



## **Covid-19 Exploratory Data Analysis using R**

Theodoros Panagiotis Vagenas

PhD student

[tpvagenas@mail.ntua.gr](mailto:tpvagenas@mail.ntua.gr)

Athens 2021

# Introduction

The following analysis utilizes up-to-date Covid-19 data from John Hopkins CSSE containing information about cumulative confirmed cases and deaths on different countries during the interval 22/01/20 - 12/01/21. At first, the dataset is preprocessed in order to be analysed. Then summarized measurements and graphs are produced for public use. The graphs are drawn for the top 20 countries with the most Covid-19 cases worldwide. Finally, conclusions are extracted for the European Union countries including Greece, to which we have special interest. All the measurements and conclusions are accompanied with graphs (implemented codes in R provided).

## 1.1 R libraries

Below there is a list of the R libraries used in the codes presented in the project:

```
1 library(data.table)
2 library(date)
3 library(lubridate)
4 library(dplyr)
5 library(ggplot2)
6 library(cowplot)
7 library(gridExtra)
8 library(ggrepel)
9 library(arules)
10 library(ggExtra)
```

## 1.2 Data preparation

In this section, methods applied in order to clean and transform the data.

```
1 # Load Data for cumulative confirmed cases and deaths:
2 cases_data <- fread('time_series_covid19_confirmed_global.csv')
3 deaths_data <- fread('time_series_covid19_deaths_global.csv')
4 #The following columns are deleted from both datasets:
5 cases_data$'Province/State' <- NULL
6 cases_data$Lat <- NULL
7 cases_data$Long <- NULL
8 deaths_data$'Province/State' <- NULL
9 deaths_data$Lat <- NULL
10 deaths_data$Long <- NULL
11 #Data conversion from wide to long format:
12 cases_long <- melt(cases_data)
13 deaths_long <- melt(deaths_data)
14 #The produced columns are renamed to confirmed and deaths, and Country.Region to
    Country:
15 setnames(cases_long, "Country/Region", "Country")
16 setnames(deaths_long, "Country/Region", "Country")
17 setnames(cases_long, "value", "confirmed")
18 setnames(deaths_long, "value", "deaths")
19 #Conversion of date variable from character type to date object:
20 cases_long$variable <- cases_long[,lapply(.SD, mdy), .SDcols="variable"]
21 setnames(cases_long, "variable", "date")
22 deaths_long$variable <- deaths_long[,lapply(.SD, mdy), .SDcols="variable"]
23 setnames(deaths_long, "variable", "date")
24 #Group by country and date:
25 cases_long2 <- cases_long[,sum(confirmed), by = .(Country, date)]
26 setnames(cases_long2, "V1", "confirmed")
27 deaths_long2 <- deaths_long[,sum(deaths), by = .(Country, date)]
28 setnames(deaths_long2, "V1", "deaths")
29 #Merge the datasets:
30 com <- merge(cases_long2, deaths_long2, by=c("Country", "date"))
31 #Calculate counts for the whole world: cases= 91595225, deaths=1962228
32 world <- com[,.(confirmed = sum(confirmed), deaths = sum(deaths)), by=date]
33 world$confirmed[length(world$confirmed)]
34 world$deaths[length(world$deaths)]
```

```

35 #Sort by country and date:
36 com <- com[order(Country,date)]
37 #Calculate daily confirmed cases and deaths:
38 day1 <-min(com$date)
39 com <- com %>% mutate(confirmed.inc = ifelse(date==day1, NA, confirmed- lag(
    confirmed, n=1)))
40 com <- com %>% mutate(deaths.inc = ifelse(date==day1, NA, deaths- lag(deaths, n=1)))

```

# Exploratory Data Analysis

## 2.1 World Covid-19 cases and deaths Analysis

First, a table with the same columns as table "com" is created with com data grouped by date with summation. As a result, confirmed cases and deaths are summed up for each date to compute worldwide results.

```

1 max_date <- max(com$date)
2 world_data_per_day <- com[,.(cases=sum(confirmed),deaths=sum(deaths),world.confirmed
    .inc=sum(confirmed.inc),world.deaths.inc=sum(deaths.inc)),by=date]
3 # Set English as language for date labels
4 Sys.setlocale("LC_TIME", "C")

```

In **Figure 2.1** we can see the cumulative cases and deaths as they are represented graphically for every date worldwide. Confirmed cases escalated very fast. There is a period after November that the number of cases and deaths have a quite stable increase. As it was expected there is a strong correlation between the number of deaths and the number of cases.

```

1 g1 <- ggplot()+
2   geom_line(data=world_data_per_day,aes(x=date,y=cases),color="Blue") +
3   labs(x = "Date", y="Cases",title=paste0("Worldwide cases - ",max_date))+
4   scale_x_date(date_breaks="1 month",date_labels = "%b")+
5   theme(legend.title=element_blank(),
6         legend.position='bottom',
7         plot.title =element_text(size = 13,hjust = 0.5))
8 g2 <- ggplot()+
9   geom_line(data=world_data_per_day,aes(x=date,y=deaths),color="Blue")+
10  labs(x = "Date", y="Deaths",title=paste0("Worldwide deaths - ",max_date))+
11  scale_x_date(date_breaks="1 month",date_labels = "%b")+
12  theme(legend.title=element_blank(),
13        legend.position='bottom',
14        plot.title = element_text(size = 13,hjust = 0.5))
15 plot_grid(g1, g2, labels=c("", "", ""), ncol = 2, nrow = 1, align = "v")

```

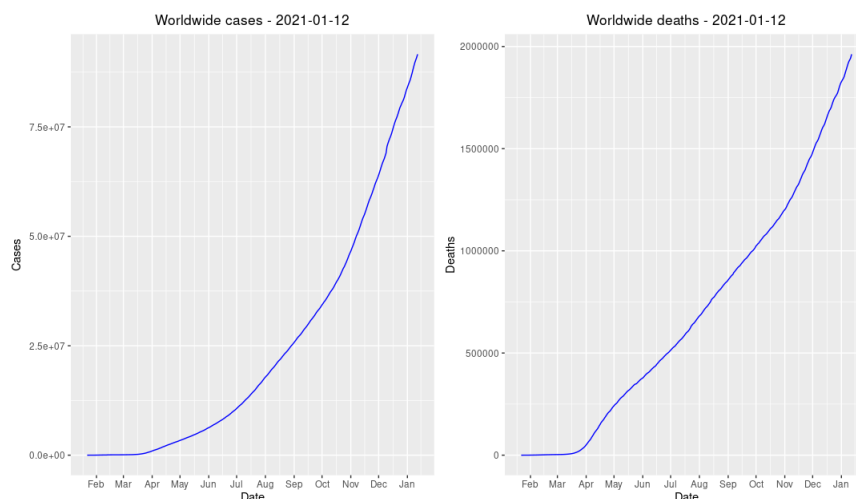


Figure 2.1: Worldwide cases and deaths

In **Figure 2.2** the same results are plotted together in the logarithmic scale. This way a number of values can be displayed without compressing down the smallest values of cases and deaths. Also with logarithmic

scale the constant percentage of the increase in the cumulative cases/deaths is shown as a constant vertical increase. So, a straight line means a constant growth rate. As a result, we can observe two large sudden increases in growth rate the first months and a more smooth continuously increasing growth rate after May.

```
1 g2 <- ggplot()+
2   geom_line(data=world_data_per_day,aes(x=date,y=cases,color='cases')) +
3   geom_line(data=world_data_per_day,aes(x=date,y=deaths,color='deaths'))+
4   labs(x = "Date", y="Cases",title="Worldwide number of cases and deaths in log
5     scale")+
6   theme(legend.title=element_blank(),
7         legend.position='right',
8         plot.title =element_text(hjust=0.5))+
9   scale_y_continuous(trans='log10')+
10  scale_x_date(date_breaks="1 month",date_labels = "%b")
11 plot(g2)
```

