

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, 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 cases with zero reported death cases since I did ln transformation where records with zero death could not have been defined. 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.

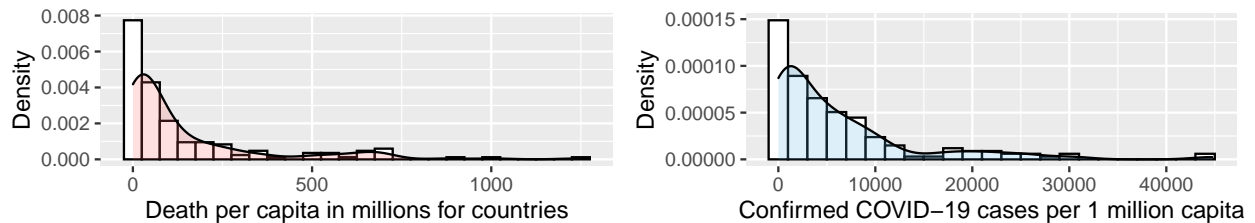


Figure 1: Distribution of death and confirmed cases in relation to Covid-19 on 15th Oct, 2020

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 'x' and 'y' variables

EDA

Figure 1 shows the distribution of variables which is skewed with a fat right tail and a few extreme values (e.g. USA etc...) 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 is an indicator that ln transformations might strengthen the results for both variables in the analysis.

Substantive reasoning for choosing log-log model:

- distribution of variables doesn't follow a normal distribution,
- something else
- something else

Statistical reasons for choosing ln-ln transformation:

- one
- and a two

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. α is not really meaningful to interpret, but it shows how many deaths can be associated in a case where we have 1 confirmed case. β is the slope parameter, and it states that in this sample $\beta\%$ increase in the number of confirmed cases is associated with 1% higher number of deceased people due to Covid-19.

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

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

The below table shows the result of the regression which resulted in a 95% CI of [3.4, 3.23], which means that I can reject the H_0 with 95% confidence. The overall conclusion is that β is significant in the level that I was interested in, therefore there is a positive pattern of association between the ln of confirmed covid-19 cases and the ln of cases resulting in deaths.

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 see 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 weakened 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

Model comparison table with scatter plot visualisation 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