

Final Covid 19

A. Ayzenberg

2025-04-28

Necessary Libraries

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.1      v tibble    3.2.1
## v lubridate  1.9.4      v tidyr     1.3.1
## v purrr      1.0.4
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(lubridate)
library(ggplot2)
```

Import the data

```
url_in <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_data"
file_names <- c("time_series_covid19_confirmed_US.csv", "time_series_covid19_confirmed_global.csv", "time_series_covid19_deaths_US.csv", "time_series_covid19_deaths_global.csv")
urls <- str_c(url_in, file_names)
global_cases <- read_csv(urls[2])
```

```
## Rows: 289 Columns: 1147
## -- Column specification -----
## Delimiter: ","
## chr (2): Province/State, Country/Region
## dbl (1145): Lat, Long, 1/22/20, 1/23/20, 1/24/20, 1/25/20, 1/26/20, 1/27/20,...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
global_deaths <- read_csv(urls[4])
```

```
## Rows: 289 Columns: 1147
## -- Column specification -----
## Delimiter: ","
## chr      (2): Province/State, Country/Region
## dbl (1145): Lat, Long, 1/22/20, 1/23/20, 1/24/20, 1/25/20, 1/26/20, 1/27/20,...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
US_cases <- US_cases <- read_csv(urls[1])
```

```
## Rows: 3342 Columns: 1154
## -- Column specification -----
## Delimiter: ","
## chr      (6): iso2, iso3, Admin2, Province_State, Country_Region, Combined_Key
## dbl (1148): UID, code3, FIPS, Lat, Long_, 1/22/20, 1/23/20, 1/24/20, 1/25/20...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
US_deaths <- read_csv(urls[3])
```

```
## Rows: 3342 Columns: 1155
## -- Column specification -----
## Delimiter: ","
## chr      (6): iso2, iso3, Admin2, Province_State, Country_Region, Combined_Key
## dbl (1149): UID, code3, FIPS, Lat, Long_, Population, 1/22/20, 1/23/20, 1/24...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

