
COronaVirus Disease 2019

A PREDICTION OF THINGS TO COME

CHAN DOMINIC
- DS13 -

01

BACKGROUND

04

EDA

02

PROBLEM STATEMENT

05

PREDICTIONS

03

THE DATA

06

CONCLUSIONS

In January 2020, a novel coronavirus, SARS-CoV-2, was identified as the cause of an outbreak of viral pneumonia in Wuhan, China. The disease, later named coronavirus disease 2019 (COVID-19), subsequently spread globally. In the first three months after COVID-19 emerged nearly 1 million people were infected and 50,000 died.

Similar to the severe acute respiratory syndrome (SARS) and Middle East respiratory syndrome (MERS) of the past.





MERS is a viral respiratory disease that was first reported in Saudi Arabia in September 2012 and has since spread to 27 countries. From its emergence through January 2020, WHO confirmed 2,519 MERS cases and 866 deaths (about 1 in 3).

Infection with SARS coronavirus (SARS-CoV) can cause a severe viral respiratory illness. SARS was first reported in Asia in February 2003, though cases subsequently were tracked to November 2002. SARS quickly spread to 26 countries before being contained after about four months. More than 8,000 people fell ill from SARS and 774 died. Since 2004, there have been no reported SARS cases.



02

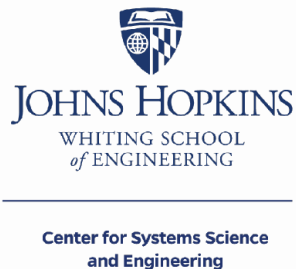
PROBLEM STATEMENT

Determining the significance of variables to accurately predict the global COVID-19 cases; **to effectively predict their eventual progression of deaths and cases by means of the Regression models.** This is to ensure that global citizens are made **aware to the situation's progression** and can make the necessary preparations should there be an extended lockdown period. Governments and private companies can leverage on this model to make **a informed decision as to growth of the economy and businesses.**

covid19 SG

SG RELATED DATA

co.vid19.sg provided detailed information of the individual cases up till 19 April 20



GLOBAL SOURCE

Dataset operated by the Johns Hopkins University Center for Systems Science and Engineering and was updated daily. This was also the source to the Kaggle dataset

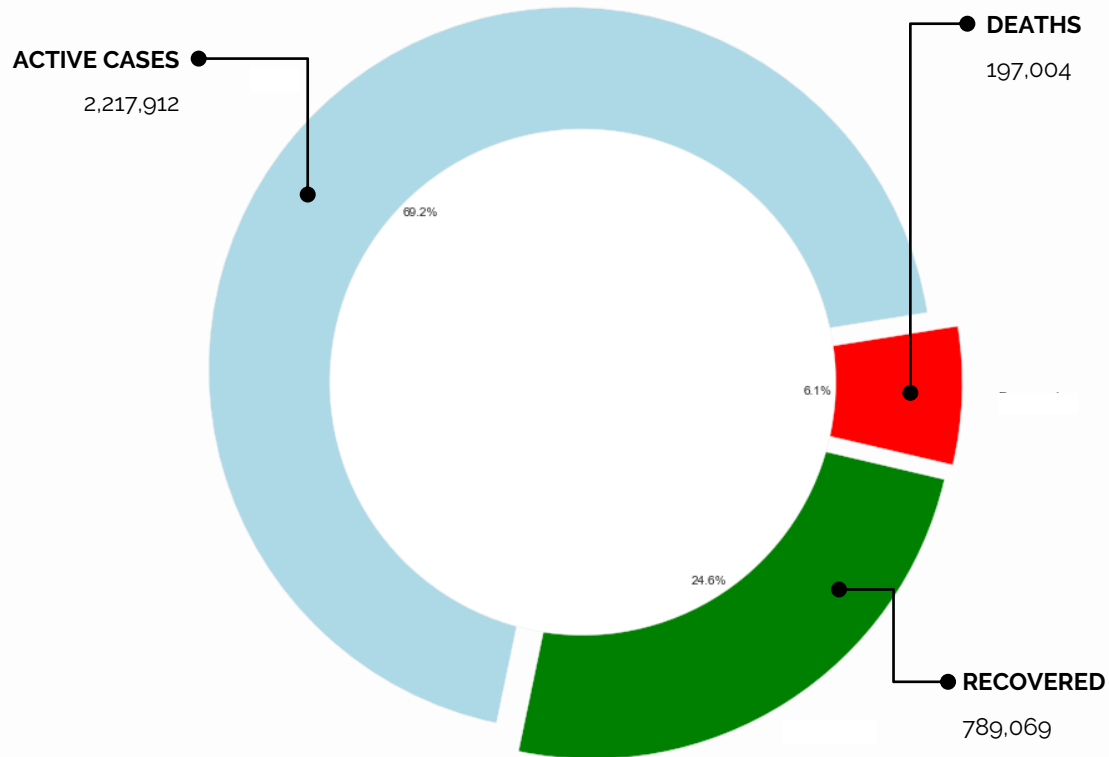
kaggle™

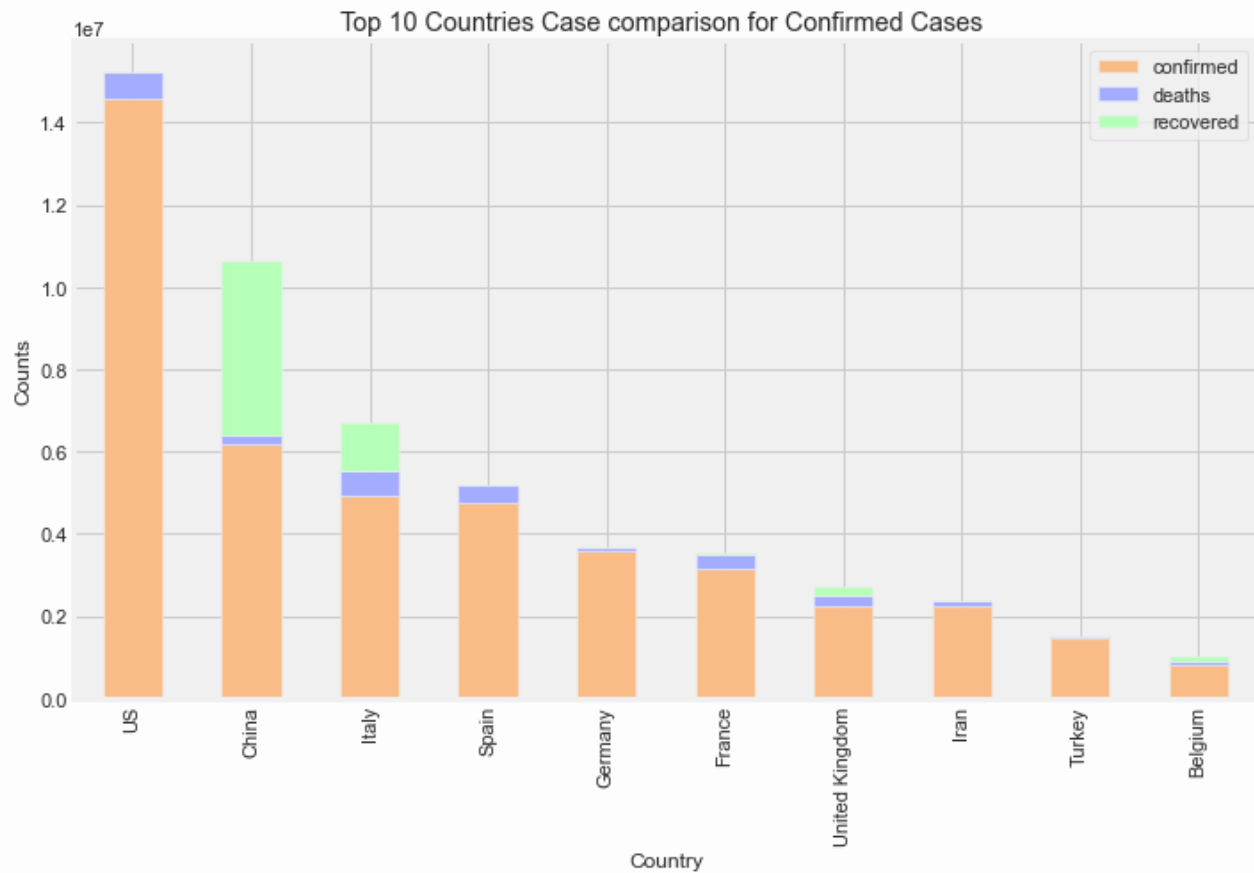
COMPETITION WEEK 4

The challenge involves forecasting confirmed cases and fatalities between April 15 and May 14. Training set contain features from prior.

