about each country and their cases, deaths and recoveries.
 - <https://www.kaggle.com/sudalairajkumar/novel-corona-virus-2019-dataset>
- Countries population by year 2020
 - This dataset provides us with the world population and top 20 countries' live clock. It contains population data for the past, present and future.
 - <https://www.kaggle.com/eng0mohamed0nabil/population-by-county-2020>
- Covid-19 Daily Vaccination
 - This data set contains vaccination data for countries showing how many people are vaccinated daily.
 - <https://github.com/owid/covid-19-data/tree/master/public/data/vaccinations>

3.2. Data mining technology

3.2.1. Programing languages and programs

- 3.2.1.1. Python - was used as the main programming language that was used to analyze and build models. As one of the most widely used data mining languages it allowed us to do everything from cleaning and data organization to applying machine learning algorithms.
- 3.2.1.2. Jupyter Notebook - was used to execute the python code and display the visualizations
- 3.2.1.3. GitHub (version control) - was used to stay organized and up to date with all the changes to the project.

3.2.2. Python libraries

- 3.2.2.1. Numpy - was used to make it easier to perform mathematical and logical operations on the data in our datasets.

- 3.2.2.2. Pandas - was used to allow us to explore and manipulate data in a very efficient manner and helped in ingesting data into python as well as display it and integrate with other visualization tools.
- 3.2.2.3. Matplotlib - allowed us to visualize the data as it allowed us to create all the needed graphs and plots for our data.
- 3.2.2.4. Seaborn - was used in visualizing the data in conjunction with matplotlib due to the fact that it provided easier implementation in some instances.
- 3.2.2.5. Sklearn - was used as our main machine learning library and used to implement the machine learning model.

4. Data exploration and data preprocessing

4.1. Data exploration

We explored our datasets using various techniques to examine relationships between attributes and discover the properties of the datasets.

4.2. Data Cleaning

Prior to processing the data, we cleaned and organized the datasets by removing unnecessary columns that we do not need. For example, in the dataset `country_vaccinations.csv` we removed columns that contained vaccine type, name, source as we are not going to analyze those values. Another thing we did was narrow down all datasets to only U.S. locations.

After selecting the columns that we are interested in working with we proceeded to make sure that our dataset was free of null and NaN values by finding all NaN and null values in the dataset and replacing them with nearest neighbor values. For example, there is no data on December 22, 2020 however vaccines were administered and we have data from the day before and after, so we decided to replace the null value with the nearest neighbor value.

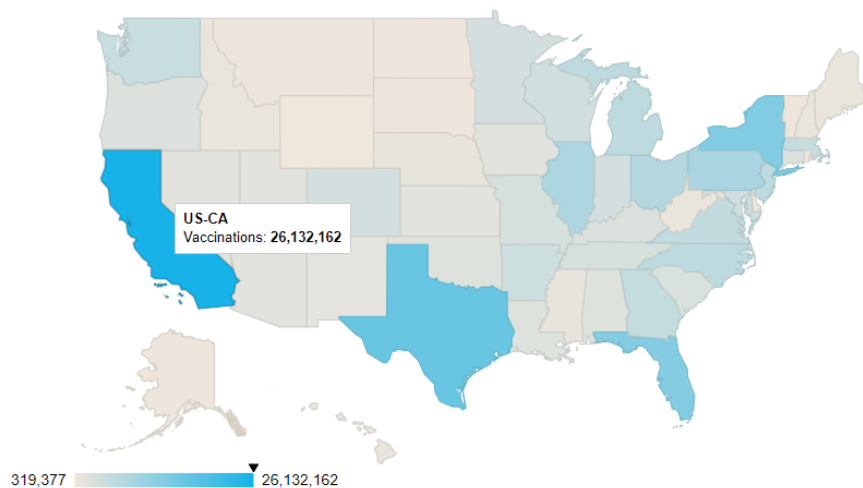
Below is the table that shows vaccination data in the U.S. [2]. Later we will use this cleaned dataset to further preprocess it by finding outliers, filling it missing values, etc.

