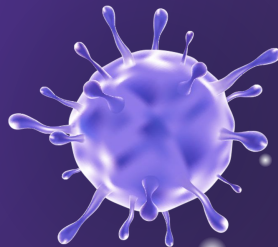# Covid-19 Infection Analysis

*Tonirose Babasoro, Elsa Bustos, Utsav Chaudhary, India Scott*

Program: Data Analytics
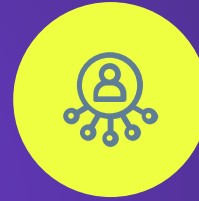Group Name: Endless Knot
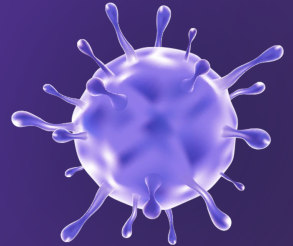Presentation Date: April 7, 2022

# INTRODUCTION

## TOPIC

What is the likelihood of being infected by Covid-19? How is infection affected by factors such as vaccination rates, gender, and ethnicity?

## PURPOSE

It is important to analyze future trends following the Covid-19 pandemic to understand the prevalence of infection within the American population.

# Questions Answered

- Are certain populations more likely to be infected than others?
- How do these factors affect the other?
- What other factors should be considered in identifying risks of infection?

# Data Sources

We gathered data from reliable organizations such as New York Times and the Center for Disease Control (CDC) which provide csv files on their findings. Our data identifies:
- vaccination rates
- gender ratios
- ethnicity statistics

# Data Exploration

- After manual review of multiple csvs, we were able to identify a primary key: "States"
- Data is filtered and aggregated via SQL server and filter null values with pandas (Python)

# Data Analysis

- Ordinary Least Squares (OLS)
- Linear regression
- SVM support vector machine
- Autoregressive Integrated Moving Average (ARIMA) for Time series model

coefficient of determination: 0.5703788862651746

## OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | Date | R-squared: | 0.607 |
| Model: | OLS | Adj. R-squared: | 0.599 |
| Method: | Least Squares | F-statistic: | 75.33 |
| Date: | Sat, 02 Apr 2022 | Prob (F-statistic): | 0.00 |
| Time: | 00:12:10 | Log-Likelihood: | -48837. |
| No. Observations: | 5320 | AIC: | 9.789e+04 |
| Df Residuals: | 5212 | BIC: | 9.860e+04 |
| Df Model: | 107 | | |
| Covariance Type: | nonrobust | | |

## SARIMAX Results

| | | | |
|---|---|---|---|
| Dep. Variable: | recovered | No. Observations: | 5320 |
| Model: | ARIMA(5, 1, 0) | Log Likelihood | -55067.188 |
| Date: | Fri, 01 Apr 2022 | AIC | 110146.376 |
| Time: | 23:52:26 | BIC | 110185.850 |
| Sample: | 0 | HQIC | 110160.168 |
| | - 5320 | | |
| Covariance Type: | opg | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| ar.L1 | -0.8148 | 0.007 | -124.881 | 0.000 | -0.828 | -0.802 |
| ar.L2 | -0.7459 | 0.008 | -97.523 | 0.000 | -0.761 | -0.731 |
| ar.L3 | -0.5328 | 0.007 | -76.355 | 0.000 | -0.546 | -0.519 |
| ar.L4 | -0.3472 | 0.006 | -58.862 | 0.000 | -0.359 | -0.336 |
| ar.L5 | -0.2130 | 0.004 | -52.946 | 0.000 | -0.221 | -0.205 |
| sigma2 | 5.937e+07 | 1.79e-10 | 3.31e+17 | 0.000 | 5.94e+07 | 5.94e+07 |

