

Covid19 PT Report

Vasco Pereira

2020-05-01

Synopsis

This is a data analysis report of the **public available** data about Covid19 in Portugal. The report demonstrates how to download and process data, ranks Portugal by cases and deaths and demonstrates the pandemic evolution in this country. Some final remarks are made about the limited availability of the data provided.

This report is also an example of reproducible research in data analysis making it possible to anyone to reproduce or adapt for any country.

Raw Data

World Data Source

The primarily data source used for this work was available by the **European Centre for Disease Prevention and Control**. By analyzing this data we can get a gist of the evolution of covid19 in the world. Some information were added trough time to this data, such as new formats. I choose the JSON format.

```
library(jsonlite)
library(dplyr)
url <- "https://opendata.ecdc.europa.eu/covid19/casedistribution/json"
dataRaw <- read_json(url, simplifyVector = TRUE)
data <- as_tibble(dataRaw$records)
str(data)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':  14242 obs. of  11 variables:
## $ dateRep      : chr  "01/05/2020" "30/04/2020" "29/04/2020" "28/04/2020" ...
## $ day          : chr  "1" "30" "29" "28" ...
## $ month        : chr  "5" "4" "4" "4" ...
## $ year         : chr  "2020" "2020" "2020" "2020" ...
## $ cases        : chr  "222" "122" "124" "172" ...
## $ deaths       : chr  "4" "0" "3" "0" ...
## $ countriesAndTerritories: chr  "Afghanistan" "Afghanistan" "Afghanistan" "Afghanistan" ...
## $ geoId        : chr  "AF" "AF" "AF" "AF" ...
## $ countryterritoryCode  : chr  "AFG" "AFG" "AFG" "AFG" ...
## $ popData2018   : chr  "37172386" "37172386" "37172386" "37172386" ...
## $ continentExp  : chr  "Asia" "Asia" "Asia" "Asia" ...
```

World Data Pre-Processing - getting Tiddy

Since all available variables are in character format I transformed the numeric values and date values into their respective formats.

```
library(lubridate)
data$dateRep <- dmy(data$dateRep)
data$day <- as.numeric(data$day)
data$month <- as.numeric(data$month)
data$year <- as.numeric(data$year)
data$cases <- as.numeric(data$cases)
data$deaths <- as.numeric(data$deaths)
data$popData2018 <- as.numeric(data$popData2018)
str(data)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':  14242 obs. of  11 variables:
## $ dateRep      : Date, format: "2020-05-01" "2020-04-30" ...
## $ day          : num  1 30 29 28 27 26 25 24 23 22 ...
## $ month        : num  5 4 4 4 4 4 4 4 4 4 ...
## $ year         : num  2020 2020 2020 2020 2020 2020 2020 2020 2020 2020 ...
## $ cases        : num  222 122 124 172 68 112 70 105 84 61 ...
## $ deaths       : num  4 0 3 0 10 4 1 2 4 1 ...
## $ countriesAndTerritories: chr  "Afghanistan" "Afghanistan" "Afghanistan" "Afghanistan" ...
## $ geoId        : chr  "AF" "AF" "AF" "AF" ...
## $ countryterritoryCode  : chr  "AFG" "AFG" "AFG" "AFG" ...
## $ popData2018   : num  37172386 37172386 37172386 37172386 37172386 ...
## $ continentExp  : chr  "Asia" "Asia" "Asia" "Asia" ...
```

Now we can work with these numbers and make some exploratory analysis. Let's check out a summary of our data:

```
summary(data)
```

```
##      dateRep      day      month      year
## Min.   :2019-12-31  Min.   : 1.00  Min.   : 1.000  Min.   :2019
## 1st Qu.:2020-02-22  1st Qu.: 9.00  1st Qu.: 2.000  1st Qu.:2020
## Median :2020-03-27  Median :17.00  Median : 3.000  Median :2020
## Mean   :2020-03-17  Mean   :16.59  Mean   : 3.073  Mean   :2020
## 3rd Qu.:2020-04-14  3rd Qu.:24.00  3rd Qu.: 4.000  3rd Qu.:2020
## Max.   :2020-05-01  Max.   :31.00  Max.   :12.000  Max.   :2020
##
##      cases      deaths  countriesAndTerritories  geoId
## Min.   : -1430.0  Min.   :  0.00  Length:14242  Length:14242
## 1st Qu.:   0.0  1st Qu.:  0.00  Class :character  Class :character
## Median :   1.0  Median :  0.00  Mode  :character  Mode  :character
## Mean   :  225.7  Mean   : 16.33
## 3rd Qu.:  29.0  3rd Qu.:  1.00
## Max.   :48529.0  Max.   :4928.00
##
## countryterritoryCode  popData2018  continentExp
## Length:14242  Min.   :1.000e+03  Length:14242
## Class :character  1st Qu.:2.790e+06  Class :character
## Mode  :character  Median :9.942e+06  Mode  :character
```

```
##           Mean      :5.520e+07
##           3rd Qu.   :3.717e+07
##           Max.      :1.393e+09
##           NA's      :141
```