| country | iso_code | date | total_vaccinations | people_vaccinated | people_fully_vaccinated | daily_vaccinations_raw | daily_vaccinations | total_vaccinations_per_hundred |
|---------------|----------|------------|--------------------|-------------------|-------------------------|------------------------|--------------------|--------------------------------|
| United States | USA | 2021-01-20 | 16525281 | 14270441 | 2161419 | 817693 | 892403 | 4.94 |
| United States | USA | 2021-01-19 | 15707588 | 13595803 | 2023124 | 1130189 | 911493 | 4.7 |
| United States | USA | 2021-01-15 | 12279180 | 10595866 | 1610524 | 1130189 | 798707 | 3.67 |
| United States | USA | 2021-01-16 | 12279180 | 10595866 | 1610524 | 1130189 | 811670 | 3.67 |
| United States | USA | 2021-01-17 | 12279180 | 10595866 | 1610524 | 1130189 | 824632 | 3.67 |
| United States | USA | 2021-01-18 | 12279180 | 10595866 | 1610524 | 1130189 | 837595 | 3.67 |
| United States | USA | 2021-01-14 | 11148991 | 9690757 | 1342086 | 870529 | 747082 | 3.33 |
| United States | USA | 2021-01-12 | 9327138 | 9327138 | 0 | 339816 | 641524 | 2.79 |

4.3. Clustering of vaccination data

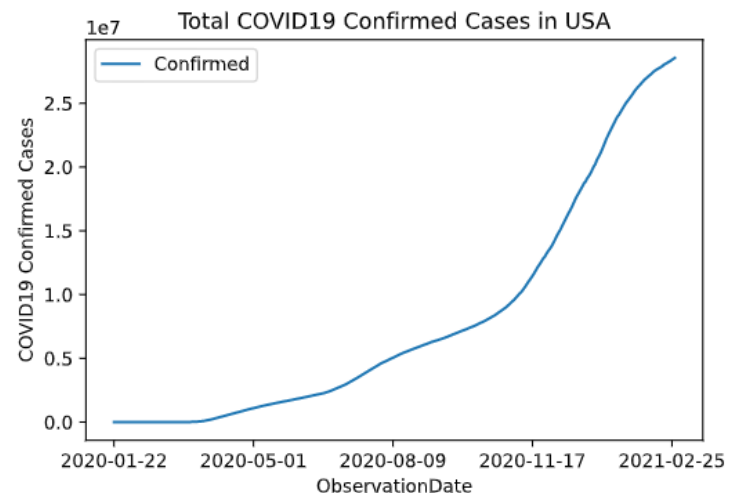
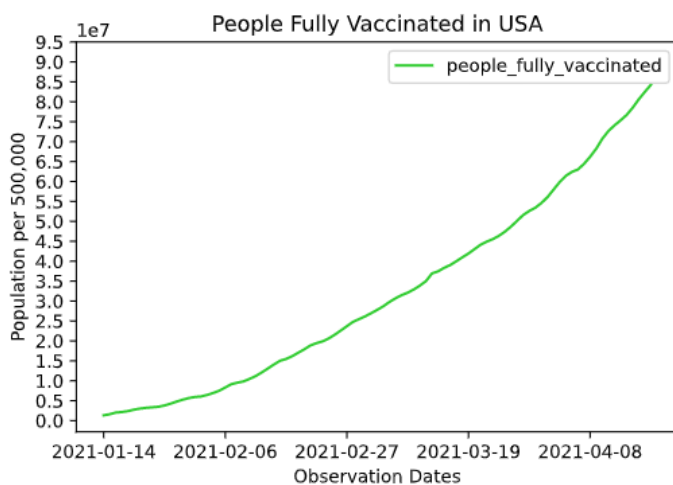
The figure above shows clustering of vaccination data per state where we can observe that California, Texas, Florida and New York have the most prominent clusters (i.e., shown in darker blue compared to the greyed out states). This clustering occurred for a few reasons, one being that these states have high populations and more vaccines are being delivered. Therefore, more people have access to vaccinations. However, to reach 70% immunity in these states, vaccination should continue to increase. For example, California's population is 39 million, and total vaccination (on 4/23/2021) was 26 million.