Transform the data for visualization and analysis

```
# transform global data
global_cases <- global_cases %>% pivot_longer(cols = -c(`Province/State`, `Country/Region`, Lat, Long),
global_deaths <- global_deaths %>% pivot_longer(cols = -c(`Province/State`, `Country/Region`, Lat, Long),
global <- global_cases %>% full_join(global_deaths) %>% rename(Country_Region = `Country/Region`, Provi
```

```
## Joining with 'by = join_by('Province/State', 'Country/Region', date)'
```

```
global <- global %>% filter(cases > 0)
summary(global)
```

```
## Province_State    Country_Region      date      cases
## Length:306827     Length:306827    Min.   :2020-01-22  Min.   :      1
## Class :character   Class :character 1st Qu.:2020-12-12 1st Qu.:    1316
## Mode  :character   Mode  :character Median :2021-09-16 Median :    20365
##                      Mean  :2021-09-11 Mean  :   1032863
##                      3rd Qu.:2022-06-15 3rd Qu.:   271281
```

```
##                                     Max.      :2023-03-09   Max.      :103802702
##      deaths
## Min.      :      0
## 1st Qu.:      7
## Median :    214
## Mean   :   14405
## 3rd Qu.:   3665
## Max.    : 1123836
```

```
# Verify that the following isn't a single outlier
global %>% filter(cases>100000000)
```

```
## # A tibble: 80 x 5
##   Province_State Country_Region date           cases  deaths
##   <chr>          <chr>      <date>        <dbl>   <dbl>
## 1 <NA>          US        2022-12-20  100050937 1088341
## 2 <NA>          US        2022-12-21  100233060 1089383
## 3 <NA>          US        2022-12-22  100329204 1089979
## 4 <NA>          US        2022-12-23  100368433 1090186
## 5 <NA>          US        2022-12-24  100374955 1090208
## 6 <NA>          US        2022-12-25  100378169 1090223
## 7 <NA>          US        2022-12-26  100390601 1090252
## 8 <NA>          US        2022-12-27  100501536 1090608
## 9 <NA>          US        2022-12-28  100614880 1091598
## 10 <NA>         US        2022-12-29  100718983 1092522
## # i 70 more rows
```

```
# transform US data
```

```
US_cases <- US_cases %>% pivot_longer(cols = -(UID:Combined_Key), names_to = "date", values_to = "cases")
US_deaths <- US_deaths %>% pivot_longer(cols = -(UID:Population), names_to = "date", values_to = "deaths")
US <- US_cases %>% full_join(US_deaths)
```

```
## Joining with `by = join_by(Admin2, Province_State, Country_Region,
## Combined_Key, date)`
```

```
global <- global %>% unite("Combined_Key", c(Province_State, Country_Region), sep = ", ", na.rm = TRUE,
# add population statistics to global data
uid_lookup_url <- "https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/
uid <- read_csv(uid_lookup_url) %>% select(-c(Lat,Long_, Combined_Key, code3, iso2, iso3, Admin2))
```

```
## Rows: 4321 Columns: 12
## -- Column specification -----
## Delimiter: ","
## chr (7): iso2, iso3, FIPS, Admin2, Province_State, Country_Region, Combined_Key
## dbl (5): UID, code3, Lat, Long_, Population
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
global <- global %>% left_join(uid, by = c("Province_State", "Country_Region")) %>% select(-c(UID, FIPS))
# Final preparations for plotting and visualizing the data
US_by_state <- US %>% group_by(Province_State, Country_Region, date) %>% summarize(cases = sum(cases),
```

```
## 'summarise()' has grouped output by 'Province_State', 'Country_Region'. You can
## override using the '.groups' argument.
```

```
US_totals <- US_by_state %>% group_by(Country_Region, date) %>% summarize(cases = sum(cases), deaths = sum(deaths))
```

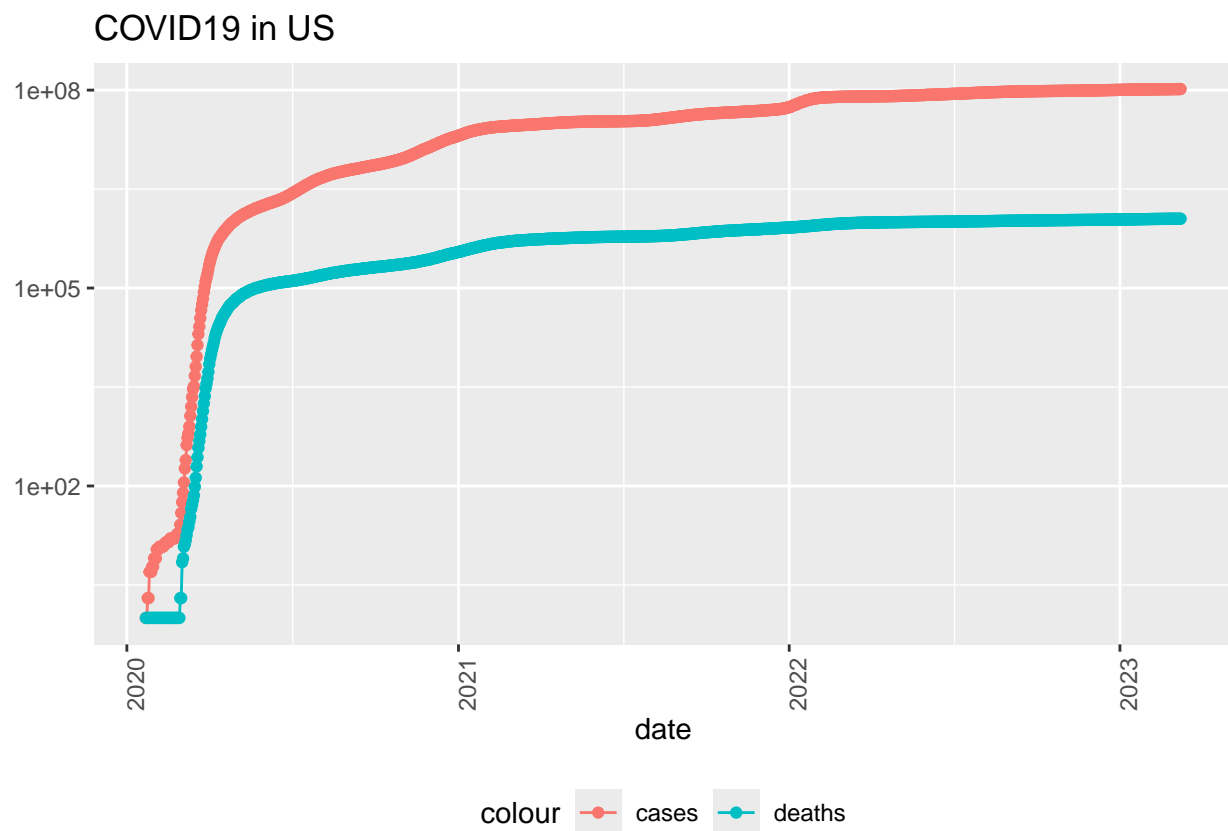
```
## 'summarise()' has grouped output by 'Country_Region'. You can override using
## the '.groups' argument.
```

