

# Using Machine Learning to Analyze COVID-19 Case Numbers

Simon Chow, Aruneshwar Venkat, Gabriel Fonseca, Jackie Vo

April 28, 2022

# 2020 ...

When you are laughing at  
all the corona memes but  
the laughing suddenly  
turns into coughing



The Anonymous Hashtagoholic  
@einfreakinstein

A Cone Like My Dog Has So I Can't Touch My Face.

[#ThingsINeedInMyCOVID911Kit](#)



11:15 AM · Mar 11, 2020 · [Twitter Web App](#)

# Introduction

- Using COVID-19 epidemiological database to analyze disease severity in different states at different times
- Predict severity of COVID-19 at a certain time at certain location.

[HTML] **Using machine learning to predict COVID-19 infection and severity risk among 4510 aged adults: a UK Biobank cohort study**

[AA Willette](#), [SA Willette](#), [Q Wang](#), [C Pappas...](#) - Scientific Reports, 2022 - nature.com

... **Older adults** with prior **infection** show exhaustion of the naïve ... to **predict COVID-19 severity** among all test cases in this **study**, as ... For **COVID-19 severity**, antibody titers to the HTLV1 virus ...

☆ Save  Cite Cited by 6 Related articles

**Forecast and prediction of COVID-19 using machine learning**

[D Painuli](#), [D Mishra](#), [S Bhardwaj](#), [M Aggarwal](#) - Data Science for COVID-19, 2021 - Elsevier

COVID-19 outbreaks only affect the lives of people, they result in a negative impact on the economy of the country. On Jan. 30, 2020, it was declared as a health emergency for the entire globe by the World Health Organization (WHO). By Apr. 28, 2020, more than 3 million people were infected by this virus and there was no vaccine to prevent. The WHO released certain guidelines for safety, but they were only precautionary measures. The use of information technology with a focus on fields such as data Science and machine learning ...

☆ Save  Cite Cited by 37 Related articles All 4 versions

# Data Source:

- Google COVID-19 Open Data Repository
- Provides general demographic information, other COVID-19 related information such as vaccination status, death etc.
- Focus only on United States data for analysis

## COVID-19 Open Data Repository



The Google Health COVID-19 Open Data Repository is one of the most comprehensive collections of up-to-date COVID-19-related information. Comprising data from more than 20,000 locations worldwide, it contains a rich variety of data types to help public health professionals, researchers, policymakers and others in understanding and managing the virus.

# Methods

1. K-means clustering
2. Time series Poisson regression
3. Time series with Random Forest
4. Time series with SVM

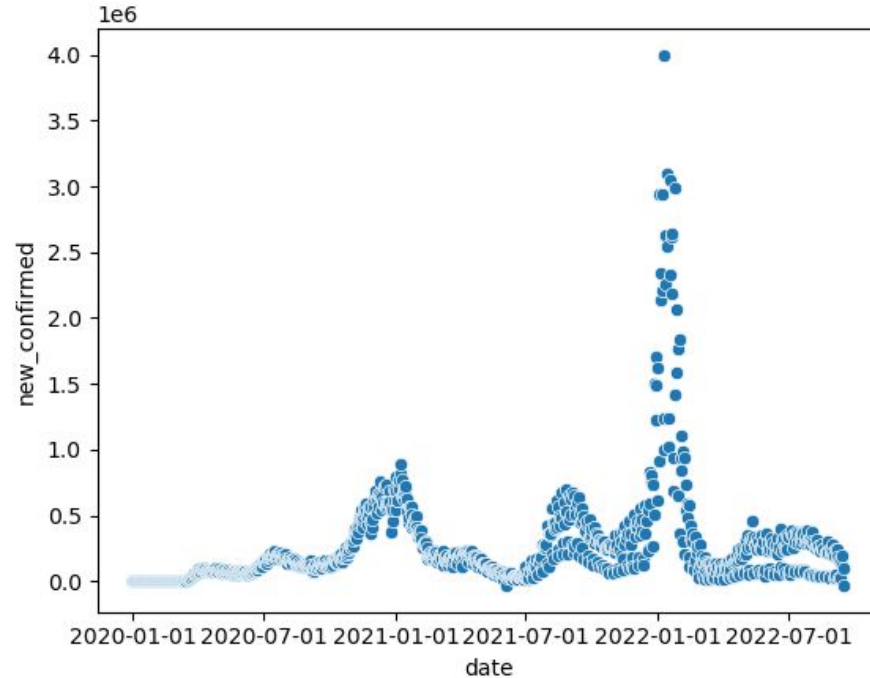
# **Method 1: K Means Clustering**

# Overview - K-means clustering

- Focused on a few features to cluster for predicting new covid cases over time
  - New deceased, new people vaccinated, new people fully vaccinated, new hospitalized patients, new intensive care patients
- $N = 640$  time points
- We clustered the actual new covid cases to generate true labels

