

# Covid Analysis and Models

Cody S

```
library("tidyverse")
```

## Goal

My main goal will be to see if I can get a general sense of how the number of deaths relative to the number of cases changed as vaccines were developed and rolled out.

## Import Data

```
# set urls

global_cases_url <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/refs/heads/master/csse_github_data/csse_github_data_global_cases.csv"
global_deaths_url <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/refs/heads/master/csse_github_data/csse_github_data_global_deaths.csv"
US_cases_url <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/refs/heads/master/csse_github_data/csse_github_data_us_cases.csv"
US_death_url <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/refs/heads/master/csse_github_data/csse_github_data_us_deaths.csv"
```

```
#import to dataframes

global_cases <- read.csv(global_cases_url)
global_deaths <- read.csv(global_deaths_url)
US_cases <- read.csv(US_cases_url)
US_deaths <- read.csv(US_death_url)

all_data <- list(global_cases, global_deaths, US_cases, US_deaths)
```

```
lapply(all_data, head)
```

## Start cleaning and organizing

The data is broken down into pretty small regions, and the regions are the rows and the dates are the columns. My main goal will be consolidated all of the global data into a single time series of cases and deaths, and all of the USA data in the same way for comparison. This will mostly involve dropping the geographic breakdown and pivoting the data so that the dates are the rows.

```
geo_rename <- function(df) {  
  df <- df %>%  
    rename(Province.State = Province_State,  
           Country.Region = Country_Region) %>%  
    select(-Admin2, -UID, -FIPS, -code3, -Combined_Key, -starts_with("iso"))  
  
  return(df)  
}  
  
US_cases_renamed <- geo_rename(US_cases)  
US_deaths_renamed <- geo_rename(US_deaths)  
  
lapply(list(US_cases_renamed, US_deaths_renamed), head)
```

```
# drop some columns and pivot to dates  
pivot_organize <- function(df) {  
  df <- df %>%  
    select(-Lat, -contains("Long"), -contains("Population")) %>%  
    pivot_longer(  
      cols = -c(Province.State, Country.Region),  
      names_to = "Date",  
      values_to = "Cases"  
    ) %>%  
    select(Date, Country.Region, Province.State, Cases)  
    %>%  
    mutate(Date = sub("^X", "", Date))  
    %>%  
    mutate(Date = as.Date(Date, format = "%m.%d.%y"))  
  
  return(df)  
}
```

```
# apply the pivot and reorganization to all data sets
global_cases_clean <- pivot_organize(global_cases)
global_deaths_clean <- pivot_organize(global_deaths)
us_cases_clean <- pivot_organize(US_cases_renamed)
us_deaths_clean <- pivot_organize(US_deaths_renamed)

# make sure to keep cases and deaths straight
global_deaths_clean <- global_deaths_clean %>%
  rename(Deaths = Cases)
us_deaths_clean <- us_deaths_clean %>%
  rename(Deaths = Cases)

all_data_clean <- list(global_cases_clean, global_deaths_clean, us_cases_clean, us_deaths_clean)

lapply(all_data_clean, head)
```