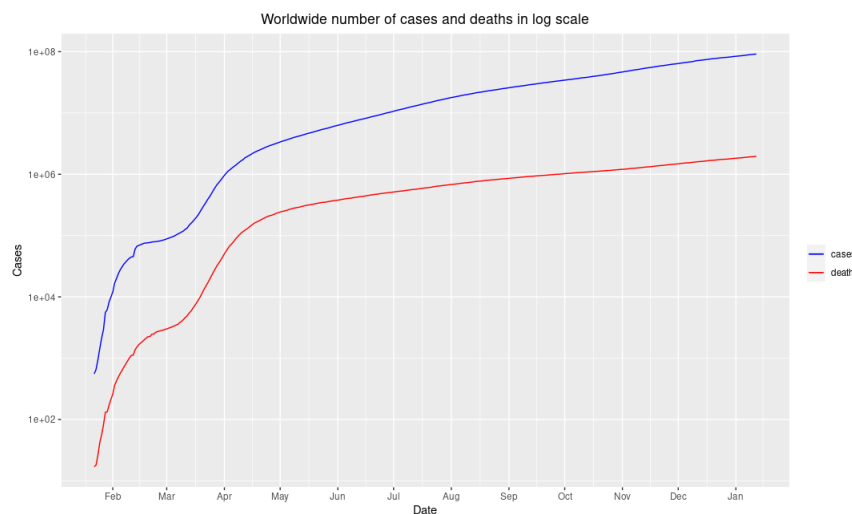


Figure 2.2: Worldwide cases and deaths (log scale)

In **Figure 2.3** we can see the daily confirmed cases and deaths for the whole world for the same interval. Confirmed cases have some spikes between periods of smaller increases in new cases per day. In April, August and November these spikes are bigger so the everyday cases became suddenly much larger. Daily deaths also shown a big increase throughout April-May and also in November keeping a stable course in-between the two periods. After November daily new cases and deaths have been increased from day to day indicating an outbreak.

```
1 g3 <- ggplot(data=world_data_per_day,aes(x=date,y=world.confirmed.inc),color='blue')
2   +
3   geom_point() +
4   geom_smooth()+
5   labs(x = "Date", y="Cases",title="Worldwide daily cases")+
6   scale_x_date(date_breaks="1 month",date_labels = "%b")+
7   theme(plot.title =element_text(hjust=0.5))
8 g4 <- ggplot(data=world_data_per_day,aes(x=date,y=world.deaths.inc),color='blue')+
9   geom_point() +
10  geom_smooth()+
11  labs(x = "Date", y="Cases",title="Worldwide daily deaths")+
12  scale_x_date(date_breaks="1 month",date_labels = "%b")+
13  theme(plot.title =element_text(hjust=0.5))
14 grid.arrange(g3, g4, ncol=2)
```



Figure 2.3: Worldwide cases and deaths

In the next step a death rate is calculated from the cumulative confirmed cases and deaths with the following equation:

$$death\_rate = 100 \frac{deaths}{cases} \quad (2.1)$$

The code for this calculation and the plot is provided below:

```
1 world_data_per_day <- world_data_per_day %>% mutate(death_rate=100*deaths/cases)
2 g5 <- ggplot()+
3   geom_line(data=world_data_per_day,aes(x=date,y=death_rate)) +
4   labs(x = "Date", y="Death rate",title="Worldwide death rate")+
5   scale_x_date(date_breaks="1 month",date_labels = "%b")+
6   theme(legend.title=element_blank(),
7         legend.position='bottom',
8         plot.title = element_text(hjust=0.5))
9 plot(g5)
10 ggsave("death_rate.jpg", dpi=300)
```

In [Figure 2.4](#) we can observe that the biggest value of death rate was between April and May which was the period after the first outbreak. The next months there was a decreasing tendency probably due to the lock-downs and restrictions that were imposed and due to the better understanding and treatment of the Covid-19.

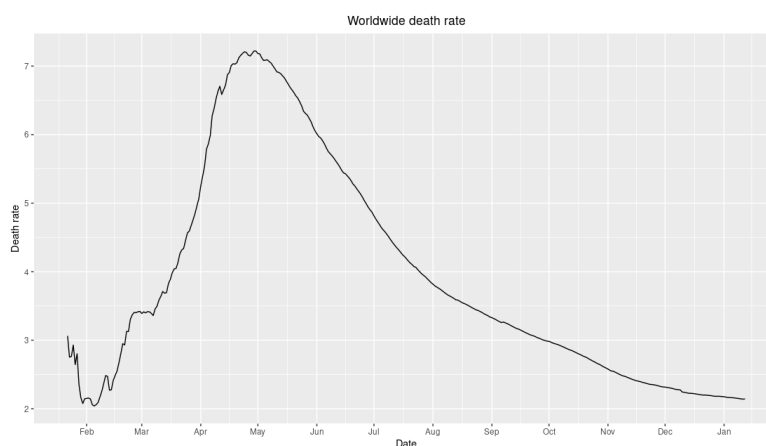


Figure 2.4: Worldwide death rate

## 2.2 Top 20 countries with the most confirmed cases

In this sections measurements and graphs are presented for the 20 countries with the most confirmed cases in the world. The countries are sorted by the number of cumulative cases. Also a column named "other" is constructed and contains the cases and deaths for the rest countries. We can see that US have more Covid-19 cases and deaths than the rest of the countries, minus the top presented in the graph. In the second row, deaths are presented with the same order as before, so it can be easily observed that there are countries with more cases but less deaths than others with less cases. For example, Mexico is 13th in

confirmed cases but 4th in deaths from the top 20. This statistics may indicate inaccuracies in confirmed cases/deaths or a health system and generally a mechanism with smaller capability of dealing with the pandemic. This is more clear in the 3rd graph presenting the death rate, in which we can observe better the variability in deaths and death rate between countries with most cases.

```

1 top_countries <- com%>% filter(date== max(date))
2 top_countries$confirmed.inc <- NULL
3 top_countries$deaths.inc <- NULL
4 top_countries$date <- NULL
5 top_countries <- top_countries[order(-confirmed)]
6 top_countries <- top_countries %>% mutate(death_rate=100*deaths/confirmed)
7 other_countries <- top_countries[21:dim(top_countries)[1],lapply(.SD, sum),.SDcols =
  c("confirmed","deaths")]
8 other_countries_death_rate <- top_countries[21:dim(top_countries)[1],lapply(.SD,
  mean),.SDcols = c("death_rate")]
9 top_countries_20 <- rbind(top_countries[1:20,],cbind(Country="other",other_countries
  ,other_countries_death_rate))
10 top_countries_20$Country <- factor(top_countries_20$Country, levels = top_countries_
  20$Country)

1 # bar plot for comparison between confirmed cases, deaths and death rate
2 g6 <- ggplot(top_countries_20,aes(x = Country,y=confirmed,group=Country))+
3   geom_bar(stat='identity')+
4   geom_text(aes(label=confirmed, y=confirmed), vjust=0,size=2)+
5   labs(x = element_blank(), y="Cases",title="Top 20 Countries with most cases: Cases
  ")+
6   theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text(size =
  8))
7 g7 <- ggplot(top_countries_20,aes(x = Country,y=deaths,group=Country))+
8   geom_bar(stat='identity')+
9   geom_text(aes(label=deaths, y=deaths), vjust=0,size=2)+
10  labs(x = element_blank(), y="Deaths",title="Top 20 Countries with most cases:
  Deaths")+
11  theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text( size =
  8))
12 g8 <- ggplot(top_countries_20,aes(x = Country,y=death_rate,group=Country))+
13   geom_bar(stat='identity')+
14   geom_text(aes(label=round(death_rate,2), y=death_rate), vjust=0,size=3)+
15   labs(x = "Country", y="Death rate",title="Top 20 Countries with most cases: Death
  rate")+
16   theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text(size =
  8))
17 #grid.arrange(g6, g7,g8, nrow=3,top="Top 20 countries with the most confirmed cases
  ")
18 plot_grid(g6, g7, g8, labels=c("", "", ""), ncol = 1, nrow = 3, align = "v")

