



# Predicting The Incidence Rate And Case Fatality Rate Of COVID-19

by  
Mike Weltevrede (SNR 1257560)

A thesis submitted in partial fulfillment of the requirements for the degree of Master in  
Econometrics and Mathematical Economics.

Tilburg School of Economics and Management  
Tilburg University

Supervised by:  
dr. Otilia Boldea

Second reader:  
dr. George Knox

Date:  
June 15, 2020



TODO

**Abstract**

## Acknowledgements

TODO

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Problem description</b>	<b>2</b>
<b>3</b>	<b>Methodology</b>	<b>4</b>
3.1	Model 1: Within-Region Spread . . . . .	4
3.2	Model 2: Weighted Within-Region Spread . . . . .	4
3.3	Model 3: Within and Between-Region Spread . . . . .	5
3.4	Model 4: Full Model . . . . .	5
<b>4</b>	<b>Dataset</b>	<b>7</b>
4.1	Coronavirus data . . . . .	7
4.2	Independent variables . . . . .	9
<b>5</b>	<b>Results</b>	<b>11</b>
<b>6</b>	<b>Conclusion</b>	<b>13</b>
<b>7</b>	<b>Future research</b>	<b>14</b>
	<b>References</b>	<b>15</b>
	<b>Appendices</b>	<b>17</b>
<b>A</b>	<b>Tables</b>	<b>17</b>

# 1 Introduction

## 2 Problem description

In this section, we will elaborate on the methods that we apply in order to model the epidemiological spread of SARS-CoV-2 and the disease it causes: COVID-19. We are basing our model on specifications as used by Adda (2016). In the paper, Adda (2016) investigates the spread of several viral diseases in the past, namely influenza, gastroenteritis, and chickenpox. The key additions made are, firstly, that a spatial spillover effect is considered and, secondly, that we allow for some sort of weighting on the parameters on the basis of region specific variables. Adda (2016) starts from the Standard Inflammatory Response (SIR) model, the most commonly used model in epidemiology (Kermack & McKendrick, 1927; Anderson & May, 1992).

The SIR model splits the total population into three groups.  $S$  denotes the fraction of individuals who are susceptible to being infected,  $I$  denotes the fraction of individuals who are currently infected, also called infectives, and  $R$  denotes the fraction of individuals who have been removed from the model, be that because they successfully recovered from the disease or that they have deceased. Adda (2016) defines  $R$  to be the group of individuals who have recovered but who are still immune, i.e. the deceased people are not included in  $R$ .

Q: All other sources except for Adda look at the group  $R$  as the removed, i.e. people who overcame the disease but also deaths. It is not clear how Adda deals with deaths, although I suspect he subtracts them from the susceptible population and adjusts the total population accordingly. I think it is likely not an issue since the number of deaths is negligible compared to the size of the large population. (Also see the next TODO note)

As such, at any point in time, we have that

$$S, I, R \in [0, 1] \text{ and } S + I + R = 1.$$

The SIR model is postulated in continuous time, i.e. the equations in (2.1), (2.2), and (2.3) depict the change in the variables  $S$ ,  $I$ , and  $R$ , respectively, for one time period ahead. This type of model is also called a stock-and-flow model because there is a certain stock (for instance the number of infectives) to which a flow is added or subtracted.

$$\frac{dS}{dt} = -\alpha SI + \lambda R \tag{2.1}$$

$$\frac{dI}{dt} = \alpha SI - \beta I \tag{2.2}$$

$$\frac{dR}{dt} = \beta I - \lambda R \tag{2.3}$$

It is important to grasp the main assumptions of the SIR model, which also tell us how these equations are constructed. The first assumption that is made, is that the population is constant, meaning that births and deaths are ignored.

TODO: I do not; I take it into account in calculating  $S$ . Should see if this matters.

Next, note that the spread of the virus is determined by the interaction between the infectives and the susceptible population. The second assumption that is made under the SIR model in this light is that there is a constant rate of change in infectives that is proportional to this interaction between the infectives and the susceptible population. This is represented by the term  $\alpha SI$  in equations (2.1) and (2.2), which is also called the transmission term (Keeling & Rohani, 2011). The third assumption that the SIR model makes is that there is a constant rate of change at which infectives recover or deacease. This relates to the term  $\beta I$  in equations (2.2) and (2.3).

TODO: Look into this, if we do use the definition that people who die are included. Reason: the fuller hospitals are, the more people will likely decess.

Finally, we assume that there is a constant rate of change at which immune individuals lose their immunity. This is denoted by the term  $\lambda R$  in equations (2.1) and (2.3). For instance, Adda (2016) mentions that  $\lambda$  is set to 0 for chickenpox as individuals acquire a lifetime immunity while  $\lambda$  will be high for gastroenteritis due to almost no immunity emerging. In the case of COVID-19, some studies show that it is likely that individuals who recovered from COVID-19 may be immune to reinfection, at least temporarily (Kirkcaldy, King, & Brooks, 2020). Nonetheless, no definitive results have been shown.

TODO: Explain more later on immunity since this is currently still researched.

One of the main measures resulting from the SIR model is the estimation of the basic reproduction number  $R_0 := \alpha/\beta$ . An epidemic is said to develop if  $R_0 > 1$ .

TODO: Explain how  $R_0 > 1$  is determined (why not 2, for instance).

This measure is widely used to indicate that an ongoing epidemic is dying out if  $R_0$  drops below 1. For instance, the Italian health ministry has posted an article on May 9, 2020 stating that the  $R_0$  reproduction rate for COVID-19 was below 1 in Italy, at between 0.5 and 0.7 (Ministero della Salute, 2020a), showing that this measure is also used communicated to citizens as a way of informing them whether the pandemic is tending to end.