Ok, something awkward is going on. The **new cases** variable have negative values, because, as one may notice, the minimum value is **-1430**. There shouldn't be negative cases, and there is not an explanation available anywhere for this observations.

>I tried to address this issue with European Centre for Disease Prevention and Control without response.

One thing to always have in mind is the source of the data: our source must be **reliable and trustworthy**.

I will consider this some sort of correction to the number of reported cases and not mess with this awkward input values for the sake of a globally cohesive information.

Now that I have a **somewhat** tidy data, I can think a bit about the variables available:

```
names(data)
```

```
## [1] "dateRep"           "day"
## [3] "month"             "year"
## [5] "cases"             "deaths"
## [7] "countriesAndTerritories" "geoId"
## [9] "countryterritoryCode" "popData2018"
## [11] "continentExp"
```

The most interesting data I will address will be the number of detected **cases** and **deaths**, by **date reported**.

Portugal in the World

How is Portugal rated in death and cases counts? There is some social media discussion about how to address this rates. Let us compare absolute numbers with percentages related to country population:

```
CasesRank <- data %>%  
  group_by(countriesAndTerritories) %>%  
  summarise("TotalCases" = sum(cases)) %>%  
  arrange(desc(TotalCases))  
  
PTrankCases <- grep("Portugal", CasesRank$countriesAndTerritories)  
  
numberOfCases <- CasesRank %>% filter(countriesAndTerritories == "Portugal")  
numberOfCases <- as.integer(numberOfCases[2])  
  
knitr::kable(head(CasesRank, 10))
```

countriesAndTerritories	TotalCases
United_States_of_America	1069826
Spain	213435
Italy	205463
United_Kingdom	171253
Germany	159119
France	129581
Turkey	120204
Russia	106498
Iran	94640
Brazil	85380

The above table reflects the ranked 10 worst countries in COVID-19 diagnosed cases. Portugal is ranked **18**, with **25056** total cases.

```
DeathsRank <- data %>%
  group_by(countriesAndTerritories) %>%
  summarise("TotalDeaths" = sum(deaths)) %>%
  arrange(desc(TotalDeaths))

PTrankDeaths <- grep("Portugal", DeathsRank$countriesAndTerritories)

numberOfDeaths <- DeathsRank %>% filter(countriesAndTerritories == "Portugal")
numberOfDeaths <- as.integer(numberOfDeaths[2])

knitr::kable(head(DeathsRank, 10))
```

countriesAndTerritories	TotalDeaths
United_States_of_America	63006
Italy	27967
United_Kingdom	26771
Spain	24543
France	24376
Belgium	7594
Germany	6288
Iran	6028
Brazil	5901
Netherlands	4795

The above table reflects the ranked 10 worst countries in COVID-19 deaths. Portugal is ranked **21**, with **989** total deaths.

Let us now check the same rates in percentage.

```
CasesPerRank <- data %>%
  group_by(countriesAndTerritories) %>%
  summarise("TotalCases" = sum(cases), "Pop" = unique(popData2018)) %>%
  mutate("PercentageCases" = TotalCases/Pop * 100) %>%
  arrange(desc(PercentageCases))

DeathsPerRank <- data %>%
  group_by(countriesAndTerritories) %>%
  summarise("TotalDeaths" = sum(deaths), "Pop" = unique(popData2018)) %>%
  mutate("PercentageDeaths" = TotalDeaths/Pop * 100) %>%
  arrange(desc(PercentageDeaths))

PTPerRankCases <- grep("Portugal", CasesPerRank$countriesAndTerritories)
PTPerRankDeaths <- grep("Portugal", DeathsPerRank$countriesAndTerritories)

knitr::kable(head(CasesPerRank, 10))
```

countriesAndTerritories	TotalCases	Pop	PercentageCases
Cases_on_an_international_conveyance_Japan	696	3000	23.2000000
San_Marino	569	33785	1.6841794
Holy_See	11	1000	1.1000000
Andorra	745	77006	0.9674571
Luxembourg	3784	607728	0.6226470
Iceland	1797	353574	0.5082387
Qatar	13409	2781677	0.4820473
Spain	213435	46723749	0.4568020
Gibraltar	144	33718	0.4270716
Belgium	48519	11422068	0.4247830

```
knitr::kable(head(DeathsPerRank, 10))
```

