# Final Project

Code ▾

Hoin Lee & Changbin Kim

13/DEC/2022

## Setting

Hide

```
# load all packages here: `mosaic`, `tidyverse`, `lubridate`, and all others used
library(tidyverse)
library(lubridate)
library(ggplot2)
library(dplyr)
library(esquisse)
library(mosaic)
library(mosaicData)
library(dcData)
```

Hide

```
primary <- Provisional_COVID_19_Deaths_by_Sex_and_Age
death <- read_csv("United_States_COVID-19_Cases_and_Deaths_by_State_over_Time_-_ARCHIVED.csv")
```

```
Rows: 60060 Columns: 15── Column specification ──────────────────────────────────────────
───────────────────────────────────────
Delimiter: ","
chr  (5): submission_date, state, created_at, consent_cases, consent_deaths
dbl (10): tot_cases, conf_cases, prob_cases, new_case, pnew_case, tot_death, conf_death, prob_dea
th, new_death, pn...
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.
```

Hide

```
condition <- Conditions_Contributing_to_COVID_19_Deaths_by_State_and_Age_Provisional_2020_2022
population <- poppulation_by_state
```

## Guiding Question

It's been several years since the pandemic hit the world. What kind of groups in US were impacted most from Covid-19?

## Inspecting Data

Hide

```
head(primary)
```

| Data As Of <chr> | Start Date <chr> | End Date <chr> | Group <chr> | Y… <dbl> | Mo… <dbl> | State <chr> | Sex <chr> | Age Group <chr> |
|---|---|---|---|---|---|---|---|---|
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | All Ages |
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | Under 1 year |
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | 0-17 years |
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | 1-4 years |
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | 5-14 years |
| 11/16/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | All Sexes | 15-24 years |

6 rows | 1-9 of 16 columns

Hide

```
str(primary)
```

```
spec_tbl_df [107,406 × 16] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Data As Of                        : chr [1:107406] "11/16/2022" "11/16/2022" "11/16/202
2" "11/16/2022" ...
 $ Start Date                        : chr [1:107406] "01/01/2020" "01/01/2020" "01/01/202
0" "01/01/2020" ...
 $ End Date                          : chr [1:107406] "11/12/2022" "11/12/2022" "11/12/202
2" "11/12/2022" ...
 $ Group                             : chr [1:107406] "By Total" "By Total" "By Total" "By
Total" ...
 $ Year                              : num [1:107406] NA NA NA NA NA NA NA NA NA NA ...
 $ Month                             : num [1:107406] NA NA NA NA NA NA NA NA NA NA ...
 $ State                             : chr [1:107406] "United States" "United States" "Unit
ed States" "United States" ...
 $ Sex                               : chr [1:107406] "All Sexes" "All Sexes" "All Sexes"
"All Sexes" ...
 $ Age Group                         : chr [1:107406] "All Ages" "Under 1 year" "0-17 year
s" "1-4 years" ...
 $ COVID-19 Deaths                   : num [1:107406] 1069807 381 1360 207 426 ...
 $ Total Deaths                      : num [1:107406] 9579442 55597 98848 10612 16602 ...
 $ Pneumonia Deaths                  : num [1:107406] 991029 745 2060 439 574 ...
 $ Pneumonia and COVID-19 Deaths     : num [1:107406] 541427 63 347 48 123 ...
 $ Influenza Deaths                  : num [1:107406] 12744 44 273 89 112 ...
 $ Pneumonia, Influenza, or COVID-19 Deaths: num [1:107406] 1530307 1105 3332 683 982 ...
 $ Footnote                          : chr [1:107406] NA NA NA NA ...
 - attr(*, "spec")=
  .. cols(
  ..   `Data As Of` = col_character(),
  ..   `Start Date` = col_character(),
  ..   `End Date` = col_character(),
  ..   Group = col_character(),
  ..   Year = col_double(),
  ..   Month = col_double(),
  ..   State = col_character(),
  ..   Sex = col_character(),
  ..   `Age Group` = col_character(),
  ..   `COVID-19 Deaths` = col_double(),
  ..   `Total Deaths` = col_double(),
  ..   `Pneumonia Deaths` = col_double(),
  ..   `Pneumonia and COVID-19 Deaths` = col_double(),
  ..   `Influenza Deaths` = col_double(),
  ..   `Pneumonia, Influenza, or COVID-19 Deaths` = col_double(),
  ..   Footnote = col_character()
  .. )
 - attr(*, "problems")=<externalptr>
```