## 3 Methodology

### 3.1 Model 1: Within-Region Spread

We start with a simple model ignoring effects across regions. First, it is important to understand the concept of an incubation period. This is defined as the period between an infection and the moment that the infected individual starts showing symptoms, at which point the infective is said to be symptomatic. Note that this is not the same as the period between an infection and the moment that the infective is infectious, which is called the latent period. For COVID-19, the latent period is estimated to be approximately 2 days shorter than the incubation period (He et al., 2020). That is, infectives are able to infect others before showing symptoms.

This may be an issue when considering policies such as self-isolation when one is sick, because an infective may have already spread the virus before feeling sick. The incubation period for COVID-19 is estimated to be above 2 and below 11.5 (Lauer et al., 2020), 12.5 (Li et al., 2020), and 14 days (Linton et al., 2020). This is a large range, but this is not rare. For instance, .

While the maximum incubation period is not agreed upon by Lauer et al. (2020) and Li et al. (2020), their results on the median are similar. Lauer et al. (2020) reports a median incubation period of 5.1 days (95% CI: 4.5 to 5.8 days), while Li et al. (2020) reports a median incubation period of 5.2 days (95% CI: 4.1 to 7.0 days). For comparison, Linton et al. (2020) give the result of a mean incubation period of 5.0 days (95% CI: 4.2 to 6.0 days) when excluding Wuhan residents and 5.6 days (95% CI: 5.0 to 6.3 days) when including Wuhan residents.

TODO: Explain more later on how this affects our choice for  $\tau$ . We will include this in this section too, as we have tested multiple lags.

The within-region model is henceforth defined as:

$$I_{r,t} = \alpha_{within} I_{r,t-\tau} S_{r,t-\tau} + X_{r,t} \delta + \eta_{r,t} \quad (3.1)$$

where the subscript  $\tau$  is a lag indicating the length of the incubation period.

The matrix  $X$  includes fixed effects for regions, as well as weekend and week of the year dummy variables.

TODO: Elaborate on the definition of  $X$  and why we choose these variables, possibly in the Dataset section.

Lastly, we include an idiosyncratic error term  $\eta$ . The model is estimated by ordinary least squares (OLS). Because fixed effects for regions are included in  $X$ , note that this means that running OLS is actually a least-squares dummy variables (LSDV) regression. In general, the main issue with LSDV regression is that there needs to be an indicator variable for each observed individual (in our case, these are regions). However, it is feasible to run an LSDV regression since we consider a relatively small number of regions with a large number of time periods. This will be explained more in section 4.

### 3.2 Model 2: Weighted Within-Region Spread

In the previous model, it has been assumed that the incidence rate within a certain region is only determined by the previous incidence rates plus some other effects. However, the transmission rate  $\alpha$  is likely influenced by other factors as well. These may include policies, such as shutting down restaurants or public transport, but also persistent regional characteristics such as metrics on the quality of hospitals or economic development. In this section, we incorporate these factors in the within-region model (3.1). After defining the between-regions model in section 3.3, we will apply the same methodology to obtain the full weighted model in section 3.4.

Let the tensor  $W$  contain  $K$  region-specific variables that may influence the transmission rate  $\alpha$ . As such, we now allow for  $\alpha_{within}$  to differ for these  $K$  variables. In section 4, we elaborate on how these

TODO: Add to this with example

variables included in  $W$  are specifically defined and selected. For instance, we include the number of rail travellers, which changes over time, but also a measure of the development of health care through the number of available hospital beds, which does not change over time. We define  $X$  and  $\eta$  in the same way as for (3.1). Taking this into account, the weighted within-region model is defined as:

$$I_{r,t} = I_{r,t-lag} S_{r,t-lag} \sum_{k=1}^K \alpha_{within}^k W_{r,t-lag}^k + X_{r,t} \delta + \eta_{r,t} \quad (3.2)$$

TODO: Possibly update later

### 3.3 Model 3: Within and Between-Region Spread

A key addition made by Adda (2016) is recognizing that there is spatial spillover between regions. That is, there may be infectives in one region that travel to another region and then infect individuals there. As such, the number of new cases would be modeled as  $\alpha_{within} SI + \alpha_{between} S\tilde{I}$  where  $\tilde{I}$  is the fraction of infectives from outside the region of interest who meet susceptible people from within the region. Clearly, this is an important addition to the model and we acknowledge and incorporate this in this thesis.

TODO: Consider the difference in definition in  $I$  between the SIR model and Adda. Possibly use the notation from Keeling and Rohani (but this includes  $X$ )

The following model is defined:

$$\begin{aligned} I_{r,t} = & \alpha_{within} I_{r,t-lag} S_{r,t-lag} \\ & + \alpha_{between} \sum_{c \in R \setminus r} I_{c,t-lag} S_{r,t-lag} \\ & + X_{r,t} \delta + \eta_{r,t} \end{aligned} \quad (3.3)$$

In (3.3), the transmission parameter  $\alpha$  is now allowed to be different within and between regions. Adda (2016) estimates (3.3) by OLS and by instrumental variable estimation (IV). Weather episodes, such as the amount of rain and temperature-related instruments, are used as instruments. There is a biological reasoning behind choosing these instruments, for instance that warmer temperatures tend to have a negative effect on the proliferation of some viruses. A social reason is also given, namely that bad weather conditions impact the amount of social interaction between people, meaning that there are less opportunities for viruses to spread. We challenge the choice of these instruments, particularly in the case of SARS-CoV-2. