

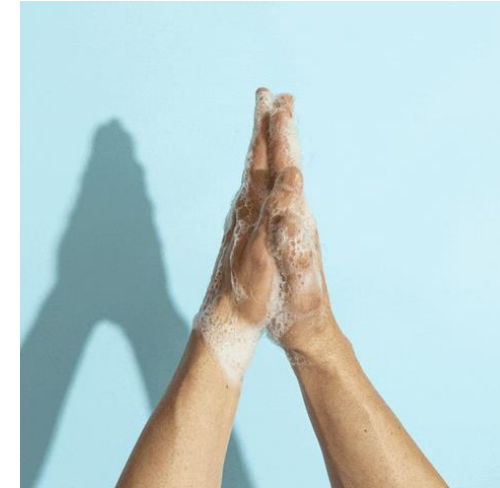


# **COVID-19 Data Analysis Project**

**Blue Team**

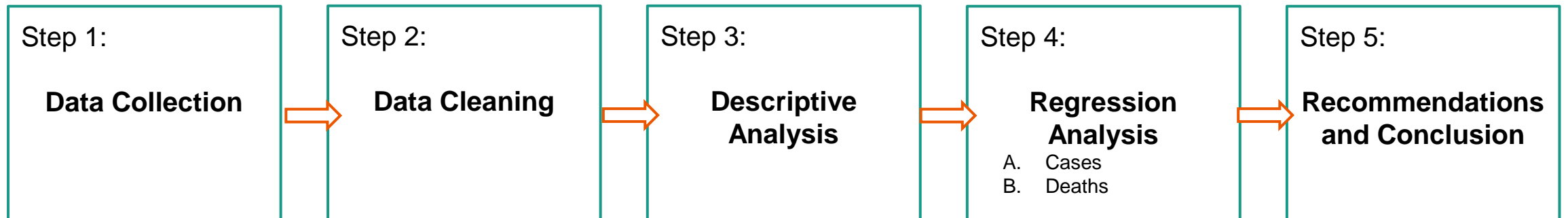
Nicole Gee, Nicholas Richmond, Matthew Clark, Swati Kohli, Joseph Carozza

# Introduction





# Project Outline



# Step 1: Data Collection

- The team brainstormed about different variables that may affect the number of COVID cases and deaths.
- Various sources of data were researched and added to a central document
- Those data files were downloaded or pulled directly from web pages to be cleaned and combined in R



**DATA COLLECTION**

## Step 2: Combining and Cleaning Data

- All of the .csv and .xlsx files were loaded into the environment in RStudio
- Country names and labels were 'cleaned' to match the primary JHU data source
- Finally, desired variables from the researched data sets were joined with the JHU data



# Final Data Set



## Finance & Sustainability

GDP/Capita (USD)  
Median Household Income  
Median Percapita Income  
Per Capita Spending on Health  
R&D Expenditure (% of GDP)



## Population & Trends

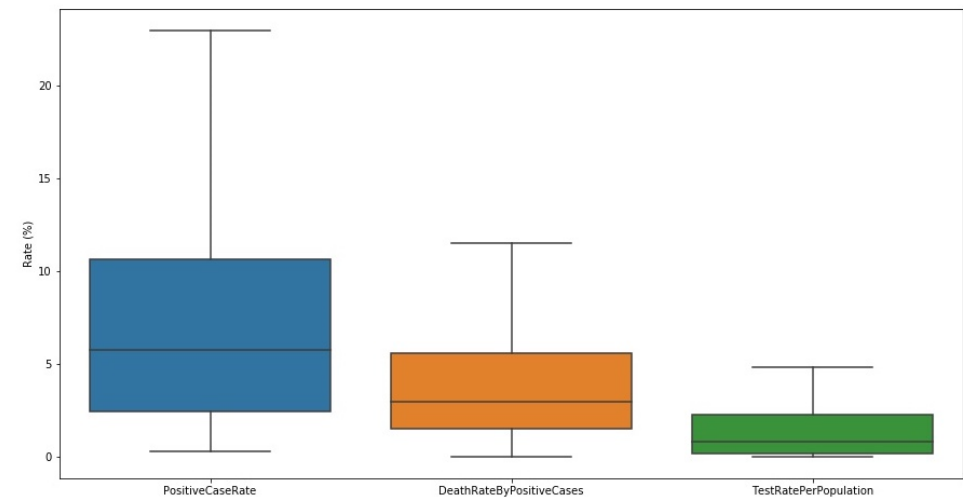
Population 2020  
% Urban Population  
People 65 & Over  
Median Age  
Proportion of Female Population  
Smoking prevalence (% of adults)



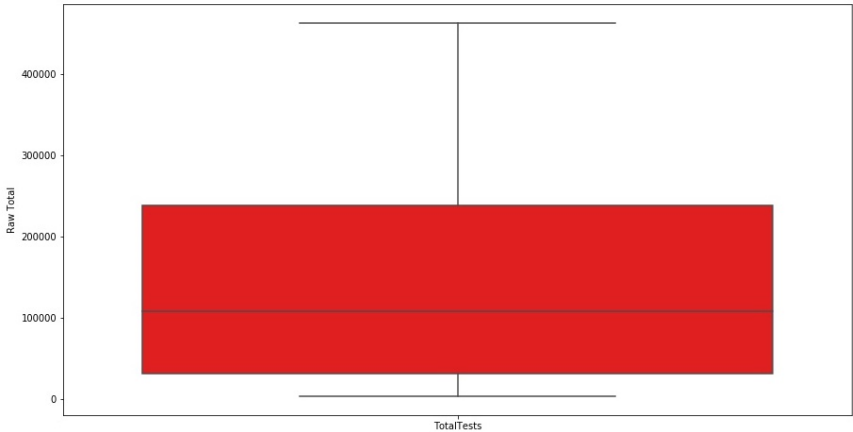
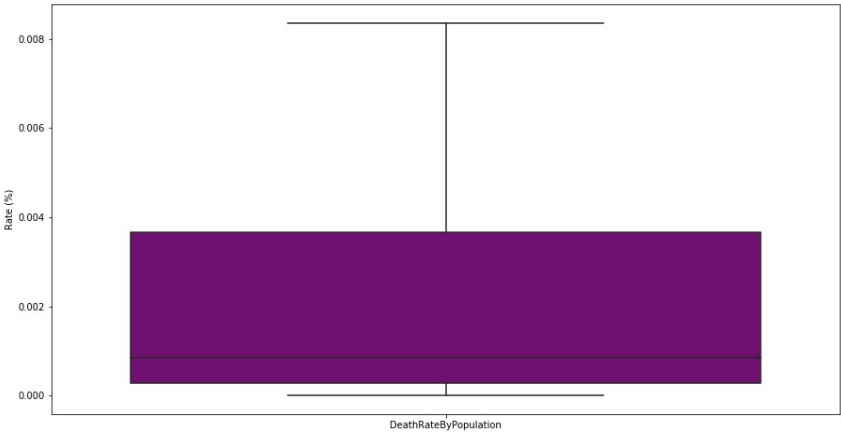
## Human & Capital Mobility

