

# How much did restricting mobility help in controlling the spread of COVID?

Satheesh Joseph, Paco Valdez, Yi Zhang

## Introduction

### Overview

The Covid-19 pandemic has taken more than 2.5 million lives worldwide. In the U.S. alone, more than half a million people died because of it. COVID-19 is thought to spread mainly through close contact from person to person, including among people who are physically near each other (within about 6 feet).[1]

As the virus spread, many cities and other areas have implemented a shutdown policy that significantly restricted the mobility of the population. Specifically, many companies have opted for a complete work-from-home policy.

As a team of Data Scientists, we're interested in assessing the effectiveness of restricting mobility in controlling the spread of Covid-19.

### Research Question

Rather than investigating the effect of social distancing on a micro and personal interaction level, we are interested in its effect on a policy level. Specifically, the Research Question we're asking is:

How much did restricting social mobility help in controlling the spread of COVID?

The basic causal theory we're working under can be expressed by the diagram below:

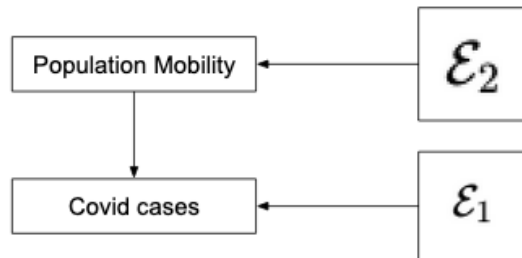


Figure 1: The Basic Causal Theory

## Operationalization

To operationalize the research question, we propose the following plan for the causal models.

**Firstly**, we will use the New York Times Covid-19 data for the number of Covid cases per State. And we will use the COVID-19 Community Mobility Report for the population mobility scores. For the base model, we will use the “residential mobility score” as a relatively stable representative of the general population’s social mobility. We will also use Mask-Wearing survey data from the New York Times as a control variable.

**Secondly**, mobility have drastically different effect on the spread of the virus depending on the population density. Thus, we believe that having the raw number of cases as the outcome variable can be misleading. Therefore, we will be using the number of cases per 100,000 residents as the outcome variable to normalize the population density. The population density data will come from the COVID-19 US State Policy Database.

More formally, our base model will assume a causal relationship of the form:

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score$$

where each data point will represent a State in the U.S., and the outcome variable will be the total number of Covid-19 cases per 100,000 residents for that state from March 2020 up until to July 2020, when the mask wearing survey was taken. . And the independent variable will be the **average** residential mobility score of the state in the time period.

The mask wearing survey data will be operationalized by transforming it into probability of people wearing mask. The probability of people wears masks is calculated by assuming that survey respondents who answered ‘Always’ were wearing masks all of the time, those who answered ‘Frequently’ were wearing masks 80 percent of the time, those who answered ‘Sometimes’ were wearing masks 50 percent of the time, those who answered ‘Rarely’ were wearing masks 20 percent of the time and those who answered ‘Never’ were wearing masks none of the time.

For a more advanced model, we propose to use some control variables to assess the difference mobility score makes in relation to other factors, notably we’ll include whether the “Stay at Home” order is in effect and whether the face mask policy is in effect and how long has been in place.

For the most inclusive model, we will explore all the other mobility scores, including to workplaces, parks, grocery & pharmacy, as well as retail & recreation.

## The models

### Data Exploration

##	state	cases_per_100k	workplaces	transit
##	Length:50	Min. : 89.48	Min. : -38.03	Min. : -58.000
##	Class :character	1st Qu.: 594.01	1st Qu.: -33.01	1st Qu.: -29.927
##	Mode :character	Median : 865.29	Median : -30.29	Median : -14.419
##		Mean : 917.10	Mean : -29.94	Mean : -12.763
##		3rd Qu.: 1220.34	3rd Qu.: -26.62	3rd Qu.: 4.145
##		Max. : 2091.42	Max. : -21.77	Max. : 38.742
##	parks	grocery	residential	stay_home
##	Min. : -34.23	Min. : -19.1290	Min. : 2.194	Min. : 0.00
##	1st Qu.: 68.67	1st Qu.: -0.1452	1st Qu.: 5.847	1st Qu.: 1.00
##	Median : 110.27	Median : 4.5806	Median : 7.484	Median : 1.00
##	Mean : 121.23	Mean : 4.9032	Mean : 7.558	Mean : 0.78
##	3rd Qu.: 176.15	3rd Qu.: 7.9758	3rd Qu.: 9.097	3rd Qu.: 1.00
##	Max. : 342.16	Max. : 30.8065	Max. : 12.097	Max. : 1.00
##	fm	mask_prob	mask_days	

##	Min.	:0.00	Min.	:0.5645	Min.	: 0.00
##	1st Qu.	:1.00	1st Qu.	:0.7095	1st Qu.	: 29.75
##	Median	:1.00	Median	:0.7736	Median	: 69.50
##	Mean	:0.82	Mean	:0.7761	Mean	: 55.92
##	3rd Qu.	:1.00	3rd Qu.	:0.8448	3rd Qu.	: 75.00
##	Max.	:1.00	Max.	:0.9327	Max.	:103.00

All the variables have their means and medians very close to each other, that suggests that there is no severe skewness.

The histogram plots below shows no major deviation from normality for the variables.