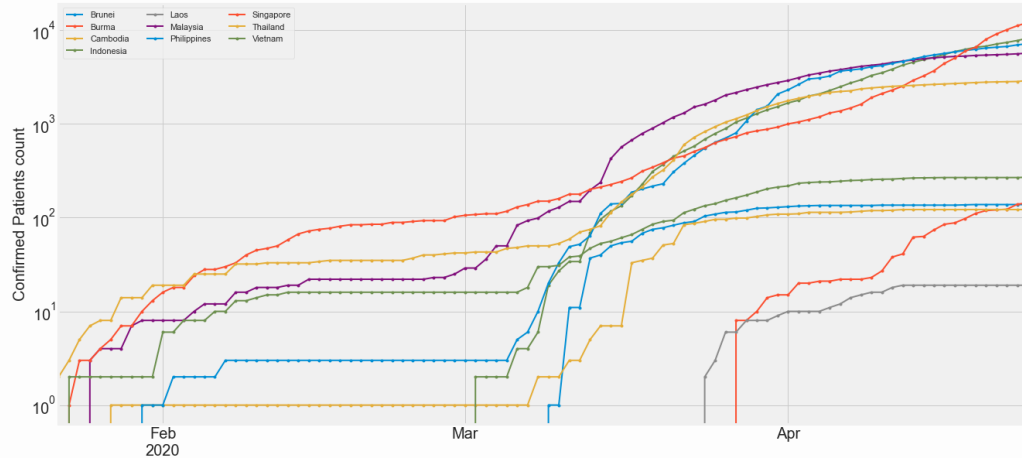


World COVID-19 Cases

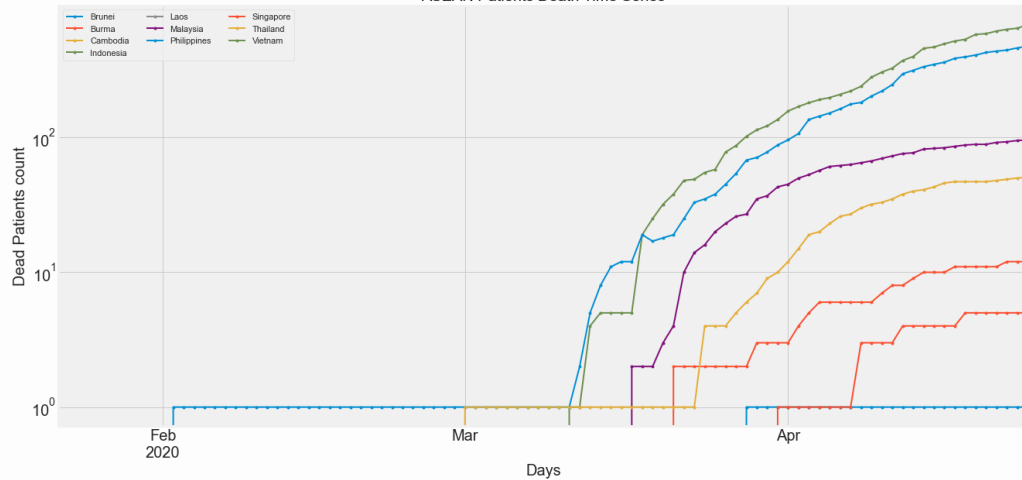




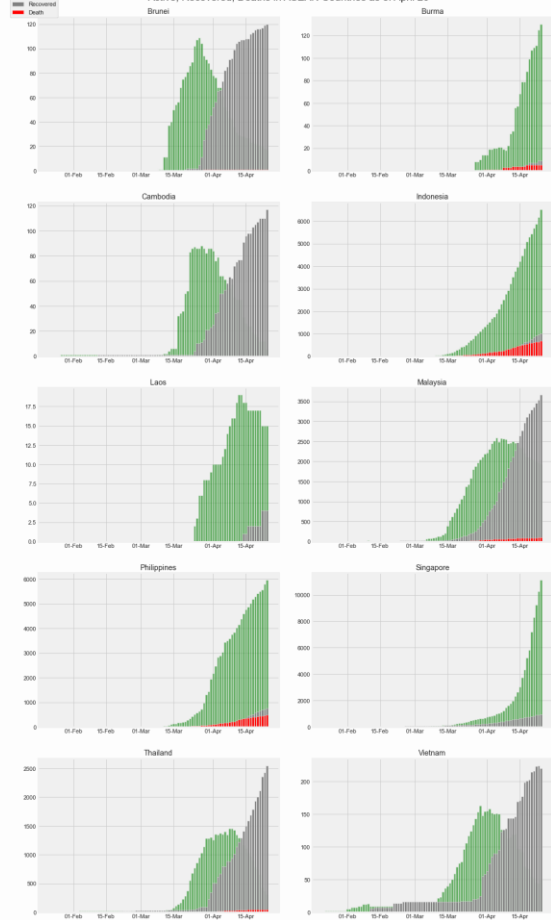
ASEAN Confirmed Patients Time Series