The data provides the number of death by week-ending date and state. The number of deaths reported in this table is the total number of deaths received and coded as of the date of analysis.
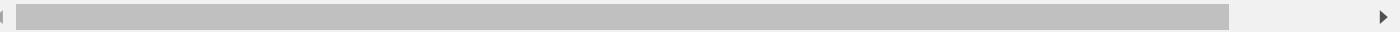
Starting from January 2022, there has been a huge inflated number of deaths from Covid-19 in the US. As a result, the number of deaths in the US has also increased. When we are investigating whether the Covid-patients are decreasing or increasing, we need to watch out for the outliers in the dataset, where there has been a sudden increasing rate of deaths.

```
head(death)
```

| submission_date <chr> | state <chr> | tot_cases <dbl> | conf_cases <dbl> | prob_cases <dbl> | new_c... <dbl> | pnew_c... <dbl> | tot_death <dbl> | cor |
|---|---|---|---|---|---|---|---|---|
| 03/11/2021 | KS | 297229 | 241035 | 56194 | 0 | 0 | 4851 | |
| 12/01/2021 | ND | 163565 | 135705 | 27860 | 589 | 220 | 1907 | |
| 01/02/2022 | AS | 11 | NA | NA | 0 | 0 | 0 | |
| 11/22/2021 | AL | 841461 | 620483 | 220978 | 703 | 357 | 16377 | |
| 05/30/2022 | AK | 251425 | NA | NA | 0 | 0 | 1252 | |
| 05/17/2020 | RMI | 0 | 0 | 0 | 0 | 0 | 0 | |

6 rows | 1-9 of 15 columns

```
str(death)
```

```
spec_tbl_df [60,060 × 15] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ submission_date: chr [1:60060] "03/11/2021" "12/01/2021" "01/02/2022" "11/22/2021" ...
 $ state          : chr [1:60060] "KS" "ND" "AS" "AL" ...
 $ tot_cases      : num [1:60060] 297229 163565 11 841461 251425 ...
 $ conf_cases     : num [1:60060] 241035 135705 NA 620483 NA ...
 $ prob_cases     : num [1:60060] 56194 27860 NA 220978 NA ...
 $ new_case       : num [1:60060] 0 589 0 703 0 ...
 $ pnew_case      : num [1:60060] 0 220 0 357 0 0 NA 274 0 0 ...
 $ tot_death      : num [1:60060] 4851 1907 0 16377 1252 ...
 $ conf_death     : num [1:60060] NA NA NA 12727 NA ...
 $ prob_death     : num [1:60060] NA NA NA 3650 NA 0 NA 429 NA NA ...
 $ new_death      : num [1:60060] 0 9 0 7 0 0 0 8 0 6 ...
 $ pnew_death     : num [1:60060] 0 0 0 3 0 0 NA 3 0 0 ...
 $ created_at     : chr [1:60060] "03/12/2021 03:20:13 PM" "12/02/2021 02:35:20 PM" "01/03/2022 0
3:18:16 PM" "11/22/2021 12:00:00 AM" ...
 $ consent_cases  : chr [1:60060] "Agree" "Agree" NA "Agree" ...
 $ consent_deaths : chr [1:60060] "N/A" "Not agree" NA "Agree" ...
 - attr(*, "spec")=
  .. cols(
  ..   submission_date = col_character(),
  ..   state = col_character(),
  ..   tot_cases = col_double(),
  ..   conf_cases = col_double(),
  ..   prob_cases = col_double(),
  ..   new_case = col_double(),
  ..   pnew_case = col_double(),
  ..   tot_death = col_double(),
  ..   conf_death = col_double(),
  ..   prob_death = col_double(),
  ..   new_death = col_double(),
  ..   pnew_death = col_double(),
  ..   created_at = col_character(),
  ..   consent_cases = col_character(),
  ..   consent_deaths = col_character()
  .. )
 - attr(*, "problems")=<externalptr>
```

The data provides the number of cases and deaths by the state over time in the US. The dataset contains archived aggregate daily counts of COVID-19 cases and death by the state.

There is a US Jurisdiction(Puerto Rico) that is officially not a US state. In order to analyze, we need to eliminate the Jusrsdictions that are not US states.

Hide

```
head(condition)
```

| Data As Of | Start Date | End Date | Group | Y... | Mo... | State | Condition Group |
|---|---|---|---|---|---|---|---|
| <chr> | <chr> | <chr> | <chr> | <dbl> | <dbl> | <chr> | <chr> |
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |

| Data As Of<br><chr> | Start Date<br><chr> | End Date<br><chr> | Group<br><chr> | Y...<br><dbl> | Mo...<br><dbl> | State<br><chr> | Condition Group<br><chr> |
|---|---|---|---|---|---|---|---|
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |
| 11/13/2022 | 01/01/2020 | 11/12/2022 | By Total | NA | NA | United States | Respiratory diseases |