Visualizing US statistics

```
# Find out maximum date (how far the data goes)
max(US_totals$date)
```

```
## [1] "2023-03-09"
```

```
US_totals %>% filter(cases > 0) %>% ggplot(aes(x = date, y = cases)) + geom_line(aes(color = "cases")) +
```

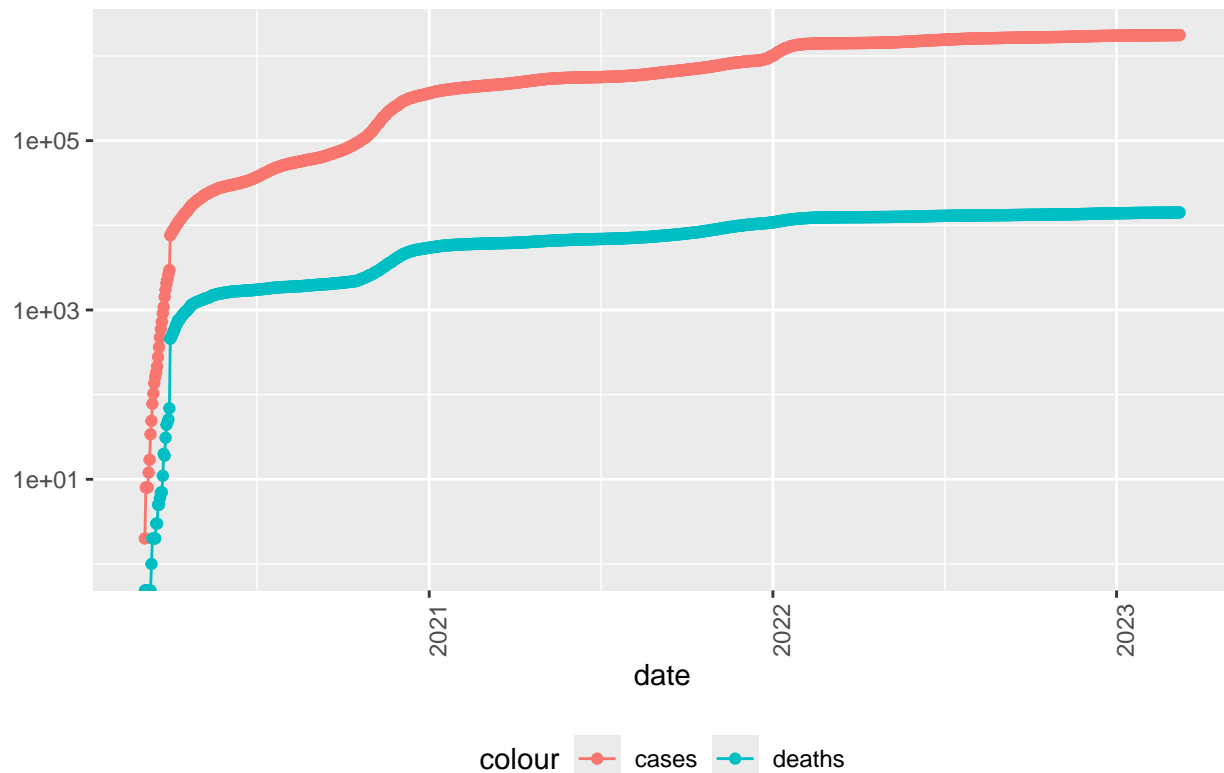


Visualizing data for Colorado

```
state <- "Colorado"
US_by_state %>% filter(Province_State == state) %>% filter(cases > 0) %>% ggplot(aes(x = date, y = cases))
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
## log-10 transformation introduced infinite values.
```

COVID19 in Colorado



Evaluate only new cases

In this case only new cases will be looked at to visualize if there is a trend for US data

```
## Transforming and looking at US new cases
US_by_state <- US_by_state %>% mutate(new_cases = cases - lag(cases), new_deaths = deaths - lag(deaths))