```

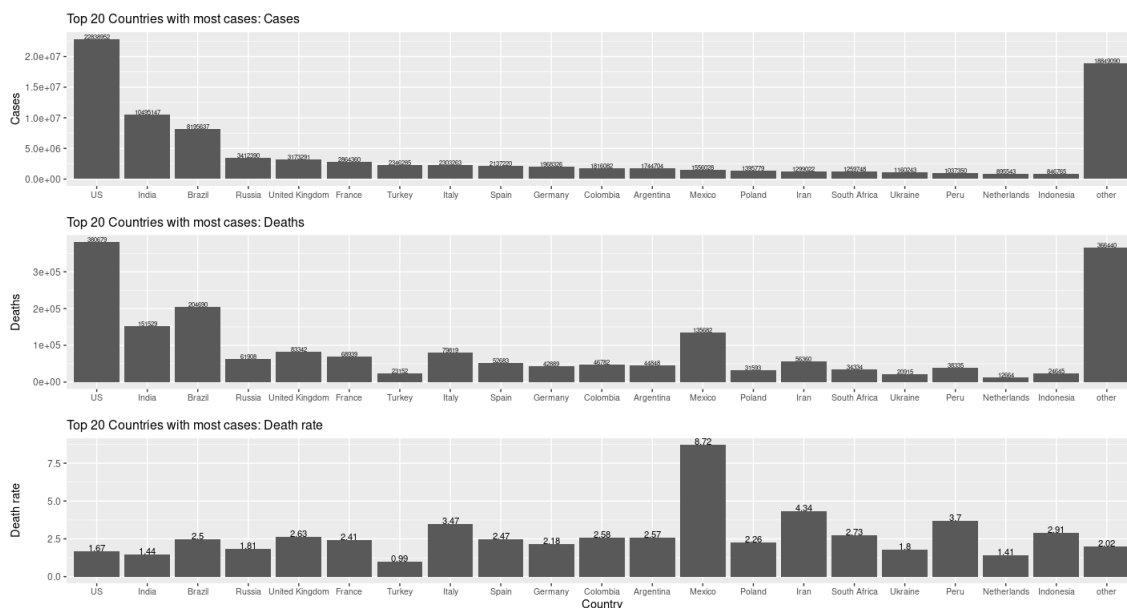


Figure 2.5: Top 20 countries cases, deaths and death rate

In **Figure 2.6** cases are plotted as opposed to deaths for the previous 20 countries. Colors show the countries and the diameter of the points shows the death rate. This display depicts clearly the relationship between cases and deaths in these countries. Especially, there is a concentration of countries below 5 millions cases and 100 thousands deaths and two other distinctive smaller clusters that differ a lot from this concentration.

```
1 g9 <- ggplot(data=top_countries_20,aes(x=confirmed,y=deaths,group=Country))+
2   geom_point(aes(color=Country,size=death_rate)) +
3   labs(x = "Confirmed cases", y="Deaths",title="Top 20 countries with most cases:
4     Deaths vs confirmed",color="Country",size="Death rate")+
5   geom_text_repel(aes(label = Country))+
6   theme(legend.title=element_text(),legend.position='bottom',plot.title =element_
7     text(hjust = 0.5,size=12),axis.text.x=element_text(size=10),axis.text.y=element_
8     text(size = 10))
9 plot(g9)
10 ggsave("top20_confirmed_deaths.jpg", dpi=300)
```