United States Vaccination Rates Per State



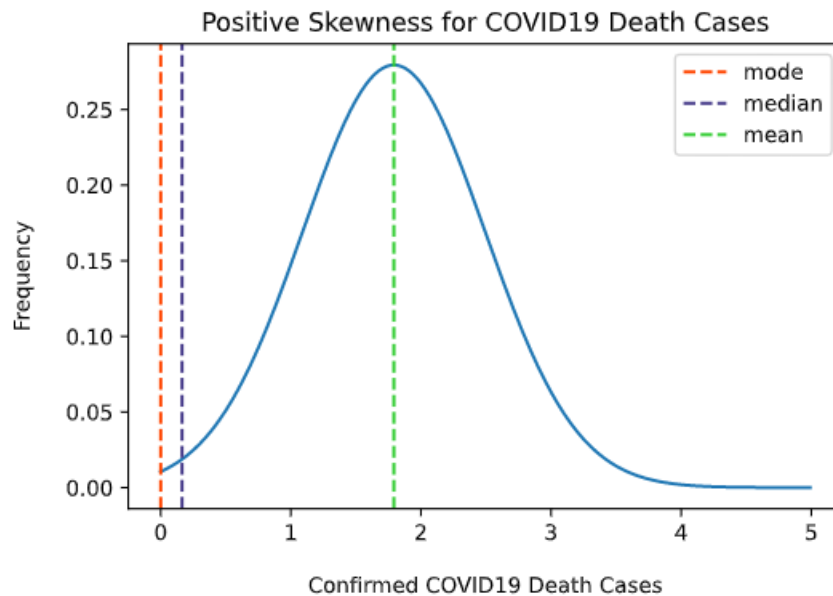
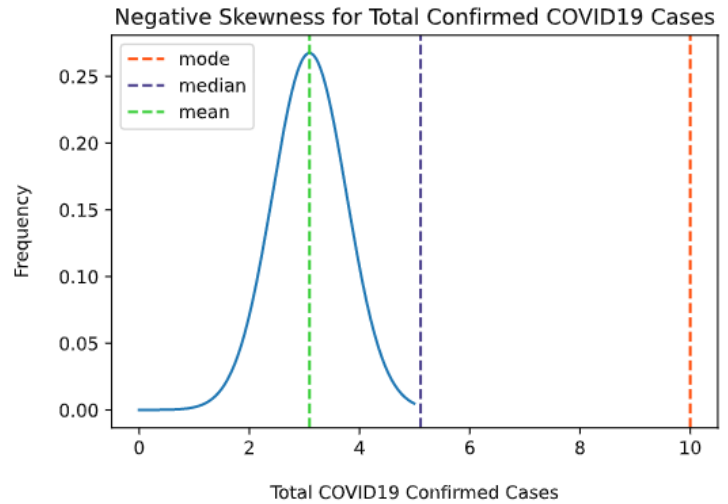
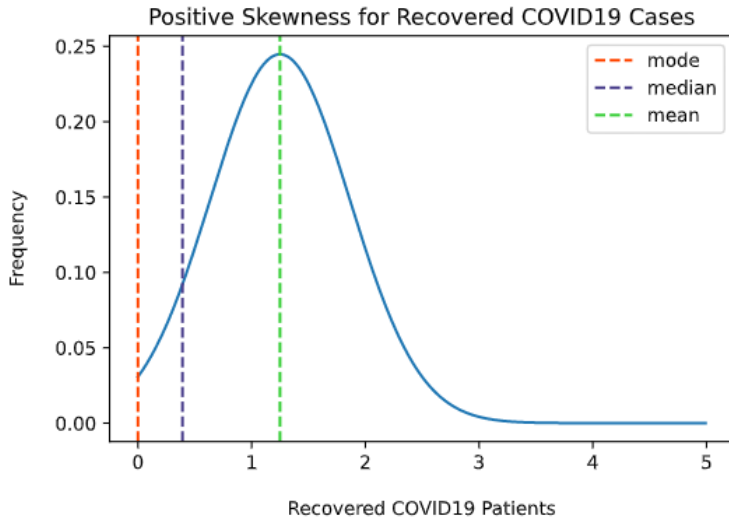
4.4. Datasets visualizations

The figures below represent visualizations of our datasets where we are able to see total covid-19 cases, death, recoveries and vaccinations. Our datasets and graphs show data from January 21, 2020 until February 25, 2021 for Covid-19 cases, as well as vaccination data ranging from January 14, 2021 until April 23, 2021.