6 rows | 1-8 of 14 columns

Hide

```
str(condition)
```

| Data As Of<br><chr> | Start Date<br><chr> | End Date<br><chr> | Group<br><chr> | Y...<br><dbl> | Mo...<br><dbl> | State<br><chr> | Condition Group<br><chr> |
|---|---|---|---|---|---|---|---|

```
spec_tbl_df [484,380 × 14] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Data As Of        : chr [1:484380] "11/13/2022" "11/13/2022" "11/13/2022" "11/13/2022" ...
 $ Start Date        : chr [1:484380] "01/01/2020" "01/01/2020" "01/01/2020" "01/01/2020" ...
 $ End Date          : chr [1:484380] "11/12/2022" "11/12/2022" "11/12/2022" "11/12/2022" ...
 $ Group             : chr [1:484380] "By Total" "By Total" "By Total" "By Total" ...
 $ Year              : num [1:484380] NA NA NA NA NA NA NA NA NA NA ...
 $ Month             : num [1:484380] NA NA NA NA NA NA NA NA NA NA ...
 $ State             : chr [1:484380] "United States" "United States" "United States" "United Sta
tes" ...
 $ Condition Group   : chr [1:484380] "Respiratory diseases" "Respiratory diseases" "Respiratory
diseases" "Respiratory diseases" ...
 $ Condition         : chr [1:484380] "Influenza and pneumonia" "Influenza and pneumonia" "Influe
nza and pneumonia" "Influenza and pneumonia" ...
 $ ICD10_codes       : chr [1:484380] "J09-J18" "J09-J18" "J09-J18" "J09-J18" ...
 $ Age Group         : chr [1:484380] "0-24" "25-34" "35-44" "45-54" ...
 $ COVID-19 Deaths   : num [1:484380] 1448 5660 14767 36761 80664 ...
 $ Number of Mentions: num [1:484380] 1513 5874 15375 38195 83580 ...
 $ Flag              : chr [1:484380] NA NA NA NA ...
 - attr(*, "spec")=
  .. cols(
  ..   `Data As Of` = col_character(),
  ..   `Start Date` = col_character(),
  ..   `End Date` = col_character(),
  ..   Group = col_character(),
  ..   Year = col_double(),
  ..   Month = col_double(),
  ..   State = col_character(),
  ..   `Condition Group` = col_character(),
  ..   Condition = col_character(),
  ..   ICD10_codes = col_character(),
  ..   `Age Group` = col_character(),
  ..   `COVID-19 Deaths` = col_double(),
  ..   `Number of Mentions` = col_double(),
  ..   Flag = col_character()
  .. )
 - attr(*, "problems")=<externalptr>
```

This data set shows the health conditions and contributing causes of deaths involving COVID-19 by age group and occurrence.

Hide

```
head(population)
```

| State | Total Resident Population | Resident Population Age 18 Years and ( |
|---|---|---|
| <chr> | <dbl> | |
| United States | 331893745 | 25832 |
| Northeast | 57159838 | 4542 |
| Midwest | 68841444 | 5338 |

| State | Total Resident Population | Resident Population Age 18 Years and ( |
|---|---|---|
| <chr> | <dbl> | |
| South | 127225329 | 9845 |
| West | 78667134 | 6106 |
| Alabama | 5039877 | 391 |

6 rows

◄ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬ ►

Hide

```
str(population)
```

```
spec_tbl_df [56 × 3] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ State                                : chr [1:56] "United States" "Northeast" "Midwes
t" "South" ...
 $ Total Resident
Population                   : num [1:56] 3.32e+08 5.72e+07 6.88e+07 1.27e+08 7.87e+07 ...
 $ Resident Population Age 18
 Years and Older: num [1:56] 2.58e+08 4.54e+07 5.34e+07 9.85e+07 6.11e+07 ...
 - attr(*, "spec")=
 .. cols(
 ..    State = col_character(),
 ..    `Total Resident
 ..  Population` = col_number(),
 ..    `Resident Population Age 18
 ..    Years and Older` = col_number()
 .. )
 - attr(*, "problems")=<externalptr>
```

The data provides the total number of resident population by state.

# Clearing data

In order to get clear output, we removed all of each overlapped values and each data set's outlier because there are some overlapped datas from in our data set and the outlier would lead to mess our output.

# 1. The number of COVID patient between 2020/01/01 and 2022/11/12

# Death by age

Hide

```
Primary <-
  primary %>%
  filter(State != "United States")

Primary <-
  Primary %>%
  filter(State != "Puerto Rico")

Primarya <-
  Primary %>%
  filter(`End Date` == "11/12/2022")

Primarya <-
  Primarya %>%
  filter(`Start Date` == "01/01/2020")

PrimaryaR <-
  Primarya %>%
  group_by(Sex,`Age Group`) %>%
  summarise(total = sum(`COVID-19 Deaths`, na.rm = TRUE)) %>%
  filter(Sex %in% c('Male','Female')) %>%
  filter(`Age Group` %in% c("0-17 years", "18-29 years", "30-39 years", "40-49 years", "50-64 yea
rs", "65-74 years", "75-84 years", "85 years and over") )
```

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

Hide

```
PrimaryaR1 <-
  PrimaryaR %>%
  relocate(`Age Group`)

Death_by_Age <-
  PrimaryaR1 %>%
  pivot_wider(names_from = Sex, values_from = total) %>%
  mutate(total = Female + Male)

Death_by_Age
```

| Age Group | Female | Male | total |
|---|---|---|---|
| <chr> | <dbl> | <dbl> | <dbl> |
| 0-17 years | 513 | 592 | 1105 |
| 18-29 years | 2523 | 4034 | 6557 |
| 30-39 years | 7272 | 11792 | 19064 |
| 40-49 years | 16684 | 28085 | 44769 |
| 50-64 years | 74193 | 121424 | 195617 |
| 65-74 years | 98748 | 144156 | 242904 |

| Age Group | Female | Male | total |
| --- | ---: | ---: | ---: |
| <chr> | <dbl> | <dbl> | <dbl> |
| 75-84 years | 121812 | 155486 | 277298 |
| 85 years and over | 157850 | 124226 | 282076 |

8 rows

Overlapping data: "United States, Puerto Rico, All sexes" and also the time period of start and end date are overlapped, so we filter the start date with "01/01/2020" and the end date with "11/12/2022". About age group, there are also overlapped data of age interval, so we filter out the specific age interval: "0-17 years", "18-29 years", "30-39 years", "40-49 years", "50-64 years", "65-74 years", "75-84 years", "85 years and over".

