

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
Advantages compared to other methods	9
Disasadvantages	9
Literature review of related work	10
Time Series	10
Implementation	10
.1 Log Scale	10
.2 Adjust Zero	10
Technical Details	11
.1 Preliminary	11
.2 Facet Grid	11
Advantages compared to other methods	11
Disasadvantages	21
Discussion on analysis results	21
Discussion on other Aspects	21
Literature review of related work	21

TODO Parallel Co-ordinates	21
Technical Details	23
Advantages compared to other methods	23
Disasadvantages	23
Discussion on analysis results	23
Discussion on other Aspects	23
Literature review of related work	23
For Each Visualisation	23
Technical Details	23
Advantages compared to other methods	23
Disasadvantages	23
Discussion on analysis results	23
Discussion on other Aspects	23
Literature review of related work	23
Appendix	ATTACH 23
References	26

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

On December 31st 2019 a viral pneumonia was reported in Wuhan, China, this was later found to be a result of a new strain of virus named *Sars-CoV2*, the disease caused by such an infection, usually resulting in viral pneumonia, is known as *Corona Virus Disease 2019 (COVID-19)*. The outbreak of this disease was declared a Public Health Emergency of International Concern on the 30th January 2020. [worldhealthorganization2020] December 2012 first cases of *COVID-19* were reported, the disease has since attributed to the *SARS-CoV2* virus.

A data set detailing the location, deaths, tests and cases related to the *COVID-19* pandemic has been made available through the website /Our World in Data/[ritchiet2020], documented in this report is a visual anylisis performed entirely using the *Free Software*¹ **R** [rcoreteam2020] primarily with the ggplot2 package [wickham2016] (see listing 20 in the appendix).

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 3.

Discussion

Worldwide

The first plot appears to show a very limited amount of difference in deaths attributable to *COVID-19* across regions other than North America and Europe.

¹Free as in Speech and beer

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 disease incidence.

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 and hence much lower tolerance for government intervention.
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 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, particularly considering the comparatively permissive borders within the *EU*, but this could be indicative of policy decisions and warrants further research.

Technique

Worldwide Map

In order to produce a choropleth map 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 listing 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 3, this removal is performed in listing 3.

```

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

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  ## Retrieve the map data
2  some_maps <- map_data("world", region = fatalprops$location)
3
4  ## Join the Data Frames Together
5  fatalmap <- left_join(fatalprops, some_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 for a natural background to be implemented, this is demonstrated in listing 5 and the output is provided in figure 1.

A bubble overlay may also be implemented in order to make clearer the spread of cases (see section for a brief literature review), it is necessary however to adjust the *USA* location to represent the mainland population centre in order to make the visualisation more effective. This is demonstrated in listing 6 and shown in figure 2

Europe Centric

The choropleth map clearly shows that the disease has caused significantly 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

```

1 wmp <- 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)
17 wmp

```

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