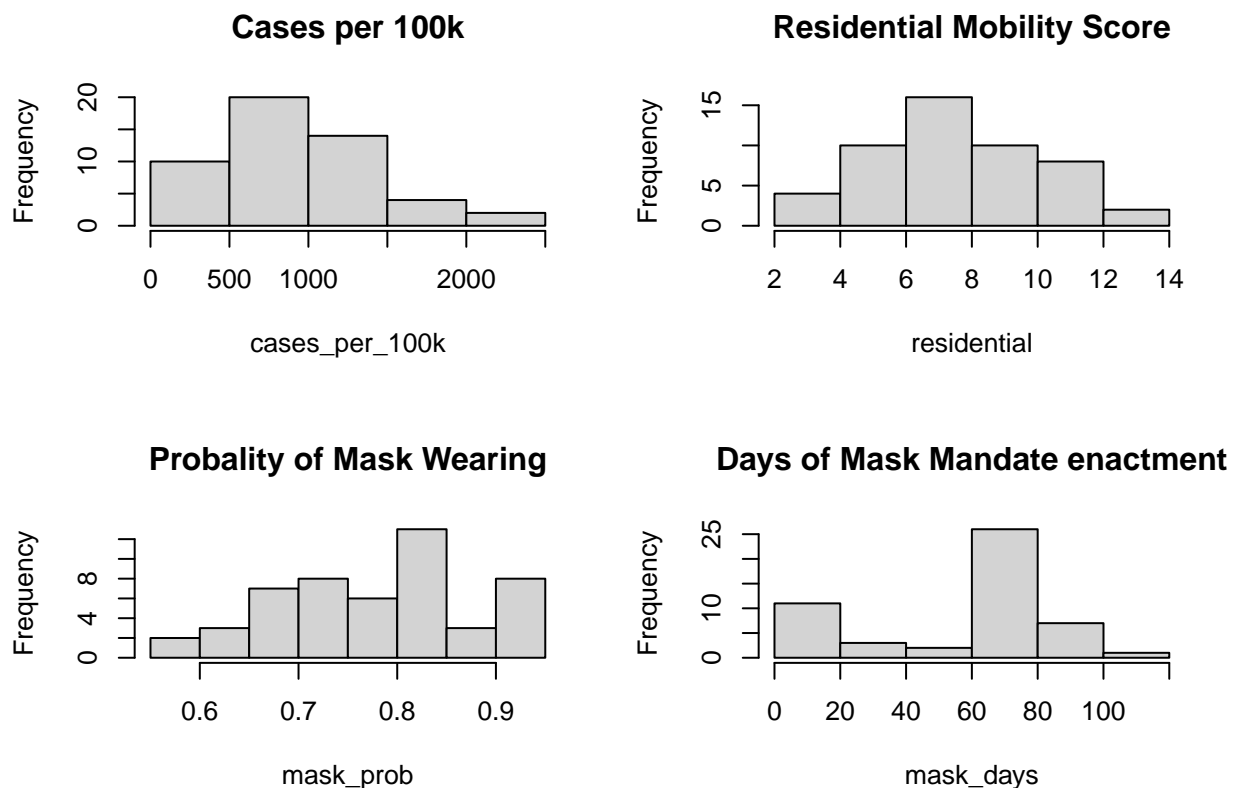


Figure 2: Histograms of the left-hand side variable and selected right-hand side variables

The other mobility variables also doesn't show major deviation from normality.

Using natural logarithm transformation on the variables doesn't show a significant improvement. We will continue the analysis without the transformation.

By looking at the scatter plots we can notice that **workplaces** is the inverse of the **residential** mobility variable. Using both in a OLS regression probably will cause the coefficients to cancel each other. The **residential** variable shows promising results, it looks like it has a very linear relationship with **cases\_per\_100k**.

## The Linear Models

Our intuition told us that the residential variable would be the best proxy to represent mobility. At this point in time during the pandemic, the U.S. was registering its second wave of Covid-19 cases and non-essential business were still closed in most of the states. One way to explain this relationship is that in the states

