

Ananlysis of COVID Data

Ryan Greenup

May 21, 2020

Contents

Preliminary	1
Load Packages and Data	1
Load the Data	1
Set Working Directory	2
Introduction	2
Chloropleth Map	2
Discussion	2
.1 Worldwide	2
.2 Europe	3
Technique	3
.1 Woldwide Map	3
.2 Europe Centric	4
Time Series	6
Technical Details	8
Advantages compared to other methods	8
Disasadvantages	9
Discussion on analysis results	9
Discussion on other Aspects	9
Literature review of related work	9
Bar Chart	9
Pie Chart	9
Spider Chart / Star Plot	9
Multiple Line Charts	9
Parallell Co-ordinates	9
3D Scatter Plot	9
Log Scaled from 100th case	ATTACH 10

Bubble Plot	ATTACH 10
Animation of 3d Chloropleth heatmap	10
Technical Details	11
Advantages compared to other methods	11
Disasadvantages	11
Discussion on analysis results	11
Discussion on other Aspects	11
Literature review of related work	11
For Each Visualisation	11
Technical Details	11
Advantages compared to other methods	11
Disasadvantages	11
Discussion on analysis results	11
Discussion on other Aspects	11
Literature review of related work	11
Apendix	ATTACH 11
References	11

Preliminary

Load Packages and Data

```

1  if (require("pacman")) {
2    library(pacman)
3  }else{
4    install.packages("pacman")
5    library(pacman)
6  }
7  pacman::p_load(xts, sp, gstat, ggplot2, rmarkdown, reshape2, ggmap,
8                 parallel, dplyr, plotly, tidyverse, reticulate,
9                 ↪ UsingR, Rmpfr,
10                 swirl, corrplot, gridExtra, mise, latex2exp,
11                 ↪ tidyverse, xts, maptools, plyr, ggplot2, maps,
12                 ↪ viridis)
13  mise()

```

Load the Data

```
1 covid <- read.csv("/home/ryan/Notes/DataSci/Visual_Analytics/Assessment_1  
  ↪ 2/owid-covid-data.csv")
```

Set Working Directory

Introduction

- in December 2012 first cases of *COVID-19* were reported, the disease has since attributed to the *SARS-CoV2* virus.
- The disease became endemic throughout China before spreading throughout Europe in an epidemic fashion and finally reaching the rest of the globe as a pandemic outbreak.

Chloropleth Map

A Chloropleth map of the number of deaths can offer an insight into the impact that the disease has had with respect to individual countries.

The Total deaths should be scaled relative to the population of the country, that way countries with a smaller and sparser population will still be represented by the visualisation (this is quite important given that many countries such as Italy have a small population compared to the US and much of Asia [2020n]).

A worldwide Chloropleth map visualising the total number of deaths attributed to *COVID-19* is shown in figure 1 and a Europe-centric visualisation is shown in figure 2.

Discussion

Worldwide

The first plot appears to show a very limited amount of difference in deaths attributable to *COVID-19* across regions other than the North America and Europe. While first-world countries such as New Zealand and Australia are somewhat insulated from the disease by virtue of geography and population density, it's striking that much of Asia and Russia have such low levels of outbreak.

This could be attributed to the fact that a more power-centric regime such as in China, Russia, North Korea, etc. may have more capacity to:

1. Diminish the spread of the disease by implementing policy decisions,
 - (a) whereas countries such as the US and Europe have a much higher expectation of civil liberties.
2. Control the spread of information for want of international reputation.

- (a) In saying that though research suggests that under-reporting has even occurred in countries such as the US [sood2020] so such under-reporting could merely be incidental.

A similar disease, *MERS*, emerged in 2012 in Middle-Eastern Regions [woodley2020] and a Korean outbreak of the *MERS* disease occurred in 2015 [serrano2015], these outbreaks likely prepared Korea, the Middle East and other Asian Regions regions for an outbreak which helps explain the dichotomous nature of the deaths attributable to *COVID-19* for those Countries.

Europe

A closer look at Europe shows that Belgium and Italy have been the most affected by this disease, it isn't very clear why those regions have been impacted so significantly but this could be indicative of policy decisions and warrants further research.