```

1 # Compute the centroid as the mean longitude and latitude
2 # Used as label coordinate for country's names
3 region_lab_df <- some.eu.maps %>%
4   group_by(region) %>%
5   summarise(long = mean(long), lat = mean(lat)) %>%
6   full_join(aggregate(total_deaths_per_million ~ region, fatalmap,
   ↪ mean))
7 # Manually Adjust US to be population Centre
8 region_lab_df[region_lab_df$region == "USA",]$long <- -92.47
9 region_lab_df[region_lab_df$region == "USA",]$lat <- 37.37
10
11
12 wmp +
13   scale_size_continuous(range = c(1, 9), name = "Total Number \n of
   ↪ Deaths") +
14   guides(size = FALSE) +
15   geom_point(data = region_lab_df, aes(y = lat, x = long, size =
   ↪ total_deaths_per_million), alpha = 0.5, col = "purple")

```

Listing 6: use ggplot2 to create a choropleth map from data, output in figure 1

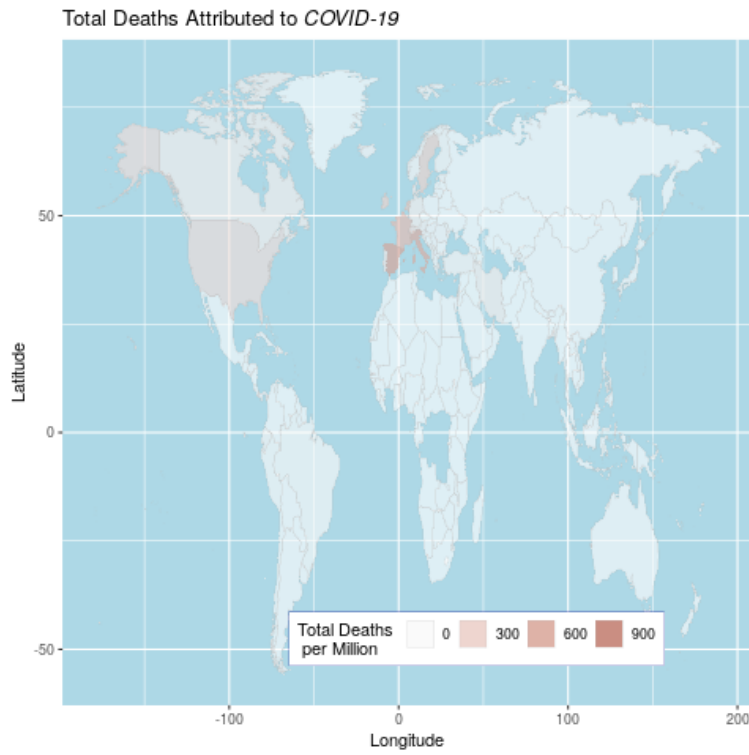


Figure 1: Choropleth map of total deaths attributed to *COVID-19* (per Million people)

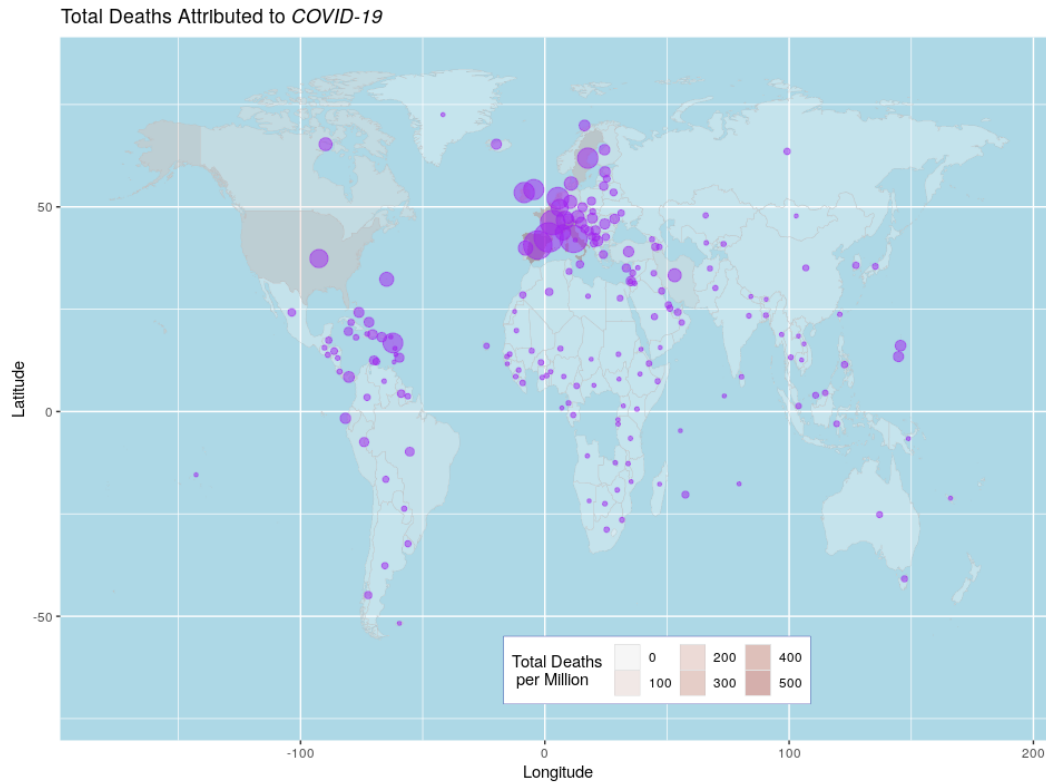


Figure 2: Choropleth map with bubble overlay to aid in case visualisation

isn't detectable in the visualisation and skews the palette, for this reason it will be removed), this is demonstrated in listing 7

```
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 7: Rename the features of the data and remove San Marino

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 dplyr to generate a second data set as shown in listing 8 which can then be used to generate a plot with the ggmap add on as shown in listing 9, this produces the output shown in figure 3, bubbles were also implemented in order to help visualise the number of relative cases.

```
1  fatalmap <- left_join(fatalprop, some.eu.maps, by = "region")
2
3  ## Filter out only Europe
4  fatalmap <- fatalmap %>%
5    filter(30 < lat & lat < 65) %>%
6    filter(-30 < long & long < 35)
7
8  ## Create Label Data Frame
9  region_lab_df <- fatalmap %>%
10    dplyr::group_by(region) %>%
11    dplyr::summarise(long = mean(long), lat = mean(lat)) %>%
12    full_join(aggregate(total_deaths_per_million ~ region, fatalmap,
13      ↪ mean))
```

Listing 8: use dplyr to reduce the plot size and create a data frame of country labels

Advantages compared to other methods

- A Choropleth map provides a very clear way to visualise the occurrence of disease in a geographical sense, in contrast to other methods such as scatter plots, heatmaps and bar charts, the choropleth map provides a clear way to distinguish the impact of the disease on individual countries.