countriesAndTerritories	TotalDeaths	Pop	PercentageDeaths
Cases_on_an_international_conveyance_Japan	7	3000	0.2333333
San_Marino	41	33785	0.1213556
Belgium	7594	11422068	0.0664853
Andorra	42	77006	0.0545412
Spain	24543	46723749	0.0525279
Italy	27967	60431283	0.0462790
United_Kingdom	26771	66488991	0.0402638
France	24376	66987244	0.0363890
Sint_Maarten	13	41486	0.0313359
Netherlands	4795	17231017	0.0278277

There are evident differences from percentage to absolute numbers in the extremes (but in fact Portugal doesn't change much at the time of this report writing). In percentage Portugal is in **22** place for cases and **22** for deaths. Is this a fair comparison? There may be missing variables to understand our data: some index of number of urban centers per country for example, and also the predominance of respiratory

diseases, atmospheric pollution and elderly people percentage. Also, the number of tests each country do highly influence the reported cases. All this summed together can cause a great impact at the political level for managing the emergency states of each countries and interpreting the results. With too many unknown factors the effectiveness of some policies may result only by chance.

Since there are just too many unknown factors in the percentage rate maybe the fairest way to compare is the absolute numbers (aware of not being ideal as well).

Evolution of the Disease in Portugal and Simulations

In the next lines I try to simulate new cases and deaths from COVID-19 in Portugal with a 15 day advance from the last data reported and to identify a peak, where we could consider a turning point for the pandemic in Portugal, meaning that the contingency politics are taking effect.

```
PTdata <- filter(data, geoId == "PT")
PTdataArranged <- arrange(PTdata, dateRep)

library(ggpmisc)

x <- 1:length(PTdataArranged$dateRep)
y <- log10(PTdataArranged$cases)

xsq <- x^2
xcub <- x^3

fit <- lm(y~x+xsq+xcub)

xv <- seq(min(x), 100, 1)
yv <- predict(fit, list(x = xv, xsq = xv^2, xcub = xv^3))

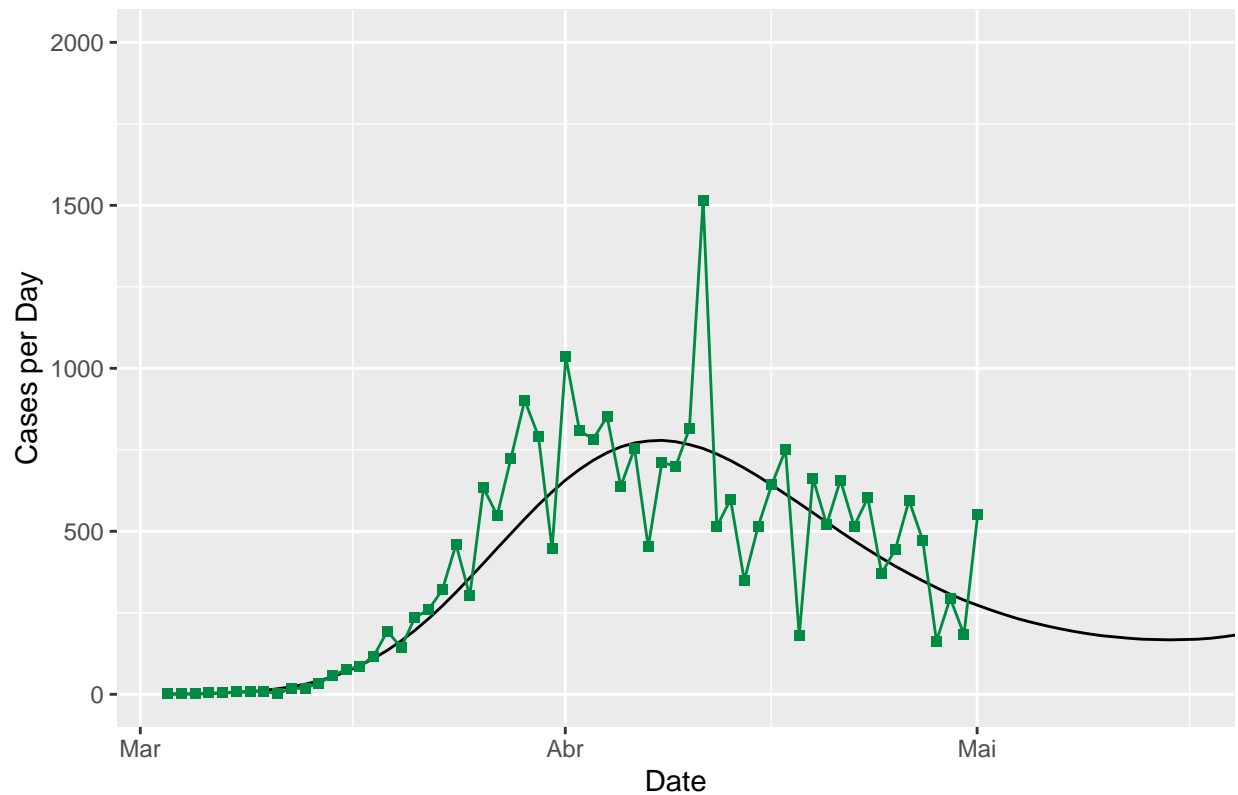
#Prediction <- tibble(Day = xv, logCases = yv)
PredictionCases <- tibble(Day = as.Date("2020-03-02")+xv,
                          SimCases = as.integer(10^yv),
                          RealCases = c(PTdataArranged$cases, rep(NA,
                                                                    100-length(PTdataArranged$cases))))
PredictionCases$Day <- as.POSIXct(PredictionCases$Day)

CasesMaxDay <- as.Date(PredictionCases$Day[
  grep(max(PredictionCases$SimCases[1:length(PTdataArranged$cases)]),
        PredictionCases$SimCases)])

curvePredict <- ggplot(PredictionCases, aes(Day, SimCases))
PTgSimCasesNEW <- curvePredict + geom_line() +
  geom_point(aes(Day, RealCases), col = "springgreen4", pch = 15) +
  geom_line(aes(Day, RealCases), col = "springgreen4") +
  labs(y = "Cases per Day", x = "Date",
        title = "Portugal new Cases Simulation") +
  stat_peaks(col = "tomato3", ignore_threshold = .9) +
  theme(plot.title = element_text(hjust = 0.5)) +
  coord_cartesian(xlim = c(PredictionCases$Day[1], as.POSIXct(Sys.Date() + 15)),
                  ylim = c(0, 2000))

PTgSimCasesNEW
```