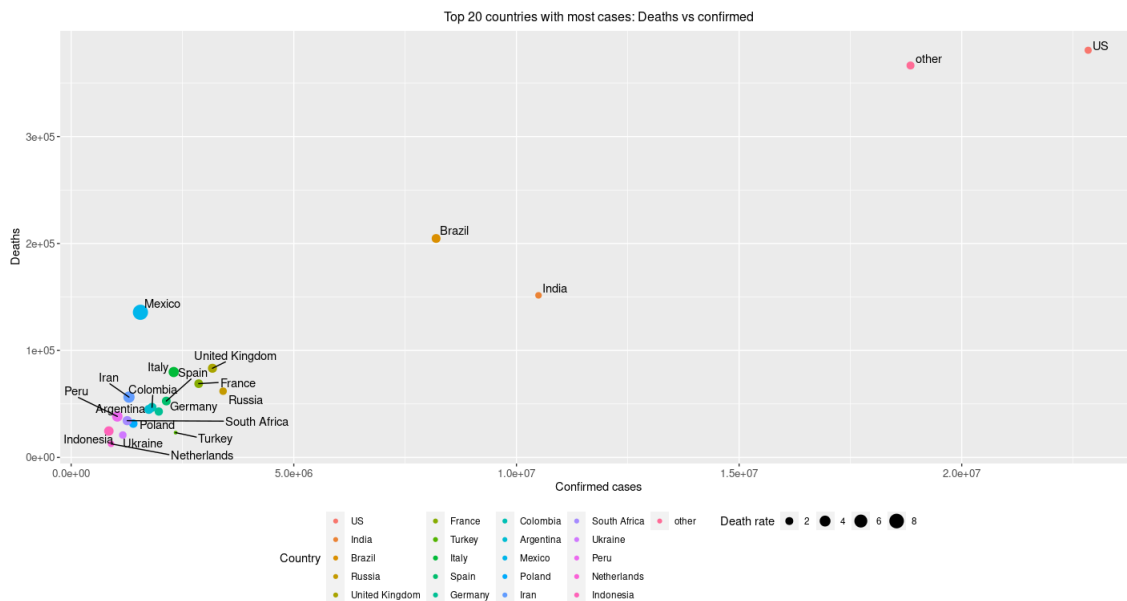


Figure 2.6: Top 20 Confirmed cases vs deaths

## World map of cases

In the **Figure 2.7** below, a map with confirmed cases is plotted. According to the color scaled graph, yellow indicates the countries with less cases and red with more cases. Also grey level is used for countries that are not in the dataset. There was also some differences in region names that were corrected manually in order to generate graph for the most of the countries. In the map, we can see the spatial distribution of the Covid-19 cases until the last date of the dataset.

```
1 # Construct world map with confirmed cases
2 # Load again dataset without dropping columns
3 full_world <- map_data("world")
4 cases_data <- fread('time_series_covid19_confirmed_global.csv')
5 deaths_data <- fread('time_series_covid19_deaths_global.csv')
6 cols_vector = c('Province/State', 'Country/Region', 'Lat', 'Long', names(cases_data)[
7   length(names(cases_data))])
8 cases_data <- cases_data[,..cols_vector]
9 cases_data$'Province/State' <- NULL
10 setnames(cases_data, "Country/Region", "region")
11 cases_data <- cases_data[,.(lat=Lat, long=Long, confirmed=sum(`1/12/21`)),by=.(region)
12   ] # here last date
13 sequence_break <- c(1,100,1000,10000,100000,1000000)
14 full_world$region[full_world$region=="USA"] <- "US"
15 full_world$region[full_world$region=="Myanmar"] <- "Burma"
16 full_world$region[full_world$region=="Macedonia"] <- "North Macedonia"
17 full_world$region[full_world$region=="South Korea"] <- "Korea, South"
18 full_world$region[full_world$region=="UK"] <- "United Kingdom"
19 full_world$region[full_world$region=="Republic of Congo"] <- "Congo (Brazzaville)"
```

```

18 full_world$region[full_world$region=="Democratic Republic of the Congo"] <- "Congo (
    Kinshasa)"
19 full_world$region[full_world$region=="Ivory Coast"] <- "Cote d'Ivoire"
20 full_world$region[full_world$region=="Taiwan"] <- "Taiwan*"
21 full_world$region[full_world$region=="Czech Republic"] <- "Czechia"
22 full_world$region[full_world$region=="French Guiana"] <- "Czechia"
23 nn <- left_join(cases_data, full_world, by = "region")
24 gworldmap <- ggplot()+
25   geom_polygon(data = full_world, aes(x=long,y=lat,group=group),alpha=0.2)+
26   geom_polygon(data = nn, aes(x=long,y=lat,y,group=group,fill=confirmed),alpha=1)+
27   scale_fill_gradientn(colours = rev(heat.colors(7)),breaks = sequence_break,trans =
    "log10")+
28   annotate(geom="text", x=150, y=-100, label="Grey countries' measurements are not
    available",color="red",alpha=0.7)+
29   theme(
30     plot.title =element_text(hjust = 0.5),
31     legend.position = "right",
32     legend.background = element_rect(fill = "#ffffff", color = NA)
33   )
34 plot(gworldmap)

```

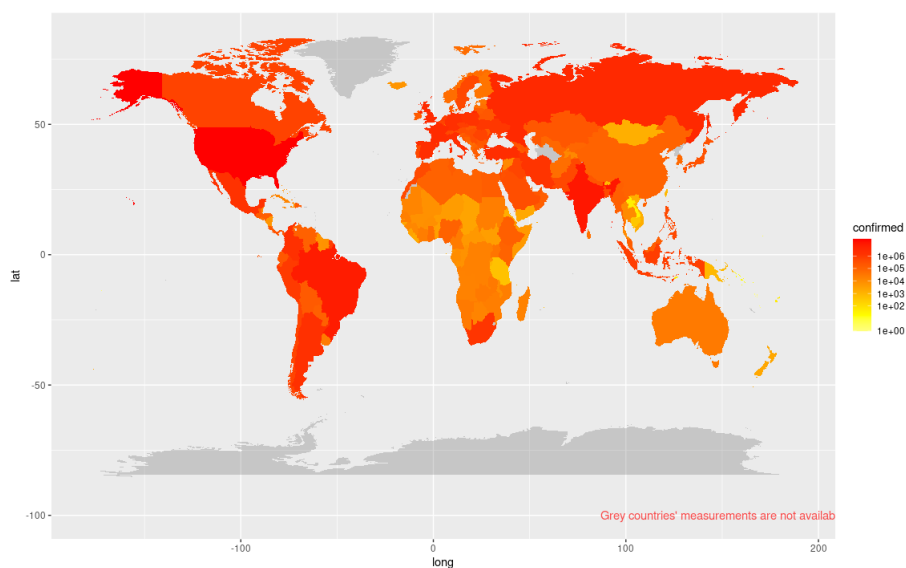


Figure 2.7: World map with Covid-19 confirmed cases

## 2.3 European Union countries analysis

The next region of interest is the European Union countries with a special analysis for Greece. It can also provide more information about the effects and the Covid-19 progression between neighbor countries. The appropriate table is constructed similar to the "Top 20 countries" case.

```

1 europeanUnion <- c("Austria","Belgium","Bulgaria","Croatia","Cyprus",
2   "Czech Rep.,","Denmark","Estonia","Finland","France",
3   "Germany","Greece","Hungary","Ireland","Italy","Latvia",
4   "Lithuania","Luxembourg","Malta","Netherlands","Poland",
5   "Portugal","Romania","Slovakia","Slovenia","Spain",
6   "Sweden","United Kingdom")
7 europe_com <- com %>% filter(com$Country %in% europeanUnion)
8 euro_countries <- europe_com%>% filter(date== max(date))
9 euro_countries$confirmed.inc <- NULL
10 euro_countries$deaths.inc <- NULL
11 euro_countries$date <- NULL
12 euro_countries <- euro_countries[order(-confirmed)]
13 euro_countries <- euro_countries %>% mutate(death_rate=100*deaths/confirmed)
14 euro_countries$Country <- factor(euro_countries$Country, levels = euro_countries$
    Country)
15 temp <- europe_com %>% filter(Country %in% euro_countries$Country)
16 euro_countries <- euro_countries %>% mutate(categ = discretize(euro_countries$death_
    rate, method = "frequency", breaks = 3))
17 # European confirmed cases,deaths and death rates ordered by cases (descent)
18 g18 <- ggplot(euro_countries,aes(x = Country,y=confirmed,group=Country))+

```



```

19 geom_bar(stat='identity')+
20 geom_text(aes(label=round(confirmed, 2), y=confirmed), size=2, vjust=0)+
21 labs(x = "", y="Cases",title="Europe: Confirmed cases")+
22 theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text(size =
  7))
23 g19 <- ggplot(euro_countries,aes(x = Country,y=deaths,group=Country))+
24 geom_bar(stat='identity')+
25 geom_text(aes(label=round(deaths, 2), y=deaths), size=2, vjust=0)+
26 labs(x = "", y="Deaths",title="Europe: Deaths")+
27 theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text(size =
  7))
28 g20 <- ggplot(euro_countries,aes(x = Country,y=death_rate,group=Country))+
29 geom_bar(stat='identity')+
30 geom_text(aes(label=round(death_rate, 2), y=death_rate), size=2, vjust=0)+
31 labs(x = "Country", y="Death rate",title="Europe: Death rate")+
32 theme(plot.title =element_text(hjust = 0,size=12),axis.text.x=element_text(size =
  7))
33 # align vertical for more clear results
34 plot_grid(g18, g19, g20, labels=c("", "", ""), ncol = 1, nrow = 3, align = "v")