International Inbound  
Tourists(thousands)  
Import/Export (% of GDP)

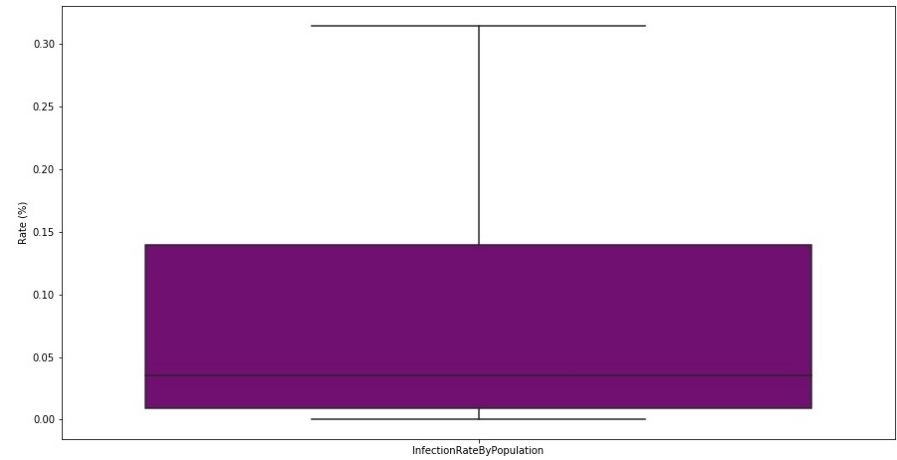
# Step 3: Descriptive Analytics



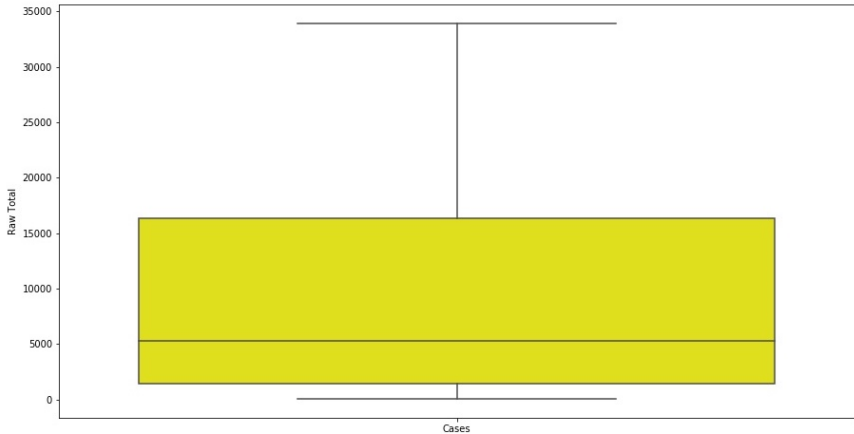
Rate Field	Calculation Method
InfectionRateByPopulation	Total Cases / Total Population
DeathRateByPopulation	Total Deaths / Total Population
DeathRateByPositiveCases	Total Deaths / Total Cases
PositiveCasesPerTestsGiven	Total Cases / Total Tests
TestRatePerPopulation	Total Tests / Total Population



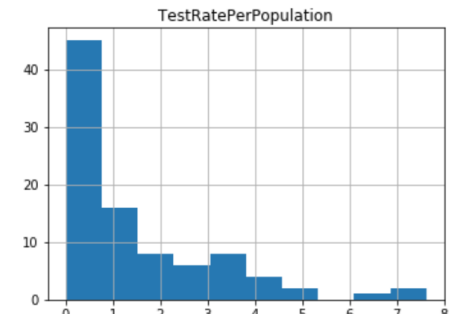
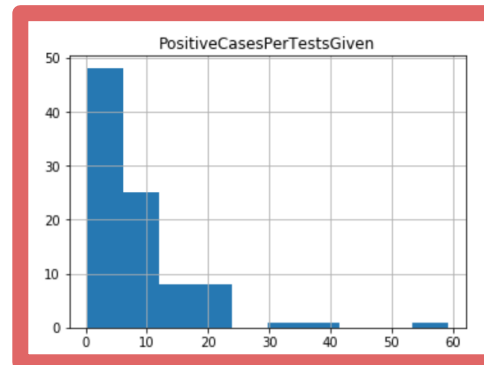
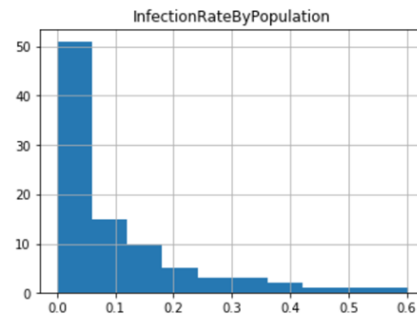
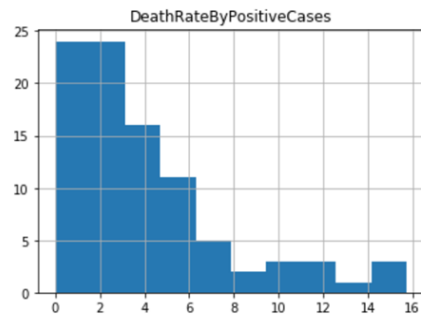
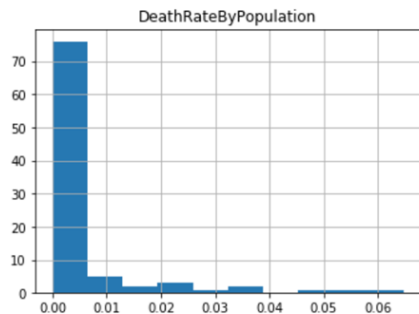
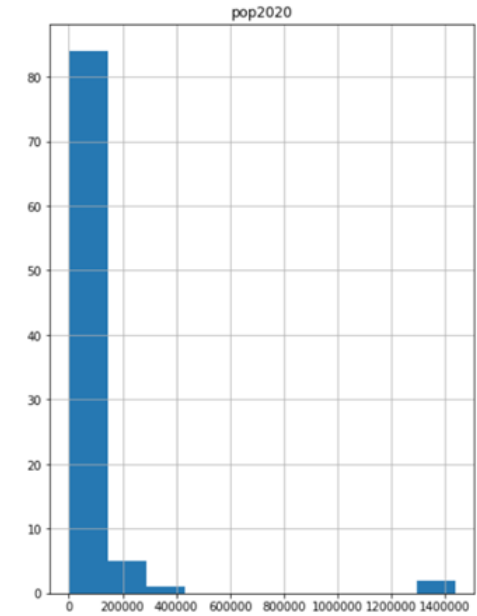
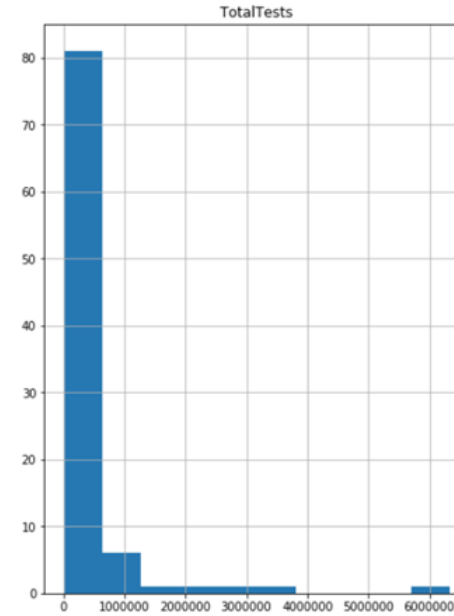
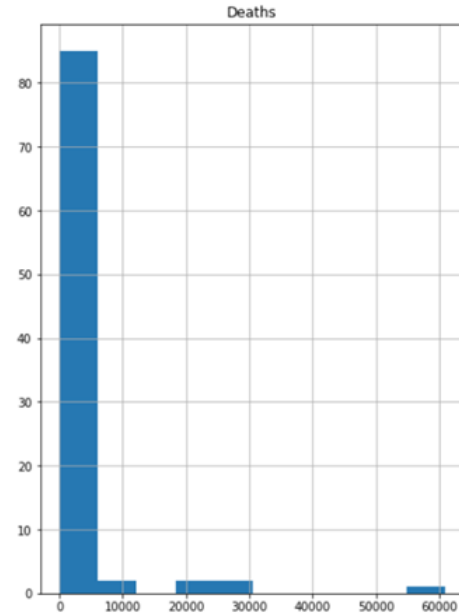
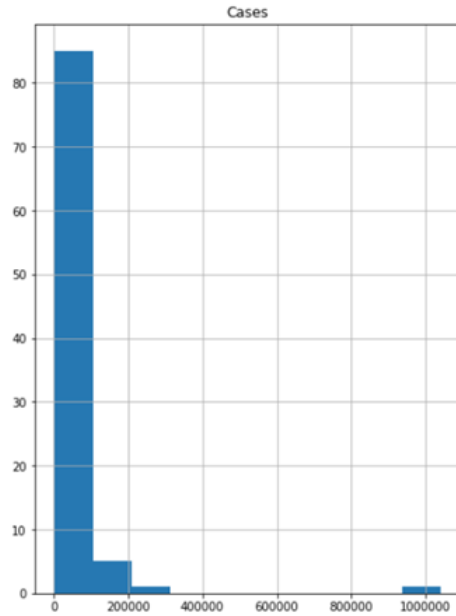
Raw (Total) Cases & Tests