US_totals <- US_totals %>% mutate(new_cases = cases - lag(cases), new_deaths = deaths - lag(deaths))
US_totals %>% ggplot(aes(x = date, y = new_cases)) + geom_line(aes(color = "new_cases")) + geom_point(aes(color = "new_deaths"))
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

```
## Warning in transformation$transform(x): NaNs produced

## Warning in scale_y_log10(): log-10 transformation introduced infinite values.

## Warning in transformation$transform(x): NaNs produced

## Warning in scale_y_log10(): log-10 transformation introduced infinite values.

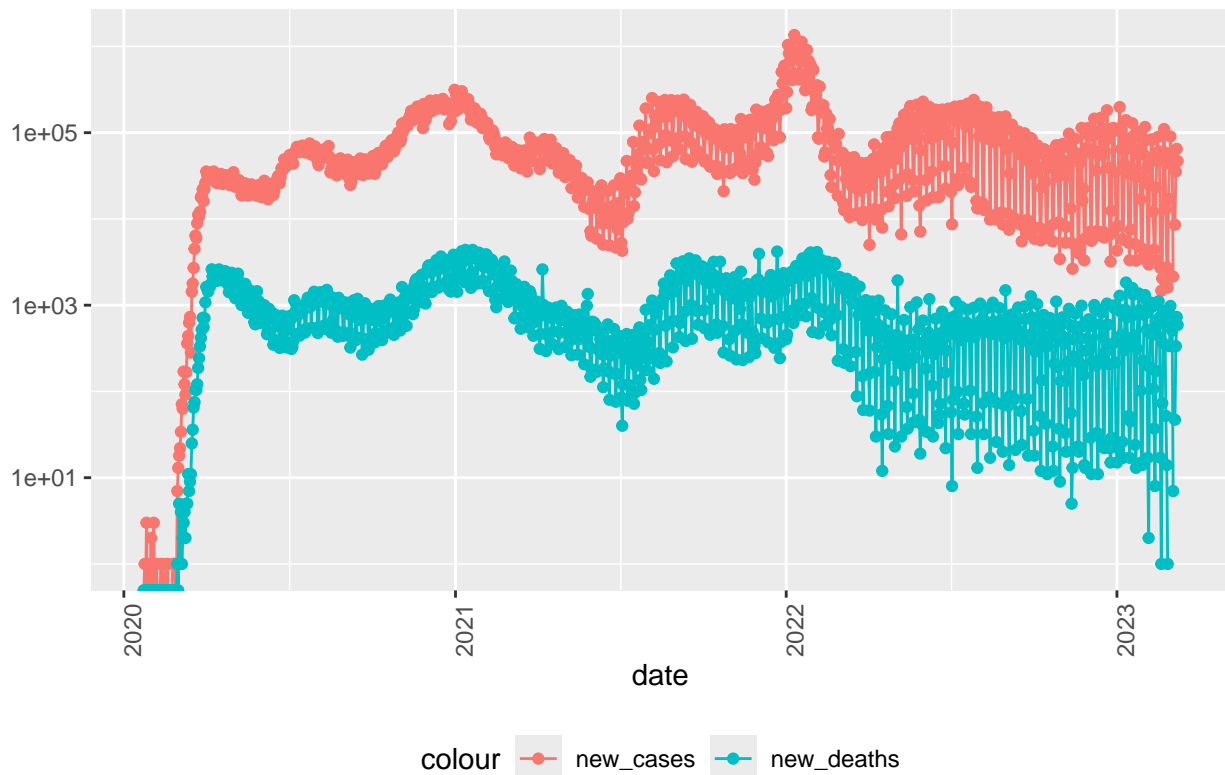
## Warning: Removed 1 row containing missing values or values outside the scale range
## ('geom_line()').