4.5. Skewness of Covid-19 cases data

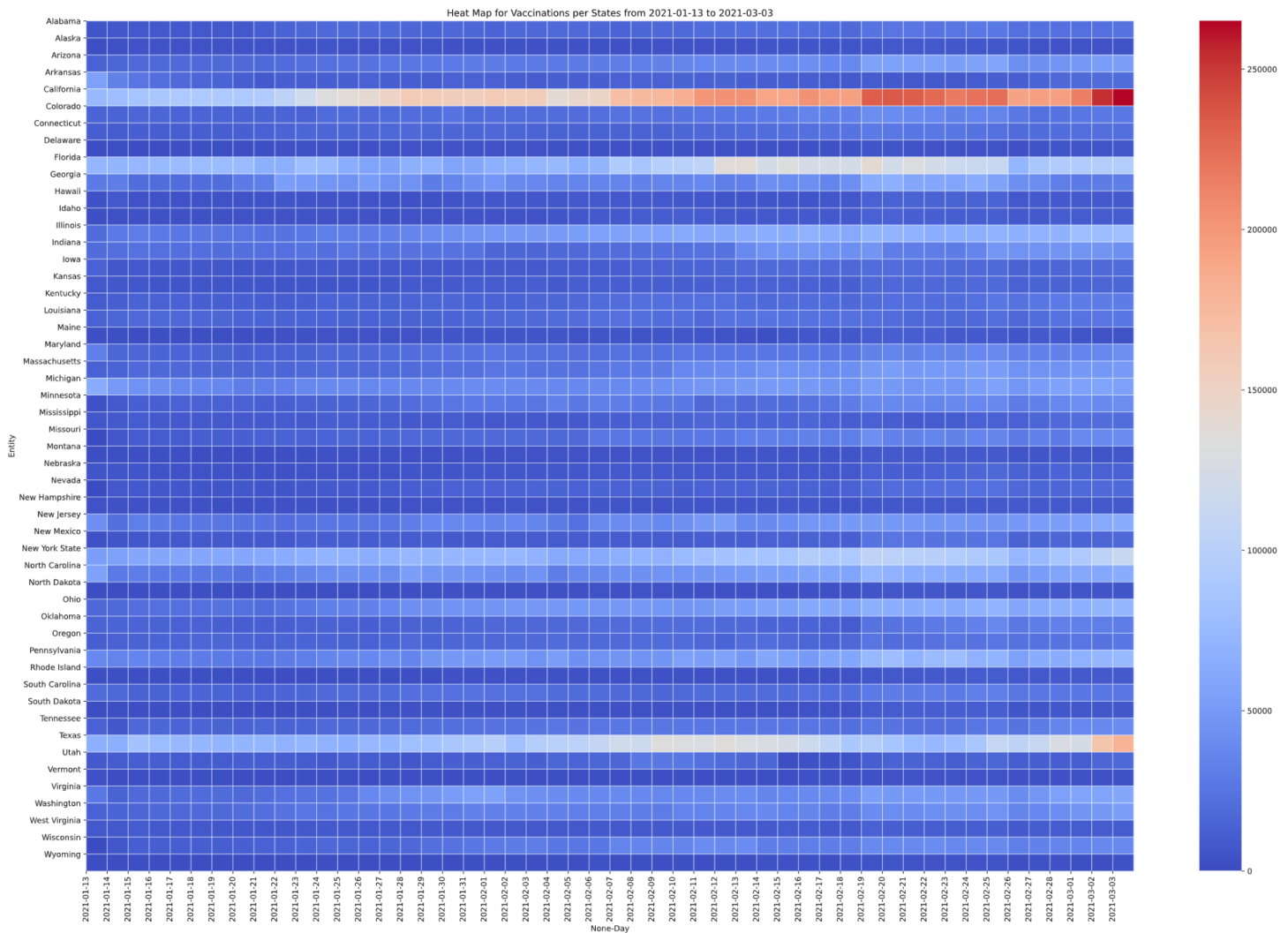
We discovered the skewness in our covid-19 cases dataset by graphing the skewness graph of each feature in the dataset and observing it in comparison to the mean, mode and median of the dataset.



4.6. Kernel Density Estimate - was used to to estimate the probability density function of the random variables to make inferences about the population based on the finite data sample that we have.