Calculated (Rate) Fields



# Histogram Plots



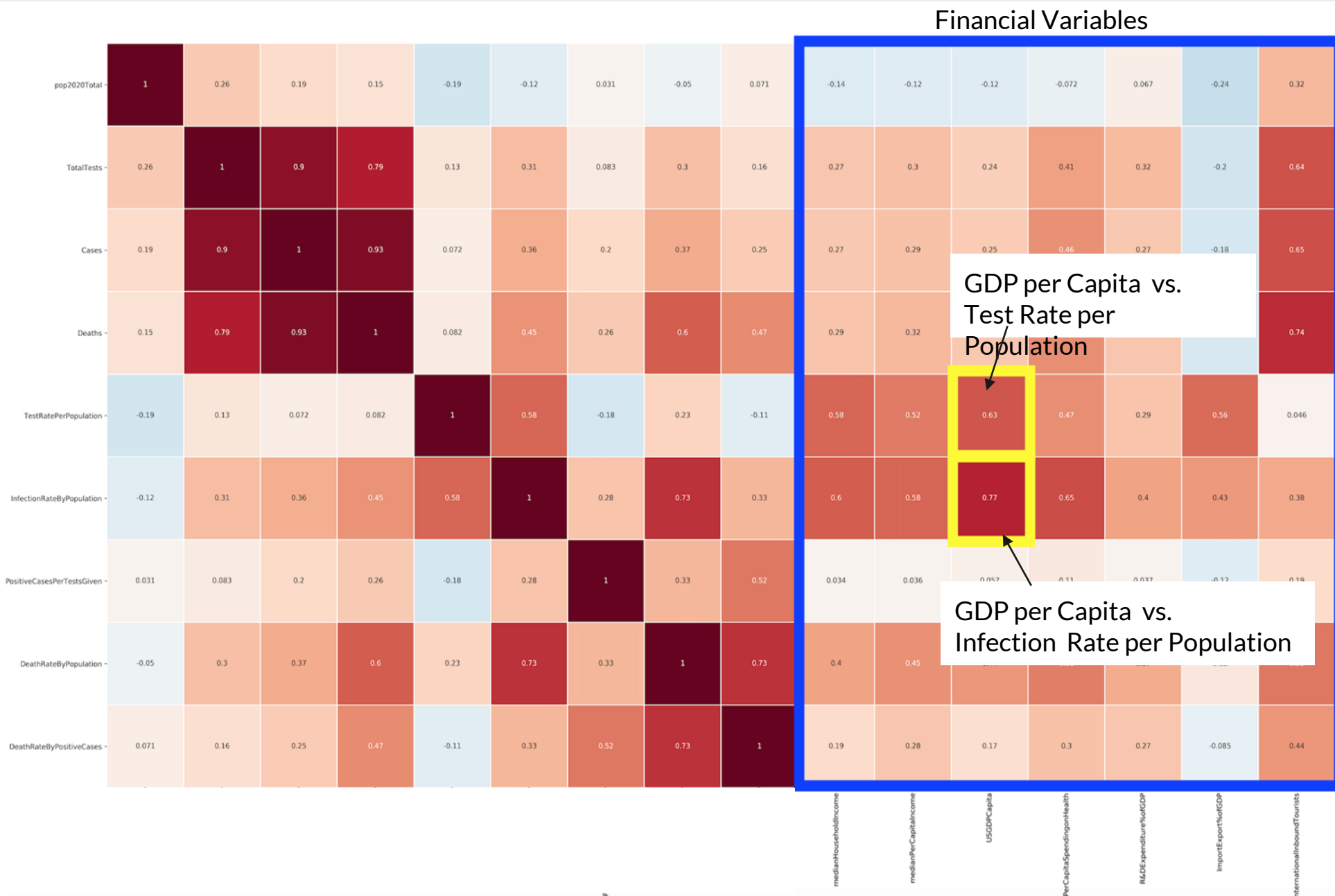


# Correlations

Pt. 1

## Financial Variables

1. Household Income
2. Per Capita Income
3. GDP
4. Healthcare Spending
5. R&D Spending
6. Import/Export %
7. International Students