# Results - K-means clustering

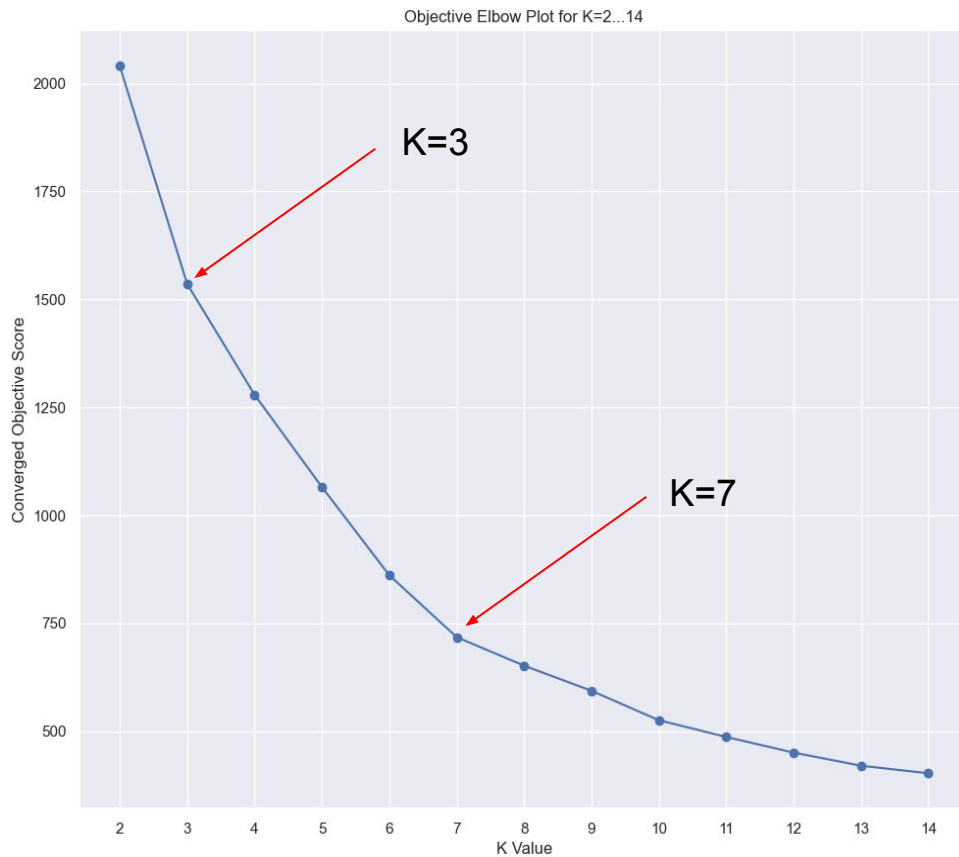
New cases over times plot  
(no k means yet)