4.7. Correlation between attributes by using heatmap

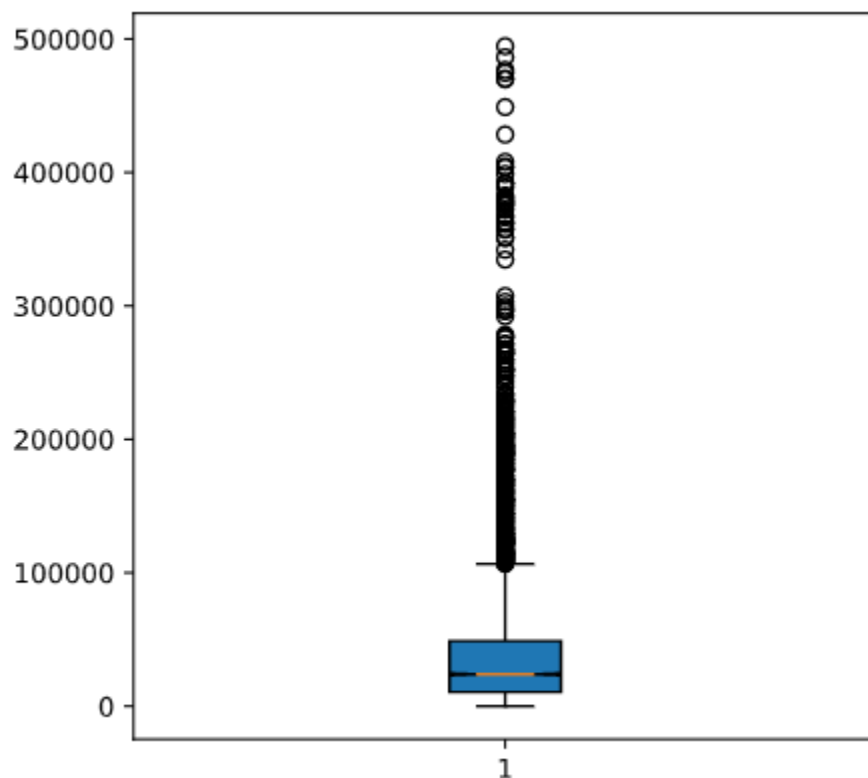
The figure below shows a correlation matrix using seaborn library between vaccinations in each state and the quantity of vaccine distribution over time. We can observe that in states like California and Texas where most cells are light and dark orange, vaccinations have been occurring a lot more compared to states like Alabama where most cells are blue. As of now, California vaccinated 26 million people, therefore there is a correlation between California vaccination and some of the recent months (March and April) where more people are vaccinated and it is continuing to increase.