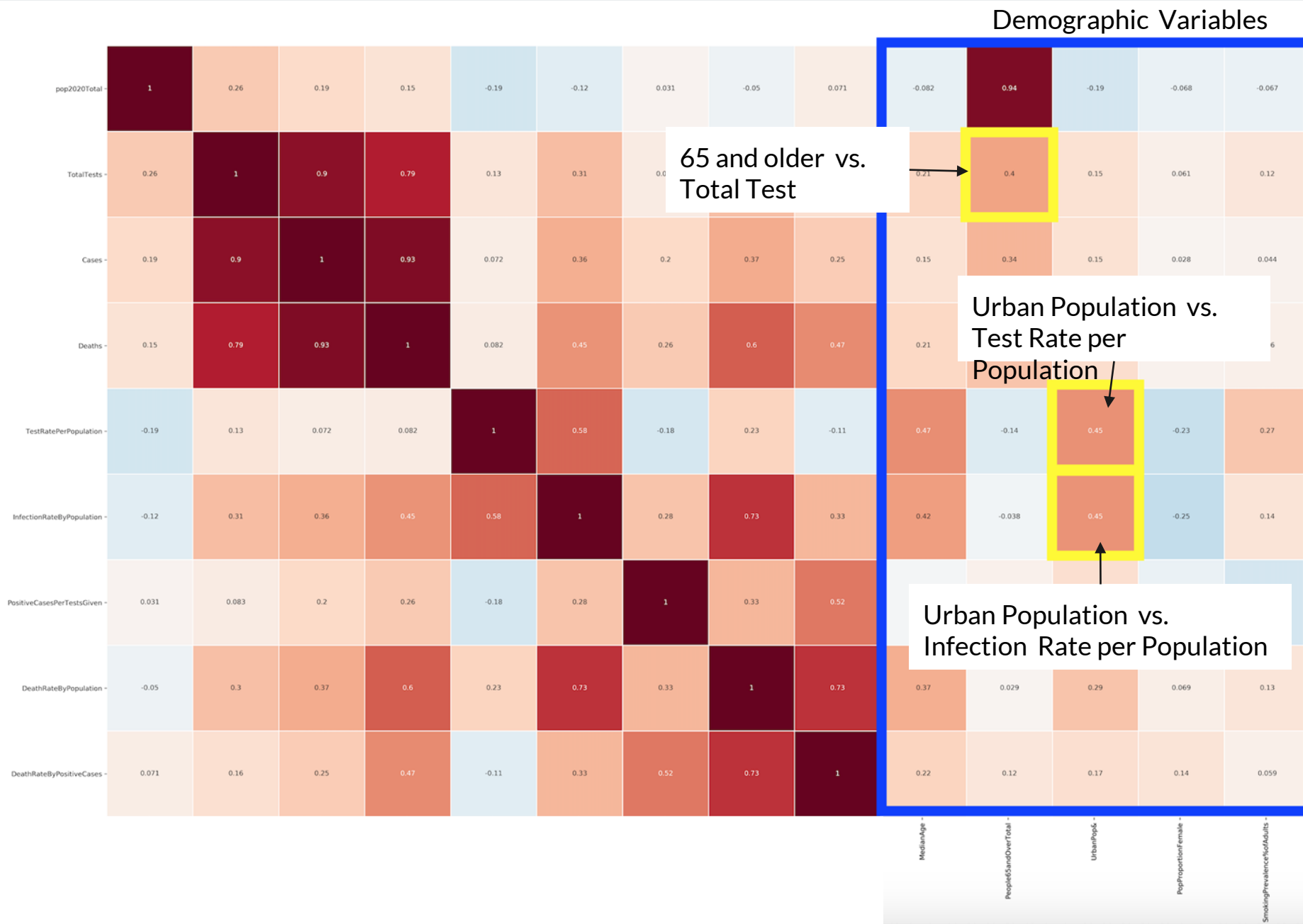


# Correlations

## Pt. 2

### Demographic Variables:

1. Median Age
2. People Over 65
3. Urban Population %
4. Female Population %
5. Smoking Prevalence





## Step 4a: Regression Analysis on Positive Cases

*Why are richer countries seemingly more impacted?*

# Analysis Prep



Country “richness” was benchmarked by GDP per Capita. Also tested household income, income per capita, and healthcare spending per capita.

To correct for size, case numbers were converted from raw count to cases per million residents.

Removed country lines without GDP, population, or testing information.  
Leaves 106 countries for analysis.

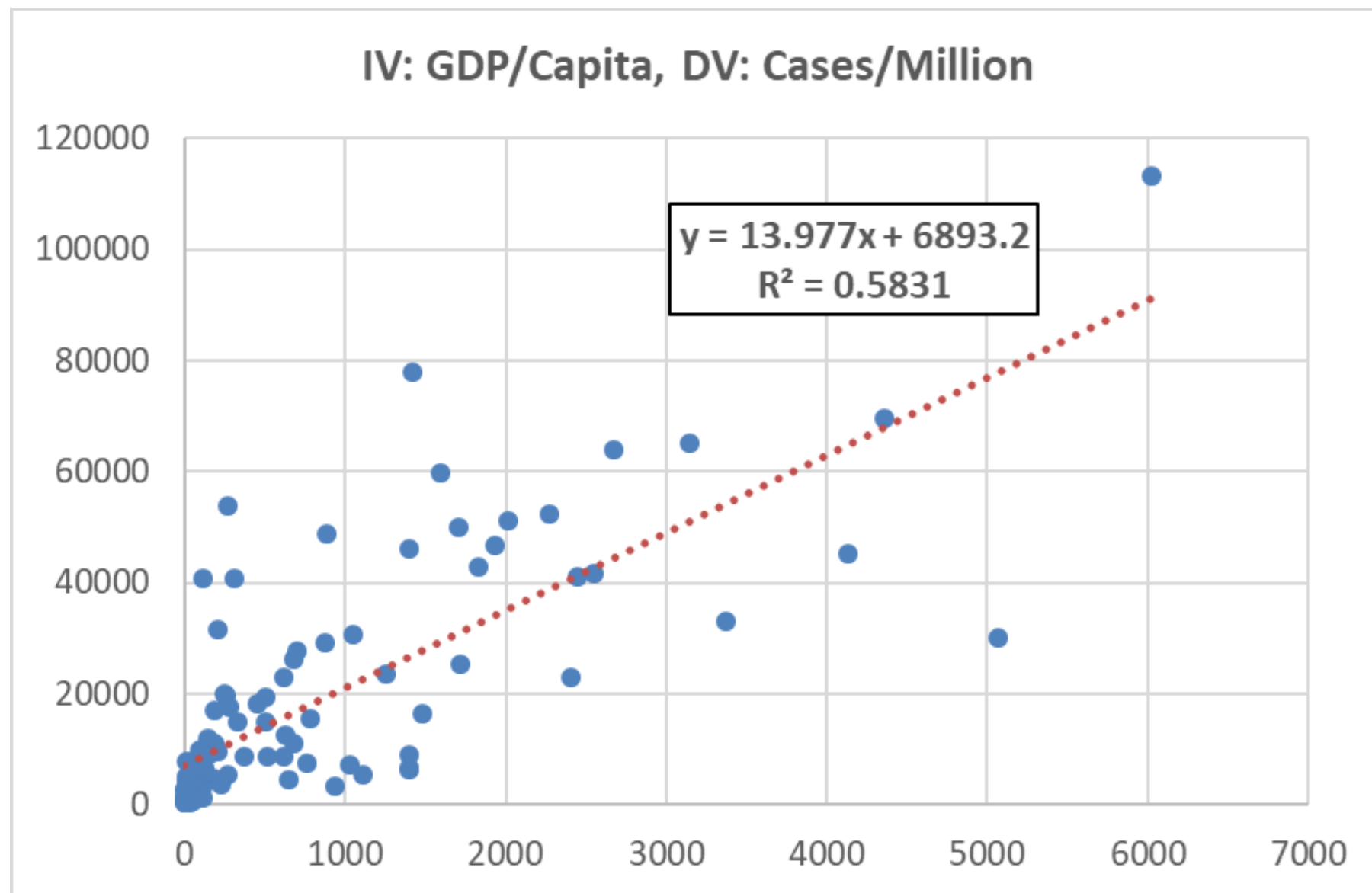
# Why are richer countries more impacted?