Hide

```
ggplot(Death_by_Age) +
  aes(x = `Age Group`, y = total, fill = `Age Group`) +
  geom_col() +
  scale_fill_hue(direction = 1) +
  labs(title = "Death by Age") +
  theme_minimal()
```



The data shows that people who are over 50 years old have a high percentage of death from COVID-19. Overall, we were able to find out the elder people are more likely to die from COVID-19; especially those over 50 years old, and its number of counts is greatly increased.

# Death by sex

```
Death_by_sex <-
  PrimaryaR1 %>%
  group_by(Sex) %>%
  summarise(Total = sum(total))

Death_by_sex
```
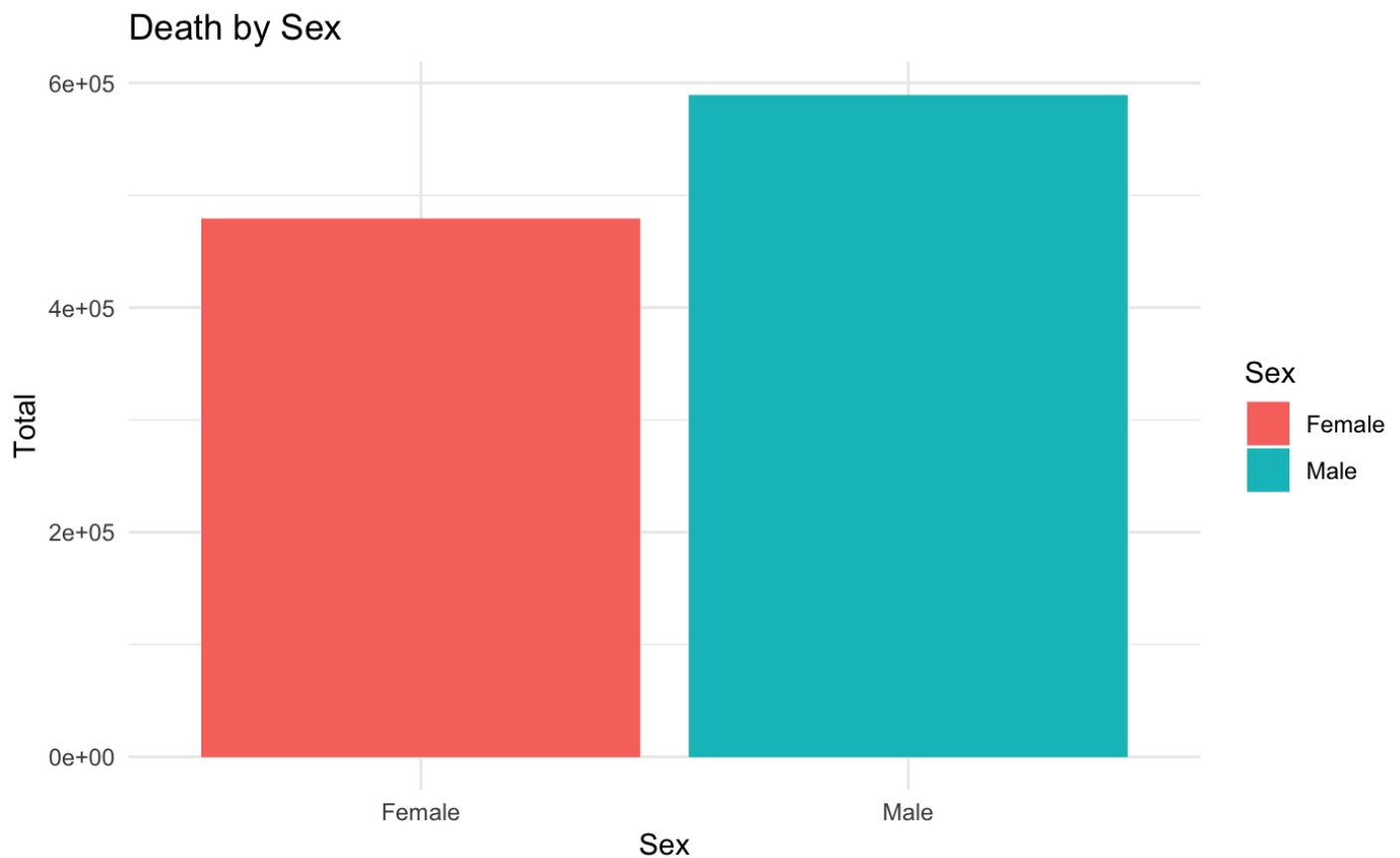
| Sex | Total |
| --- | --- |
| <chr> | <dbl> |
| Female | 479595 |
| Male | 589795 |

2 rows

NA

```
ggplot(Death_by_sex) +
  aes(x = Sex, y = Total, fill = Sex) +
  geom_col() +
  scale_fill_hue(direction = 1) +
  labs(title = "Death by Sex") +
  theme_minimal()
```

## Death by Sex



According to the US census("https://www.statista.com/statistics/737923/us-population-by-gender/#
(https://www.statista.com/statistics/737923/us-population-by-
gender/#):~:text=Projection%20estimates%20calculated%20using%20the,US%20Census%20data%20for%202021."),
they calculated each gender's population, and population ratio is about 49.5% male and 50.47 female. the number
of males and females in the United States is most likely the same. However, the graph shows that a number of
males had died than females. Based on this graph, we may conclude that COVID-19 lead male to die at a higher
percentage compared to female.