Technique

Worldwide Map

First the data must be aggregated in order to retrieve the total number of deaths, this can be achieved by taking the maximum of the total deaths across countries (the total number of death rates will be a strictly positive and monotone trend, otherwise the outbreak would be an entirely different type of pandemic!), this can be performed by using the `aggregate` function as demonstrated in figure 1.

```
1 fatalprop <- aggregate(total_deaths_per_million ~ location, covid, max)
2 ## Order the Values in Descending Order
3 fatalprop <- fatalprop[order(-fatalprop$total_deaths_per_million),]
4 ## Rename USA
5 covid$location[covid$location=="United States"] <- "USA"
```

Listing 1: Use Aggregate to aggregate total number of deaths

It is next necessary to rename `location` to `region` so map data will be consistent with the provided data set, this is shown in listing 2.

```
1 ## Rename to facilitate joining with map
2 names(fatalprop) <- c("region", "total_deaths_per_million")
```

Listing 2: Rename Features for consistency

For a broad overview of the data, small regions such as San Marino and Belgium will not be visible and will skew the colour palette, so instead they should be removed and instead a separate plot of Europe will be created as shown in figure 2, this removal is performed in listing 3.

Next it is necessary to retrieve map data, this can be done using the `map_data` function, this data may then be combined by region with the provided data set using the `left_join` function, this is shown in listing 4.

```

1  ## San Marino will be shown by Italy and this skews the results
2  ## Belgium and San Marino are very hard to visualise from above
3  ## They skew the results and so will be removed.
4  fatalprops <- fatalprop %>% filter(region!="San Marino")
5  fatalprops <- fatalprop %>% filter(region!="Belgium")

```

Listing 3: Filter out small dense regions to prevent scale issues

```

1  ## Retrieve the map data
2  some.eu.maps <- map_data("world", region = fatalprops$location)
3
4  ## Join the Data Frames Together
5  fatalmap <- left_join(fatalprops, some.eu.maps, by = "region")

```

Listing 4: Combine Map Data with Provided Data

Finally this data frame can be plotted by using `ggplot2` and the `geom_map` layer, modifying the theme layer will allow to provide a natural background, this is demonstrated in listing 5 and the output is provided in figure 1.

Europe Centric

The choropleth map clearly shows that the disease has caused more fatalities per capita in Europe and so the plot will be adjusted central to Europe.

As before it is necessary to rename the features of the dataset, however in this instance small European countries such as Belgium should be retained (San Marino is a very small Italian province that isn't detectable in the visualisation and skews the palette, for this reason it will be removed), this is demonstrated in figure 6

In this map it will be desirable to have labels for the European countries (whereas this would have made the worldwide map too busy), so this will be implemented by using `dyplr` to generate a second data set as shown in listing 7 which can then be used to generate a plot with the `ggrepel` add on as shown in listing 8, this produces the output shown in figure 2, for this plot bubbles were also implemented in order to help visualise the number of relative cases. The inspiration for the use of bubbles was the *John Hopkins Coronavirus Dashboard [2020a]* where a similar strategy was implemented to visualise the number of cases, a screenshot of this is provided in the appendix at figure 3.

Time Series

The spread of disease over time can often be modelled by an exponential model as demonstrated in equations (1) and (2), for this reason the use of a \log -scale will linearise trends and so the use of a \log -scale will make it easier to compare the rates of population change between different countries.

```

1  ggplot(fatalmap, aes(map_id = region)) +
2  geom_map(map = fatalmap, color = "grey", aes(fill =
   ↪ total_deaths_per_million), lwd = 0.1, alpha = 0.6)+
3  expand_limits(x = fatalmap$long, y = fatalmap$lat)+
4  scale_fill_gradient(high = "darkred", low = "white") +
5  guides(fill = guide_legend("Total Deaths \n per Million")) +
6  # Change the colors of background
7  # and the color of grid lines to white
8  theme(
9    panel.background = element_rect(fill = "lightblue",
10                                     colour = "lightblue",
11                                     size = 0.5, linetype = "solid"),
12    legend.position = c(0.6, 0.1),
13    legend.direction = "horizontal",
14    legend.background = element_rect(fill = "white", size = 0.1,
   ↪ colour = "darkblue", linetype = "solid")) +
15  labs(x = "Longitude", y = "Latitude", title = TeX("Total Deaths
   ↪ Attributed to \\textit{COVID-19}"))
16 # geom_text(data = region_lab_df, aes(y = lat, x = long, label =
   ↪ region), size = 1)