```

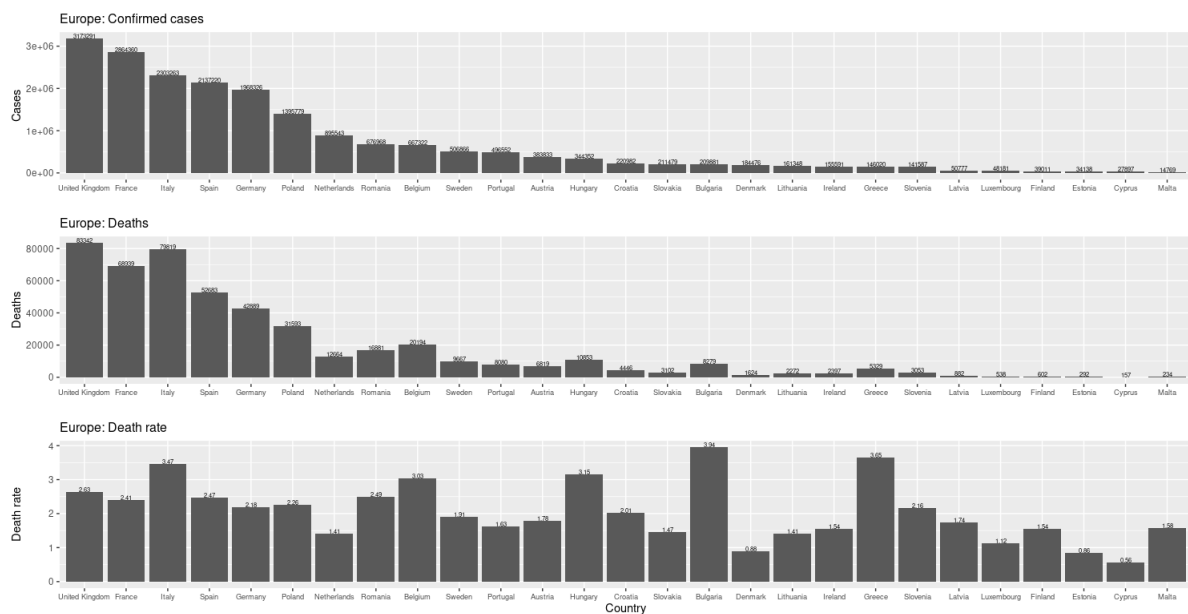


Figure 2.8: European Union countries cases, deaths and death rate

In **Figure 2.8** above European Union members are sorted by confirmed cases. There is a vertical alignment so they can be compared to each other. An interesting point of view the study of the reasons that lead to less cases in them so that can help the other countries to protect themselves from another outbreak. Death rate can also be an indication that some countries confronted Covid-19 better. Biggest death rate here is observed in Italy and UK where there was a very fast increase in cases, a situation that caused problems in the health system. Death rate also has a different behavior from cases as we can see from the fact that countries are sorted by the cases.

The next diagram **Figure 2.9** plots deaths versus cases. The form of the graph between these two variables for these countries is very similar but has different scales. This fact was also observed in the previous plots that also show the death rates.

```

1 g21 <- ggplot(data=temp,aes(x=confirmed,y=deaths,group=Country))+
2   geom_line(aes(color=Country)) +
3   labs(x = "Confirmed cases", y="Deaths",title="Europe: Deaths vs Confirmed cases")+
4   theme(legend.title=element_blank(), legend.position='right',plot.title =element_
  text(hjust = 0.5))
5 plot(g21)
6
7
8 g22 <- ggplot(data=euro_countries,aes(x=confirmed,y=deaths,group=Country))+
9   geom_point(aes(color=Country,size=death_rate)) +
10  labs(x = "Confirmed cases", y="Deaths",title="Europe: Deaths vs Confirmed cases",
  color="Country",size="Death rate")+

```

```

11 geom_text_repel(aes(label = Country))+
12 theme(legend.title=element_text(),legend.position='bottom',plot.title =element_
    text(hjust = 0.5,size=12),axis.text.x=element_text(size=10),axis.text.y=element_
    text(size = 10))
13 plot(g22)

```

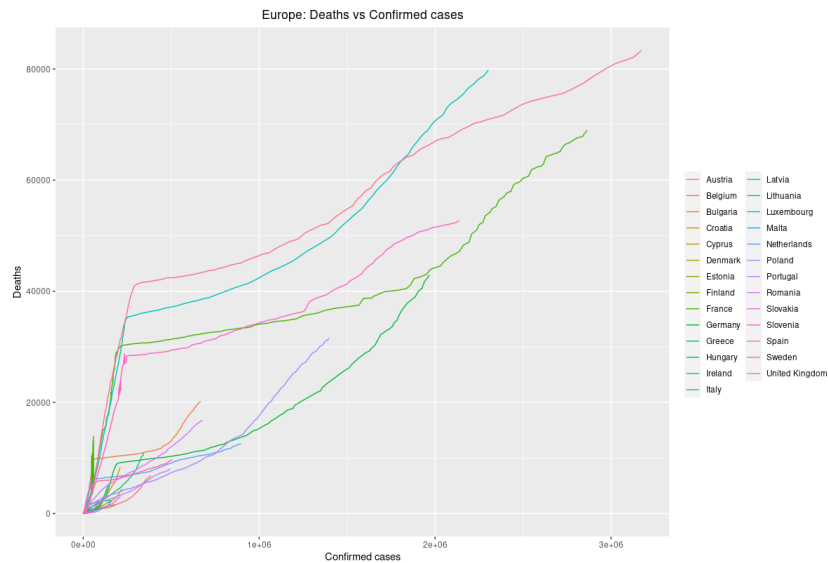


Figure 2.9: European Union countries deaths vs confirmed

In the next diagram **Figure 2.10** the same plot is presented but with points and death rates, in correspondence with the "Top 20 countries" occasion.

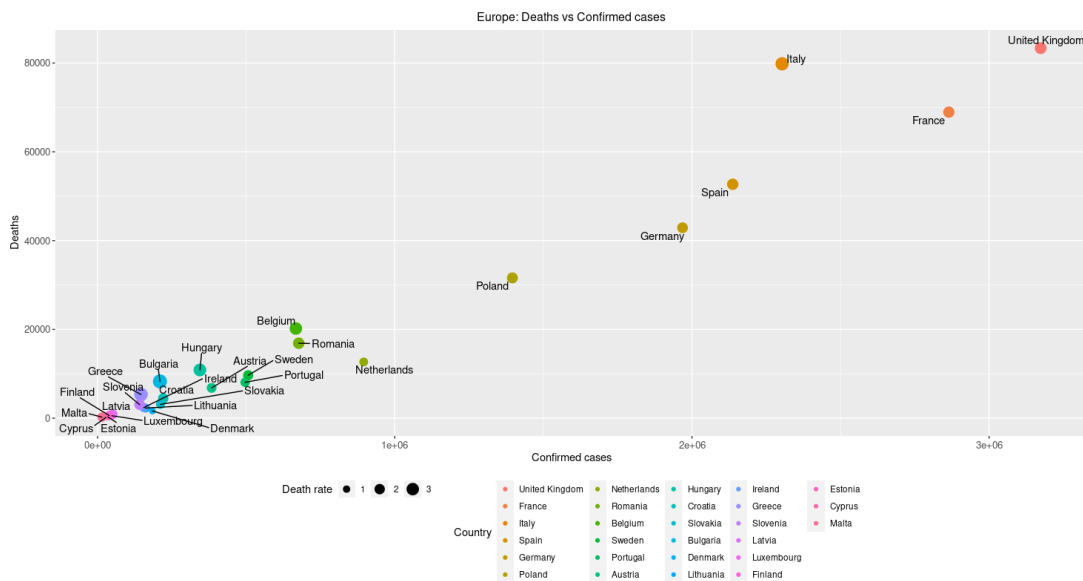


Figure 2.10: European Union countries deaths vs confirmed

With the diverging bars below (**Figure 2.11**) we can observe the deviations between EU countries according to cases and deaths. It shows how much the countries differ from the average. So we can see that countries with cases and deaths below the average are more close to the average/midpoint. This fact indicates that Covid-19 crisis is lower in the most EU countries. On the contrary, countries above the average are further than midpoint, which is an indication that these countries have a much difficult situation than the others.

```

1
2 # Diverging bars Europe confirmed and deaths
3 euro_countries$confirmed_norm <- round((euro_countries$confirmed - mean(euro_
    countries$confirmed))/sd(euro_countries$confirmed), 2)
4 euro_countries$confirmed_type <- ifelse(euro_countries$confirmed_norm < 0, "below",
    "above")
5 euro_countries <- euro_countries[order(euro_countries$confirmed_norm), ] # sort