# Death by state

Hide

```
Primary3 <-
  Primary2 %>%
  group_by(`State`) %>%
  summarise(total = sum(`COVID-19 Deaths`, na.rm = TRUE)) %>%
  arrange(desc(total))

Primary3
```
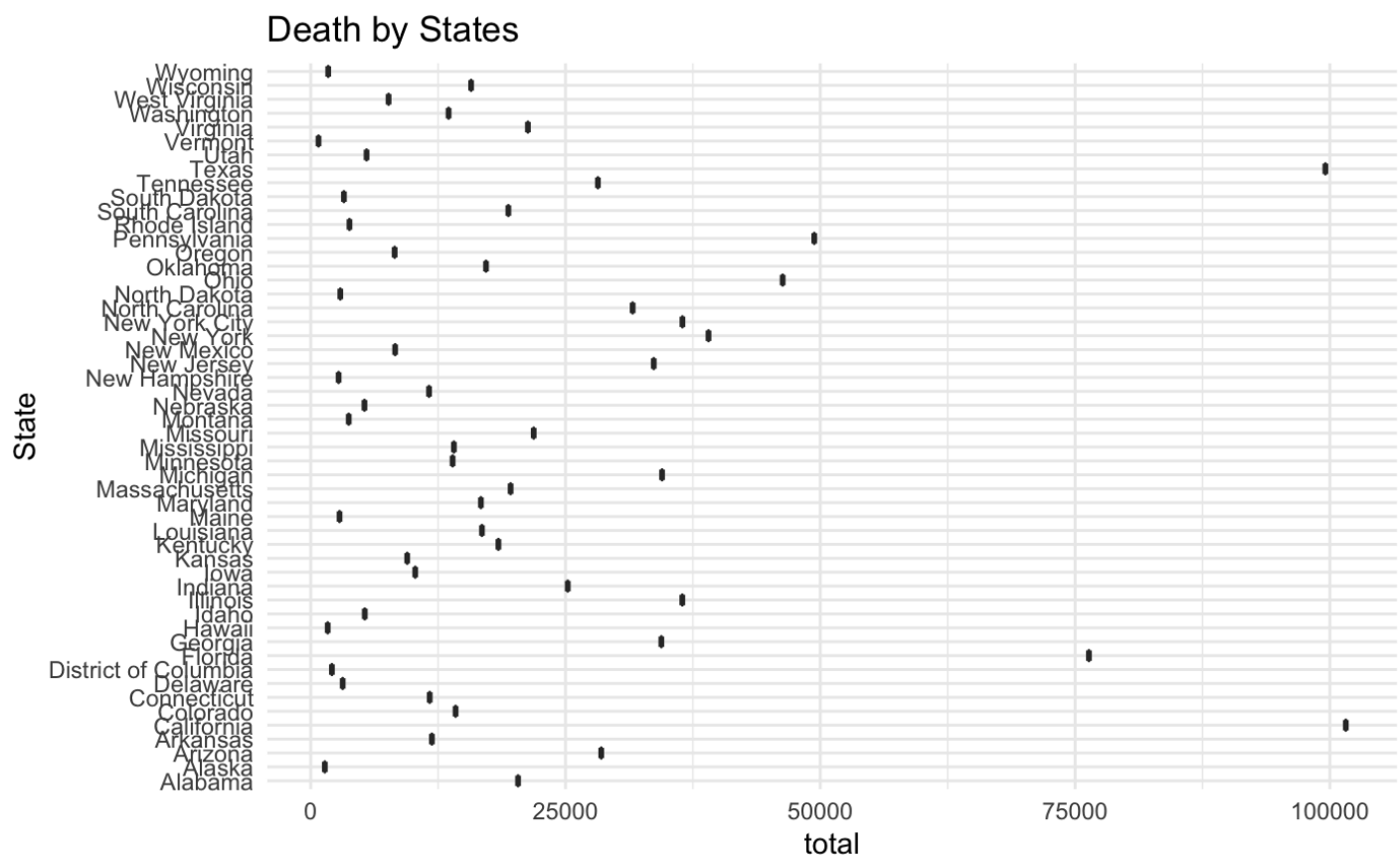
| State | total |
|---|---|
| <chr> | <dbl> |
| California | 101561 |
| Texas | 99555 |
| Florida | 76356 |

| State | total |
|---|---|
| <chr> | <dbl> |
| Pennsylvania | 49408 |
| Ohio | 46306 |
| New York | 39029 |
| New York City | 36472 |
| Illinois | 36460 |
| Michigan | 34472 |
| Georgia | 34403 |

1-10 of 52 rows       Previous  **1**  2  3  4  5  6  Next

Hide

```
ggplot(Primary3) +
  aes(x = total, y = State) +
  geom_boxplot(fill = "#112446") +
  labs(title = "Death by States") +
  theme_minimal()
```