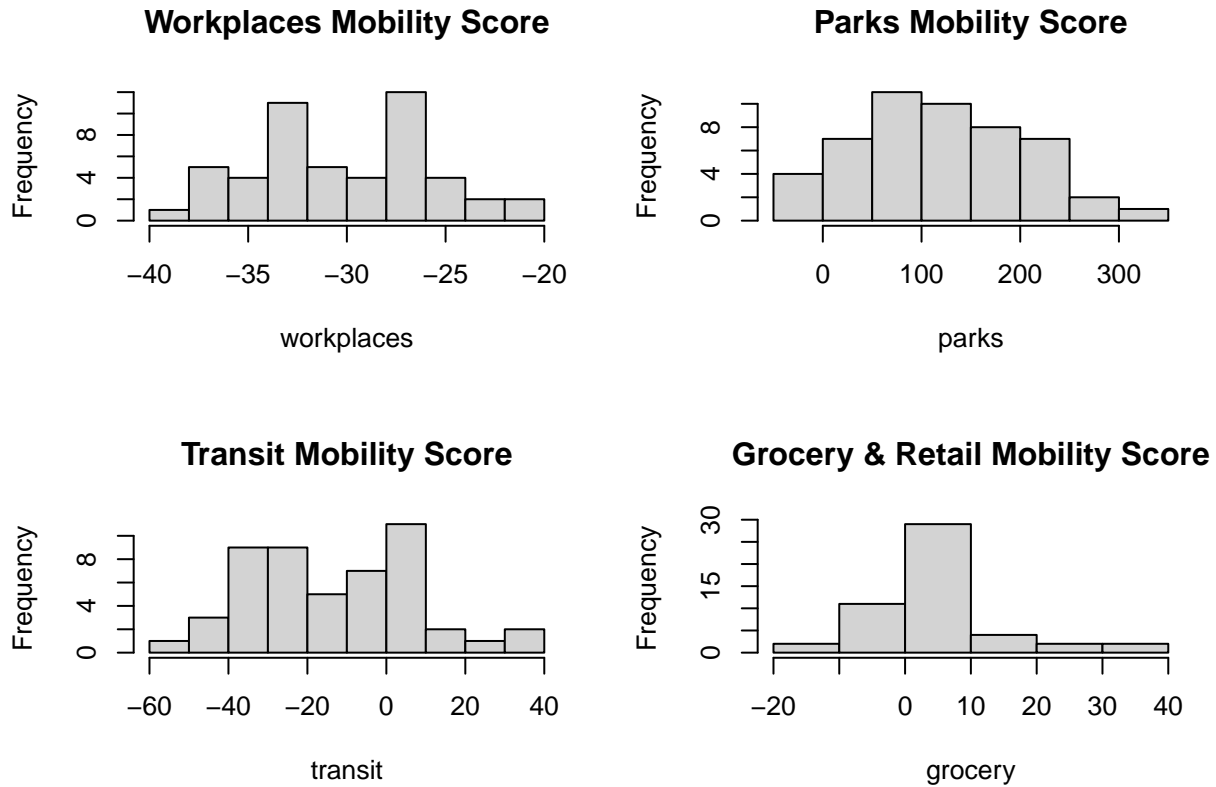


Figure 3: Histograms of the mobility variables

where the residential mobility score increased, people were not sheltering at their homes. Instead people were going out with friends and family, thus spreading the virus.

## Model One

In our initial model we will use the residential mobility score and cases per 100,000 inhabitants:

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1261.67  -283.09   -38.77   281.55   859.63
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   167.67     186.18   0.901   0.372
## residential    99.16      23.36   4.245 9.96e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 417.9 on 48 degrees of freedom
```

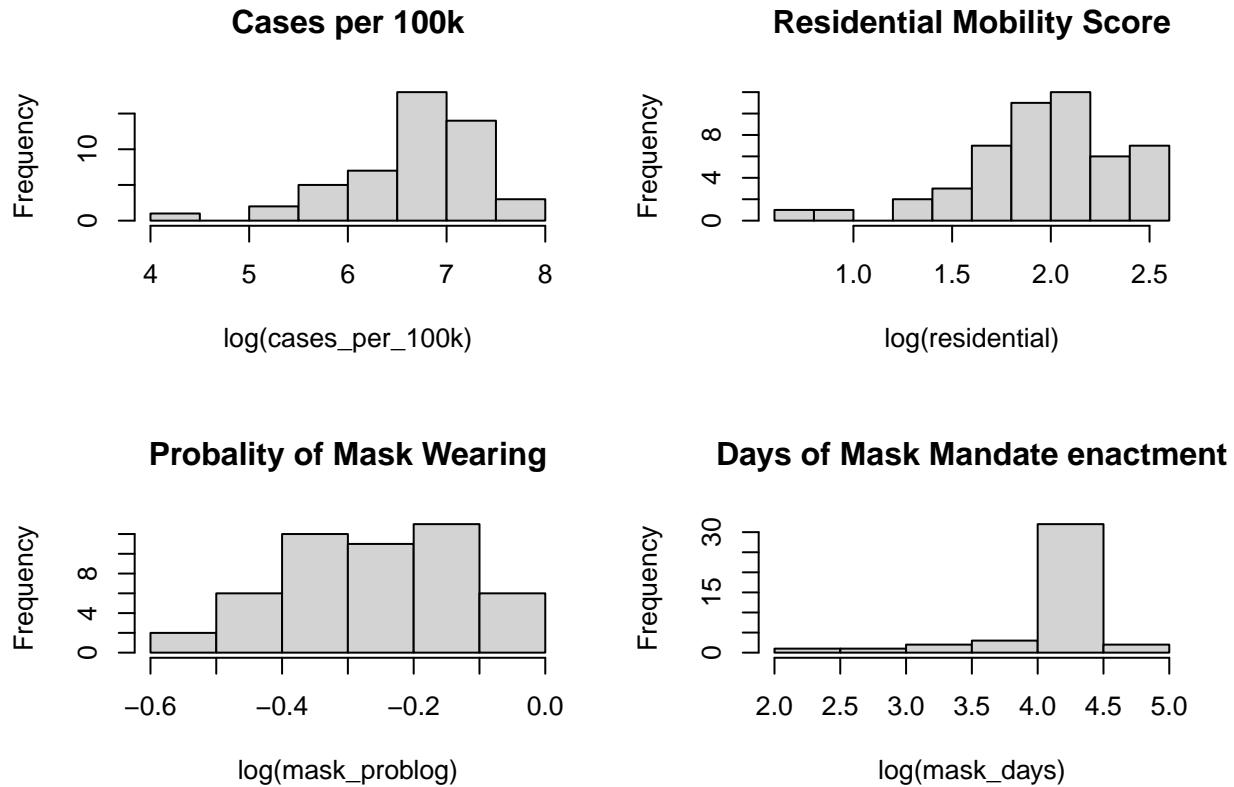


Figure 4: Histograms of selected log-transformed variables