- Of the top 10 countries with the most cases/million:
  - 4 are within the top 10% of GDP/capita
  - 8 are within the top 20% of GDP/capita
  - All are within the top 30% of GDP/capita
- Of the 10 countries with the least cases/million:
  - 6 are within the lowest 10% of GDP/capita
  - 8 are within the lowest 20% of GDP/capita
  - All are within the lowest 30% of GDP/capita
- But **WHY?** Are richer countries more infectious?

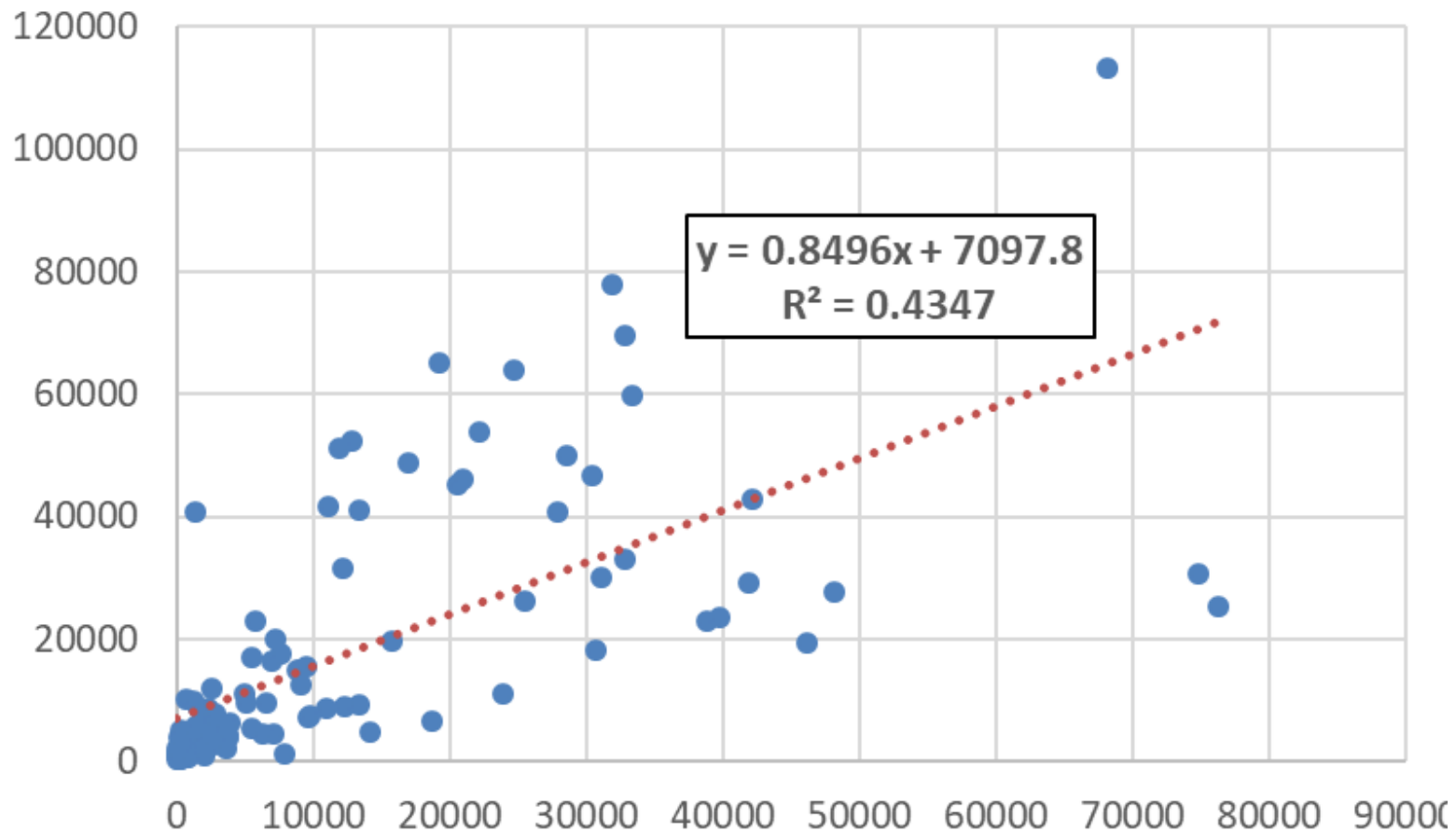
	Country	Cases/M Ranking	GDP/Capita Ranking
↑ Top 10	Luxembourg	1	1
	Spain	2	23
	Qatar	3	3
	Belgium	4	14
	Italy	5	20
	United States	6	4
	Singapore	7	5
	France	8	16
	United Kingdom	9	17
	Portugal	10	29
↓ Bottom 10	Haiti	97	100
	Zambia	98	95
	Madagascar	99	104
	Vietnam	100	83
	Laos	101	84
	Nepal	102	97
	Uganda	103	101
	Mauritania	104	93
	Burundi	105	106
	Yemen	106	98

## Regression Analysis

Are cases/million correlated with GDP/capita?