```

Listing 5: use ggplot2 to create a choropleth map from data, output in figure 1

```

1  ## Rename to facilitate joining with map
2  names(fatalprop) <- c("region", "total_deaths_per_million")
3
4  ## San Marino will be shown by italy
5  fatalprop <- fatalprop %>% filter(region!="San Marino")

```

Listing 6: Rename the features of the data and remove San Marino

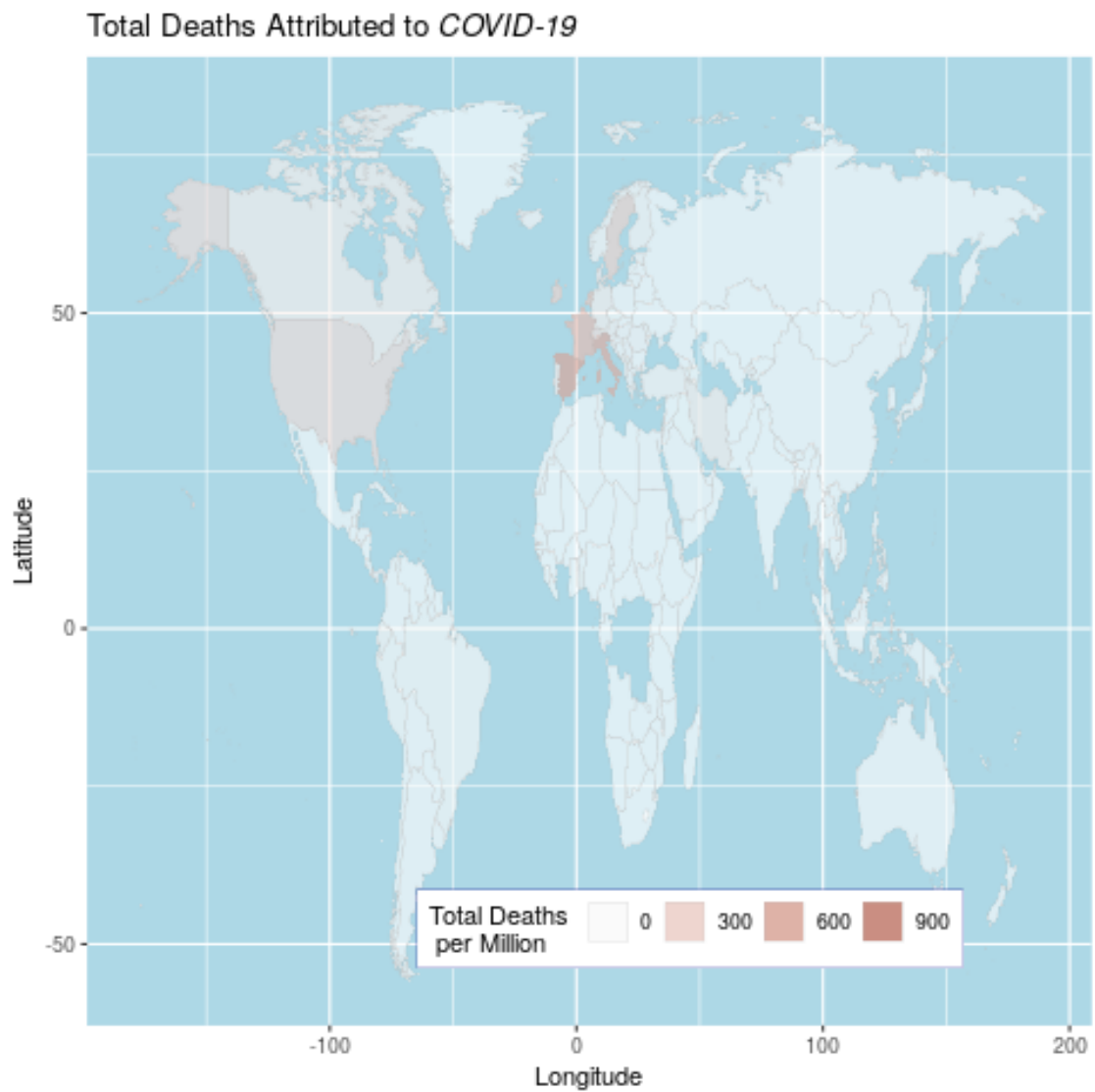


Figure 1: Choropleth map of total deaths attributed to *COVID-19* (per Million people)


```

1 library(ggrepel)
2 ggplot(fatalmap, aes(map_id = region, label = region)) +
3   geom_map(map = fatalmap,
4             aes(fill = total_deaths_per_million),
5             color = "white") +
6   geom_point(data = region_lab_df, aes(y = lat, x = long, size =
7     ↪ total_deaths_per_million), alpha = 0.45, colour = "blue", stroke
8     ↪ = 1, fill = "white", shape = 21) + scale_size_continuous(range =
9     ↪ c(1, 25), name = "Total Number \n of Deaths") +
10  guides(size = FALSE) +
11  expand_limits(x = fatalmap$long, y = fatalmap$lat) +
12  scale_fill_viridis_c(option = "C") +
13  scale_fill_gradient(high = "darkred", low = "white") +
14  guides(fill = guide_legend("Total Deaths \n per Million")) +
15  # Change the colors of plot panel background to lightblue
16  # and the color of grid lines to white
17  theme(
18    panel.background = element_rect(
19      fill = "lightblue",
20      colour = "lightblue",
21      size = 0.5,
22      linetype = "solid"
23    ),
24    legend.position = c(0.1, 0.6),
25    legend.direction = "vertical",
26    legend.background = element_rect(
27      fill = "white",
28      size =
29        1.1,
30        colour = "darkblue",
31        linetype = "solid"
32      )
33  ) +
34  labs(
35    x = "Longitude",
36    y = "Latitude",