Portugal new Cases Simulation



```

z <- log10(PTdataArranged$deaths)
z <- sub("-Inf", "0", z)

fitD <- lm(z~x+xsq+xcub)

xv <- seq(min(x), 100, 1)
zv <- predict(fitD, list(x = xv, xsq = xv^2, xcub = xv^3))

PredictionDeaths <- tibble(Day = as.Date("2020-03-02")+xv,
                           SimDeaths = as.integer(10^zv),
                           RealDeaths = c(PTdataArranged$deaths, rep(NA,
                                                                       100-length(PTdataArranged$deaths))))

PredictionDeaths$Day <- as.POSIXct(PredictionDeaths$Day)

DeathsMaxDay <- as.Date(PredictionDeaths$Day[
  grep(max(PredictionDeaths$SimDeaths[1:length(PTdataArranged$deaths)]),
        PredictionDeaths$SimDeaths)])

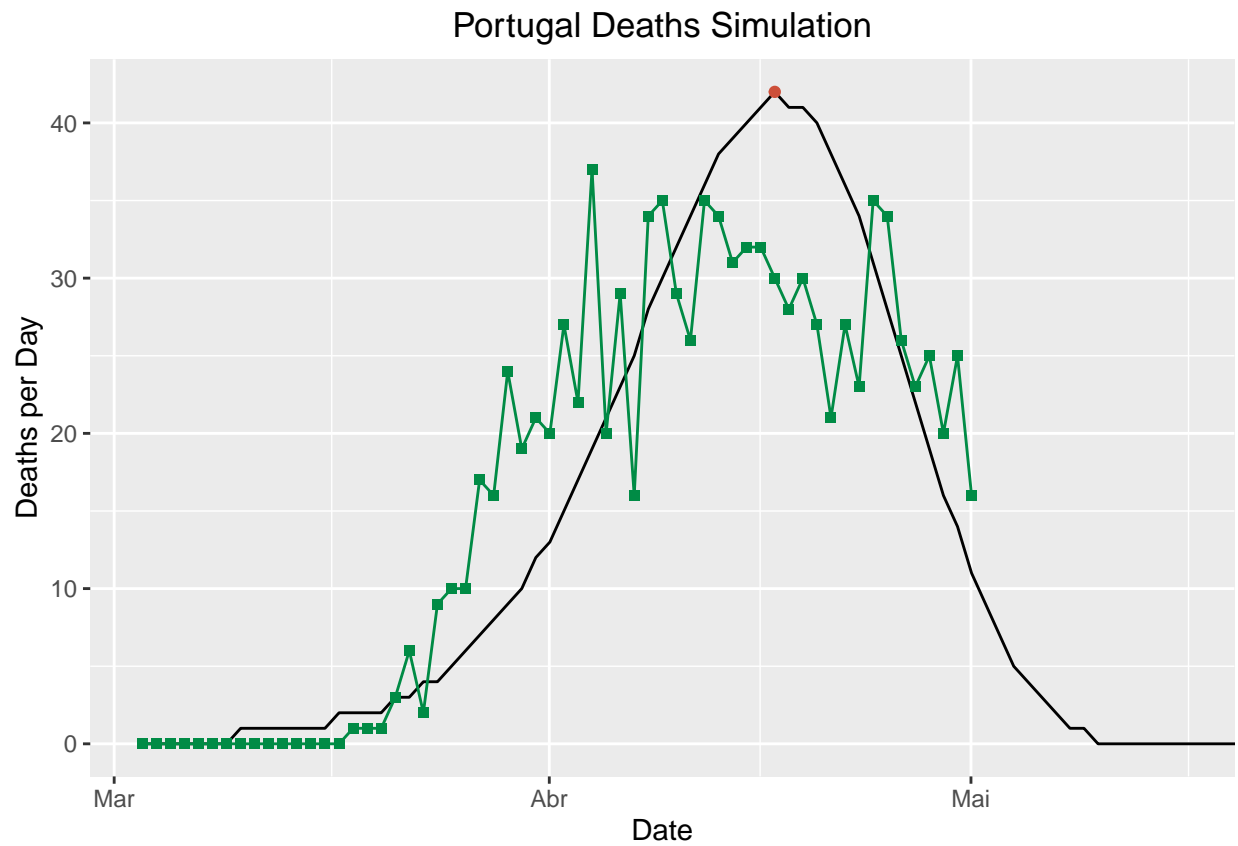
curvePredictDeaths <- ggplot(PredictionDeaths, aes(Day, SimDeaths))
PTgSimDeaths <- curvePredictDeaths + geom_line() +
  geom_point(aes(Day, RealDeaths), col = "springgreen4", pch = 15) +
  geom_line(aes(Day, RealDeaths), col = "springgreen4") +
  labs(y = "Deaths per Day", x = "Date",
       title = "Portugal Deaths Simulation") +
  stat_peaks(col = "tomato3", ignore_threshold = .7) +