## Warning: Removed 2 rows containing missing values or values outside the scale range
## ('geom_point()').

## Warning: Removed 1 row containing missing values or values outside the scale range
## ('geom_line()').

## Warning: Removed 4 rows containing missing values or values outside the scale range
## ('geom_point()').
```

COVID19 in US



```
US_state_totals <- US_by_state %>% group_by(Province_State) %>% summarize(deaths = max(deaths), cases =
US_state_totals %>% slice_min(deaths_per_thou, n = 10)
```

```
## # A tibble: 10 x 6
##   Province_State    deaths    cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl>   <dbl>      <dbl>         <dbl>         <dbl>
## 1 American Samoa      34 8.32e3    55641         150.          0.611
## 2 Northern Mariana Isl~  41 1.37e4    55144         248.          0.744
## 3 Virgin Islands     130 2.48e4   107268         231.          1.21
## 4 Hawaii            1841 3.81e5   1415872         269.          1.30
## 5 Vermont             929 1.53e5    623989         245.          1.49
## 6 Puerto Rico        5823 1.10e6   3754939         293.          1.55
## 7 Utah              5298 1.09e6   3205958         340.          1.65
## 8 Alaska             1486 3.08e5    740995         415.          2.01
## 9 District of Columbia  1432 1.78e5    705749         252.          2.03
## 10 Washington       15683 1.93e6   7614893         253.          2.06
```

Evaluate only new cases (state)

In this case only new cases will be looked at to visualize if there is a trend for CO data

```
## Transforming and looking at Colorado new cases
US_by_state %>% filter(Province_State == state) %>% ggplot(aes(x = date, y = new_cases)) + geom_line(aes
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

```
## Warning in transformation$transform(x): NaNs produced
```

```
## Warning in scale_y_log10(): log-10 transformation introduced infinite values.
```