IV: GDP/Capita, DV: Tests/Million

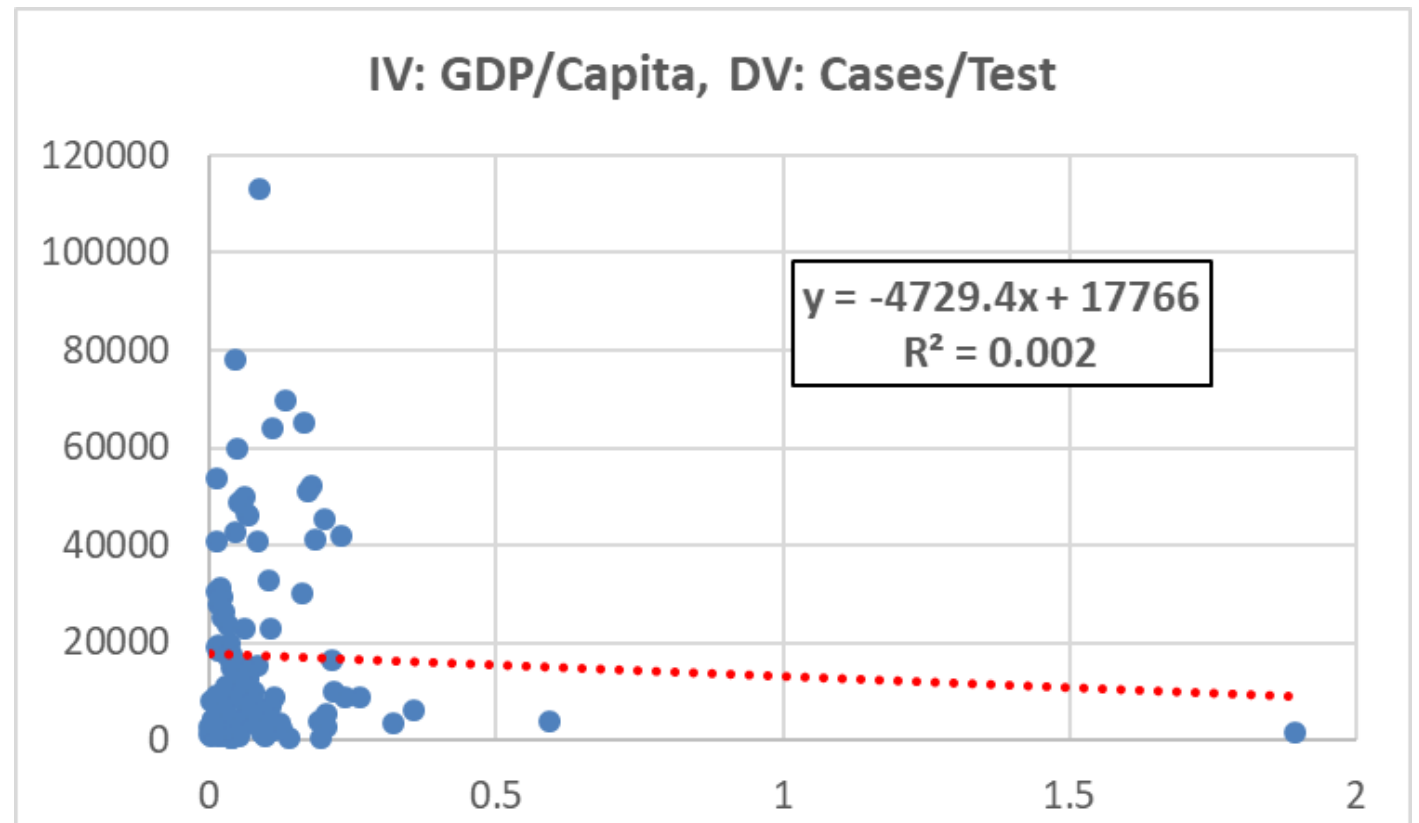


## Is testing the confounding variable?

- Will we get a similar correlation if we replace cases per million with tests per million?
- Not quite as strong, but still some correlation.

# Is Testing the Confounding Variable?

- Reminder: Original regression was IV GDP vs. DV cases/M.
- If we run the regression again, correcting population for number of tests, there is no correlation with GDP.
- Why? It is likely that the number of tests is the confounding variable.
- Therefore, richer countries are NOT more infectious, but the correlation MAY be due to the rate of testing.





# Questions asked...

7.2% Deaths in total positive cases!



Why are nations spending more on health or R&D showing more deaths?



Does Covid hit women and men equally? Why do men appear to be dying at higher rates?

Am I safe because I am young?



More Tourism more deaths??



## Step 4b- Regression Analysis on Deaths



*Can relevant financial, demographic and mobility variables predict the rate of Covid-19 deaths ?*

# Analysis Prep



## Cleaning

Countries with Null Values

China	<ul style="list-style-type: none"><li>• Total Tests=1000,000 (random assumption)</li></ul>
South Korea	<ul style="list-style-type: none"><li>• Female proportion= 50 (equal male : female)</li></ul>
Switzerland	<ul style="list-style-type: none"><li>• Median Per Capita Income(derived from median per capita Household income/ avg household size <i>assuming</i> 3.2)</li></ul>
Czechia	<ul style="list-style-type: none"><li>• Female proportion 50 (equal male : female)</li><li>• smoking prevalence (% of adults) 21.6 (average across the world)</li></ul>
Ireland	<ul style="list-style-type: none"><li>• Median Per Capita Income</li></ul>



## Assumptions

Key Assumptions for 5 significant countries



## Standardize

Scale variables to avoid overlooking of predictors with smaller numbers and ease of comparison