37    title = TeX("Total Deaths Attributed to \\textit{COVID-19}")
38  ) +
39  geom_text_repel(
40    data = region_lab_df,
41    aes(y = lat, x = long, label = region),
42    size = 2,
43    col = "black",
44    nudge_y = 0.7,
45    nudge_x = -0.5,
46    min.segment.length = 0.6,
47    force = 2
48  )

```

$$\frac{dp}{dt} \propto p \implies p = Ce^{kt} \quad \exists k, c \in \mathbb{R} \quad (1)$$

$$\frac{dp}{dt} \propto p \wedge \frac{dp}{dt} \propto (N - p) \implies p = \frac{ke^{Nt}}{1 - ke^{Nt}} \quad \exists k \in \mathbb{R}, N \in \mathbb{R}^+ \quad (2)$$

In addition to a log – scale, scaling the data to be relative to the number of days since the first case can allow the trends of the data to be compared, this was implemented by *John Hopkins University* in a visualisation published in the *Guardian* [gutierrez2020].

Technical Details

Advantages compared to other methods

- The advantage to a log-scaled plot is that it allows rates of change to be compared between countries
- Making the Data Relative to the day of the first infection allows individual countries to be compared in terms of there response

Disasadvantages

- A log-scaled plot can be misleading if it is not made clear, his particularly true for readers who have limited mathematical training.
 - For this reason a plot without log-scaling was included and the axis were labelled accordingly
- Making Data relative to the day of the first infection may not make clear that certain countries had /forewarning of the disease by virtue of the delay.

Discussion on analysis results

Discussion on other Aspects

Literature review of related work

As mentioned in section the use of the log-scaled and date-adjusted plot was implemented by *John Hopkins University* in a visualisation published in *The Guardian* newspaper [gutierrez2020].