```
## Multiple R-squared:  0.2729, Adjusted R-squared:  0.2578
## F-statistic: 18.02 on 1 and 48 DF,  p-value: 9.955e-05
```

We got a very significant p-value, which highlights the string relationship between these two variables. The first control variable we will introduce is the face-mask mandate (fm). It did not reduce the significance of the mobility variable and it showed no significant effect on the dependent variable.

$$\text{cases\_per\_100k\_residents} = \beta_0 + \beta_1 \text{residential\_mobility\_score} + \beta_2 \text{face\_mask\_mandate}$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + fm, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1276.44  -279.17   14.51   249.18   881.52
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    212.09     194.91   1.088 0.282081
## residential     108.18      26.01   4.160 0.000134 ***
## fm             -137.31     171.25  -0.802 0.426677
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 419.4 on 47 degrees of freedom
```

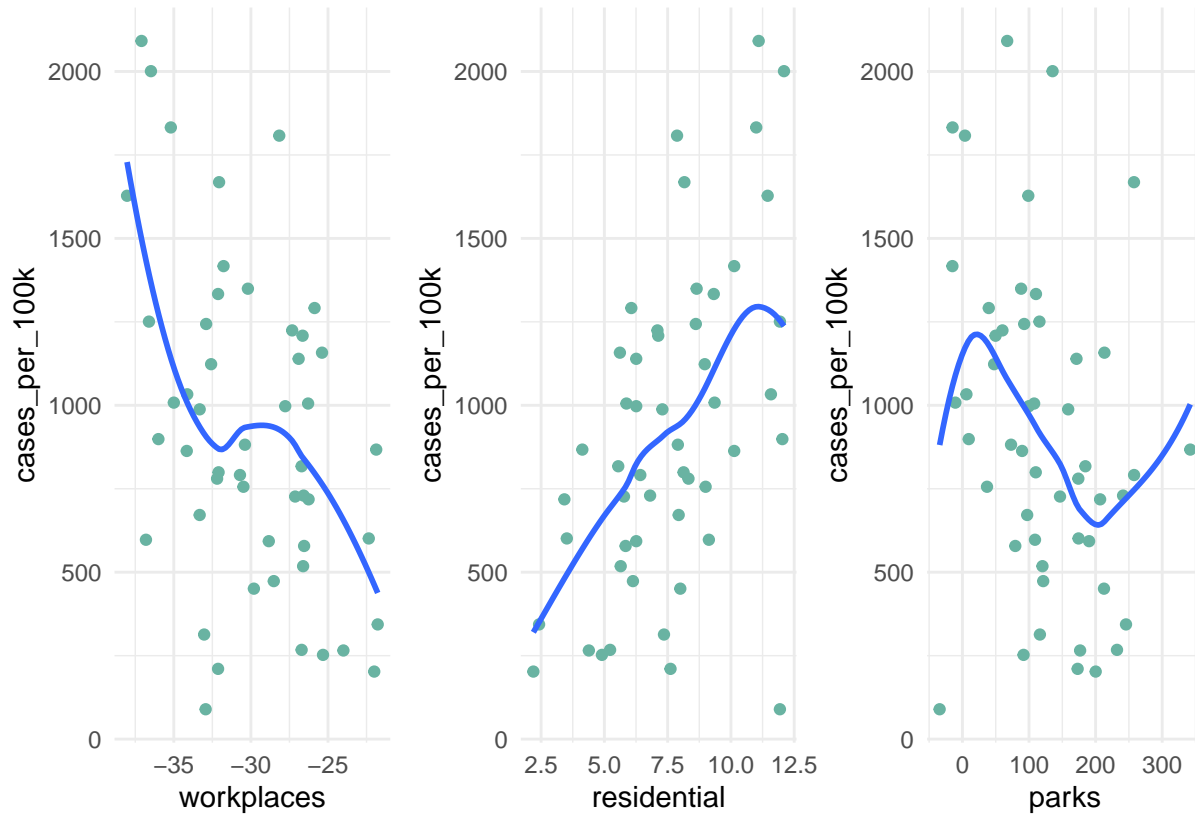


Figure 5: Scatter Plots for Cases per 100k and selected mobility variables

```
## Multiple R-squared:  0.2827, Adjusted R-squared:  0.2522
## F-statistic: 9.264 on 2 and 47 DF,  p-value: 0.0004058
```

The second control variable is if the state had a stay at home order effective on the same time frame. This control variable again did not reduce the significance of the residential mobility score and also is not significant.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 face\_mask\_mandate + \beta_3 stay\_at\_home\_order$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + fm + stay_home, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1274.91  -272.70    4.78   257.20   905.81
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   262.00     205.63   1.274  0.209012
## residential   113.77      27.05   4.206  0.000119 ***
## fm           -136.85     171.93  -0.796  0.430151
## stay_home     -118.72     150.20  -0.790  0.433355
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