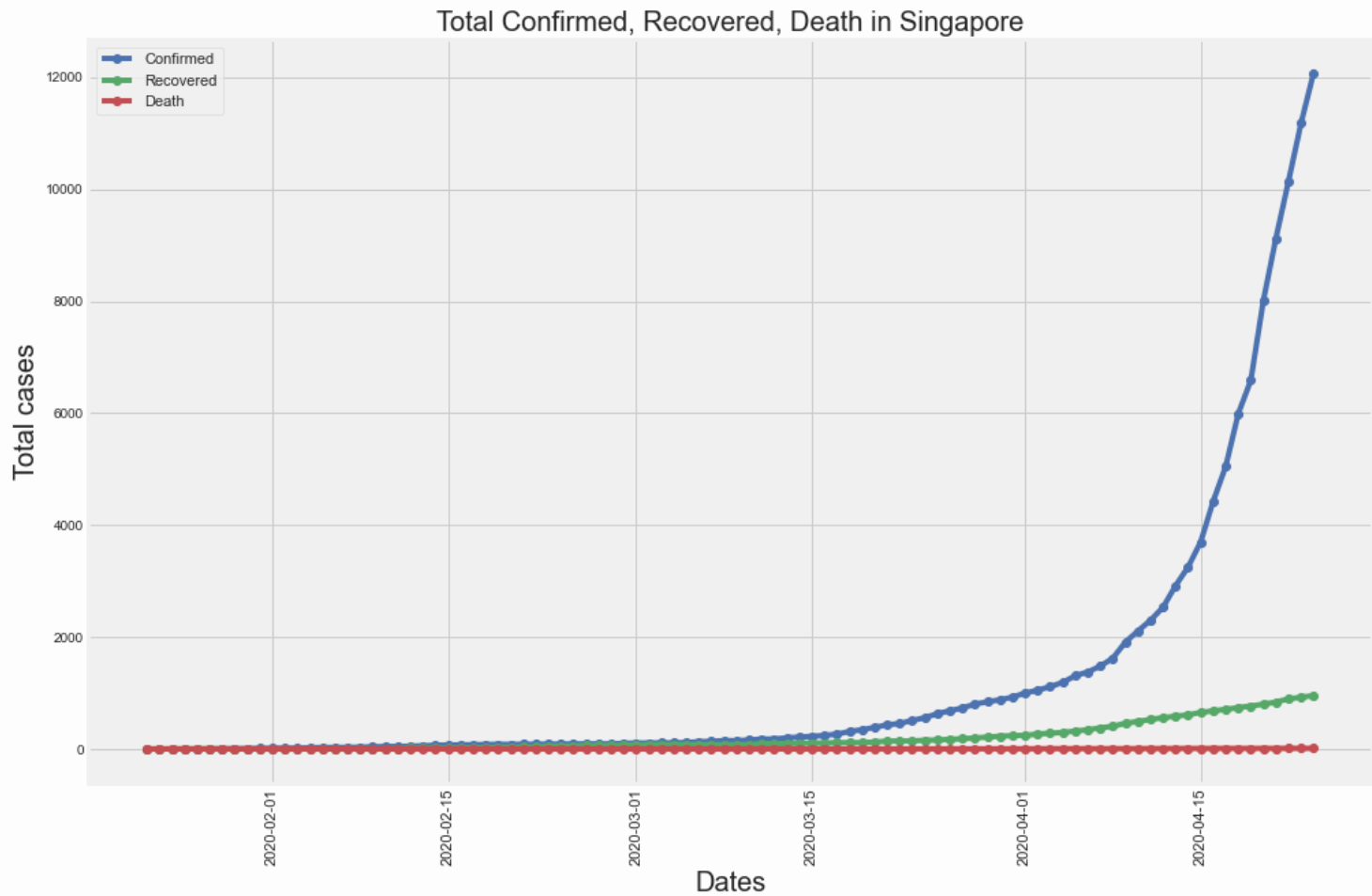


Days
ASEAN Patients Death Time Series

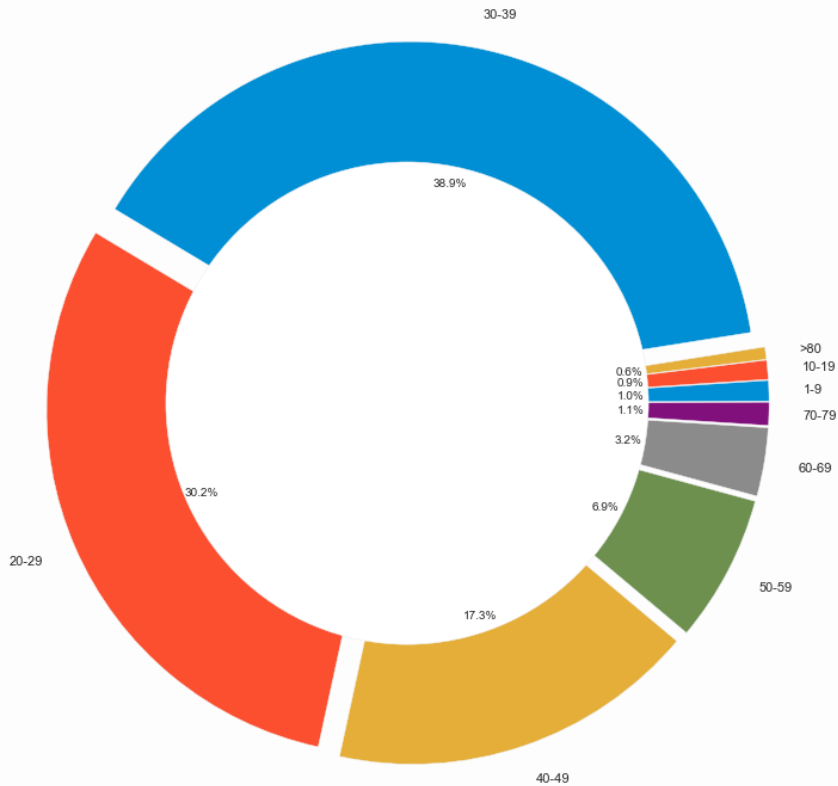


Active, Recovered, Deaths in ASEAN Countries as of April 20