Now we need to consolidate all of the sub-region data into overall totals. The next few steps will be complicated so I'm going to do it with global data first.

```
# make a simple global daily cases df with a weekly column as well
global_daily_cases <- global_cases_clean %>%
  group_by(Date) %>%
  summarise(Total_cases = sum(Cases)) %>%
  mutate(New_daily_cases = Total_cases - lag(Total_cases, n = 1, default = NA)) %>%
  mutate(Week = floor_date(Date, unit = "week", week_start = 1))

global_daily_deaths <- global_deaths_clean %>%
  group_by(Date) %>%
  summarise(Total_deaths = sum(Deaths)) %>%
  mutate(New_daily_deaths = Total_deaths - lag(Total_deaths, n = 1, default = NA)) %>%
  mutate(Week = floor_date(Date, unit = "week", week_start = 1))

tail(global_daily_cases)
tail(global_daily_deaths)

# combine daily deaths and cases
daily_global_combo <- left_join(global_daily_cases, global_daily_deaths, by = "Date")

daily_global_combo <- daily_global_combo %>%
  select(-Week.x) %>%
  rename(Week = Week.y)
```

```
# create a weekly global death/cases df
weekly_global_combo <- daily_global_combo %>%
  group_by(Week) %>%
  summarise(
    g_Weekly_total_cases = last(Total_cases),
    g_Weekly_total_deaths = last(Total_deaths),
    g_Weekly_new_cases = sum(New_daily_cases, na.rm = TRUE),
    g_Weekly_new_deaths = sum(New_daily_deaths, na.rm = TRUE)
  )

tail(daily_global_combo)
tail(weekly_global_combo)
```

```
# start plotting soon for testing, so this will set global fig dimensions
options(repr.plot.width = 10, repr.plot.height = 6)
```

```
ggplot(data = weekly_global_combo, aes(x = Week, y = g_Weekly_new_cases)) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(
    title = "Global Weekly Cases",
    x = "Week",
    y = "New Cases"
  ) +
  theme_minimal()

ggplot(data = weekly_global_combo, aes(x = Week, y = g_Weekly_new_deaths)) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(
    title = "Global Weekly Deaths",
    x = "Week",
    y = "New Deaths"
  ) +
  theme_minimal()
```

That looks good, and I can see some general trends in deaths over time, so I think I can move forward with calculating the ratio of deaths to cases, which will also allow me to compare it on the same scale to the US data later.

```
weekly_global_combo <- weekly_global_combo %>%
  mutate(g_Weekly_death_to_case_ratio = (g_Weekly_new_deaths / g_Weekly_new_cases)* 100)
```

```
tail(weekly_global_combo)
```

```
ggplot(data = weekly_global_combo, aes(x = Week, y = g_Weekly_death_to_case_ratio)) +  
  geom_bar(stat = "identity", fill = "blue") +  
  labs(  
    title = "Weekly New Deaths to Cases Ratio (%)",  
    x = "Week",  
    y = "Death/Case Ratio (%)"  
  ) +  
  theme_minimal()
```

I'm using pretty general vaccine availability dates and using as somewhat arbitrary one year span as the “rollout” duration. The post- to pre-vaccine comparison will be the most important in the end.

```
# set phase dates  
vac_rollout_start <- as.Date("2020-12-01")  
vac_rollout_end <- as.Date("2021-12-01")  
  
# add vaccine phase code  
weekly_global_combo <- weekly_global_combo %>%  
  mutate(Phase = case_when(  
    Week < vac_rollout_start ~ "1. Pre",  
    Week >= vac_rollout_start & Week < vac_rollout_end ~ "2. Rollout",  
    Week >= vac_rollout_end ~ "3. Post"  
  ))  
  
tail(weekly_global_combo)
```

That looks good, so now I'll just calculate some rudimentary death/case ratio averages for each phase.

```
global_phase_means <- weekly_global_combo %>%  
  group_by(Phase) %>%  
  summarise(global_phase_avg = mean(g_Weekly_death_to_case_ratio, na.rm = TRUE))  
  
global_phase_means
```

```
global_phase_means_plot <- ggplot(data = global_phase_means, aes(x = Phase, y = global_phase_avg)) +  
  geom_bar(stat = "identity", fill = "blue") +
```

```

labs(
  title = "Global Phase Average Ratio",
  x = "Phase",
  y = "Average Death to Case ratio (%)"
) +
theme_minimal()

global_phase_means_plot

```

## Model development

The visuals and simple checks all indicate that the ratio of deaths to cases were much lower after the vaccine rollout. Now I'll model that with a multi-part (or interrupted) linear analysis for each phase.