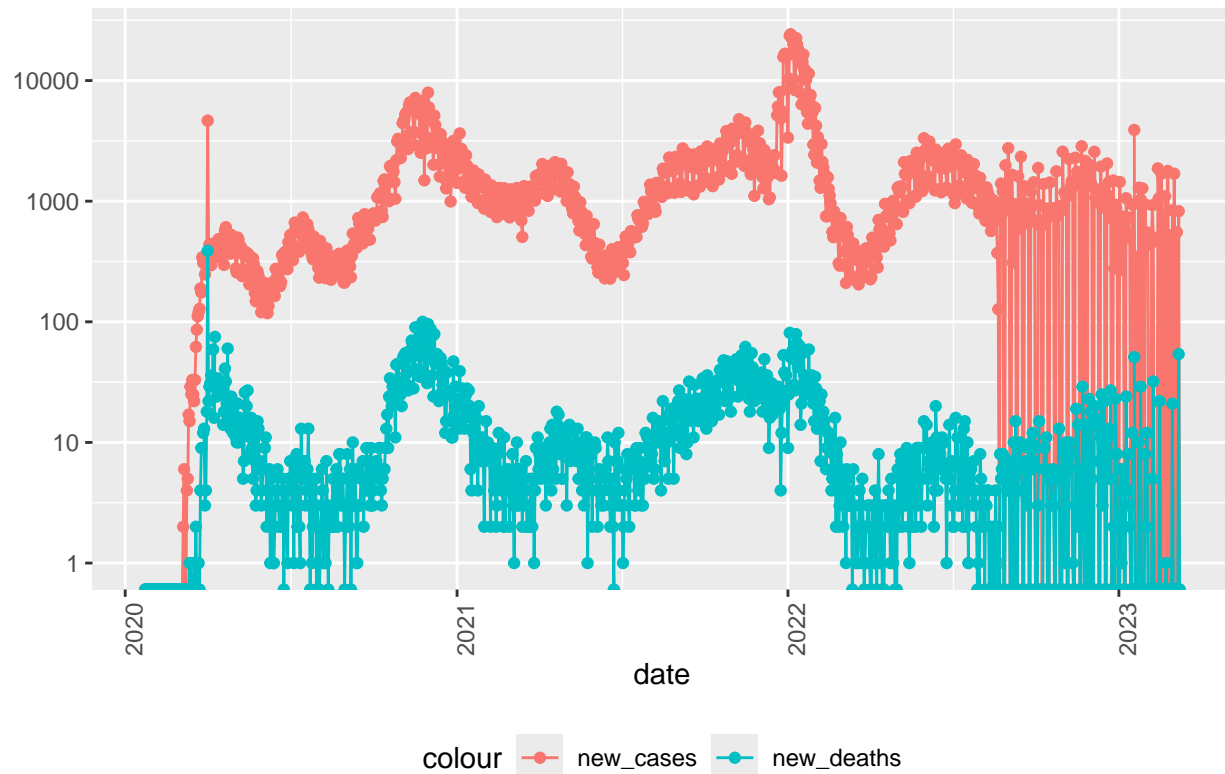
```
## Warning: Removed 1 row containing missing values or values outside the scale range
## ('geom_line()').
```

```
## Warning: Removed 2 rows containing missing values or values outside the scale range
## ('geom_point()').
```

```
## Warning: Removed 1 row containing missing values or values outside the scale range
## ('geom_line()').
```

```
## Warning: Removed 5 rows containing missing values or values outside the scale range
## ('geom_point()').
```

COVID19 in Colorado



```
US_state_totals <- US_by_state %>% group_by(Province_State) %>% summarize(deaths = max(deaths), cases =
US_state_totals %>% slice_min(deaths_per_thou, n = 10)
```

```
## # A tibble: 10 x 6
##   Province_State      deaths    cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl>   <dbl>      <dbl>         <dbl>         <dbl>
## 1 American Samoa      34 8.32e3    55641         150.           0.611
## 2 Northern Mariana Isl~  41 1.37e4    55144         248.           0.744
## 3 Virgin Islands     130 2.48e4   107268         231.           1.21
## 4 Hawaii            1841 3.81e5   1415872         269.           1.30
## 5 Vermont             929 1.53e5    623989         245.           1.49
## 6 Puerto Rico        5823 1.10e6   3754939         293.           1.55
## 7 Utah              5298 1.09e6   3205958         340.           1.65
## 8 Alaska            1486 3.08e5    740995         415.           2.01
## 9 District of Columbia 1432 1.78e5    705749         252.           2.03
## 10 Washington       15683 1.93e6   7614893         253.           2.06
```