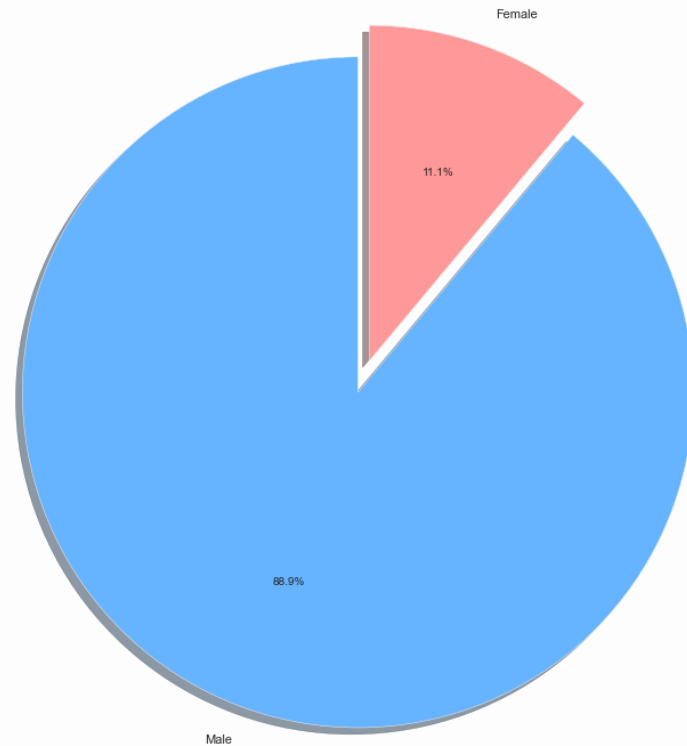


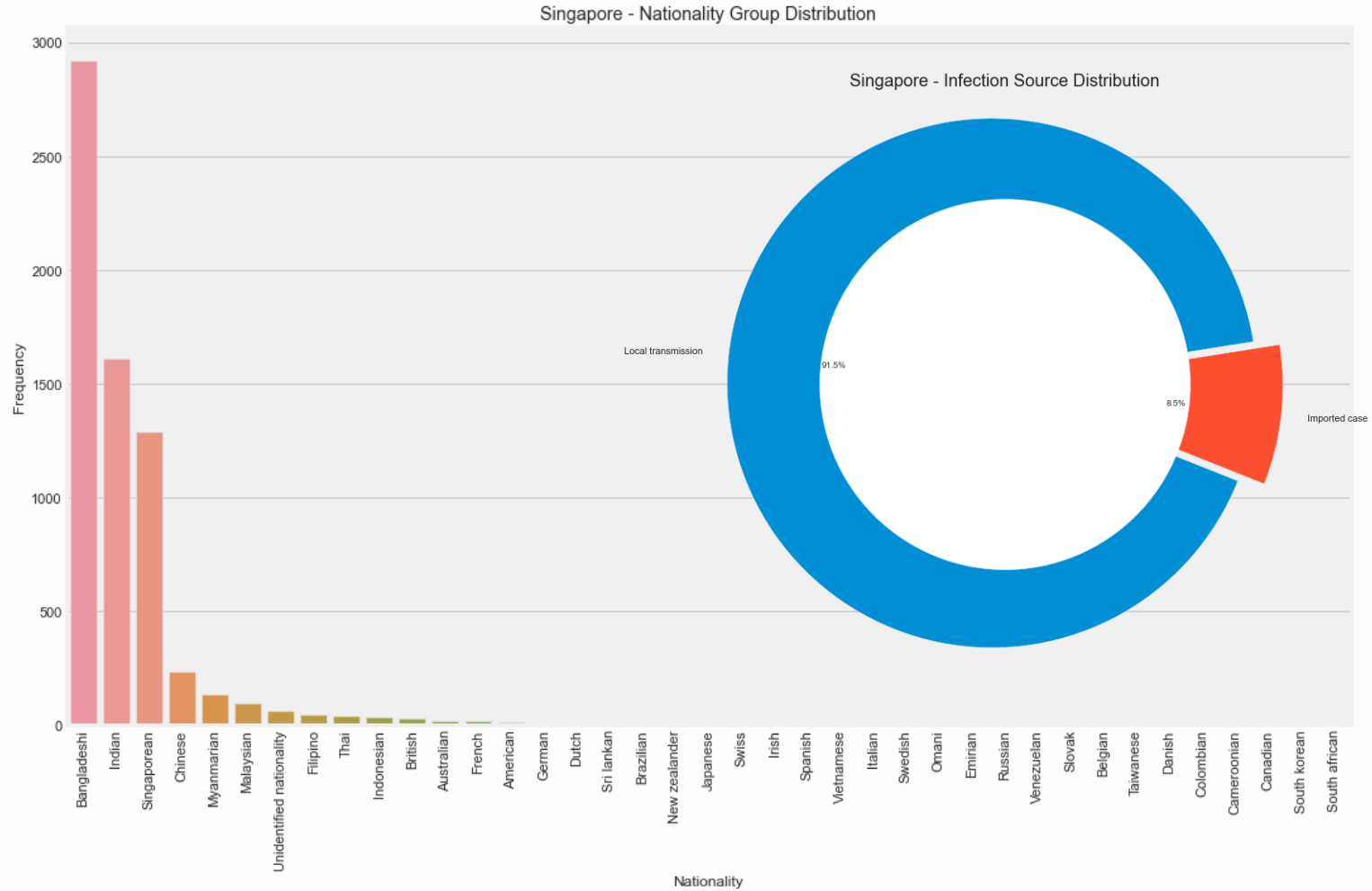


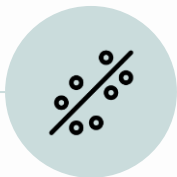
Singapore - Age Group Distribution



COVID Cases Percentage by Gender



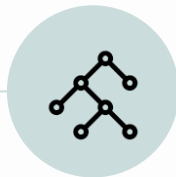




For scraped data, the linear regression model performed the worst on the training set,, though it performed the best on the validation set.