Using the "arrange" function, we easily find out which state has the most number of death from COVID-19. Obviously, the population by the state would affect the number of death by state because the population by the state is mostly proportional to the number of death by the state. The data shows that the major states have a high

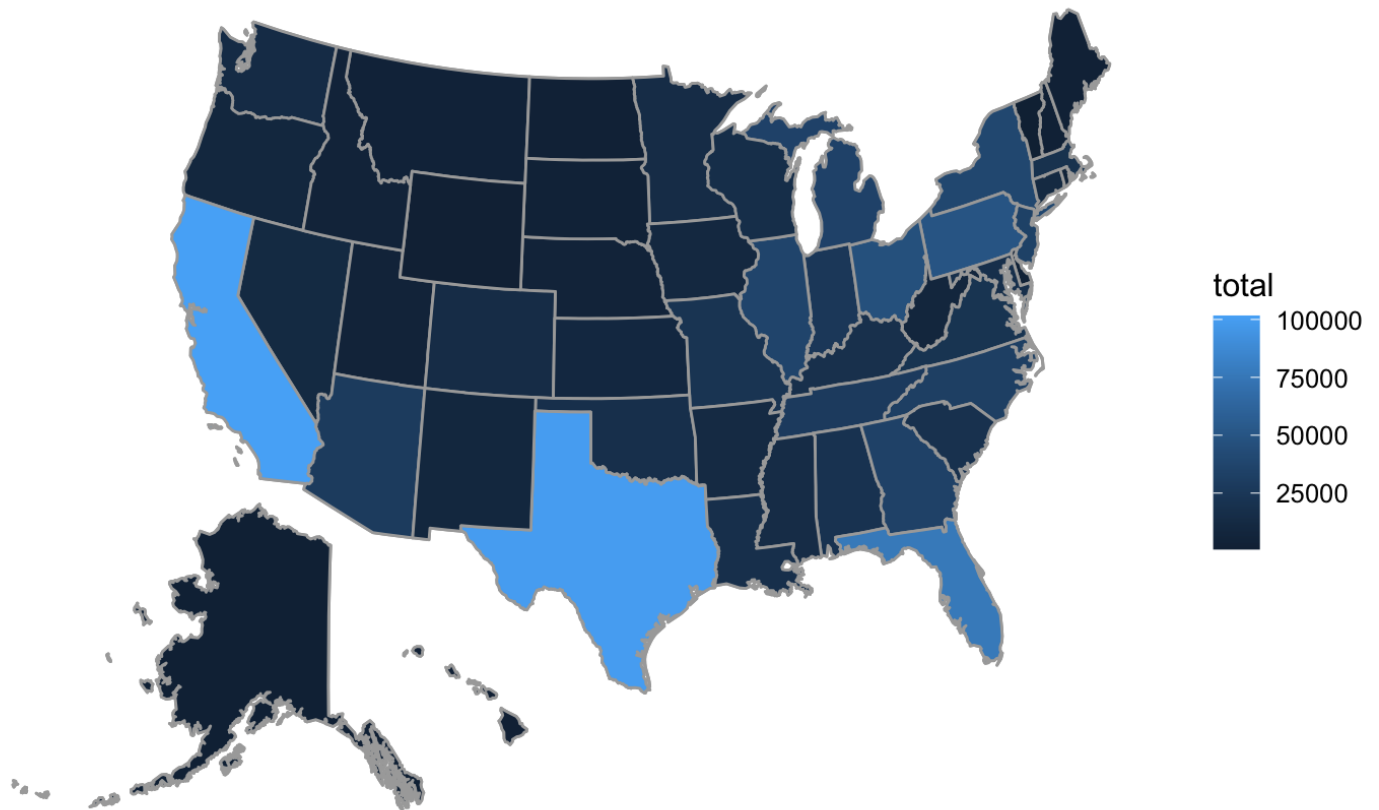number of deaths by COVID-19 considering their dense populations.

# 2. Number of deaths by each States.

## Geographically distribution

```
mUSMap(Primary3, key="State", fill = "total")
```

```
Mapping API still under development and may change in future releases.
Warning: 1 items were not translated
```



The first U.S. map illustrates the distribution of the number of death by each state. The obvious fact is that the number of death is directly proportional to each state's population. If the color of the state is closer to light blue, the number of death is higher, and if the color of the state is closer to dark blue, the number of death is lower. Based on the above map, we are able to figure out the distribution of the number of death by each state easily and emphasize the big states, such as CA, TX, and FL, are vulnerable condition with COVID-19.