Worst and Least impacted

What if the goal is to see which states suffered the least deaths, the most?

```
## Who faired the best, the worst?
US_state_totals <- US_by_state %>% group_by(Province_State) %>% summarize(deaths = max(deaths), cases =
US_state_totals %>% slice_min(deaths_per_thou, n = 10)
```



```
## # A tibble: 10 x 6
##   Province_State deaths cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl> <dbl>      <dbl>          <dbl>          <dbl>
## 1 American Samoa      34 8.32e3      55641          150.          0.611
## 2 Northern Mariana Isl~  41 1.37e4      55144          248.          0.744
## 3 Virgin Islands     130 2.48e4     107268          231.          1.21
## 4 Hawaii             1841 3.81e5     1415872          269.          1.30
## 5 Vermont             929 1.53e5      623989          245.          1.49
## 6 Puerto Rico        5823 1.10e6     3754939          293.          1.55
## 7 Utah               5298 1.09e6     3205958          340.          1.65
## 8 Alaska             1486 3.08e5      740995          415.          2.01
## 9 District of Columbia 1432 1.78e5      705749          252.          2.03
## 10 Washington        15683 1.93e6     7614893          253.          2.06
```

```
US_state_totals %>% slice_max(deaths_per_thou, n = 10)
```

```
## # A tibble: 10 x 6
##   Province_State deaths cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl> <dbl>      <dbl>          <dbl>          <dbl>
## 1 Arizona      33102 2443514     7278717          336.          4.55
## 2 Oklahoma     17972 1290929     3956971          326.          4.54
## 3 Mississippi  13370 990756     2976149          333.          4.49
## 4 West Virginia  7960 642760     1792147          359.          4.44
## 5 New Mexico    9061 670929     2096829          320.          4.32
## 6 Arkansas     13020 1006883     3017804          334.          4.31
## 7 Alabama      21032 1644533     4903185          335.          4.29
## 8 Tennessee    29263 2515130     6829174          368.          4.28
## 9 Michigan     42205 3064125     9986857          307.          4.23
## 10 Kentucky    18130 1718471     4467673          385.          4.06
```

Modeling the data

Below the data will be plotted to determine a relation between cases per thousand and deaths per thousand, and also visualize to see which points faired better than expected and which did worse.

```
mod <- lm(deaths_per_thou ~ cases_per_thou, data = US_state_totals)
summary(mod)
```

```
##
## Call:
## lm(formula = deaths_per_thou ~ cases_per_thou, data = US_state_totals)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.3352 -0.5978  0.1491  0.6535  1.2086
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.36167    0.72480  -0.499    0.62
## cases_per_thou  0.01133    0.00232   4.881 9.76e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.8615 on 54 degrees of freedom
## Multiple R-squared:  0.3061, Adjusted R-squared:  0.2933
## F-statistic: 23.82 on 1 and 54 DF,  p-value: 9.763e-06
```