For Kaggle dataset, linear regression was significantly the worst performing model.

LINEAR REGRESSION



For scraped data, the RF model performed average on the training set, though it performed worst on the validation set.

For Kaggle dataset, RF was the best perform model.

RANDOM FOREST



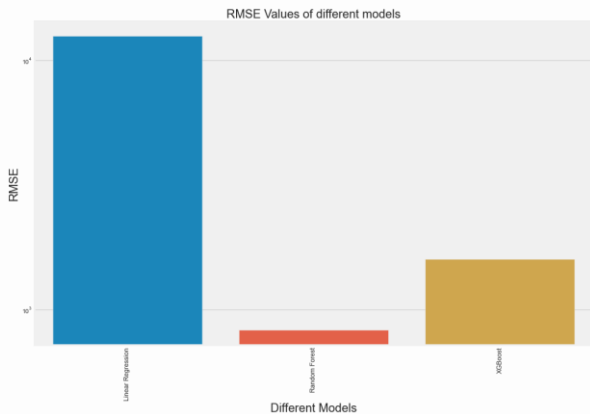
For scraped data, the XGB model performed best on the training set, though it displayed average performance on the validation set.

For Kaggle dataset, XGB was the 2nd best performing model.

XGBOOST

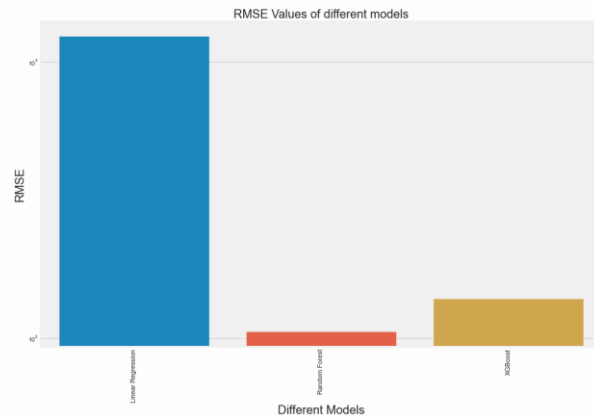


ConfirmedCases



model	mse	rmse
Linear Regression	157470458.21	12548.72
Random Forest	687457.82	829.13
XGBoost	2543786.22	1594.93

Fatalities



model	mse	rmse
Linear Regression	1529577.74	1236.76
Random Forest	11173.06	105.7
XGBoost	19406.37	139.31



Research Code Competition

COVID19 Global Forecasting (Week 4)

Forecast daily COVID-19 spread in regions around world

472 teams · 12 days ago

[Overview](#)[Data](#)[Notebooks](#)[Discussion](#)[Leaderboard](#)[Rules](#)[Team](#)[My Submissions](#)[Late Submission](#)

⚠ This competition has completed. This leaderboard reflects preliminary final standings. The result will become final after the competition organizers verify the results.

Your most recent submission

Name

submission.csv

Submitted

just now

Wait time

0 seconds

Execution time

0 seconds

Score

0.68772

Complete

[Jump to your position on the leaderboard](#) ▾

Recommendation

The model can serve as a **prediction model for countries and businesses**. Additionally, citizens can be more discerning of the potential for number of cases and deaths to rise (or fall) with time.

Notwithstanding this, more data, as mentioned above, will serve to optimise the model. **The model can also serve as a starting point, to be extrapolated to individual countries with additional research** to be done marrying the features together.



WHAT YOU SHOULD DO

- The current datasets have very little features to work with which led to me breaking up the datetime aspects of it to be used as engineered features.
- The data could be supplemented with additional features like weather info, public transport info, global events with participation rates and locations (i.e. music festivals/olympics).

THANKS!

Does anyone have any questions?

domchanjl@gmail.com

+65 8112 7772

yourcompany.com

RESOURCES

<https://github.com/CSSEGISandData/COVID-19>

<https://co.vid19.sg/singapore/dashboard>.

<https://www.niaid.nih.gov/diseases-conditions/covid-19>