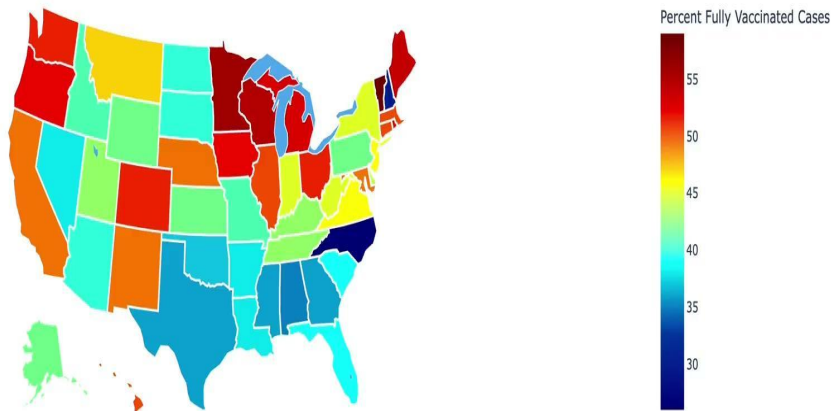| | | | |
|---|---|---|---|
| Ljung-Box (L1) (Q): | 2.19 | Jarque-Bera (JB): | 5558677.83 |
| Prob(Q): | 0.14 | Prob(JB): | 0.00 |
| Heteroskedasticity (H): | 0.12 | Skew: | 6.56 |
| Prob(H) (two-sided): | 0.00 | Kurtosis: | 160.83 |

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
[2] Covariance matrix is singular or near-singular, with condition number 5.21e+32. Standard errors may be unstable.
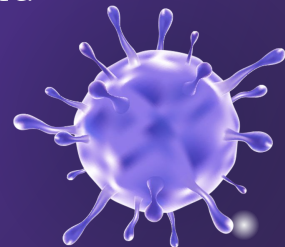
# Interactive Element