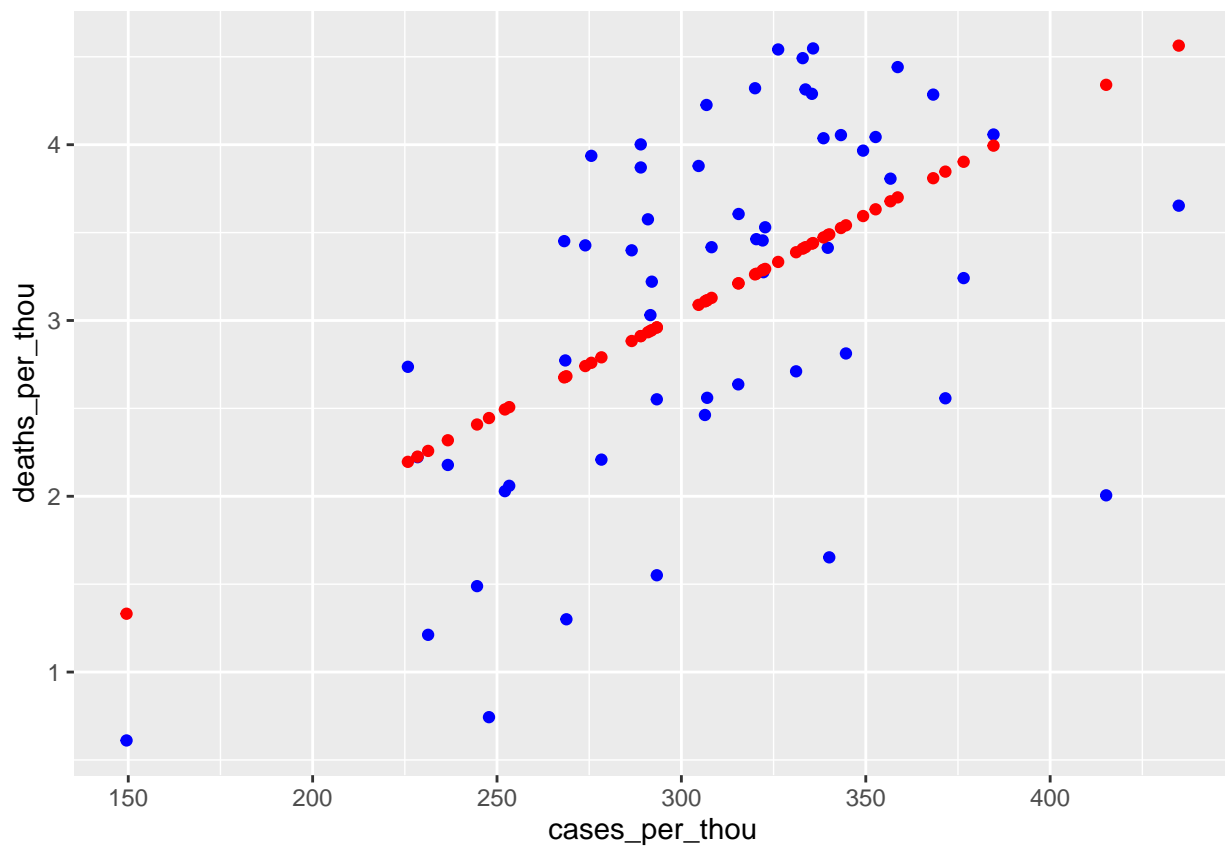
```
US_state_totals %>% slice_min(cases_per_thou)
```

```
## # A tibble: 1 x 6
##   Province_State deaths cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl> <dbl>      <dbl>          <dbl>          <dbl>
## 1 American Samoa      34  8320      55641          150.           0.611
```

```
US_state_totals %>% slice_max(cases_per_thou)
```

```
## # A tibble: 1 x 6
##   Province_State deaths cases population cases_per_thou deaths_per_thou
##   <chr>          <dbl> <dbl>      <dbl>          <dbl>          <dbl>
## 1 Rhode Island    3870 460697    1059361          435.           3.65
```

```
x_grid <- seq(150, 450)
US_tot_w_pred <- US_state_totals %>% mutate(pred = predict(mod))
US_tot_w_pred %>% ggplot() + geom_point(aes(x = cases_per_thou, y = deaths_per_thou), color = "blue") +
```



Potential Biases

The biggest source of bias in this case would be reporting or rather under reporting in some parts of the world, especially in the early days of covid, most notably due to lack of funding, insufficient numbers of tests and potential for remote areas. The same can be said for deaths as one factor not put into consideration is a country or region's wealth, meaning a more accurate comparison would be between 2 regions of comparable wealth (i.e. GDP per capita) as opposed to already striggling countries being labeled as handling the epidemic worse than regions with higher fund availability.