

# Assignment 1

Zsombor Hegedus

2020 november 25

## Introduction

In this analysis my focus is uncovering a pattern of association between the number of confirmed COVID-19 cases and the number of deaths caused by the virus. In this analysis I will use admin data, where variables, such as the number of deceased and infected people are from the github repo of Johns Hopkins University, and population data is collected from World Development Indicators maintained by the World Bank. The population of the analysis covers every country where such data was available (each country is one observation in my dataset), while the sample used for the analysis is narrowed to one date - the 15th of Oct, 2020. Potential data quality issues could arise from the fact that:

- Joining the two dataframes required entity resolution which is a manual process that always leaves room for data quality errors
- There might be a bias in how countries collected or disclosed their information (e.g. some countries might not conduct enough tests hence the number of confirmed cases are low - or they might be using different criterion to identify a COVID-19 death case)

After entity resolution, I ended up with 168 observations in my data, where I dropped out observations with zero reported death cases since I did  $\ln$  transformation where records with zero deaths could not have been used. I used *per capita* variables so divided deaths and confirmed cases by the population - expressed in millions - of given country. My dependent variable  $y$  was the number of deaths per capita caused by COVID-19 and my explanatory variable  $x$  was the number of confirmed cases per capita in relation with the disease.

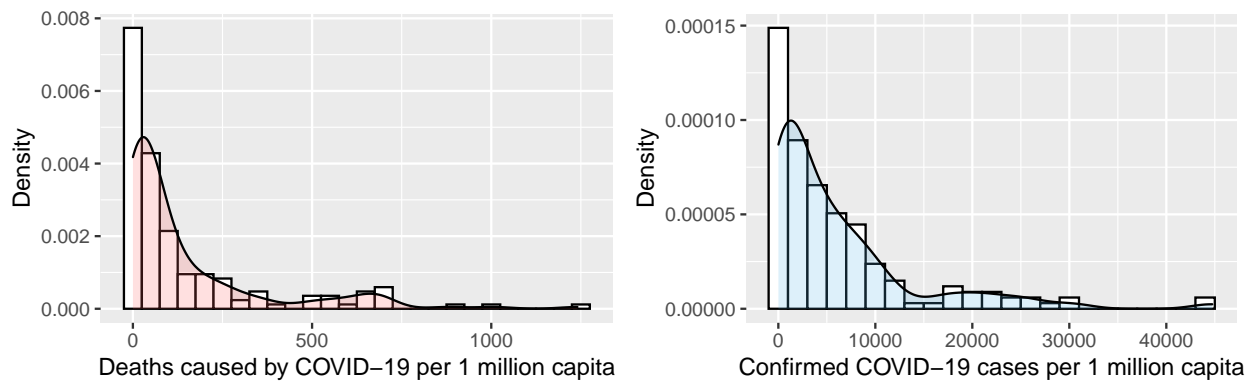


Figure 1: Distribution of deaths and confirmed cases in relation to COVID-19 for the examined countries on 15th Oct, 2020

variable	mean	median	std	iq_range	min	max	skew	numObs
Death per capita	142.92	51.38	218.83	148.00	0.09	1240.40	2.31	168
Confirmed cases per capita	6058.74	3209.04	8009.23	7371.35	8.78	44734.82	2.31	168

Table 1: Summary statistics of examined variables

## EDA

Figure 1 shows the distribution of variables which is skewed with a fat right tail and a few extreme values for both. The number of confirmed cases (visible in Table 1) far exceed the deaths - in my sample, ca.2.4% of the people that contracted COVID-19 died from the disease. Since the density plots resemble a lognormal distribution, it indicates that ln transformations for both variables might strengthen the results of the analysis.

A substantive reason for choosing log-log model is that my variables are affected in multiplicative ways - we care about percentage change when talking about the number of infected and fallen people. Maybe we could argue that e.g. the absolute number of people died is an important measure in case a given country doesn't really have a lot of reported cases, but to have a hollistic view on this phenomenon, it makes more sense to look at percentages.