4.8. Data preprocessing and transformation

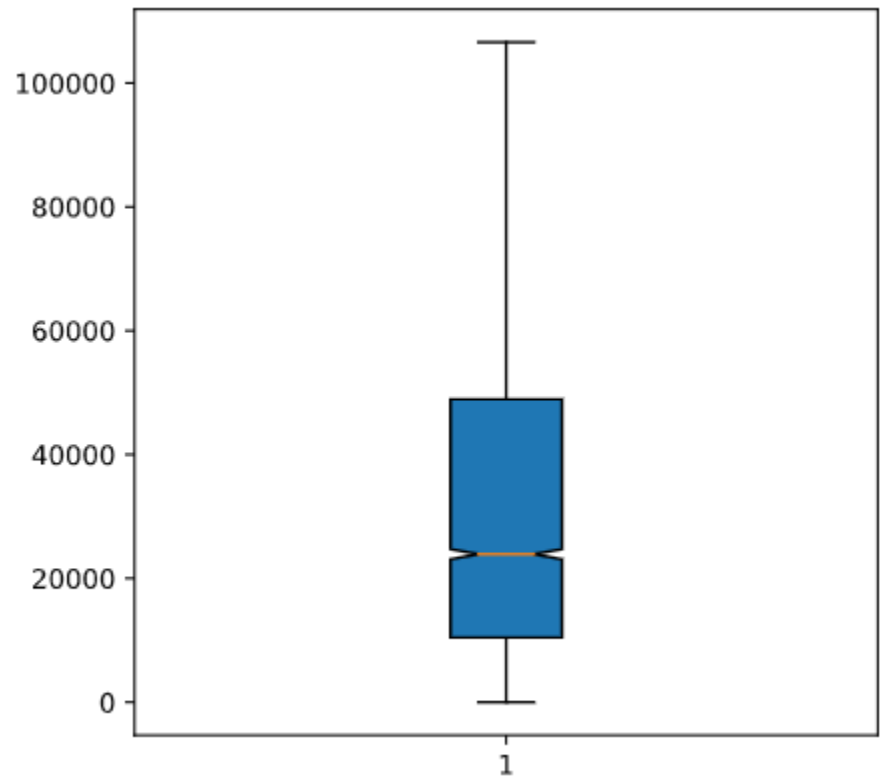
A. Outliers

After cleaning the datasets, we graphed points using boxplot in matplotlib to identify if there are any outliers. After seeing the outliers in our datasets we used the IQR scores to remove outliers.



```
#find Q1, Q3, and interquartile range for each column
Q1 = us_daily_vaccines["daily_vaccinations"].quantile(q=.25)
Q3 = us_daily_vaccines["daily_vaccinations"].quantile(q=.75)
IQR = us_daily_vaccines["daily_vaccinations"].apply(iqr)
print(Q1)
print(Q3)
#only keep rows in dataframe that have values within 1.5*IQR of Q1 and Q3
us_daily_vaccines = us_daily_vaccines[~((us_daily_vaccines["daily_vaccinations"] > (Q3+1.5*IQR)))]
```

After using IQR test, this is the new output:



B. Discretize data

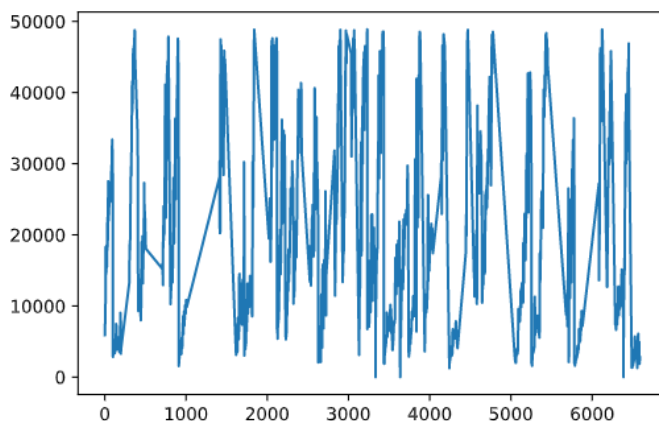
We performed discretization on our datasets because most of them contained dates that needed to be broken down into sections for easier evaluation. For example, our vaccination and covid19 cases datasets have data collected from each day, so we combined each month, and created a new csv with monthly covid19 cases instead.

C. Normalize data

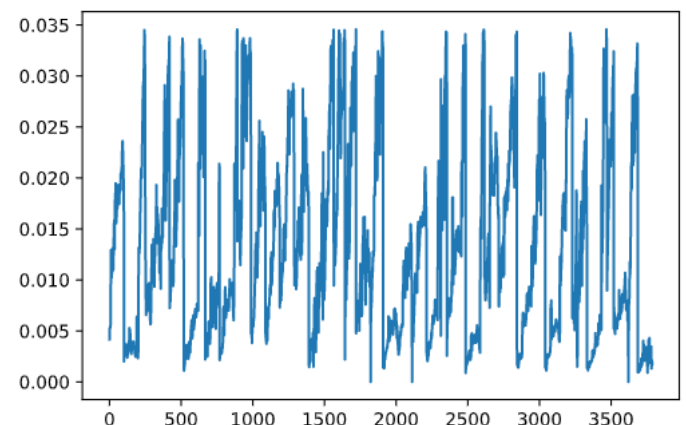
We performed data normalization to make our model faster and make calculations faster.

Below is the graph of the datasets before and after normalization.

Before:



After:



4.9. Data Quality Issues

A. Data is not up to date (lacks few months)

Some of our datasets are not up to date and are missing a couple of months. Also, different datasets contain information that is from inconsistent dates, so we had to match all datasets to be up to a certain time and month.