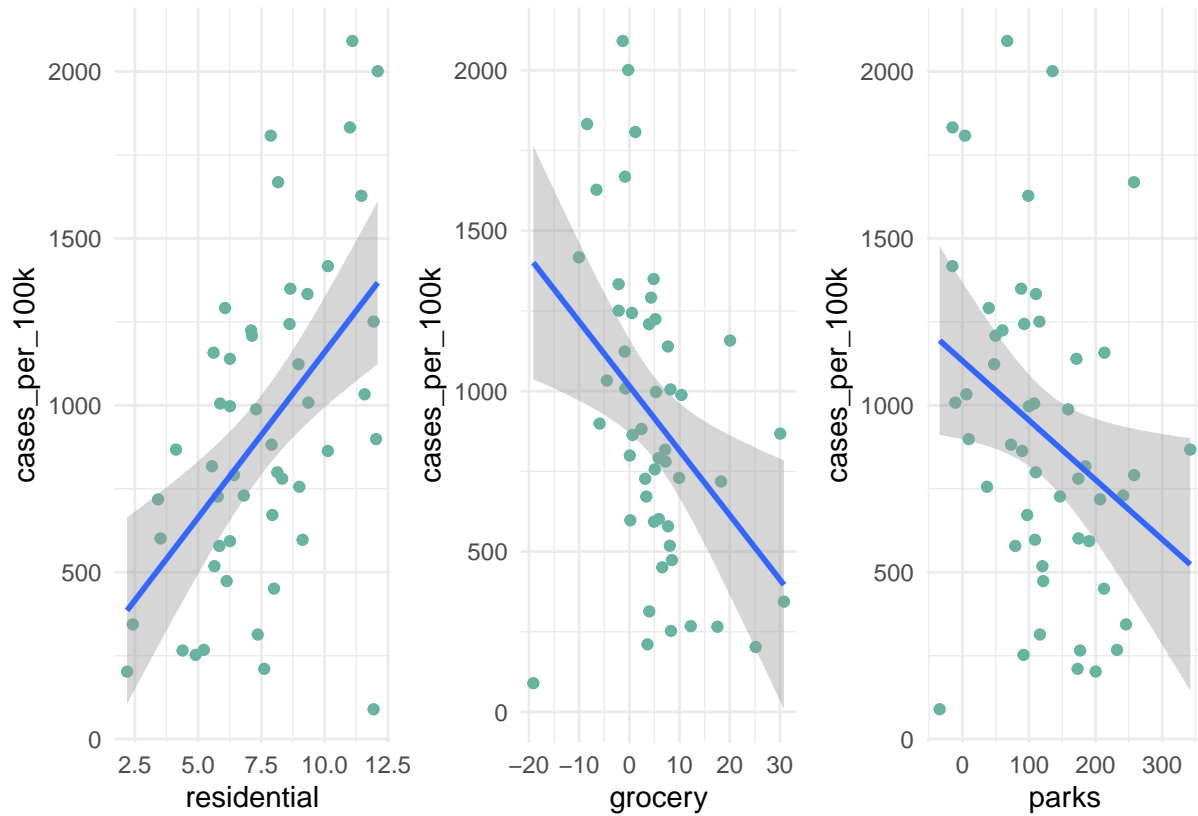


Figure 6: Scatter Plots for Cases per 100k and selected mobility variables with OLS regression line and standard errors

```
##
## Residual standard error: 421.1 on 46 degrees of freedom
## Multiple R-squared:  0.2924, Adjusted R-squared:  0.2462
## F-statistic: 6.335 on 3 and 46 DF,  p-value: 0.001095
```

## Model Two

In the second model we wanted to test if the probability of wearing a mask would have a significant impact on the amount of Covid-19 cases. Although by looking at the scatter plot of Cases per 100k and mask wearing probability there is no clear relationship, and the slope of the OLS line has the opposite direction to what our intuition tell us.

So, our second model contains the residential mobility score and the face mask wearing probability. The coefficient for Mask wearing probability (`mask_prob`), at least has the correct direction but fails to achieve significance. The residential mobility score coefficient is still significant, which give us more trust in that variable.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 mask\_wearing\_probability$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + mask_prob, data = data)
##
## Residuals:
```

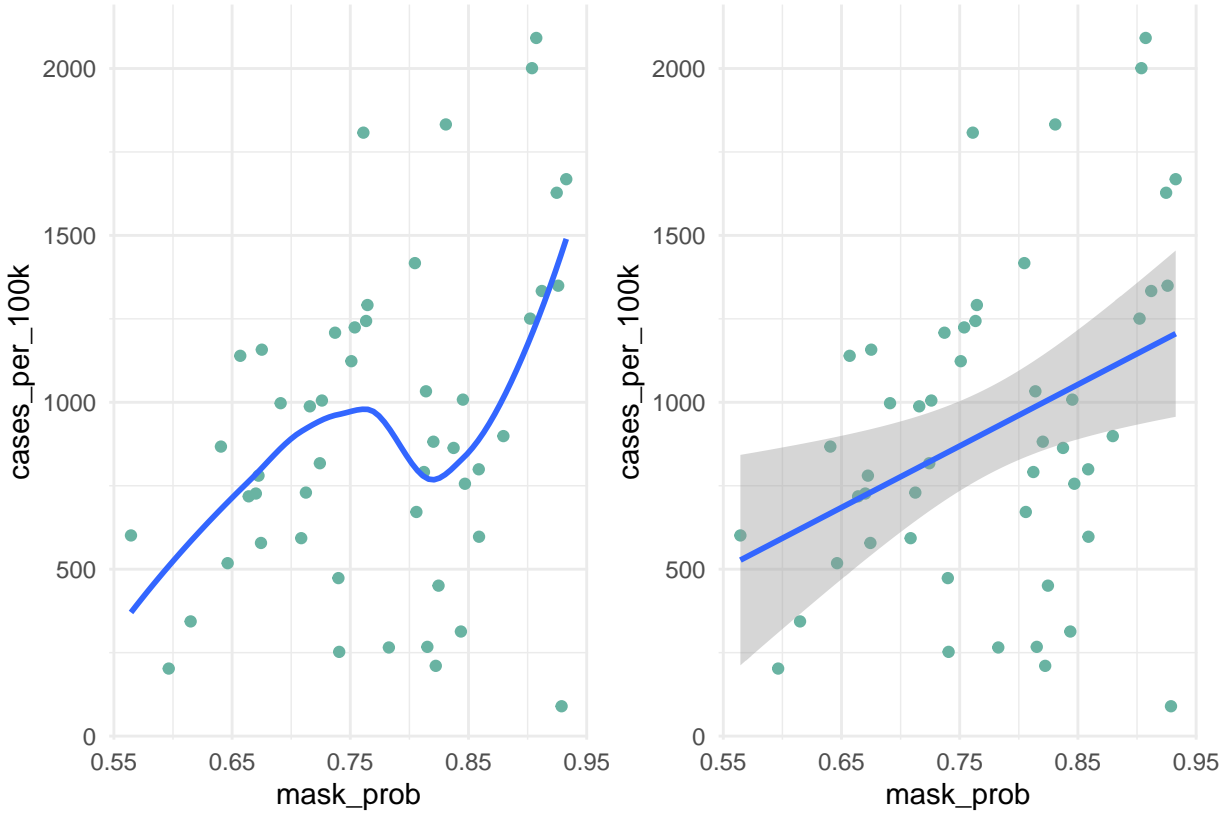


Figure 7: Scatter Plots for Cases per 100k and mask wearing probability