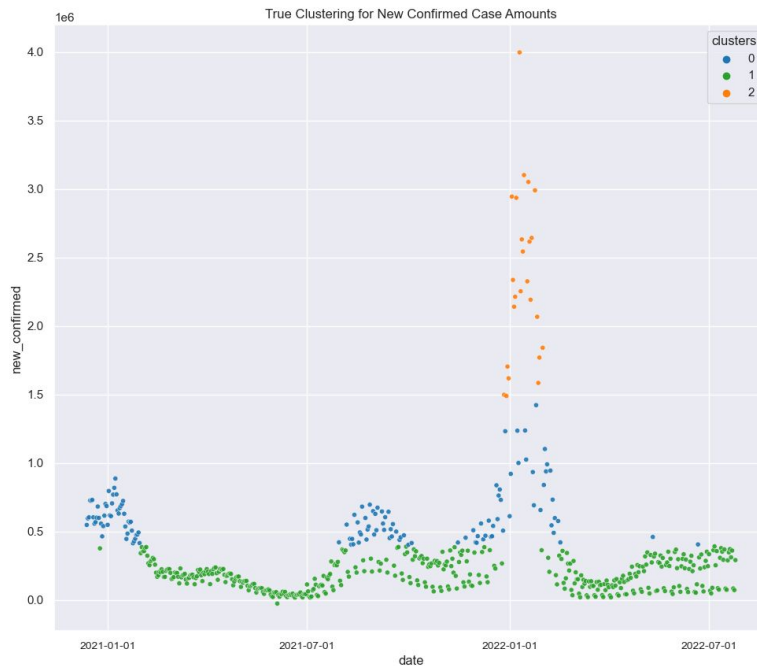
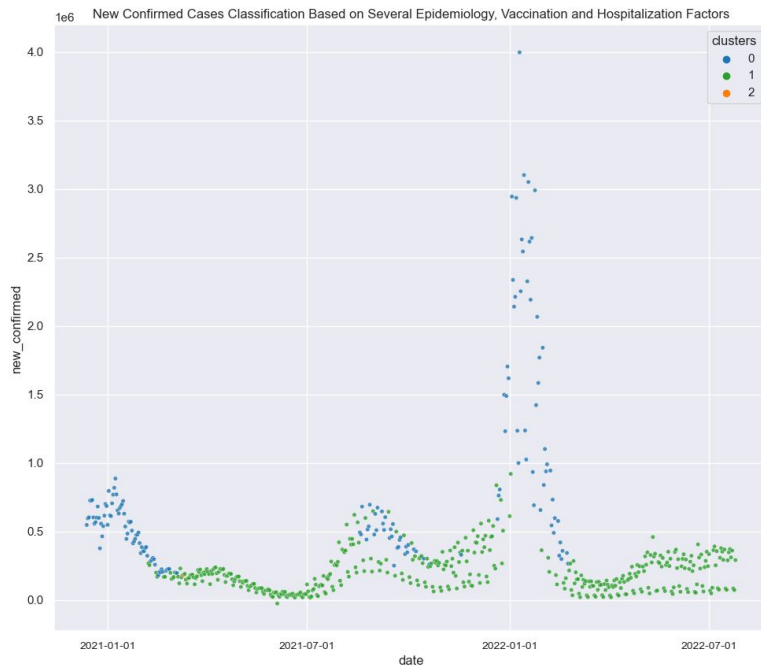
# Results - K-means clustering

**Objectives  
Elbow Plot**



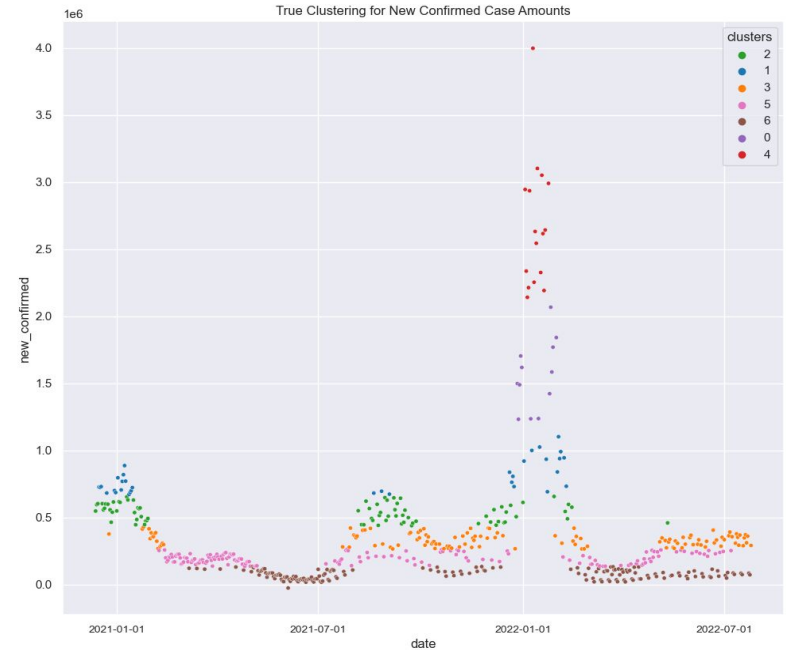
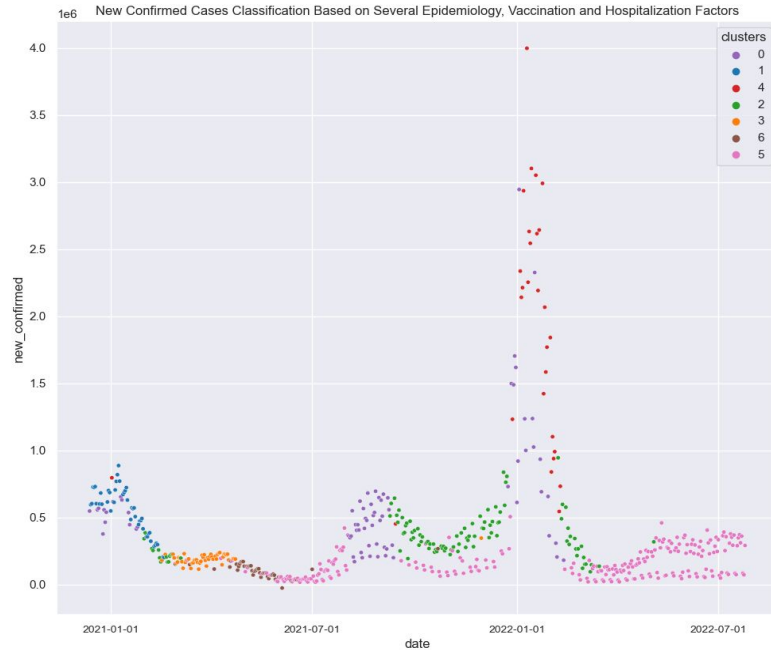
# Results - K-means clustering