Statistical reasons for choosing ln-ln transformation is that:

- the distributions of both variables are skewed with a long right tail which resemles a log-normal distirbution
- comparing scatter plots with and without ln transformation resulted in a much better fit for the model ( $R^2$  was about twice as good)

## Model Choice

For the analysis I chose a simple linear regression between the log transformed y and x variables:

$$\ln(y) = \alpha + \ln(x)$$

The estimated parameters can be found in the Appendix.  $\alpha$  is not really meaningful as  $\ln(y)$  is difficult to interpret, but it shows how much is the ln of death per 1 million capita in case we have 1 confirmed case per 1 million capita.  $\beta$  is the slope parameter, and my model indicates that in this sample the number of death per 1m capital is 0.93% higher on average for observations having 1% higher confirmed cases per 1m capita.

I conducted hypothesis testing as I was interested to see if my estimated  $\beta$  parameter is significant at 5%:

$$H_0 : \beta = 0, \quad H_A : \beta \neq 0$$

The below table shows the result of the hypothesis test which resulted in a 95% CI of [0.84, 1.02], which means that I can reject the  $H_0$  with 95% confidence. The overall conclusion is that  $\beta$  is significant in the level that I was interested in (even more, as we can see from the p value), therefore there is a positive pattern of association between the variables subject to my analysis.

variable	estimate	std.error	statistic	p.value	conf.low	conf.high
ln_conf_pc	0.93	0.04	20.76	0.00	0.84	1.02

Table 2: Hypothesis testing for the slope of the regression

## Analysing residuals

**Countries with largest negative error:** add 2-3 sentences on these guys

Countries with lowest negative error : add 2-3 sentences on these guys

## Executive summary (3-5 sentences)

I set out to analyse the pattern of association between the number of confirmed COVID cases and the number of deaths caused by COVID-19 in all countries where this data is available. I used log-log transformation and scaled the number of confirmed and fallen people by the population of each country to arrive to 'per capita' variables. I used regression analysis to uncover these tendencies, and found in my sample for the 15th Oct, 2020 that there is a positive association between the two variables. The main message my model wish to convey is that understanding the number of people contracted by the virus is crucially important for governments, so that they can have an expectation as to how many people might die due to the disease. On the one hand my results could have been strengthened by including more variables such as the age of people (or a categorical variable representing age groups) to see whether the virus is deadlier for the older population. My results, on the other hand could have been weakend by using level instead of log variables, as it was quite apparent in the EDA phase that without ln transformation my R squared would have dropped approximately to the third of what I was able to achieve with my final model.

## Appendix

### Transforming variables

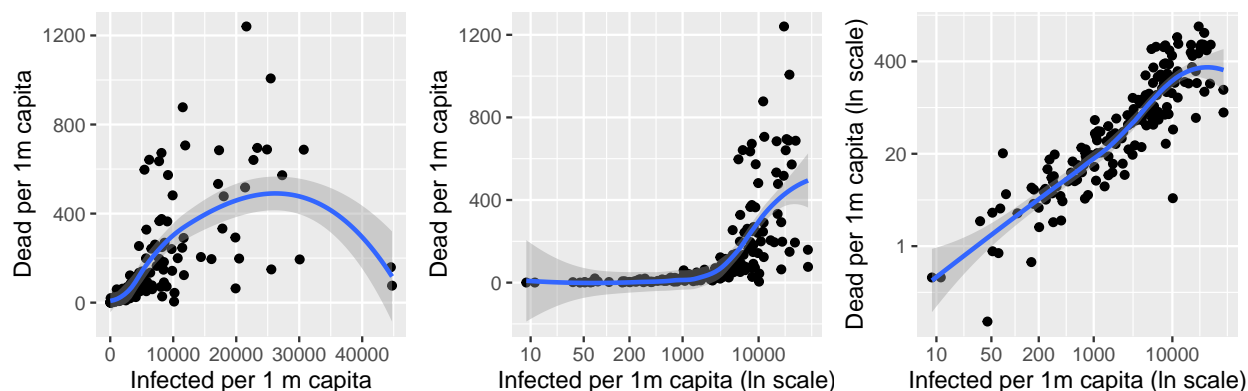


Figure 2: Scatter plots between deaths and confirmed cases in relation to COVID-19 for the examined countries on 15th Oct, 2020 - lowess functions

### Model Comparison

Model comparison table with scatter plot visaulisation with explanation on what you can see in the table 5-8 sentences - 4 model descr Stating choice of model - substantive and statistical reasoning 3-5 sentences

```
##
## =====
##           Model 1           Model 2           Model 3           Model 4
## -----
## (Intercept)          -3.42 *           -3.06 *           -2.67 *           -2.91 *
##                   [-4.11; -2.73]      [-4.95; -1.16]      [-4.36; -0.97]      [-4.20; -1.62]
## ln_conf_pc           0.93 *             0.82 *
##                   [ 0.84;  1.02]      [ 0.29;  1.35]
## ln_conf_pc_sq
##                   0.01
##                   [-0.03;  0.04]
## lspline(ln_conf_pc, ln_cutoff)1
##                               0.77 *
```

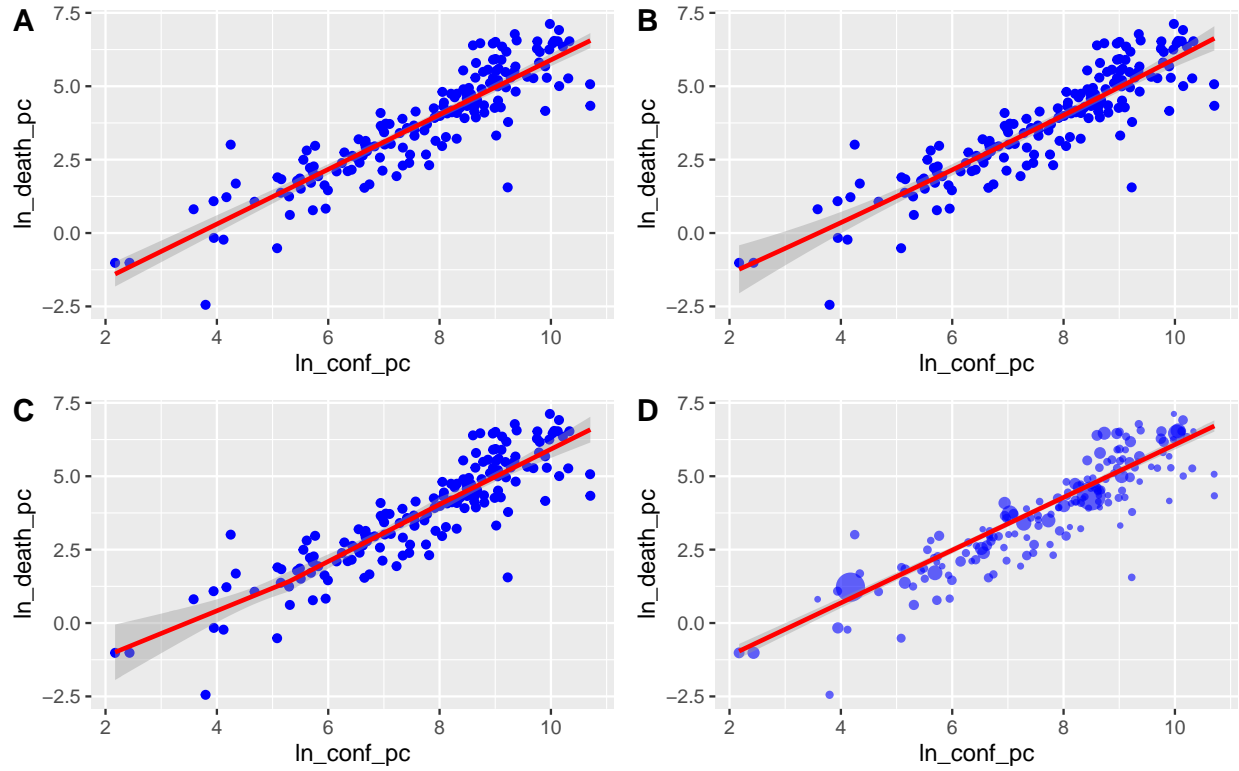


Figure 3: Comparing regression models for confirmed and died COVID-19 cases for 15th Oct, 2020

```
## [ 0.45; 1.09]
## lspline(ln_conf_pc, ln_cutoff)2 0.97 *
## [ 0.81; 1.13]
## lspline(ln_conf_pc, ln_cutoff)3 0.94 *
## [ 0.68; 1.20]
## -----
## R^2 0.79 0.79 0.79 0.90
## Adj. R^2 0.79 0.79 0.79 0.90
## Num. obs. 168 168 168 168
## RMSE 0.83 0.83 0.83 4176.82
## =====
## * Null hypothesis value outside the confidence interval.
```