B. Inconsistency between datasets

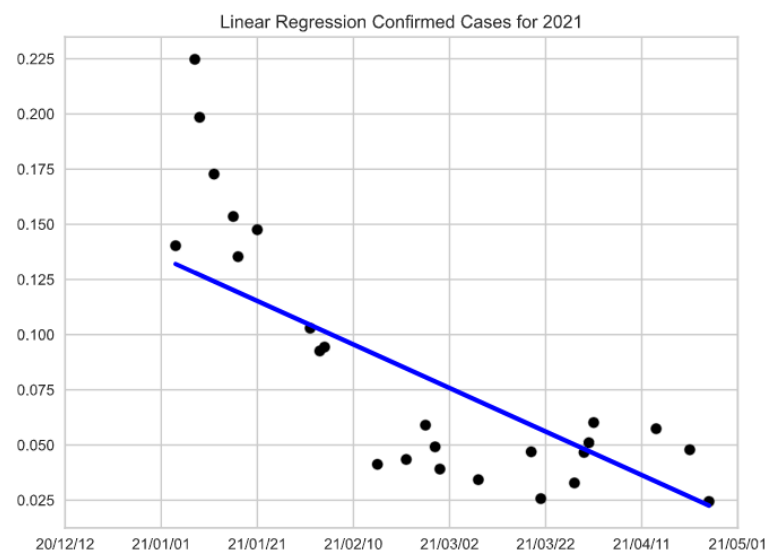
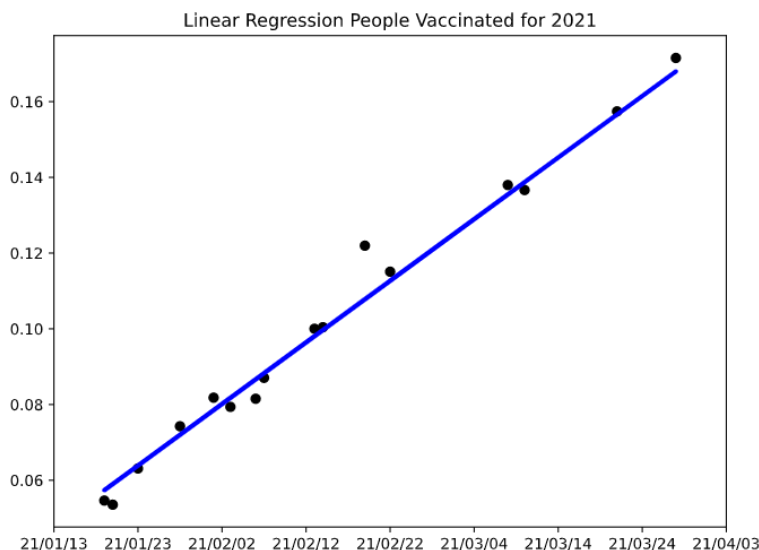
We have 2 vaccination datasets that have inconsistent data points due to the different way of collecting the data. That might have affected our results and might have caused some errors. We calculated error rates and divided our usable data to see how our model works on different data.