```
##      Min      1Q   Median      3Q      Max
## -1246.65 -279.95  -48.37   284.55   841.52
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    582.09     601.10   0.968  0.33781
## residential    122.15      39.44   3.097  0.00329 **
## mask_prob     -757.90    1044.70  -0.725  0.47176
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 419.9 on 47 degrees of freedom
## Multiple R-squared:  0.281, Adjusted R-squared:  0.2504
## F-statistic: 9.184 on 2 and 47 DF,  p-value: 0.0004298
```

Again, we can test our previous control variables. Even with the presence of the other control variables the residential mobility score coefficient is still significant.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 mask\_wearing\_probability + \beta_3 face\_mask\_mandate$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + mask_prob + fm +
```



```
##      stay_home, data = data)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -1267.2   -280.1        3.0    240.6    892.8
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    423.41     634.36   0.667  0.50789
## residential    121.67      40.08   3.035  0.00398 **
## mask_prob     -314.79    1169.01  -0.269  0.78894
## fm            -121.53     182.78  -0.665  0.50950
## stay_home     -105.04     160.01  -0.656  0.51488
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 425.4 on 45 degrees of freedom
## Multiple R-squared:  0.2935, Adjusted R-squared:  0.2307
## F-statistic: 4.673 on 4 and 45 DF,  p-value: 0.003063
```

### Model Three

Our secondary question was if the other mobility variables could have a significant impact on the cases per 100,000 variable. In our data exploration, workplaces also looked promising but it looked that it acted in the opposite direction than residential. When we tested this variable alone we confirmed our intuition, the workplaces mobility score is also significant but acts in the opposite direction than residential.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 workplaces\_mobility\_score$$

```
##
## Call:
## lm(formula = cases_per_100k ~ workplaces, data = data)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -972.46   -267.53   -55.45    335.49    976.91
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -532.23     436.14  -1.220  0.22831
## workplaces     -48.40      14.42  -3.358  0.00155 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 441 on 48 degrees of freedom
## Multiple R-squared:  0.1902, Adjusted R-squared:  0.1733
## F-statistic: 11.27 on 1 and 48 DF,  p-value: 0.001546
```

When both are present residential still has a more significant effect. And the coefficient for workplaces mobility score is significantly lower and in the opposite direction than the previous model.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 workplaces\_mobility\_score$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + workplaces, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1357.76  -266.15   -22.43   264.53   864.94
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    625.46     629.14   0.994  0.3252
## residential    138.31      56.49   2.449  0.0181 *
## workplaces     25.17      33.03   0.762  0.4498
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 419.7 on 47 degrees of freedom
## Multiple R-squared:  0.2818, Adjusted R-squared:  0.2512
## F-statistic: 9.221 on 2 and 47 DF,  p-value: 0.0004184
```

Our definitive third model is the comparison of all the mobility score variables, with the exception of the workplaces variable since it's the opposite of residential. The residential mobility score captures the information contained in the other mobility variables.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 parks\_mobility\_score + \beta_3 grocery\_mobility\_score +$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + parks + grocery +
##      transit, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1126.90  -271.39   -22.67   269.73   879.91
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -248.8399   438.7873  -0.567  0.57346
## residential   161.9561    56.7735   2.853  0.00653 **
## parks         -0.1151     1.0914  -0.105  0.91647
## grocery        6.5432     15.7669   0.415  0.68012
## transit        5.9756      6.9167   0.864  0.39220
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 424.1 on 45 degrees of freedom
## Multiple R-squared:  0.2977, Adjusted R-squared:  0.2353
## F-statistic: 4.769 on 4 and 45 DF,  p-value: 0.002709
```