## Deathratio by state

```
DeathRatioState <-
  Primary3 %>%
  inner_join(pop, by = c("State" = "table with row headers in column A and column headers in rows
3 through 5. (leading dots indicate sub-parts)")) %>%
  arrange(desc(total))
DeathRatioState
```

| State | total | ...2 | ...3 | ...4 |
|-------|------:|------|------|------|
| <chr> | <dbl> | <chr> | <chr> | <lgl> |
| California | 101561 | 39,237,836 | 30,465,205 | NA |
| Texas | 99555 | 29,527,941 | 22,052,508 | NA |
| Florida | 76356 | 21,781,128 | 17,491,848 | NA |
| Pennsylvania | 49408 | 12,964,056 | 10,290,047 | NA |
| Ohio | 46306 | 11,780,017 | 9,174,388 | NA |
| New York | 39029 | 19,835,913 | 15,722,590 | NA |
| Illinois | 36460 | 12,671,469 | 9,868,245 | NA |
| Michigan | 34472 | 10,050,811 | 7,897,432 | NA |
| Georgia | 34403 | 10,799,566 | 8,275,264 | NA |
| New Jersey | 33657 | 9,267,130 | 7,244,002 | NA |

1-10 of 51 rows                                   Previous  **1**  2  3  4  5  6  Next

# Geographically Distribution of Death Ratio by State

Hide

```
DeathRatioState <-
  Primary3 %>%
  inner_join(population)
```
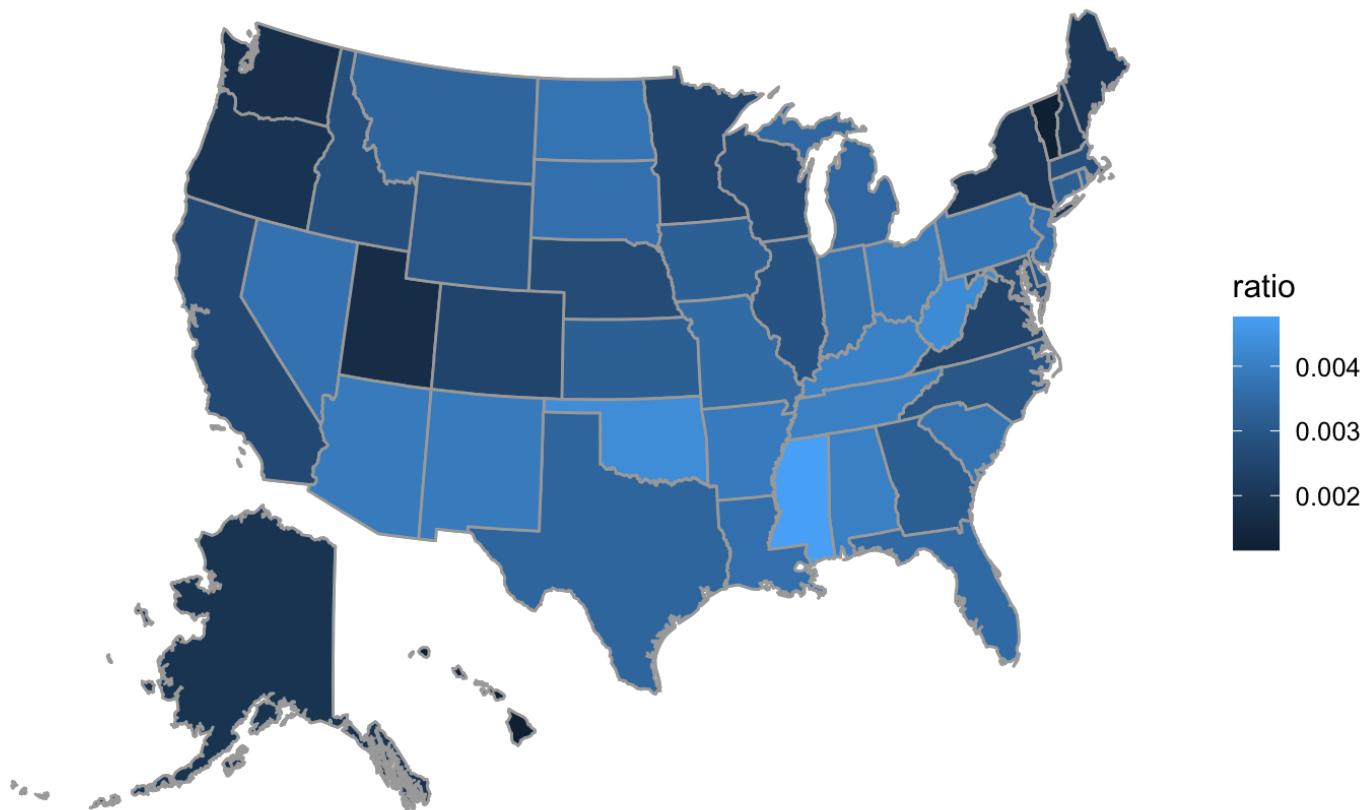
```
Joining, by = "State"
```

Hide

```
DeathRatioState1 <-
  DeathRatioState %>%
  mutate(ratio = total / `Total Resident
 Population`)

mUSMap(DeathRatioState1, key = "State", fill = "ratio")
```

```
Mapping API still under development and may change in future releases.
```

```
DeathRatioState1 %>%
  select('State', 'ratio') %>%
  arrange(desc(ratio))
```

| State | ratio |
|---|---:|
| <chr> | <dbl> |
| Mississippi | 0.004762090 |
| Oklahoma | 0.004314662 |
| West Virginia | 0.004283890 |
| Kentucky | 0.004083697 |
| Tennessee | 0.004039587 |
| Alabama | 0.004036606 |
| Ohio | 0.003930894 |
| Arkansas | 0.003929421 |
| Arizona | 0.003917367 |
| New Mexico | 0.003914216 |

1-10 of 51 rows                                    Previous  **1**  2  3  4  5  6  Next

The above map illustrates the geographically distribution of death ratio by state. If the color of the state is closer to light blue, the death ratio is higher, and if the color of the state is closer to dark blue, the death ratio is lower. The purpose of illustrating the map is to express which state has the most number of death ratio that would decide where is the most vulnerable condition with COVID-19 because each state's death ratio is calculated from the number of death by state divided by population of state. Surprisingly, contrary to expectations, MS, OK, WV, and the other small states have the high death ratio. Based on this table, perhaps, the small states have more vulnerable condition with COVID-19 than big states.