The discrete distinction between countries, a fundamental component of a choropleth map, is desirable because it is consistent with the independent legislatures across countries, this allows for a comparison of the impact that policy decisions may or may not have on a region.


```

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  )

```

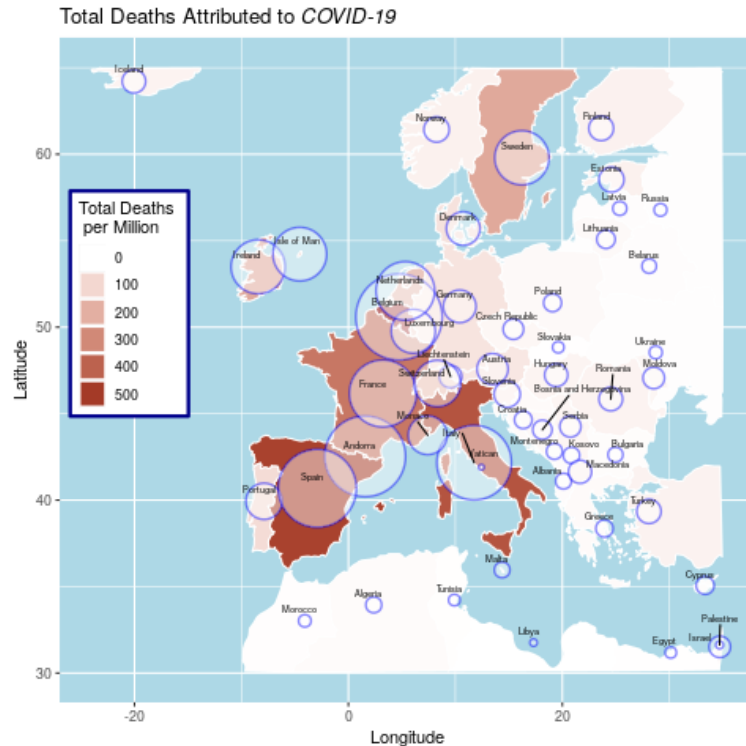


Figure 3: Europe Centred Chloropleth of Deaths Attributed to *COVID-19*

Chloropleth maps also allow trends across regions to be easily identified, e.g. figure 3 shows how severe the outbreak is in *Europe* relative to other regions, this might be lost in abstraction when using other visualization methods.

Disadvantages

When maps are projected into a 2D plane they are necessarily distorted, this distortion can impact how spread the data appears to be.

A chloropleth map can make it hard to compare metrics between to regions in any specific sense, for this a more appropriate visualization would be a bar chart.

Literature review of related work

The *John Hopkins Coronavirus Dashboard* [2020a] implemented bubbles to visualise the number of cases, a screenshot of this is provided in the appendix at figure 7, this was a part of the motivation for implementing bubbles in the chloropleth map because the visualization was so much more *striking* and promoted pre-attentive processing of the information.

In his blog, Kenneth Field produced chloropleth and bubble-map charts detailing the spread of *COVID-19*, with however, a focus on China, [field2020] these plots were very similar to those produced in this report, however the legend for the bubble plot was very nicely implemented and can be seen in figure 8 of the appendix. He also produced an example illustrating why the use of a heatmap or contour

map can make for a poor visualisation of cases due to the difficulty in interpreting the visualization compared to a bubble chart, for this reason a bubble chart was used in this report and a heatmap was not implemented.

A paper in the publication *Environment & Planning A* suggested using a cartogram to visualise the spread of disease, there example is provided in figure 9 of the appendix. [gao2020] Although the cartogram is visually quite appealing and easy to read, it is difficult to interpret quickly, the visualisation does not promote pre-attentive processing, for this reason the visualisation strategy was not implemented.

Time Series

Implementation

Time series charts can be an effective way to visualise the behaviour of a value over time, for this dataset however, two modifications will be implemented in order to make the trends more distinct.

Log Scale

The spread of disease over time can often be described 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.

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

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

Adjust Zero

In addition to a \log – scale, *sliding* 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

Preliminary

In order to log scale the data the `mutate` function from the `dplyr` package was used on data transformed into *wide* format by using the `pivot_wider` function, this is shown in listing 10.

Sliding the date back to the number of cases however was a little more difficult and required the use of a `for` loop to iterate the `lead` function over each column (where each column, after transformation with `dplyr`, represented the value for a country), this is demonstrated in listing 10 with an example of the produced *tidy* data provided in table 1; the code to produce the plot is demonstrated in listing 11, the output of which is provided in figure 4.

Rather than using a line plot or a scatter plot, a `loess` model was placed ontop of semi-opaque points, this is to enhance the continuity of the visualisation. The *Gestalt Laws* provide that continuous

shapes are easier for readers to interpret [staudinger2011] and for this reason the the overlay was implemented, to aid the reader in delineating between the different countries in a plot.