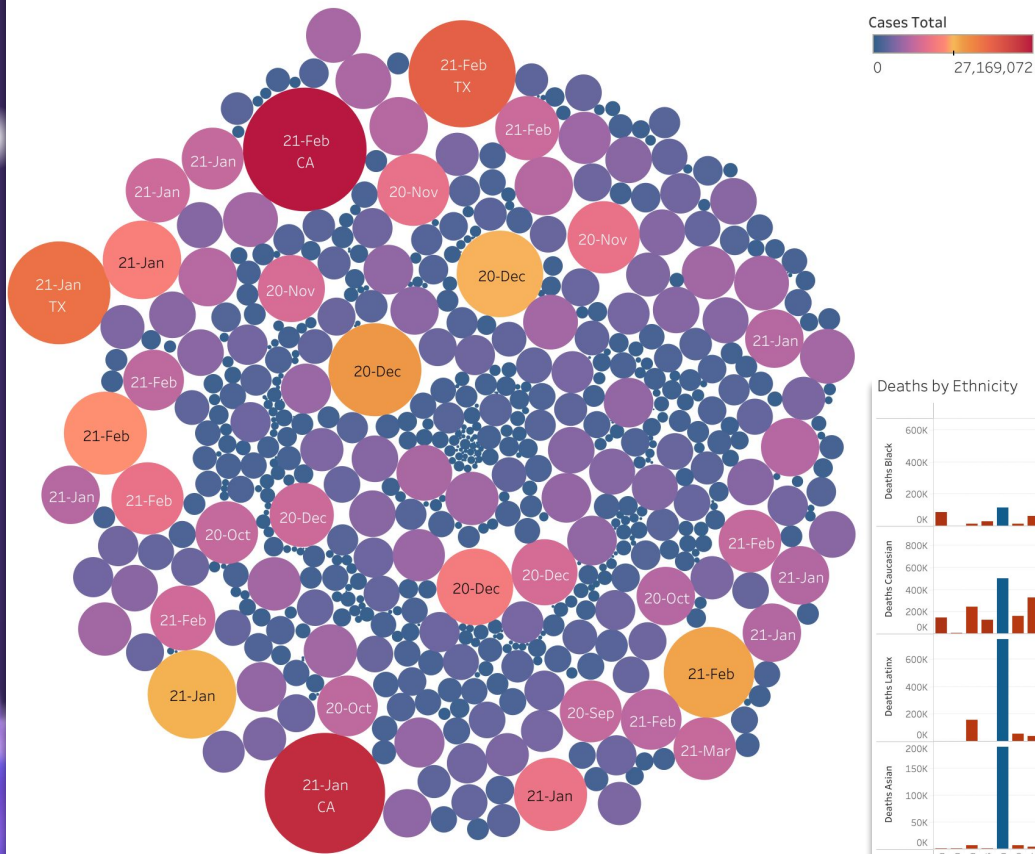US Vaccination Rates

Percent Fully Vaccinated Cases

*Vaccination Rates by State*

Plotly
- Map with options to toggle different factors of each state such as:
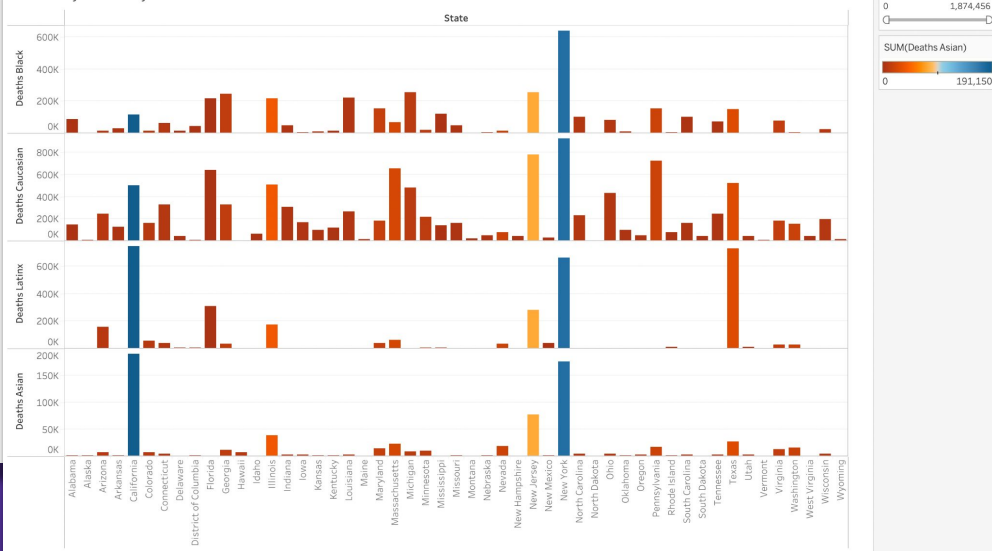  - Vaccination rates
  - Infections
  - Gender data
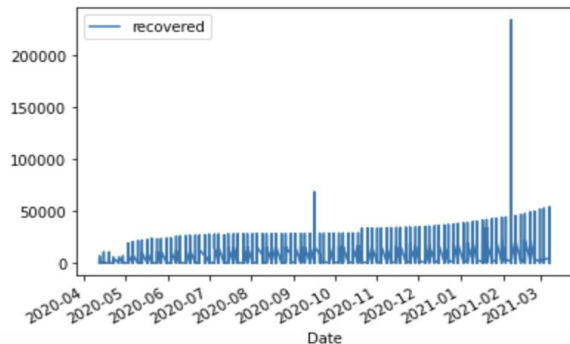
# Visualization

- Tableau
  - Correlation graphs

# Results of Analysis

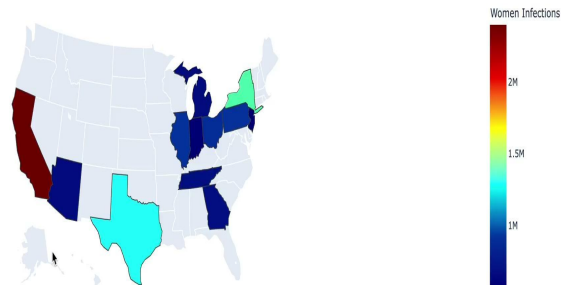- No significant differences between gender infection rates
- More vaccinations → lower cases and death rates
- People who were vaccinated and boosted had lower amounts of cases and fatalities
- On the time series analysis, the line plot is created, showing the expected values (blue) compared to the rolling forecast predictions (red). Values show some trends and are in the correct scale.

# What We Would Do Differently

- Before deciding on a topic, look into data resources
- Compare more factors from different data sets
  - Identify more primary keys to connect data
- Added other factors such as economic impact

# Recommendations for Future Analysis

- Statistical analysis can be used with bigger data sets
- Explore other factors that may have a correlation with the data we already have
- Expand to global research
  - countries that are not as ethnically homogenous
- Project future peak infections

# Gender

https://www.genderscilab.org/gender-and-sex-in-covid19/#CaseDeathRatebySex

# Ethnicity

https://covidtracking.com/race/dashboard

# Vaccination

d.cdc.gov/covid-data-tracker/#vaccinations_vacc-total-admin-rate-total

Check out our project at:
https://github.com/antirose/
CovidInfectionAnalysis

Covid-19
Infection Analysis

Tonirose Babasoro, Elsa Bustos, Utsav Chaudhary, India Scott

Program: Data Analytics
Group Name: Endless Knot
Presentation Date: April 7, 2022