## Death by body condition

<div style="text-align:right">Hide</div>

```
condition1 <-
  condition %>%
  filter(`Age Group` == 'All Ages') %>%
  filter(State == 'United States') %>%
  filter(row_number() <= 23)
condition2 <-
  condition1 %>%
  group_by(`Condition`) %>%
  summarise(total = sum(`COVID-19 Deaths`, na.rm = TRUE))%>%
  arrange(desc(total))
condition2
```

| Condition<br><chr> | tota<br><dbl |
|---|---:|
| COVID-19 | 106805 |
| Influenza and pneumonia | 50638 |
| All other conditions and causes (residual) | 43624 |
| Respiratory failure | 41203 |
| Hypertensive diseases | 19758 |
| Diabetes | 15923 |
| Cardiac arrest | 12662 |
| Renal failure | 12042 |
| Ischemic heart disease | 11661 |
| Sepsis | 10995 |

1-10 of 23 rows                                   Previous  **1**  2  3  Next

<div style="text-align:right">Hide</div>

```
#filter(`Condition Group` != 'All other conditions and causes (residual)') %>%
#filter(`Condition Group` != 'Alzheimer disease') %>%
```

The table illustrates the number of death by the condition in descending order. The reason why the conditions besides COVID-19 is higher than the number of COVID-19 is there are people who have multiple diseases; there is no guarantee people have only one disease. The table shows who have influenza and pneumonia are the most contributing to Death by COVID-19.

# Which conditions(age, sex, state, body's condition, ratio) are mostly lead to die through COVID-19

Hide

Primary1R

| Age Group<br><chr> | total<br><dbl> |
|---|---|
| 0-17 years | 1282 |
| 18-29 years | 6653 |
| 30-39 years | 19086 |
| 40-49 years | 44784 |
| 50-64 years | 195617 |
| 65-74 years | 242904 |
| 75-84 years | 277298 |
| 85 years and over | 282076 |

8 rows

Hide

Primary2R

| Sex<br><chr> | total<br><dbl> |
|---|---|
| Female | 479595 |
| Male | 589795 |

2 rows

Hide

Primary3

| State<br><chr> | total<br><dbl> |
|---|---|
| California | 101561 |

| State | total |
|-------|------:|
| <chr> | <dbl> |
| Texas | 99555 |
| Florida | 76356 |
| Pennsylvania | 49408 |
| Ohio | 46306 |
| New York | 39029 |
| New York City | 36472 |
| Illinois | 36460 |
| Michigan | 34472 |
| Georgia | 34403 |

1-10 of 52 rows                    Previous  **1**  2  3  4  5  6  Next

Hide

```
condition2 %>%
  arrange(desc(total))
```

| Condition | total |
|-----------|------:|
| <chr> | <dbl |
| COVID-19 | 106805 |
| Influenza and pneumonia | 50638 |
| All other conditions and causes (residual) | 43624 |
| Respiratory failure | 41203 |
| Hypertensive diseases | 19758 |
| Diabetes | 15923 |
| Cardiac arrest | 12662 |
| Renal failure | 12042 |
| Ischemic heart disease | 11661 |
| Sepsis | 10995 |

1-10 of 23 rows                    Previous  **1**  2  3  Next

Hide

```
DeathRatioState1 %>%
  select('State', 'ratio') %>%
  arrange(desc(ratio))
```

| State | ratio |
| :--- | ---: |
| <chr> | <dbl> |
| Mississippi | 0.004762090 |
| Oklahoma | 0.004314662 |
| West Virginia | 0.004283890 |
| Kentucky | 0.004083697 |
| Tennessee | 0.004039587 |
| Alabama | 0.004036606 |
| Ohio | 0.003930894 |
| Arkansas | 0.003929421 |
| Arizona | 0.003917367 |
| New Mexico | 0.003914216 |

1-10 of 51 rows    Previous **1** 2 3 4 5 6 Next

Hide

NA

# Conclusion

It is been several years since the pandemic hit the world. There were numerous casualties and infections around our life and the entire globe. From this motivation, Our research question is "What kind of groups in the US were impacted most by Covid-19?" According to these data outcomes, we are able to find out, in United states, the condition of '85 years and over' or 'male','who has influenza and pneumonia' or 'who lives in Mississippi' would have been mostly leading to death. Additionally, based on this evidence from the data sets, we are able to emphasize that COVID-19 is one of the greatest prevalent diseases in any disease that humanity has experienced. For now, we may feel COVID-19 is a common disease, such as influenza. Based from our research, we learned that we still need a steady concern about Covid-19.