Plots with many colours mapped to categorical variables can be difficult to interpret [wilson2017, rost2018], for this reason less than 10 countries were compared on the same plot.

```

1  cv <- as_tibble(covid)
2  cv <- cv %>%
3    mutate(date = as.Date(date))
4  cv <- cv[order(cv$date),]
5
6  # interested_locations <- c("Australia", "USA", "Italy", "Germany",
7    ↪ "Belgium", "United Kingdom", "New Zealand", "Japan", "China")
8  interested_locations <- c("Australia", "USA", "Italy", "Germany",
9    ↪ "Russia", "South Korea", "United Kingdom")
10
11 cv <- cv %>%
12   filter(location %in% interested_locations) %>%
13   filter(total_cases_per_million > 1) %>%
14   mutate(total_cases_per_million = log10(total_cases_per_million)) %>%
15   dplyr::select(date, total_cases_per_million, location) %>%
16   pivot_wider(names_from = location, values_from =
17     ↪ total_cases_per_million)
18
19 for (i in 2:ncol(cv)) {
20   ## Slide the Columns up and put the NA at the end
21   cv[,i] <- pull(cv, i) %>%
22     lead(cv[,i] %>%
23       is.na() %>%
24       sum())
25   ## Replace the date with the number of days
26   cv$date <- seq_len(nrow(cv))
27 }
28
29 cv <- cv %>%
30   pivot_longer(names(cv)[-1], names_to = "location", values_to =
31     ↪ "total_cases_per_million")

```

Listing 10: Use = dplyr= to transform the data as shown in table 1, this can then be passed to ggplot as shown in listing 11

Facet Grid

This plot however does not show all the data made available, the data set also includes information on the number of tests,cases and deaths resulting from *COVID-19*, in order to visualise this the `fact_grid` layer can be used to create a multi-scatterplot. first it is necessary to create a data frame, this can be

Table 1: Top few rows of the *tidy* data set created from listing 10.

<i>Date</i>	<i>Location</i>	<i>Total Cases Per Million</i>
1	South Korea	0.193
1	Italy	0.116
1	Australia	0.00860
1	Germany	0.122
1	United Kingdom	0.0976
1	USA	0.00903
1	Russia	0.00303
2	South Korea	0.480
2	Italy	0.339
2	Australia	0.0558

```

1 ggplot(cv , aes(y = total_cases_per_million, x = date, col = location,
  ↪ group = location)) +
2   geom_point(alpha = 0.3) +
3   geom_smooth() +
4   theme_bw() +
5   labs(y = "Total Number of Cases (Log-10 Scale)", title = "Log Scaled
  ↪ Total COVID-19 Cases per Million", x = TeX("Days since Case
  ↪ \\textit{#100}")) +
6   guides(col = guide_legend("Location"))
7 # geom_smooth()
```

Listing 11: Use *dplyr* to transform the data before plotting with *ggplot*

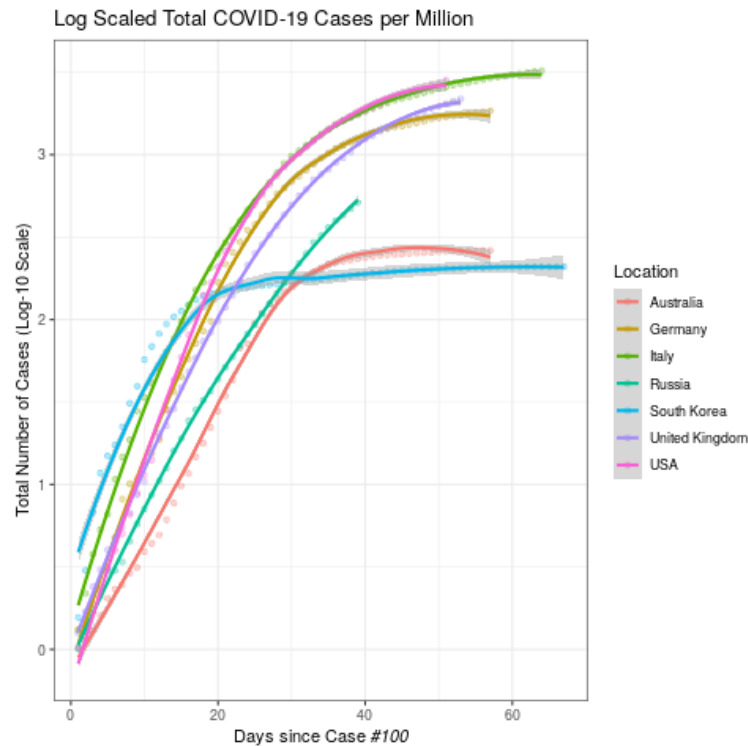


Figure 4: Chloropleth map of total deaths attributed to *COVID-19* (per Million people)

implemented by repeating the process in listing 10 for each different metric but it will also be necessary to add a feature corresponding to that metric's description, we will also create non-log scaled data as well, this is demonstrated in listings 12 through 17, finally the dataframes are merged in listing 18, the corresponding plot is shown in figure 5.