We can also add our control variables to the previous model.

$$cases\_per\_100k\_residents = \beta_0 + \beta_1 residential\_mobility\_score + \beta_2 parks\_mobility\_score + \beta_3 grocery\_mobility\_score +$$

```
##
## Call:
## lm(formula = cases_per_100k ~ residential + parks + grocery +
##      transit + fm + stay_home, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1174.95  -288.55    9.24   257.06   908.43
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -80.27862   506.56516  -0.158   0.8748
## residential   158.49786    58.19890   2.723   0.0093 **
## parks         -0.05455    1.11436  -0.049   0.9612
## grocery        3.43911    16.63787   0.207   0.8372
## transit        5.23591     7.13510   0.734   0.4670
## fm          -100.55881   183.83877  -0.547   0.5872
## stay_home    -78.88081   160.81790  -0.490   0.6263
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 431.4 on 43 degrees of freedom
## Multiple R-squared:  0.3058, Adjusted R-squared:  0.209
## F-statistic: 3.158 on 6 and 43 DF,  p-value: 0.01178
```

## CLM assumptions

**1. IID Sampling** None of the variables seem to violate IID.

**2. Linear Conditional Expectation** By looking at the predicted vs. residuals of the model it looks that relationship is very linear.

## 3. No Perfect Collinearity

```
##      (Intercept)  residential      parks      grocery      transit
## -80.27861521  158.49786207  -0.05454753   3.43910725   5.23590799
##              fm      stay_home
## -100.55881181  -78.88081299
```

None of the coefficients was dropped, that means there is no Perfect Collinearity.

**4. Homoskedastic Errors** The scale-location plot is very close to a flat line, which lead us to believe that there is no major issues with heteroskedasticity.

**5. Normally Distributed Errors** The histogram looks does not show a significant deviation from a normal distribution and the qqplot also shows very low deviation from normality.

In general we can say that all the 5 CLM Assumptions are met.