```

```

6 g25 <- ggplot(euro_countries, aes(x=Country, y=confirmed_norm, label=confirmed_norm))
7   +
8   geom_bar(stat='identity', aes(fill=confirmed_type), width=.5) +
9   scale_fill_manual(name="Confirmed cases", labels = c("Above Average", "Below
10     Average"),
11     values = c("above"="#00ba38", "below"="#f8766d")) +
12   labs(subtitle="Normalised confirmed cases",
13     title= "Diverging Bars for Europe") +
14   theme( legend.position='bottom')+
15   coord_flip()
16 euro_countries$deaths_norm <- round((euro_countries$deaths - mean(euro_countries$
17   deaths))/sd(euro_countries$deaths), 2)
18 euro_countries$deaths_type <- ifelse(euro_countries$deaths_norm < 0, "below","above"
19   )
20 euro_countries <- euro_countries[order(euro_countries$deaths_norm), ] # sort
21 g26 <- ggplot(euro_countries, aes(x=Country, y=deaths_norm, label=deaths_norm)) +
22   geom_bar(stat='identity', aes(fill=deaths_type), width=.5) +
23   scale_fill_manual(name="Deaths", labels = c("Above Average", "Below Average"),
24     values = c("above"="#00ba38", "below"="#f8766d")) +
25   labs(subtitle="Normalised deaths",
26     title= "Diverging Bars for Europe") +
27   theme( legend.position='bottom')+
28   coord_flip()
29 grid.arrange(g25,g26, ncol=2)

```

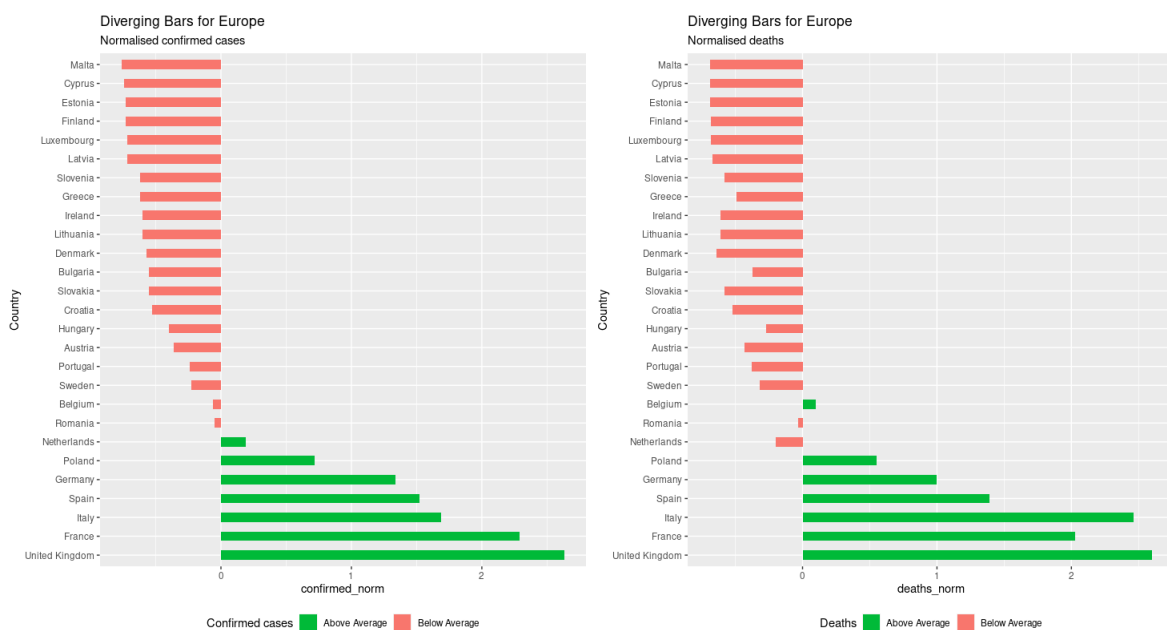


Figure 2.11: Diverging bars for European Union deaths and confirmed cases

In Dot plot [Figure 2.12](#) distances between EU countries according to cases and deaths are better visualized. As it was already mentioned there is one large group of countries with small number of cases and deaths and some countries with larger numbers. This diagram emphasizes the fact that countries with many cases/deaths present numbers significant larger than countries with less cases/deaths. Greece is in the countries with very few cases as opposed to the others.

```

1 # Ranking for Europe
2 g27 <- ggplot(euro_countries, aes(y=Country, x=confirmed)) +
3   geom_point(col="tomato2", size=3) + # Draw points
4   geom_segment(aes(y=Country,
5     yend=Country,
6     x=min(confirmed),
7     xend=max(confirmed)),
8     linetype="dashed",
9     size=0.1)+
10   labs(title="Dot Plot for Europe", subtitle="Country vs Confirmed")
11
12 g28 <- ggplot(euro_countries, aes(y=Country, x=deaths)) +

```

```

13 geom_point(col="tomato2", size=3) + # Draw points
14 geom_segment(aes(y=Country,
15                 yend=Country,
16                 x=min(deaths),
17                 xend=max(deaths)),
18              linetype="dashed",
19              size=0.1)+
20 labs(title="Dot Plot for Europe", subtitle="Country vs Deaths")
21 grid.arrange(g27,g28, ncol=2)

```

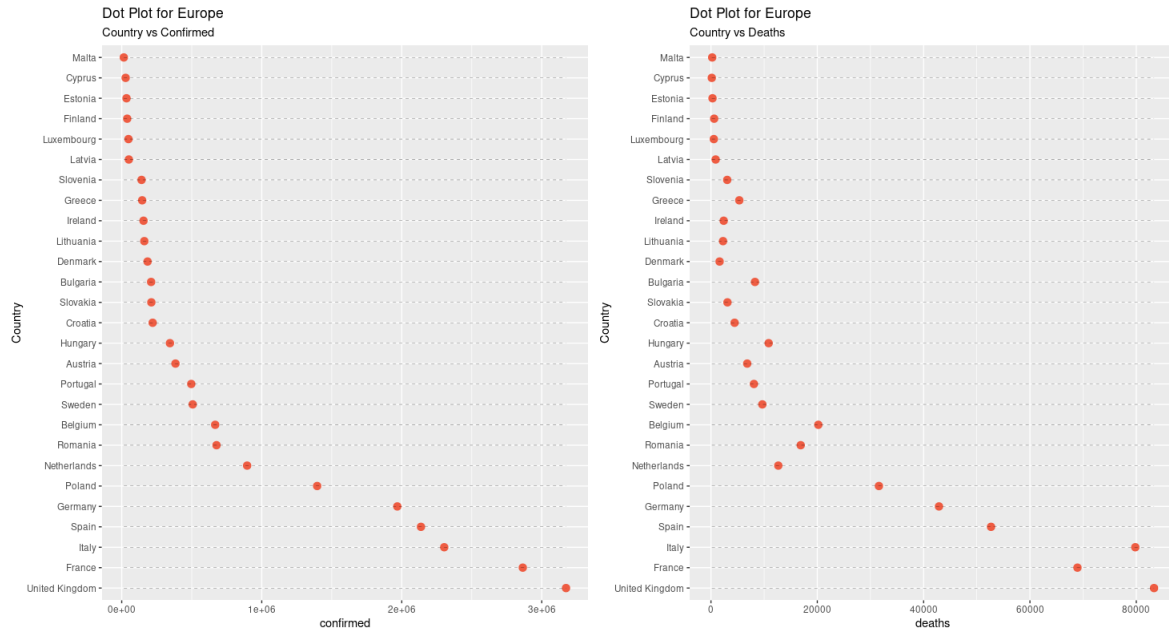


Figure 2.12: Ranking for European Union deaths and confirmed cases

## 2.4 Greece Covid-19 Analysis

In this section, graphs show the evolution of the Covid-19 pandemic in Greece. The results can underline the spread of confirmed cases during the past months and useful observations for the analysis of the situation in our country in comparison with the rest of the European Union countries. Daily cases are also calculated in order to get a better knowledge of the speed of the spreading. Also, conclusions about the distribution of daily cases and deaths are extracted.