With  $k = 3$ , accuracy =  $\sim 83.1\%$



# Results - K-means clustering

With  $k = 7$ , accuracy =  $\sim 24.9\%$  (overfitting)

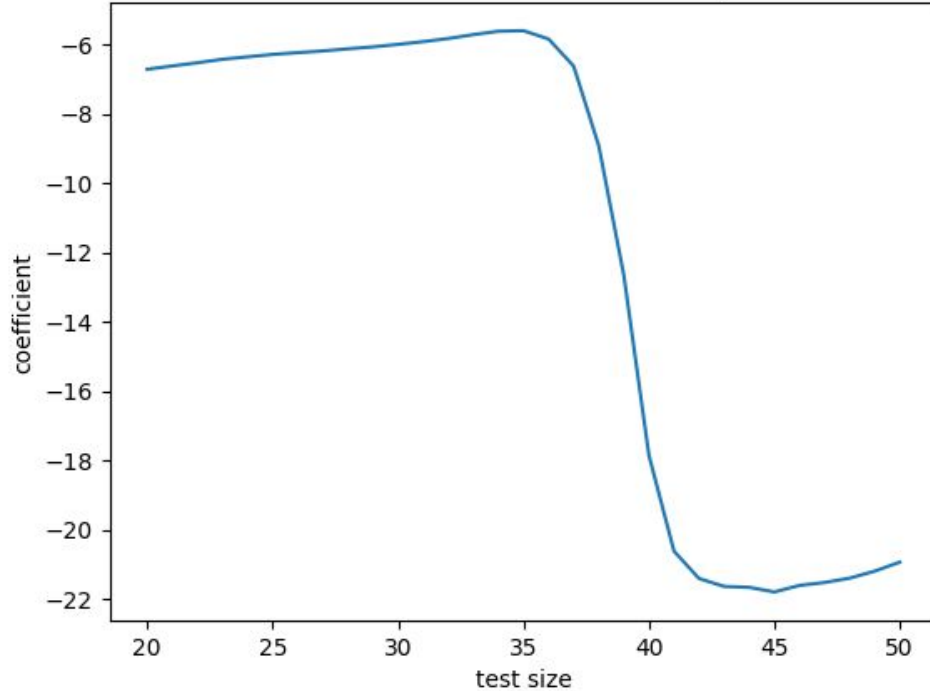




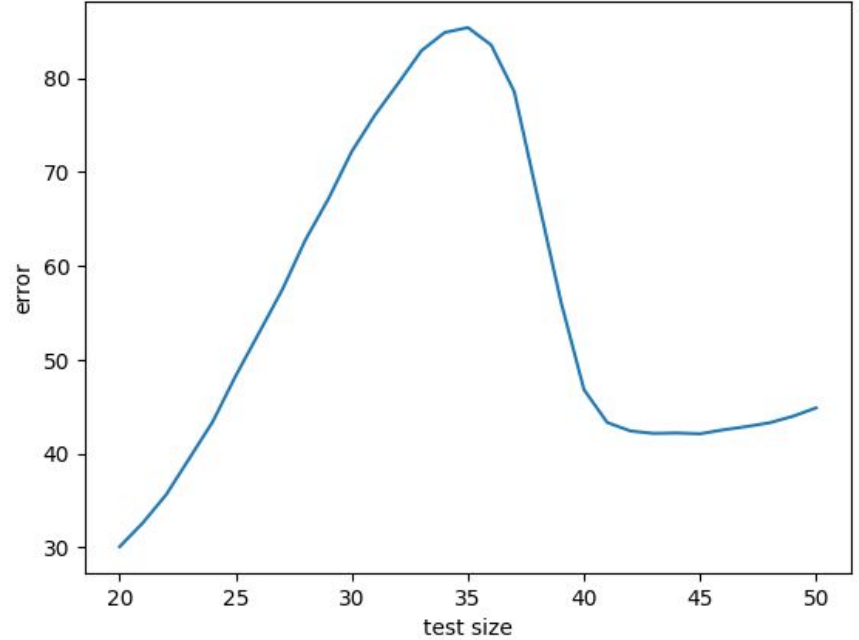
## Method 2: Poisson Regression

# Results - Poisson Regression

Poisson Coefficients for Different Train/Test Sizes



Poisson Error for Different Train/Test Sizes

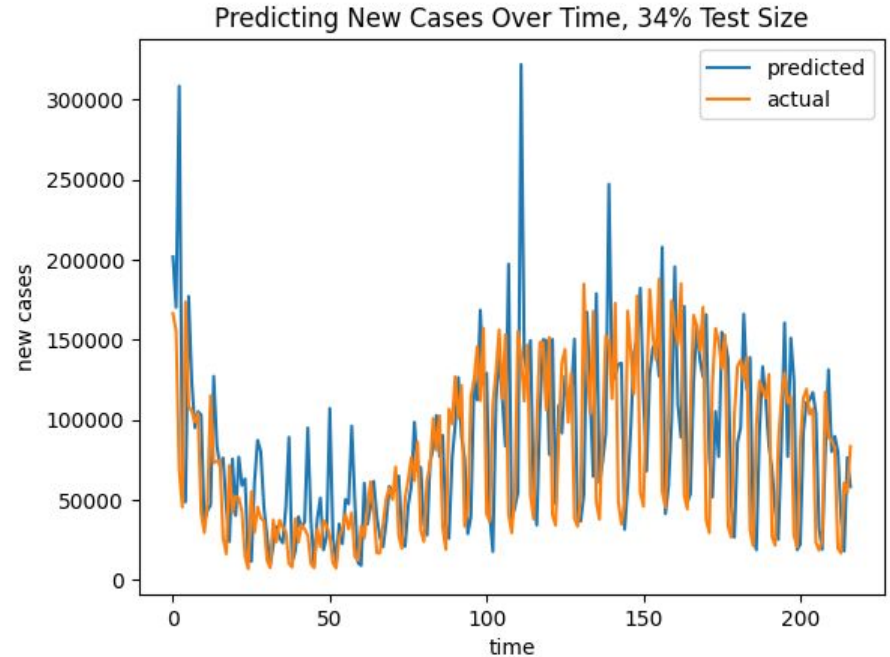


# **Methods 3 and 4: Time Series with Random Forest and SVM**

# Overview - Time Series

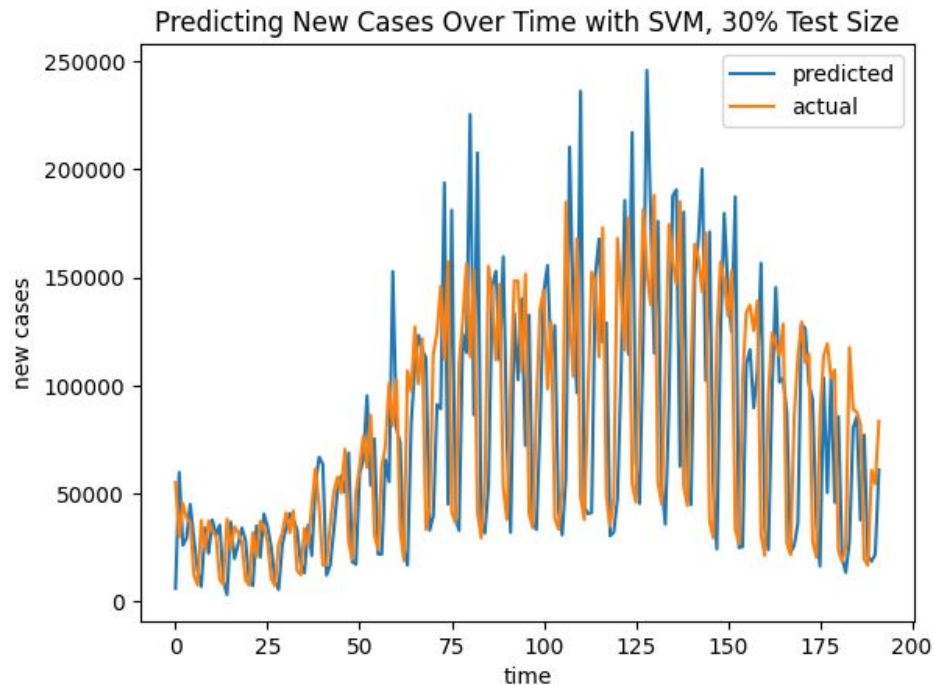
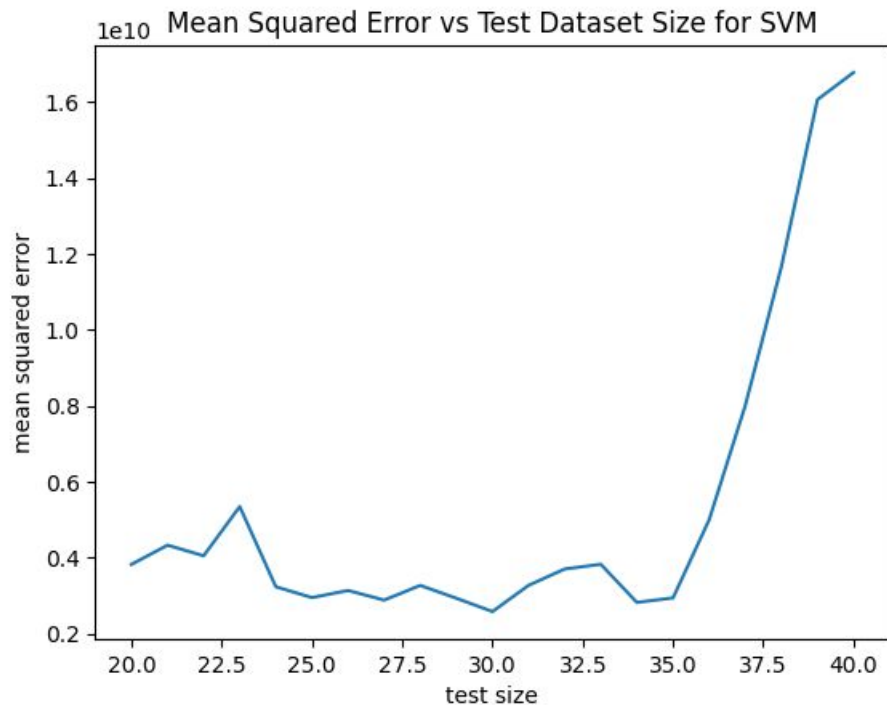
- $N = 640$  time points
- Converted data to sliding window
  - First time point is  $x$ , second time point  $y$