5. Modeling

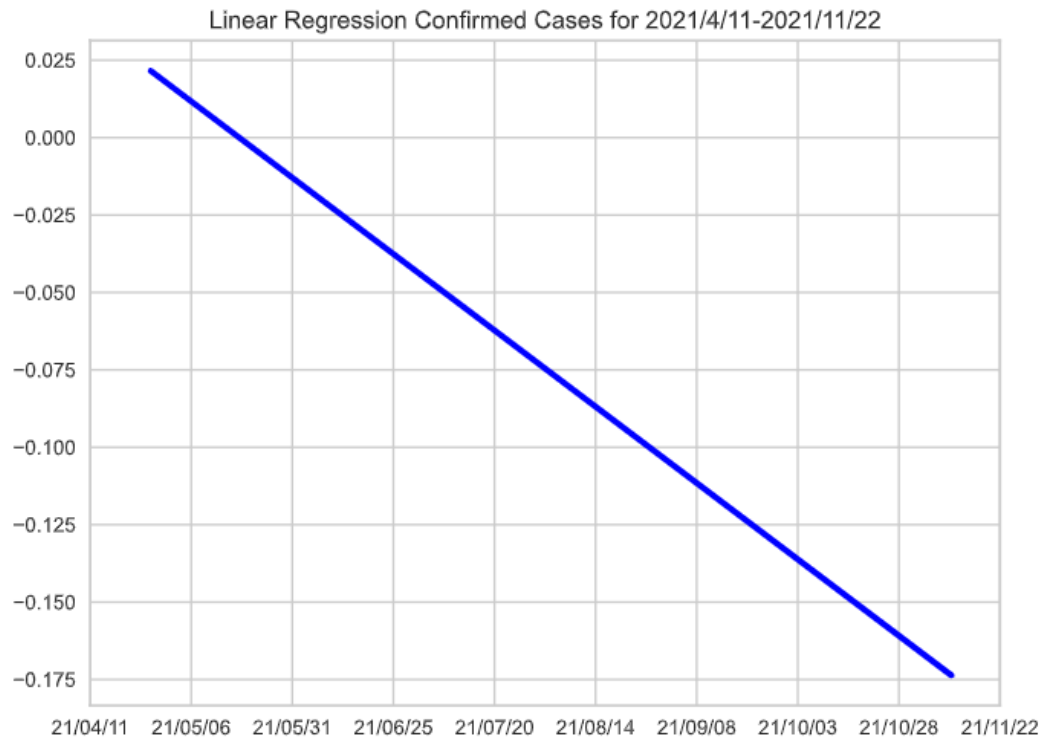
- Linear Regression

We used linear regression as a modeling technique to analyze the relationship between people getting fully vaccinated and confirmed COVID-19 cases. Linear regression is a linear approach to modelling the relationship between a scalar response and independent and dependent variables [3]. It is used to predict the value of a certain variable depending on the value of the other variable (hence, independent and dependent). Our independent variable is time, and dependent variables are people fully vaccinated, and confirmed COVID-19 cases. After independent and dependent variables have been identified, we attempted to make a model that fits linear equation relationship between the data.

We made linear regression models to see the relationship between time and vaccinations, and time and confirmed cases.



We used the regression model to predict the future of confirmed cases as you can see in the figure below. The linear regression shows us that by November of 2022, confirmed COVID-19 cases will significantly drop and vaccination rates will stabilize which means that we can safely come back in person. For example, CSUN plans to safely open more in-person classes by Spring 2022 which supports our results.



- Implementation issues and solutions

We encountered a few problems with building the model that was solved by tweaking the datasets and normalizing it and matching the time intervals.

6. Assessment of results

Results from the linear regression model show that by November 2022, we will not have confirmed COVID-19 cases given that the rate of decline of new cases remains the same. In addition we can also find out that we will reach 70% of the population vaccinated in 13.7 months.

- Interpret models according to domain knowledge

6.1.1. Data mining success criteria

Our linear regression model aligns with the business goal and confirms our hypothesis. We wanted to see when we will reach 70% of immunity and our model predicted 13.7 months from today (May 10th). Data sources we selected for this project match the goals we have set as in our datasets containing valuable information about COVID-19 cases and vaccinations. The datasets had complete data, all missing values were filled by the nearest neighbor, null and NaN values were removed.

6.1.2. Test result

The result of the model predicted that in 13.7 months, the U.S. population will reach 70% immunity and according to the experts that is an acceptable percentage to get back safely in person. Also, CSUN's plans to open campus next year align with the results of our model as it confirms that we also got the same date. Another source that we found, also concluded similar results (source [4]).

7. Summary of learning experiences

- 7.1. Overall, we learned about data mining techniques such as CRISP-DM and how to apply it to projects. CRISP-DM model uses hierarchical breakdown that consists of different sets of tasks [5] and in our case CRISP-DM was applied into all 6 parts of our project. Lifecycle of our project was as follows: 1) we developed business understanding and goals of our idea, 2) we found datasets and tried to understand their contents, 3) we cleaned and prepared our datasets for further usage, 4) we built statistical models, 5) we evaluated our models. In our learning experience, we went through all stages of the CRISP-DM lifecycle while working on this project. During the implementation stage, we learned a lot about different statistical models and how each of them differs. After our final presentation, during the QA session, we learned that for our datasets and problems we were trying to solve we could also use the Lasso regression model.

8. References

- [1] Patrícia Gonzalez-Dias et al. 2020. Methods for predicting vaccine immunogenicity and reactogenicity. (2020). Retrieved May 21, 2021 from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7062420/>
- [2] Edouard Mathieu. State-by-state data on COVID-19 vaccinations in the United States. Retrieved May 21, 2021 from <https://ourworldindata.org/us-states-vaccinations>
- [3] Shubham Jain I am currently pursuing my B.Tech in Ceramic Engineering from IIT (B.H.U) Varanasi. I am an aspiring data scientist and a ML enthusiast. I am really passionate about changing the world by using artificial intelligence. (2020, October 18). Linear, Ridge and Lasso Regression comprehensive guide for beginners. Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2017/06/a-comprehensive-guide-for-linear-ridge-and-lasso-regression/>.
- [4] Anon. Path to Normality - COVID-19 Vaccine Projections. Retrieved May 21, 2021 from <https://covid19-projections.com/path-to-herd-immunity/>
- [5] Chapman, P. (n.d.). CRISP-DM 1.0 (Tech.).

9. Appendix

Attached is the code used during the project for cleaning processing data and creating the model and validating it.