In [Figure 2.13](#) an exponential increase of cases and deaths is presented. From this plot and also from the [Figure 2.14](#) we can observe that the behavior of cases and deaths is quite similar but of course the number of deaths is much smaller than the number of cases.

```

1 greece <- com %>% filter(Country=="Greece")
2 g11 <- ggplot(greece, aes(x=date, y=confirmed)) + geom_bar(stat="identity", width
3                 =0.1) +
4   theme_classic() +
5   labs(title = "Covid-19 Confirmed Cases in Greece", x= "Date", y= "Confirmed cases"
6         ) +
7   theme(plot.title = element_text(hjust = 0.5))
8 g12 <- ggplot(greece, aes(x=date, y=deaths)) + geom_bar(stat="identity", width=0.1)
9   +
10  theme_classic() +
11  labs(title = "Covid-19 Deaths in Greece", x= "Date", y= "Deaths") +
12  theme(plot.title = element_text(hjust = 0.5))
13 grid.arrange(g11, g12, ncol=2)

```

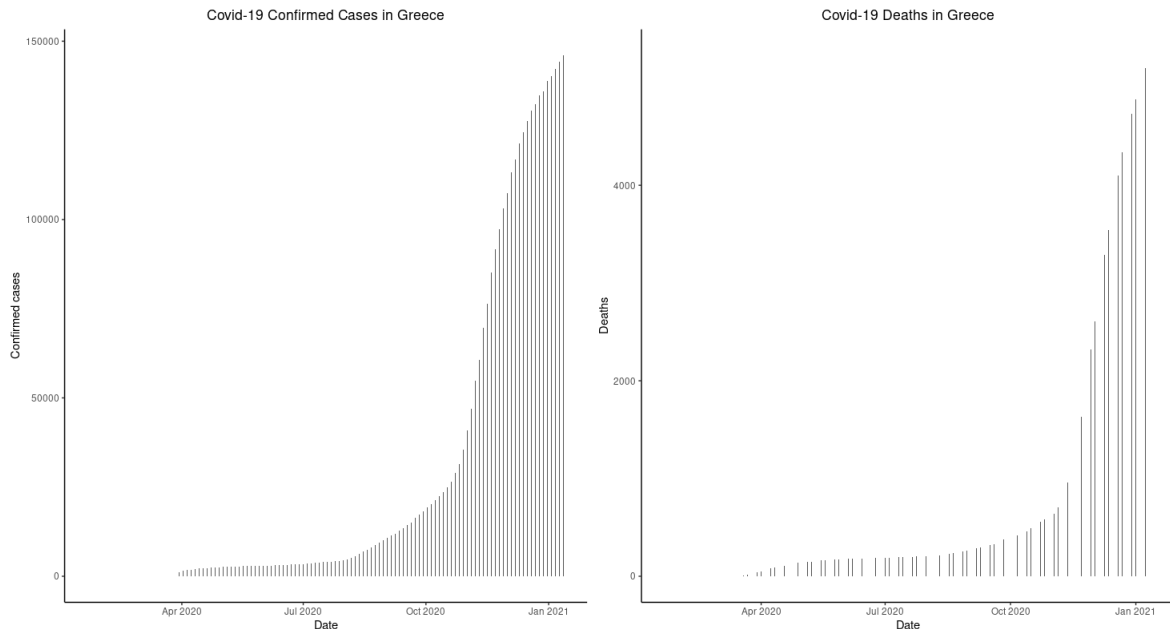


Figure 2.13: Greece confirmed cases and deaths versus date

```

1 # Cases-deaths relationships
2 g13 <- ggplot(data=greece, aes(x=confirmed, y=deaths)) +
3   geom_line() +
4   labs(x = "Confirmed cases", y="Deaths", title="Greece: Confirmed cases vs Deaths") +
5   theme(legend.title=element_blank(), legend.position='bottom', plot.title = element_
6     text(hjust = 0.5))
7 plot(g13)

```

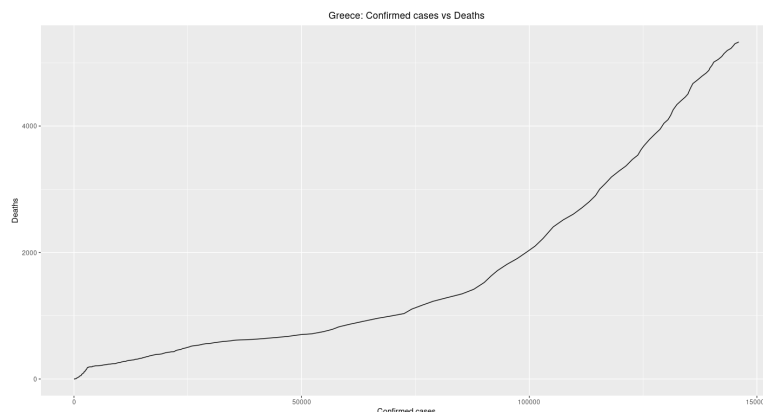


Figure 2.14: Greece confirmed cases versus deaths

In [Figure 2.15](#) the new confirmed cases and deaths for each day are plotted. An initial increase is observed in September. After that in November there is a very big spike. This was the second outbreak in Greece when the second lockdown was imposed. In December the decrease is probably the result of the restrictions of the previous days.

```

1 gg_daily1 <- ggplot(greece, aes(x=date, y=confirmed.inc)) + geom_line() +
2   theme_classic() +
3   labs(title = "Covid-19 Daily Confirmed Cases in Greece", x= "Date", y= "Confirmed
4     cases") +
5   theme(plot.title = element_text(hjust = 0.5)) +
6   scale_x_date(date_breaks="1 month", date_labels = "%b")
7 gg_daily2 <- ggplot(greece, aes(x=date, y=deaths.inc)) + geom_line() +
8   theme_classic() +
9   labs(title = "Covid-19 Daily Confirmed Deaths in Greece", x= "Date", y= "Confirmed
10    deaths") +
11   theme(plot.title = element_text(hjust = 0.5)) +
12   scale_x_date(date_breaks="1 month", date_labels = "%b")
13 grid.arrange(gg_daily1, gg_daily2, ncol=2)