### Test model with global data

```

# do a simple phased linear regression
global_phase_model <- lm(g_Weekly_death_to_case_ratio ~ Week * factor(Phase), data = weekly_g
summary(global_phase_model)

```

```

weekly_global_combo$g_Modeled <- predict(global_phase_model)

```

```

global_ratio_data_model <- ggplot(weekly_global_combo, aes(x = Week)) +
  geom_line(aes(y = g_Weekly_death_to_case_ratio), color = "blue") +
  geom_line(aes(y = g_Modeled), color = "red") +
  labs(title = "Global Data and Phase Model",
       x = "Week", y = "Ratio") +
  geom_vline(xintercept = vac_rollout_start, linetype = "dashed", color = "black") +
  geom_vline(xintercept = vac_rollout_end, linetype = "dashed", color = "black") +
  theme_minimal()

global_ratio_data_model

```

## US Data clean and model prep

Now I want to do the same thing for the USA data, then I'll replot everything at the end

```

# make a simple us daily cases df with a weekly column as well
us_daily_cases <- us_cases_clean %>%
  group_by(Date) %>%
  summarise(Total_cases = sum(Cases)) %>%
  mutate(New_daily_cases = Total_cases - lag(Total_cases, n = 1, default = NA)) %>%
  mutate(Week = floor_date(Date, unit = "week", week_start = 1))

us_daily_deaths <- us_deaths_clean %>%
  group_by(Date) %>%
  summarise(Total_deaths = sum(Deaths)) %>%
  mutate(New_daily_deaths = Total_deaths - lag(Total_deaths, n = 1, default = NA)) %>%
  mutate(Week = floor_date(Date, unit = "week", week_start = 1))

# combine daily deaths and cases
daily_us_combo <- left_join(us_daily_cases, us_daily_deaths, by = "Date")

daily_us_combo <- daily_us_combo %>%
  select(-Week.x) %>%
  rename(Week = Week.y)

# create a weekly US death/cases df
weekly_us_combo <- daily_us_combo %>%
  group_by(Week) %>%
  summarise(
    us_Weekly_total_cases = last(Total_cases),
    us_Weekly_total_deaths = last(Total_deaths),
    us_Weekly_new_cases = sum(New_daily_cases, na.rm = TRUE),
    us_Weekly_new_deaths = sum(New_daily_deaths, na.rm = TRUE)
  )

ggplot(data = weekly_us_combo, aes(x = Week, y = us_Weekly_new_cases)) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(
    title = "USA Weekly Cases",
    x = "Week",
    y = "New Cases"
  ) +
  theme_minimal()

ggplot(data = weekly_us_combo, aes(x = Week, y = us_Weekly_new_deaths)) +
  geom_bar(stat = "identity", fill = "blue") +

```

```

labs(
  title = "USA Weekly Deaths",
  x = "Week",
  y = "New Deaths"
) +
theme_minimal()

```

```

weekly_us_combo <- weekly_us_combo %>%
  mutate(us_Weekly_death_to_case_ratio = (us_Weekly_new_deaths / us_Weekly_new_cases)* 100)

ggplot(data = weekly_us_combo, aes(x = Week, y = us_Weekly_death_to_case_ratio)) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(
    title = "Weekly New Deaths to Cases Ratio (%)",
    x = "Week",
    y = "USA Death/Case Ratio (%)"
  ) +
  theme_minimal()

```

```

weekly_us_combo <- weekly_us_combo %>%
  mutate(Phase = case_when(
    Week < vac_rollout_start ~ "1. Pre",
    Week >= vac_rollout_start & Week < vac_rollout_end ~ "2. Rollout",
    Week >= vac_rollout_end ~ "3. Post"
  ))

# check that the phase and ratio calcs all worked correctly
tail(weekly_us_combo)

```

```

# calculate phase means
us_phase_means <- weekly_us_combo %>%
  group_by(Phase) %>%
  summarise(us_phase_avg = mean(us_Weekly_death_to_case_ratio, na.rm = TRUE))

# display the phase means
us_phase_means

# plot phase means
us_phase_means_plot <- ggplot(data = us_phase_means, aes(x = Phase, y = us_phase_avg)) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(

```



```

    title = "USA Phase Average Ratio",
    x = "Phase",
    y = "Average Death to Case ratio (%)"
  ) +
  theme_minimal()

us_phase_means_plot

```

```

# do a simple phased linear regression
us_phase_model <- lm(us_Weekly_death_to_case_ratio ~ Week * factor(Phase), data = weekly_us_

summary(us_phase_model)

weekly_us_combo$us_Modeled <- predict(us_phase_model)

us_ratio_data_model <- ggplot(weekly_us_combo, aes(x = Week)) +
  geom_line(aes(y = us_Weekly_death_to_case_ratio), color = "blue") +
  geom_line(aes(y = us_Modeled), color = "red") +
  labs(title = "USA Data and Phase Model",
       x = "Week", y = "Ratio") +
  geom_vline(xintercept = vac_rollout_start, linetype = "dashed", color = "black") +
  geom_vline(xintercept = vac_rollout_end, linetype = "dashed", color = "black") +
  theme_minimal()

us_ratio_data_model

```

```

# join the global and USA data to see if I can plot them together nicely
combined_means <- inner_join(global_phase_means, us_phase_means, by = "Phase")
combined_means

```

## Final Models and Visuals

### Combined phase averages

```

# Reshape the data to long format for plotting
long_combined_means <- combined_means %>%
  pivot_longer(cols = c(global_phase_avg, us_phase_avg),
               names_to = "Region",
               values_to = "Phase_avg") %>%

```

```

mutate(Region = ifelse(Region == "global_phase_avg", "Global", "USA"))

long_combined_means

combined_means_plot <- ggplot(data = long_combined_means,
  aes(x = Phase, y = Phase_avg, fill = Region)) +
  geom_bar(stat = "identity", position = "dodge") +
  scale_fill_manual(values = c("Global" = "blue", "USA" = "red")) +
  labs(
    title = "Global and USA Phase Average Ratio",
    x = "Phase",
    y = "Average Death to Case Ratio") +
  theme_minimal()

print(combined_means_plot)

combined_model_df <- inner_join(weekly_global_combo, weekly_us_combo, by = "Week")

tail(combined_model_df)

```

## Final Combined Models and Visuals

```

combined_model_plot <- ggplot(combined_model_df, aes(x = Week)) +
  # Global lines
  geom_line(aes(y = g_Weekly_death_to_case_ratio, color = "Global")) +
  geom_line(aes(y = g_Modeled, color = "Global", linetype = "Global")) +
  # USA lines
  geom_line(aes(y = us_Weekly_death_to_case_ratio, color = "USA")) +
  geom_line(aes(y = us_Modeled, color = "USA", linetype = "USA"))
  ) +
  labs(
    title = "Global and USA Death to Case Ratio and Phase Models",
    x = "Date",
    y = "Death to Case Ratio"
  ) +
  scale_color_manual(name = "Ratio Data", values = c("Global" = "blue", "USA" = "red")) +
  scale_linetype_manual(name = "Models", values = c("Global" = "dashed", "USA" = "dotted")) +
  theme_minimal()

```

```
print(combined_means_plot)
print(combined_model_plot)
```

I think that accomplishes the goal. With these two figures we can see the changes in the death to cases ratios before and after (and during) the vaccine rollout for the global and USA datasets, as well as the general linear trend during each phase.

## Bias statement

There are a few potential sources of bias or inaccuracy in this data and the analysis.

- Collection bias
  - This data was taken from the Johns Hopkins Covid project source without modification. Johns Hopkins is assumed to be a credible source with good data handling practices, but the accuracy of this data is still dependent on the accuracy of the data recorded and reporting by each region. Countries or sub-regions with less developed public health infrastructure may be missing or under-reported.
- Analysis bias
  - A single start and end date were used to define the vaccine rollout phase. In reality, several countries (mostly more wealthy) had access to the vaccines months before poorer countries and were able to distribute them to the general population much more rapidly.
  - Death counts are likely to causally lag case counts on the order of a few weeks. That lag wasn't considered when calculating the death to case ratios. It was assumed that over periods as long as the vaccine phases that the ratios would stabilize so that the averages would be meaningful representations of those phases. More granular analysis would require account for the case to death time lag.
  - Specific vaccine rates were not considered as they weren't included in the data. The demonstrated decrease in the death to case ratios can be correlated to the general availability of vaccines, but further analysis of those ratios for specific regions or over smaller time or any quantitative assessment of vaccine effectiveness would require adjusting for specific vaccination rates.

```
sessionInfo()
```