

# COVID-19 Data Analysis

Trevor Carpenter, Christina De Cesaris, Michelle Tran  
STA 135: Multivariate Data Analysis  
Professor Maxime Pouokam  
University of California, Davis  
June 7, 2021

# I. Introduction

The COVID-19 pandemic is an event that will have long-lasting effects in our society. While no country was entirely prepared for the chaos that ensued, there were definitely some countries that had a better handle on the situation than others. Countries who were quick on putting out preventative measures and were able to produce masks quickly were able to keep their population's total cases to a minimum, whereas countries who lagged behind on safety measures quickly saw the exponent rise in cases and the overcapacity of hospitals.

As we are nearing the end of pandemic, much data has been collected and analyzed. In our paper we have decided to look at data involving the mortality recovery ratio in regards to the beginning of the pandemic. Mortality recovery ratio refers the ratio of total deaths to total recovered patients. Our question of interest is: what relationship did the variables have with the mortality recovery ratio for the COVID-19 dataset? To answer our question, we decided to use linear regression and use the anova function to conduct a hypothesis testing to see if the selected variables are significant.

## II. Data, Models, and Methods

The COVID-19 data used here is publicly and available from Worldometer website <https://www.worldometers.info/coronavirus/> for March 30, April 15, and April 25, 2020. Data were captured on the next day to these specified dates. Countries with COVID-19 total cases less than 500 or countries with missing data were omitted from the analysis to keep good representability of each variable. Number of countries included in the analysis was 56 countries on March 30, 82 countries on April 15, and 91 countries on April 25.

The variables included; in any given country, total cases refers to total cases confirmed with COVID-19; active cases refers to total number of open cases (mild, serious, or critical); total deaths refers to total deaths with COVID-19; critically ill cases refers to number of serious/critically ill cases; mortality recovery ratio refers to the ratio between total deaths to total recovered patients.

[justification]

## III. Results of our Linear Model Testing

Below shows a table of the full model and its beta values.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.3038600	0.1280644	2.3727127	0.0390958
TotalCases	0.0000426	0.0000830	0.5136285	0.6186686
TotalDeaths	0.0051524	0.0001694	30.4150162	0.0000000
ActiveCases	-0.0008670	0.0001057	-8.2021353	0.0000095
Critical	0.0244141	0.0008389	29.1038920	0.0000000

After conducting the step method and looking at the full linear regression model, we see that it is likely that the Total Cases variable is not a good predictor of the Mortality Recovery Ratio. However this could be misleading due to correlation between variables, so we can check with a full and reduced ANOVA model.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.3487820	0.0903634	3.85977	0.0026544
TotalDeaths	0.0052368	0.0000398	131.66213	0.0000000
ActiveCases	-0.0008128	0.0000041	-196.44007	0.0000000
Critical	0.0239874	0.0001121	214.04372	0.0000000

We build the model based on the following variables:

$$\begin{aligned}
Y &= MortalityRecoveryRatio \\
X_1 &= TotalCases \\
X_2 &= TotalDeaths \\
X_3 &= ActiveCases \\
X_4 &= Critical
\end{aligned}$$

Our full model is  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$  and our reduced model is  $Y = \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$ . Running the anova model we get the following.

$$\begin{aligned}
H_0: & \beta_2 = \beta_3 = \beta_4 \\
H_A: & \beta_i \neq 0 \text{ for at least one } i
\end{aligned}$$

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
10	1.044377	NA	NA	NA	NA
11	1.071930	-1	-0.0275522	0.2638143	0.6186686

Interestingly the P-value for the ANOVA model is quite high and the F-value is very low, so we cannot conclude that there is any statistically significant difference between the reduced and full models in predicting the  $Y$  value. We can conclude that the total cases does not statistically contribute significantly to the mortality recovery ratio. The most statistically significant variables in determining the ratio are total deaths, active cases, and critical cases.

## IV. Conclusion & Future Work

Overall, the dataset shows that the most significant variables in predicting the mortality recovery ratio are total deaths, active case, and critical cases. We learned this from looking at the original full model's p-values as well as doing the anova test in order to check for the variables significance.

Despite nearing the end of the pandemic, there is still a lot of information to be garnered about how spread occurs quickly and countermeasures to slow down the spread of a virus. Those who lived in a poor environment and jobs that required close contact, were definitely areas that worsened the spread of the virus. Future possible research can see how socioeconomic factors and occupation played a huge role in the mortality rate ratio. How much of a significance difference would the ratios between those who have higher incomes compared to lower incomes differ? Overall, studies about the pandemic will allow for more rules and infrastructure to take place, thereby enriching our society as a whole.

## Appendix: R Script

```
knitr::opts_chunk$set(echo = F, warning = F, message = F)
library(knitr)
library(dplyr)
library(readr)
covid = as.data.frame(read_delim("data/COVID.csv", delim = ','))
covidpt = covid[-c(4, 10, 16), 3:7]

fmodel = lm(MortalityRecoveryRatio ~ ., data = covidpt)
summary(fmodel)

step(fmodel)

rmodel = lm(MortalityRecoveryRatio ~ TotalDeaths + ActiveCases
            + Critical, data = covidpt)
summary(rmodel)

anova(fmodel, rmodel)
summary(fmodel)$coefficients %>% kable()
summary(rmodel)$coefficients %>% kable()
anova(fmodel, rmodel) %>% kable()
```