Bar Chart

Pie Chart

Spider Chart / Star Plot

Multiple Line Charts

Parallell Co-ordinates

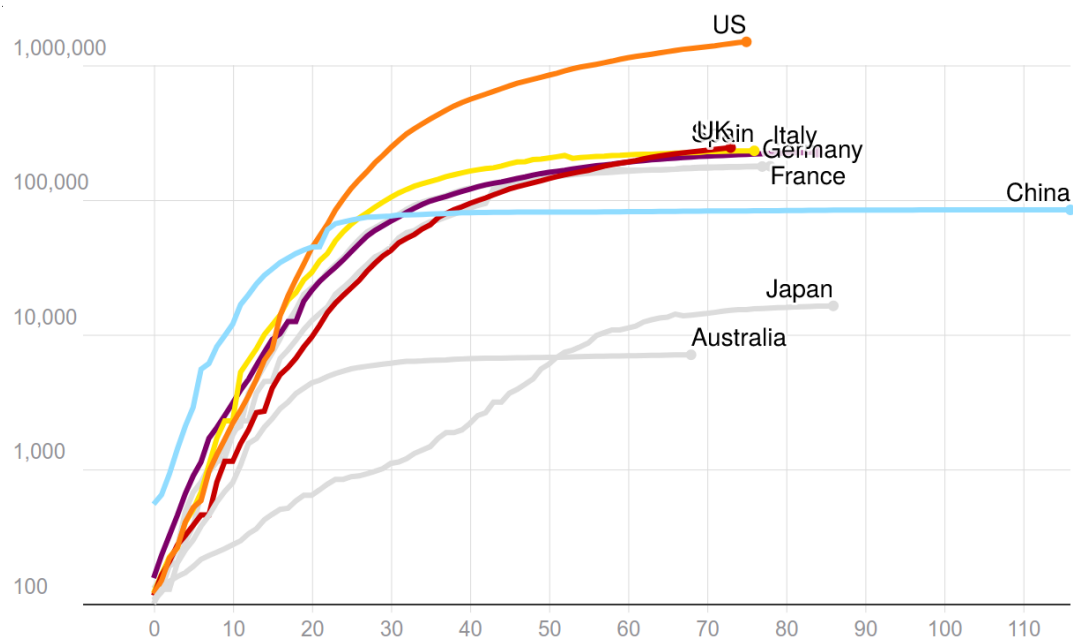
each line is a country each column is a feature like testing, death and cases.

[This Stack Post](#) shows how to make them curvy

3D Scatter Plot

Log Scaled from 100th case

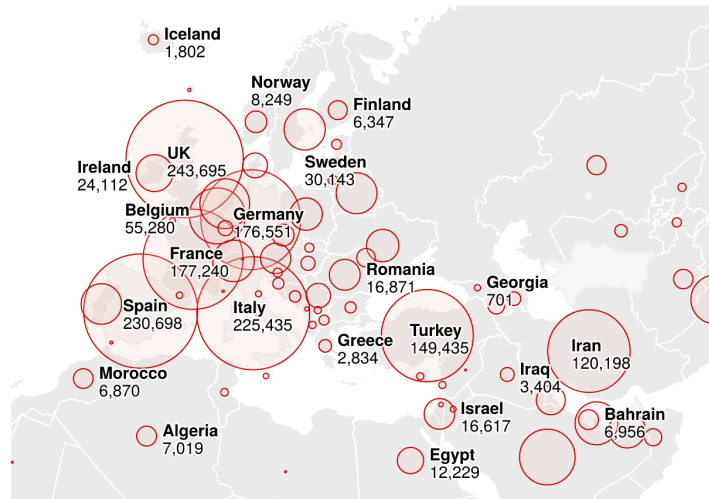
ATTACH



Bubble Plot

ATTACH

[Guardian](#)



Animation of 3d Chloropleth heatmap

visualisation

The total number of deaths per country can be analysed using

Technical Details

Advantages compared to other methods

Disasadvantages

Discussion on analysis results

Discussion on other Aspects

Literature review of related work

For Each Visualisation

Technical Details

Advantages compared to other methods

Disasadvantages

Discussion on analysis results

Discussion on other Aspects

Literature review of related work

Apendix

ATTACH



Figure 3: John Hopkins Bubble Chart [2020o]

References

../../../../Studies/Papers/references