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

Disadvantages

- 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.

```

1 interested_locations <- c("Australia", "USA", "Italy", "Germany",
  ↪ "Russia", "South Korea", "United Kingdom")
2
3 ##### Number of Cases
4 cv <- as_tibble(covid)
5 cv <- cv %>%
6   mutate(date = as.Date(date))
7 cv <- cv[order(cv$date),]
8
9 cv <- cv %>%
10   filter(location %in% interested_locations) %>%
11   filter(total_cases > 1) %>%
12   mutate(total_cases_per_million = log10(total_cases_per_million)) %>%
13   dplyr::select(date, total_cases_per_million, location) %>%
14   pivot_wider(names_from = location, values_from =
  ↪ total_cases_per_million)
15
16 for (i in 2:ncol(cv)) {
17   ## Slide the Columns up and put the NA at the end
18   cv[,i] <- pull(cv, i) %>%
19     lead(cv[,i] %>%
20       is.na() %>%
21       sum())
22   ## Replace the date with the number of days
23   cv$date <- seq_len(nrow(cv))
24 }
25
26 cv_cases_log <- cv %>%
27   pivot_longer(names(cv)[-1], names_to = "location", values_to =
  ↪ "value") %>%
28   add_column(subject = "No. of Cases") %>%
29   add_column(scale = "Log-10 Scale")

```

Listing 12: Use dplyr to create a data frame of log scaled cases

```

1  ### Number of deaths
2
3  cv <- as_tibble(covid)
4  cv <- cv %>%
5    mutate(date = as.Date(date))
6  cv <- cv[order(cv$date),]
7
8  cv <- cv %>%
9    filter(location %in% interested_locations) %>%
10   filter(total_cases > 1) %>%
11   mutate(total_deaths_per_million = log10(total_deaths_per_million))
12   ↪ %>%
13   dplyr::select(date, total_deaths_per_million, location) %>%
14   pivot_wider(names_from = location, values_from =
15     ↪ total_deaths_per_million)
16
17 for (i in 2:ncol(cv)) {
18   ## Slide the Columns up and put the NA at the end
19   cv[,i] <- pull(cv, i) %>%
20     lead(cv[,i] %>%
21       is.na() %>%
22       sum())
23   ## Replace the date with the number of days
24   cv$date <- seq_len(nrow(cv))
25 }
26
27 cv_deaths_log <- cv %>%
28   pivot_longer(names(cv)[-1], names_to = "location", values_to =
29     ↪ "value") %>%
30   add_column(subject = "No. of Deaths") %>%
31   add_column(scale = "Log-10 Scale")

```

Listing 13: Use dplyr to create a data frame of log scaled deaths


```

1  ### Number of Tests
2  cv <- as_tibble(covid)
3  cv <- cv %>%
4    mutate(date = as.Date(date))
5  cv <- cv[order(cv$date),]
6  cv <- cv %>%
7    filter(location %in% interested_locations) %>%
8    filter(total_cases > 1) %>%
9    mutate(total_tests_per_thousand = log10(total_tests_per_thousand)) %>%
10   dplyr::select(date, total_tests_per_thousand, location) %>%
11   pivot_wider(names_from = location, values_from =
12     ↪ total_tests_per_thousand)
13
14 for (i in 2:ncol(cv)) {
15   ## Slide the Columns up and put the NA at the end
16   cv[,i] <- pull(cv, i) %>%
17     lead(cv[,i] %>%
18       is.na() %>%
19       sum())
20   ## Replace the date with the number of days
21   cv$date <- seq_len(nrow(cv))
22 }
23 cv_tests_log <- cv %>%
24   pivot_longer(names(cv)[-1], names_to = "location", values_to =
25     ↪ "value") %>%
26   add_column(subject = "No. of Tests") %>%
27   add_column(scale = "Log-10")
28
29 cv <- rbind(cv_cases_log, cv_deaths_log, cv_tests_log)
30 cv %>%
31   filter(subject == "deaths")
32
33 p_per_cap <- ggplot(cv , aes(y = value, x = date)) +
34   geom_point(alpha = 0.3, aes(col = location)) +
35   geom_smooth(aes(col = location), size = 0.5) +
36   theme_bw() +
37   labs(y = TeX("Count (log_{10} Scale)"), title = TeX("log_{10} Scale;
38     ↪ Value of \\textit{COVID-19} Statistics over Time"), x = TeX("Days
39     ↪ since Case \\textit{#1}"), subtitle = "Counts Per Million of
40     ↪ population") +
41   guides(col = guide_legend("Location")) +
42   facet_grid(rows = vars(subject), scales = "free_y")
43 p_per_cap