## Dataset

Dataset with 92 countries

Response Variables

1. Deaths/Million
2. Deaths/ +ve Cases

# Approach

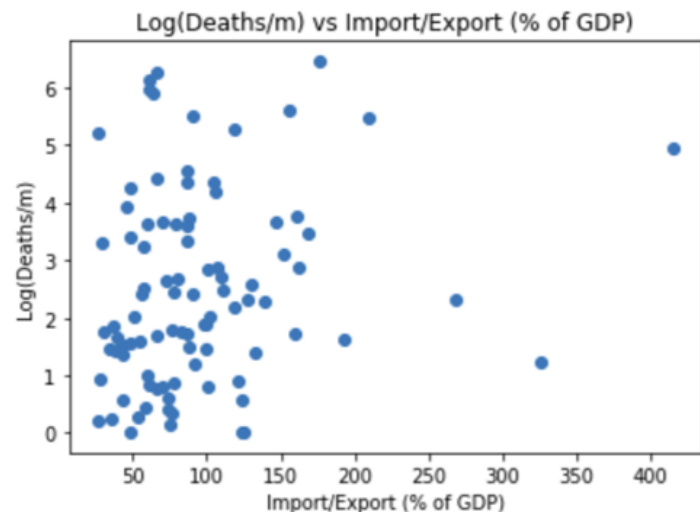
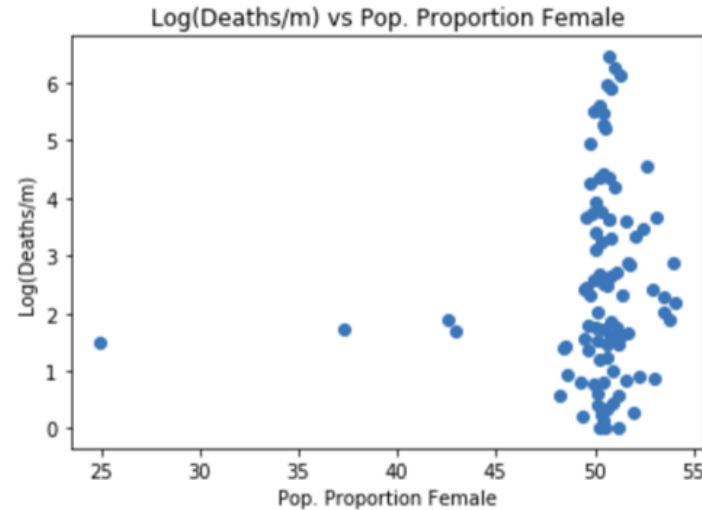
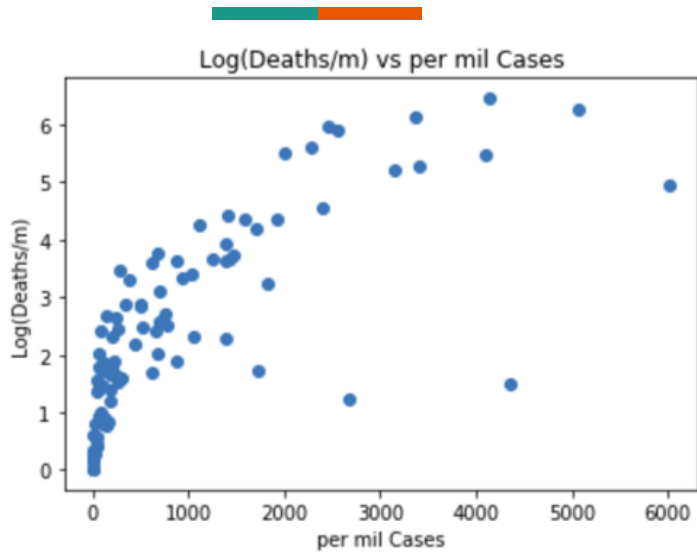
## Regression Analysis



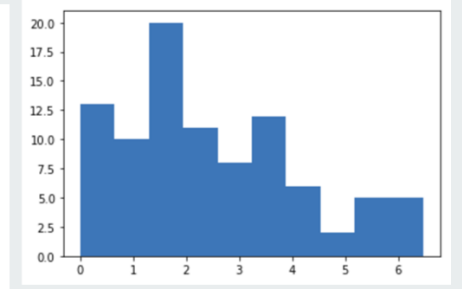
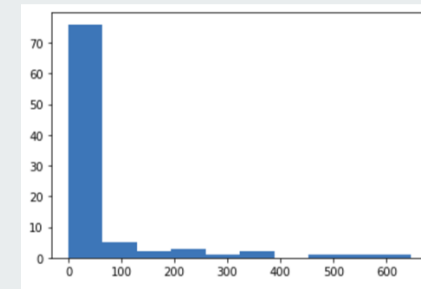
- Supervised Machine Learning for training the models
- Lasso Regression for backward regression to remove insignificant variables
- Multiple Linear Regression with
- Hypothesis Testing for Significance at 10% significance level or P value = 0.1 (increase the range of probability due to dynamics.



# Analysis, Results & Challenges



Deaths ---> deaths per million from raw data



**Uneven distribution of deaths/million**

Transformation  $y = \text{Log}(\text{Death}/\text{Million})$   
more normal or symmetric moving the big countries closer together and space out the smaller ones

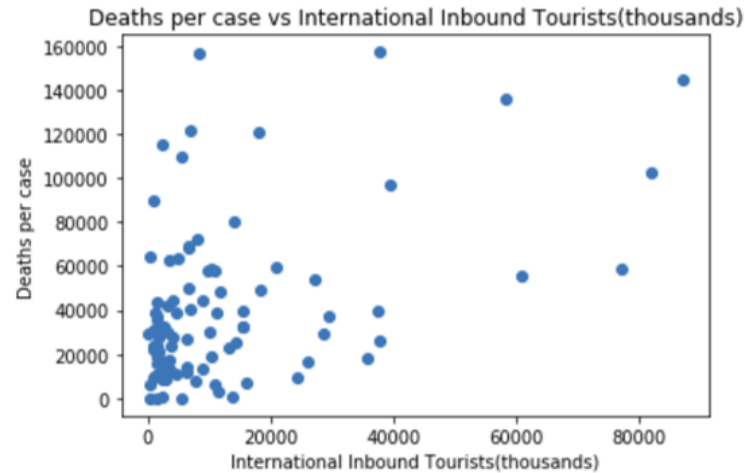
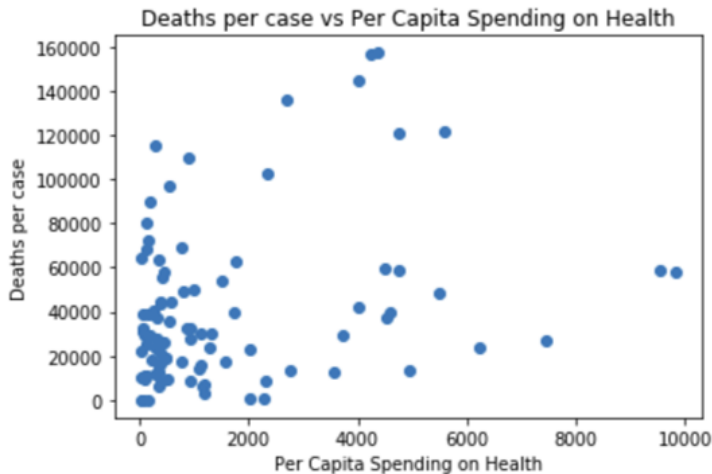
**$R^2 = 0.82$**

probable overfitting scenario

## Significant Predictors

- Positive Cases (+)
- Female proportion from population (+)
- Median Age (+)
- Import Export (% of GDP) (-)

# Analysis, Results & Challenges



**Challenge** - Finding right data due to missing values of countries

**MICE** technique can be used to find missing values



Deaths ---> deaths per Cases from raw data

$$R^2 = 0.36$$

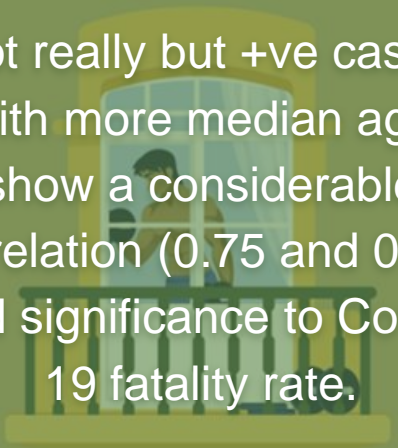
Underfitting scenario

## Significant Predictors

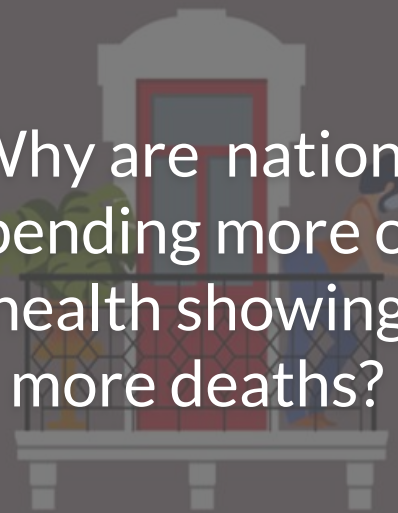
- Per Capita spending on Health (+)
- International Inbound Tourists (+)