```

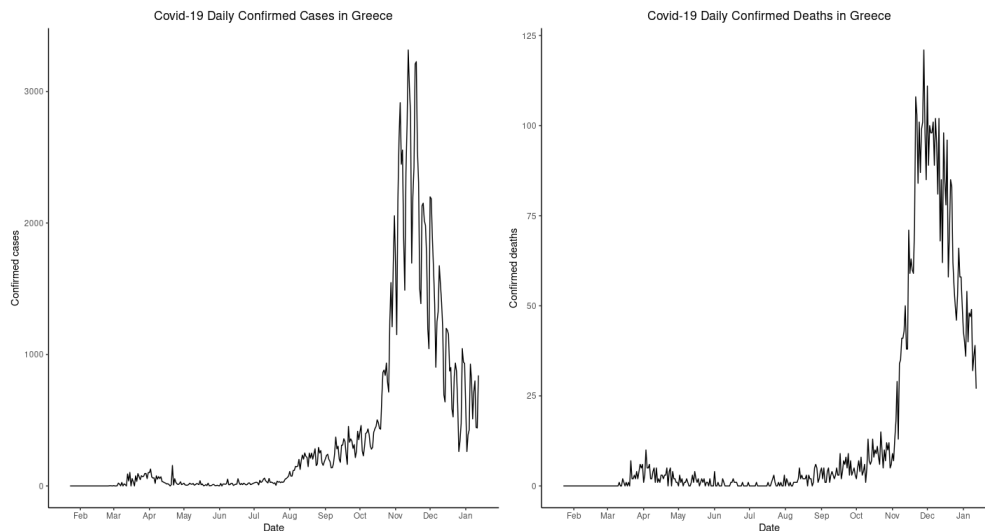


Figure 2.15: Greece: daily confirmed cases and deaths

In the following bar plot [Figure 2.16](#) mean cases and deaths according to each day of week are presented. It is easily observed that Monday and Sunday are the days with the least confirmed cases probably due to the smaller number of tests that took place in the previous days (Saturday and Sunday accordingly). But this tendency is not followed by deaths that are not dependent to the tests conducted.

```
1 # Greece results per day of week
2 days_all <- greece %>% mutate(daysofweek=weekdays(date))
3 days_all <- na.omit(days_all)
4 days_all <- days_all[,.(confirmed = mean(confirmed.inc),deaths=mean(deaths.inc)),by=
  daysofweek]
5 days_all <- days_all %>% mutate(death_rate=100*deaths/confirmed)
6 days_all$daysofweek <- factor(days_all$daysofweek, levels = c("Monday","Tuesday","
  Wednesday","Thursday","Friday","Saturday","Sunday"))
7 g14 <- ggplot(days_all,aes(x = daysofweek,y=confirmed,fill=confirmed))+
8   geom_bar(stat='identity')+
9   geom_text(aes(label=round(confirmed, 2), y=confirmed), size=2, vjust=0)+
10  labs(x = "Day", y="Confirmed cases",title="Greece: Mean of confirmed cases per day
  of week")+
11  scale_fill_gradientn(colours = rev(heat.colors(7)))+
12  theme(plot.title =element_text(hjust=0.5),axis.text.x=element_text(hjust=0.5))
13 g15 <- ggplot(days_all,aes(x = daysofweek,y=deaths,fill=deaths))+
14   geom_bar(stat='identity')+
15   geom_text(aes(label=round(deaths, 2), y=deaths), size=2, vjust=0)+
16   labs(x = "Day", y="Deaths",title="Greece: Mean of deaths per day of week")+
17   scale_fill_gradientn(colours = rev(heat.colors(7)))+
18   theme(plot.title =element_text(hjust=0.5),axis.text.x=element_text(hjust=0.5))
19
20 grid.arrange(g14,g15, nrow=2)
```

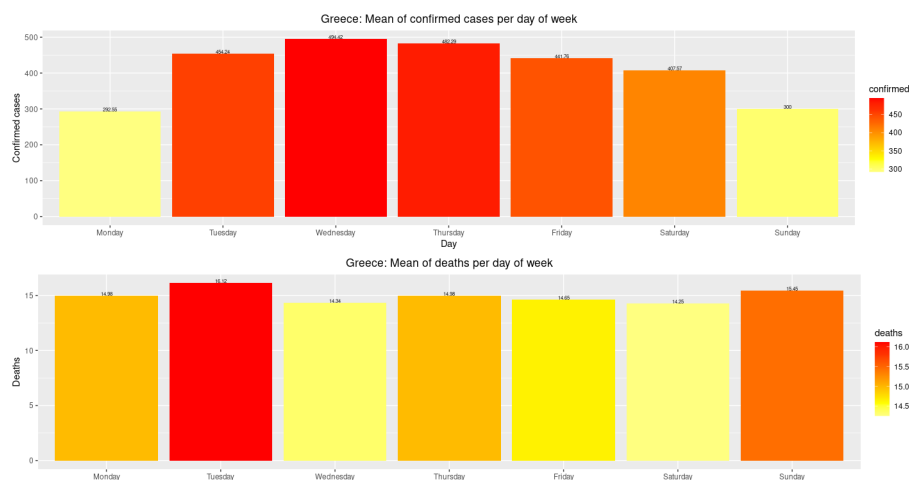


Figure 2.16: Greece confirmed cases and deaths versus day of week

Next there are two box plots [Figure 2.17](#) for the distribution of the number of new cases and deaths in Greece every day. The median value for daily confirmed cases is somewhere lower than 100 and for deaths 0-5. This is affected by the fact that in the first outbreak of the pandemic Greece had taken restriction measurements immediately and for the first months the number of cases was very small. Also the interquartile range (IQR) for cases is 0-430 which is small according to last numbers of cases. Bigger numbers of cases/deaths are presented outside IQR because they are much less than the observations with high values due to the first part of the pandemic where Greece had small numbers of cases.

```

1 # Distribution for Greece
2 g29 <- ggplot(greece, aes(Country, confirmed.inc))+
3   geom_boxplot(varwidth=T, fill="plum") +
4   labs(title="Box plot for Greece",
5        subtitle="Confirmed cases per day",
6        x="",
7        y="Confirmed cases")+
8   background_grid()+
9   theme(panel.grid.minor = element_line(colour="white", size=0.2))+
10  scale_y_continuous(minor_breaks = seq(0, 4000, 20), breaks = seq(0, 4000, 100))
11
12 g30 <- ggplot(greece, aes(Country, deaths.inc))+
13   geom_boxplot(varwidth=T, fill="plum") +
14   labs(title="Box plot for Greece",
15        subtitle="Confirmed deaths per day",
16        x="",
17        y="Confirmed deaths")+
18   background_grid()+
19   scale_y_continuous(minor_breaks = seq(0, 130, 5), breaks = seq(0, 130, 10))
20 grid.arrange(g29,g30, ncol=2)

```

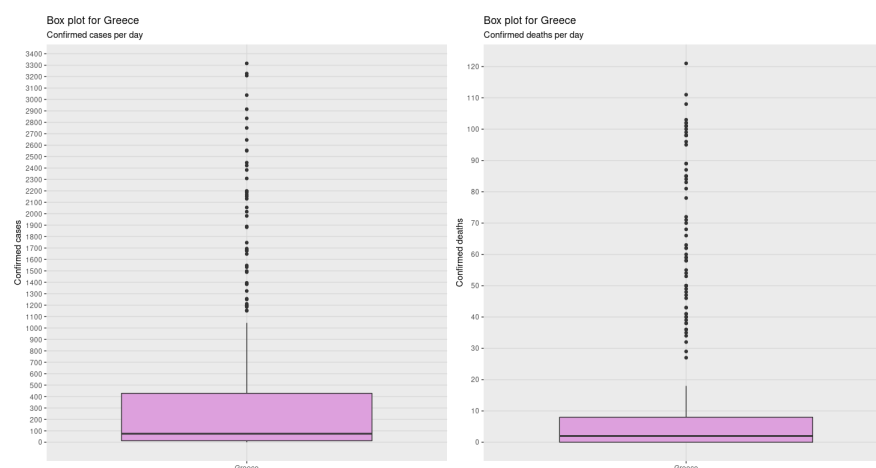


Figure 2.17: Boxplot Greece deaths and confirmed cases