```

Listing 14: Use dplyr to create a data frame of log scaled deaths

```

1 interested_locations <- c("Australia", "USA", "Italy", "Germany",
  ↪ "Russia", "South Korea", "United Kingdom")
2
3 ##### Number of Cases
4 cv <- as_tibble(covid)
5 cv <- cv %>%
6   mutate(date = as.Date(date))
7 cv <- cv[order(cv$date),]
8
9 cv <- cv %>%
10   filter(location %in% interested_locations) %>%
11   filter(total_cases > 1) %>%
12   # mutate(total_cases = log10(total_cases)) %>%
13   dplyr::select(date, total_cases, location) %>%
14   pivot_wider(names_from = location, values_from = total_cases)
15
16 for (i in 2:ncol(cv)) {
17   ## Slide the Columns up and put the NA at the end
18   cv[,i] <- pull(cv, i) %>%
19     lead(cv[,i] %>%
20       is.na() %>%
21       sum())
22   ## Replace the date with the number of days
23   cv$date <- seq_len(nrow(cv))
24 }
25
26 cv_cases_raw <- cv %>%
27   pivot_longer(names(cv)[-1], names_to = "location", values_to =
  ↪ "value") %>%
28   add_column(subject = "No. of Cases") %>%
29   add_column(scale = "Count")

```

Listing 15: use dplyr to create a data frame of non-log scaled cases

```

1  ### Number of deaths
2
3  cv <- as_tibble(covid)
4  cv <- cv %>%
5    mutate(date = as.Date(date))
6  cv <- cv[order(cv$date),]
7
8  cv <- cv %>%
9    filter(location %in% interested_locations) %>%
10   filter(total_cases > 1) %>%
11   # mutate(total_deaths = log10(total_deaths_)) %>%
12   dplyr::select(date, total_deaths, location) %>%
13   pivot_wider(names_from = location, values_from = total_deaths)
14
15  for (i in 2:ncol(cv)) {
16    ## Slide the Columns up and put the NA at the end
17    cv[,i] <- pull(cv, i) %>%
18      lead(cv[,i] %>%
19        is.na() %>%
20        sum())
21    ## Replace the date with the number of days
22    cv$date <- seq_len(nrow(cv))
23  }
24
25  cv_deaths_raw <- cv %>%
26    pivot_longer(names(cv)[-1], names_to = "location", values_to =
27      ↪ "value") %>%
27    add_column(subject = "No. of Deaths") %>%
28    add_column(scale = "Count")

```

Listing 16: use dplyr to create a data frame of non-log scaled deaths

```

1  ### Number of Tests
2  cv <- as_tibble(covid)
3  cv <- cv %>%
4    mutate(date = as.Date(date))
5  cv <- cv[order(cv$date),]
6  cv <- cv %>%
7    filter(location %in% interested_locations) %>%
8    filter(total_cases > 1) %>%
9    # mutate(total_testsd = log10(total_testsd)) %>%
10   dplyr::select(date, total_tests, location) %>%
11   pivot_wider(names_from = location, values_from = total_tests)
12
13  for (i in 2:ncol(cv)) {
14    ## Slide the Columns up and put the NA at the end
15    cv[,i] <- pull(cv, i) %>%
16      lead(cv[,i] %>%
17        is.na() %>%
18        sum())
19    ## Replace the date with the number of days
20    cv$date <- seq_len(nrow(cv))
21  }
22  cv_tests_raw <- cv %>%
23    pivot_longer(names(cv)[-1], names_to = "location", values_to =
24      ↪ "value") %>%
25    add_column(subject = "No. of Tests") %>%
26    add_column(scale = "Count")
27  cv <- rbind(cv_cases_raw, cv_deaths_raw, cv_tests_raw)
28  cv %>%
29    filter(subject == "deaths")
30
31  p_total <- ggplot(cv , aes(y = value, x = date)) +
32    geom_point(alpha = 0.3, aes(col = location)) +
33    geom_smooth(aes(col = location), size = 0.5) +
34    theme_bw() +
35    labs(y = TeX("Total Count"), title = TeX("Total Count of
36      ↪ \\textit{COVID-19} Statistics over Time"), x = TeX("Days since
37      ↪ Case \\textit{#1}")) +
38    guides(col = guide_legend("Location"), subtitle = "Per Million of
39      ↪ Population") +
40    facet_grid(rows = vars(subject), scales = "free_y")
41  p_total

```

Listing 17: use dplyr to create a data frame of non-log scaled tests

```

1 plots <- list(p_per_cap + guides(col = FALSE), p_total+
  ↳ theme(legend.position="bottom") )
2 # plots <- list(p_per_cap + theme(legend.position="bottom"), p_total+
  ↳ theme(legend.position="bottom") )
3 library(gridExtra)
4
5 gridExtra::grid.arrange(grobs = plots, layout_matrix = matrix(1:2, nrow
  ↳ = 1))