  - Second time point is  $x$ , third is  $y$ , etc.
- Used walk-forward validation for accuracy

# Results - Time Series w. Random Forest





# Results - SVM



# Limitations

- Time series: not enough data points
- No consistent trend makes it difficult to predict
- Possible confounding variables affecting case rates
  - Hard to fit linear models

# Future Directions

- Extend the analysis to worldwide
- Adjust parameters for SVM, Random Forest
- Incorporate more features into regression



# References

Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time [published correction appears in Lancet Infect Dis. 2020 Sep;20(9):e215]. *Lancet Infect Dis*. 2020;20(5):533-534. doi:10.1016/S1473-3099(20)30120-1

O. Wahltinez, et al. "COVID-19 Open-Data: curating a fine-grained, global-scale data repository for SARS-CoV-2". In: (2020). URL: <https://goo.gle/covid-19-open-data>.

Painuli D, Mishra D, Bhardwaj S, Aggarwal M. Forecast and prediction of COVID-19 using machine learning. *Data Science for COVID-19*. 2021;381-397. doi:10.1016/B978-0-12-824536-1.00027-7

Willette AA, Willette SA, Wang Q, Pappas C, Klinedinst BS, Le S, Larsen B, Pollpeter A, Li T, Mochel JP, Allenspach K, Brenner N, Waterboer T. Using machine learning to predict COVID-19 infection and severity risk among 4510 aged adults: a UK Biobank cohort study. *Sci Rep*. 2022 May 11;12(1):7736. doi: 10.1038/s41598-022-07307-z. PMID: 35545624; PMCID: PMC9092926.