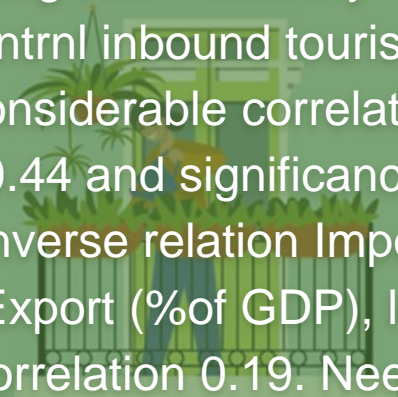
# Questions we asked and what we found



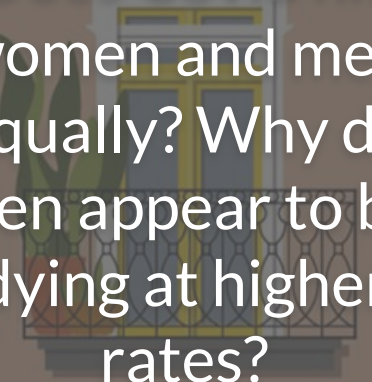
Not really but +ve cases with more median age show a considerable correlation (0.75 and 0.62) and significance to Covid-19 fatality rate.



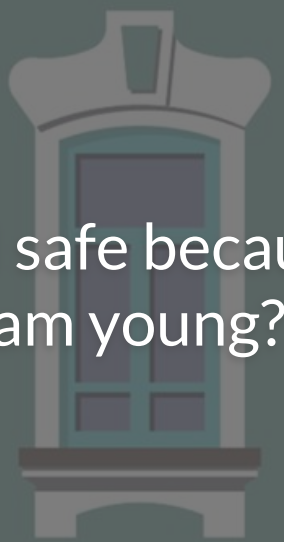
Why are nations spending more on health showing more deaths?




Contagious-> Mobility->Risk  
Intrnl inbound tourists considerable correlation 0.44 and significance. Inverse relation Import /Export (%of GDP), low correlation 0.19. Needs more research.



Does Covid hit women and men equally? Why do men appear to be dying at higher rates?



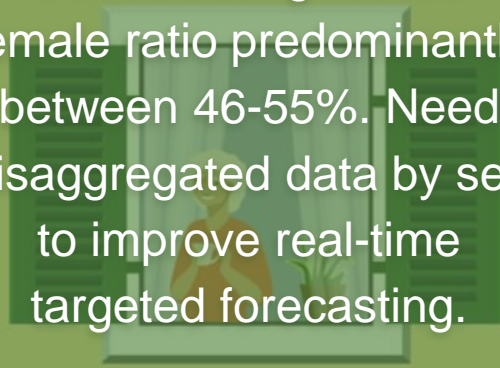
Am I safe because I am young?



Analysis shows more testing more detection implies some percentage of those dying.



More Tourism more deaths??



Contradicting results female ratio predominantly between 46-55%. Need disaggregated data by sex to improve real-time targeted forecasting.

# Recommendations & Conclusion

- Case/ death rates must be taken with a grain of salt.
- Richer countries aren't more infectious, but likely have better means for testing and tracking.
- To better understand case rates, it is imperative that countries:
  - Test more
  - Find common definitions for testing numbers
- To assist poorer countries with testing, the UN should develop a fund devoted to testing.







**Adam Kucharski** ✓  
@AdamJKucharski



Look at the Wuhan line on this new graph from [@jburnmurdoch](#). The lockdown was introduced there on 23rd Jan – 69 days ago – which means this entire Wuhan curve has happened since then. It shows how long it can take to see the effect of control measures on the number of deaths.

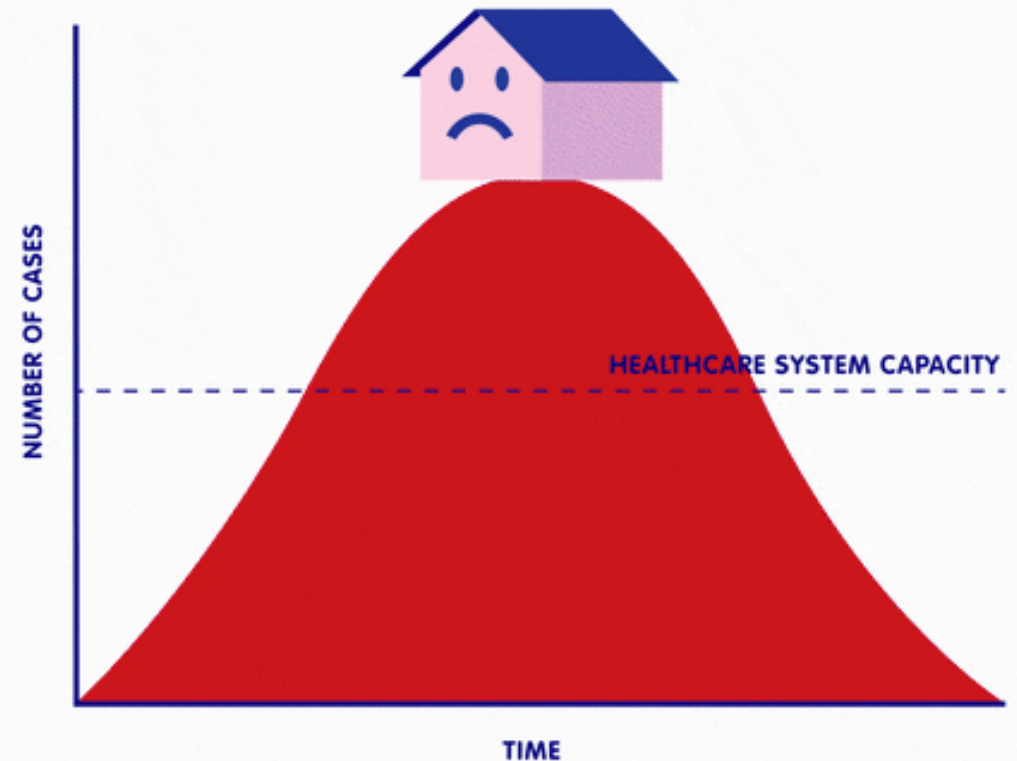


3:10 PM · Apr 1, 2020



3.6K 2.4K people are talking about this

# STAY HOME FLATTEN THE CURVE





# THANK YOU

*In addition, we would like to acknowledge Professor Rex for his help.*