```

Listing 18: Merge the plots in order to create a single visualisation

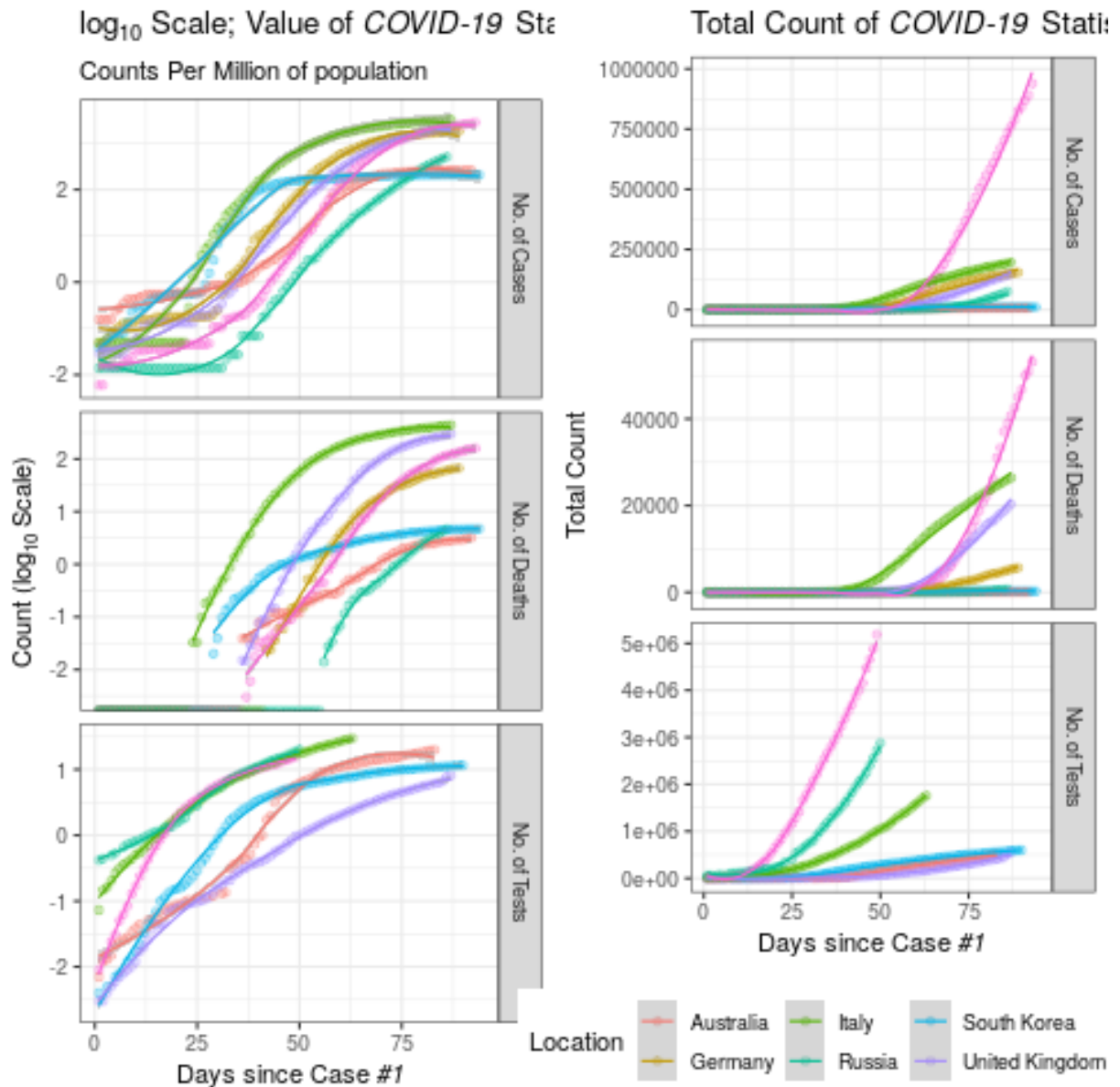


Figure 5: Multi Scatter Plot of *COVID-19* Metrics.

Discussion on analysis results

This plot demonstrates that

Discussion on other Aspects

- A potential improvement to this plot would be to plot many countries, say 30 but greyscale those countries and only apply colour to countries of interest, this would provide background information relative to those observations but not overwhelm the reader, this is a suggestion made by Andy Kirk in his *Visualising Data* blog [kirk2015].

Literature review of related work

As mentioned in section .2 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].