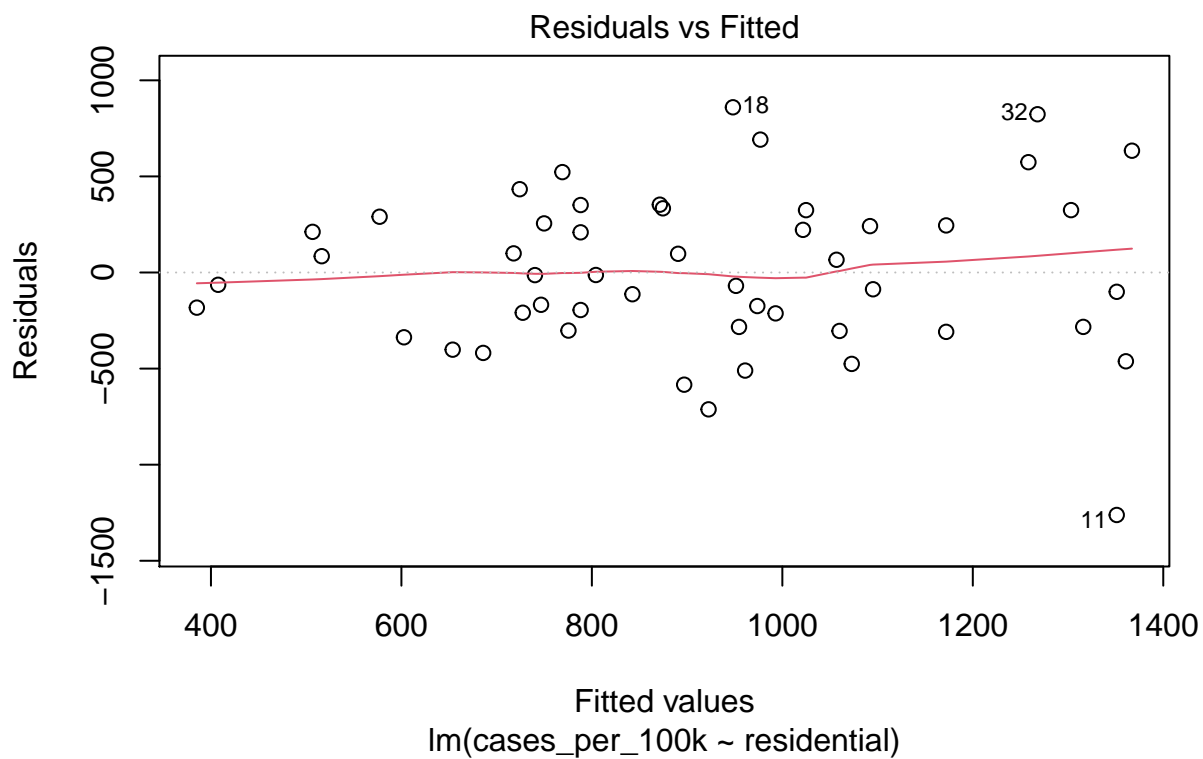


Figure 8: Residuals vs Fitted plot for Model One

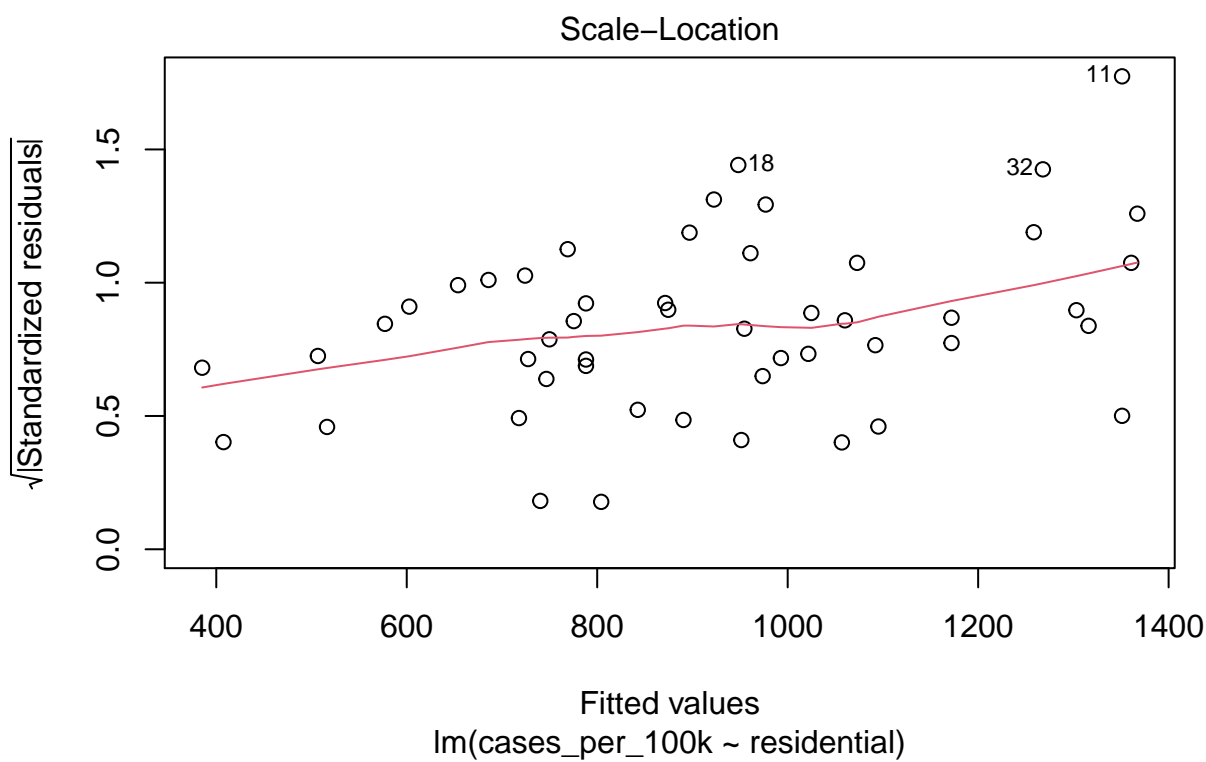


Figure 9: Scale-Location plot for Model One

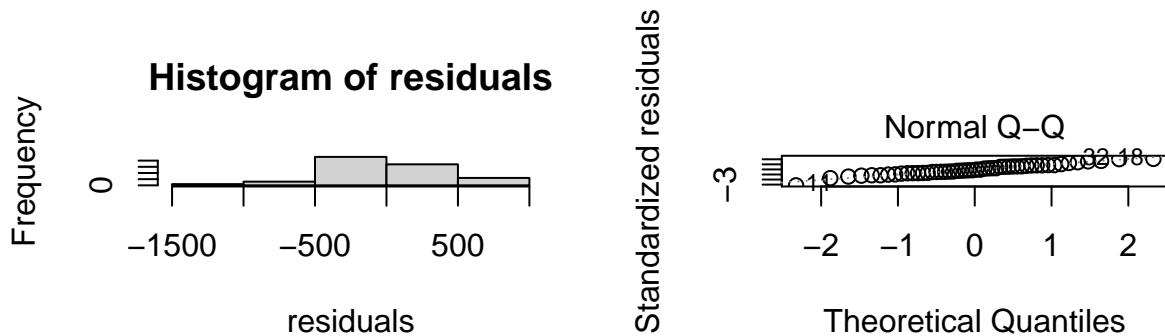


Figure 10: Histogram of residuals for model one and qq-plot for model one

## The Regression Table

## Model Limitations

The Research Question under study is well suited for a time series analysis. The initial discussion and work, were poised to well take off in including a weekly average and correlation with the cases offset by couple of weeks. As per the instruction and the limited machinery in place, the team has decided to take an average for the states from a time window where all the variables were produced and use it in the model.

### Classic Linear Model assessment

#### Omitted Variables

Diligently practicing mandates:

Though, State and County put restrictions, including Mask mandate, Curfew, Social distance, Quarantine, park, school and restaurant closures, how much general public abide to these mandates is a very good [Omitted] indicator that influence the total cases.

International/Non-Local Mobility:

As we all know Port of Entries like New York, San Francisco, Chicago and Seattle are the first affected places. State wide policies restricting mobility has varying effects on how close a location is to the airport or other transit points. Thus, closeness to transit stations and airport is an omitted variable.

All Percent drops are not equal:

A 10% drop in mobility is really good for sparsely populated regions to go below a certain threshold in transmitting the virus. The same cannot be applied to densely populated cities. Thus, the non-linearity in the measures taken vs impact seen is omitted.

## Conclusion

## References:

[1]: <https://www.cdc.gov/coronavirus/2019-ncov/faq.html>