```

```

theme(plot.title = element_text(hjust = 0.5)) +
coord_cartesian(xlim = c(PredictionDeaths$Day[1], as.POSIXct(Sys.Date() + 15)))
PTgSimDeaths

```



As we can observe from the graphs, the worst seems to have passed. The peak for new cases detected was in **2020-04-08** and the peak for deaths in Portugal was in **2020-04-17**.

Conclusions

- Cases: Portugal is ranked **18** in 207 countries, with **25056** total cases. Portugal had the peak of cases in **2020-04-08**.
- Deaths: Portugal is ranked **21** in 207 countries, with **989** total deaths. Portugal had the peak of deaths in **2020-04-17**.

The fact that the “peaks” have passed for the cases and deaths in Portugal doesn’t mean that the problem is over. For example, the peaks are always moving around because every day the simulations are iterated with the new information available, and new policies and the end of the confinement may produce more peaks for new cases and deaths.

From comparing our country with the others, we are quite better than most countries. This is due to several factors, not only policies but also with our peripheric geo-location. Also, it must be observed the predominance of the disease in the north hemisphere of the planet - so our rank may improve greatly for the next months (our rank in the end of April is about 18th to 19th place), if there is some influence with the flu-season.

Final Considerations

There is a lack of trustworthy data sources that would complement this data and give us detailed information about the **how’s** and **why’s** of SARS-CoV-2 behavior. Interesting variables that I would consider worth studying would be: **Active Cases**, **Recovered Cases**, **ICU Cases**, **Non-ICU Hospitalized Cases**, and also some information about patients, like: **Known Diseases**, a logical or more informative **IfSmoker** variable and **Air Pollution Exposure**.

Some of this information is retained by the “Sistema nacional de vigilância epidemiológica” **SINAVE**

I’ve applied for access in 2020-04-27 from an institutional e-mail but not a reply until today.