NSW Health created a visualisation of cases acquired over time using a barchart in a way that resembles a histogram, [nswhealth2020] this plot is very easy to interpret and clearly demonstrates the success of NSW in *flattening the curve*, this visualisation could have been implemented for this data as demonstrated in listing 19 shown in figure 6 for different countries in a similar fashion, this however was not effective for comparing countries and so was not pursued.

```
1  #+begin_src
2  interested_locations <- c("Australia", "USA", "Italy", "Germany",
   ↪  "Russia", "South Korea", "United Kingdom")
3  cv <- covid %>%
4    dplyr::filter(location %in% interested_locations)
5
6  ggplot(fortify(cv), aes(x = as.Date(date), y = new_cases_per_million,
   ↪  fill = location)) +
7    geom_col(col = "grey") +
8    labs(x = "Date", y = "New Cases Per Million") +
9    theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
10   theme_bw()
```

Listing 19: Use ggplot to create a bar chart

TODO 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

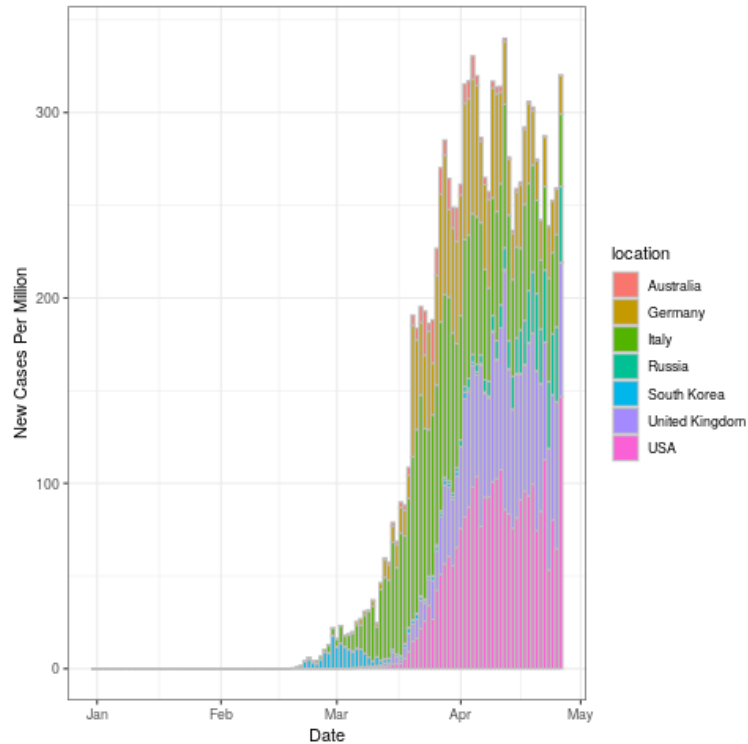


Figure 6: Bar Chart of cases over time for various locations

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

Appendix

ATTACH



Figure 7: John Hopkins Bubble Chart [2020o]

```
1 citation()  
2 citation("ggplot2")
```

Listing 20: Generate Citation for **R** programming Language

To cite R in publications use:

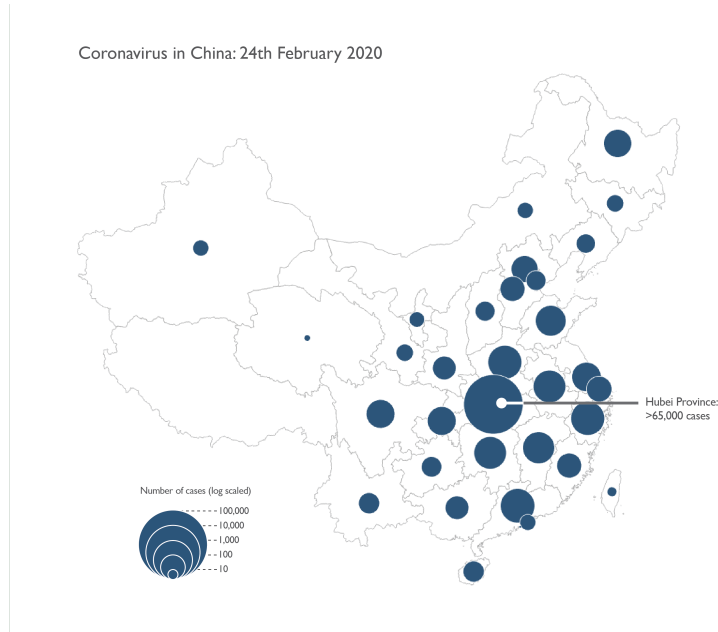


Figure 8: Bubble Plot Chart produced by Field in his blog [field2020]

R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
URL <https://www.R-project.org/>.

A BibTeX entry for LaTeX users is

```
@Manual{,
  title = {R: A Language and Environment for Statistical Computing},
  author = {{R Core Team}},
  organization = {R Foundation for Statistical Computing},
  address = {Vienna, Austria},
  year = {2020},
  url = {https://www.R-project.org/},
}
```

We have invested a lot of time and effort in creating R, please cite it when using it for data analysis. See also `citation("pkgname")` for citing R packages.

References

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

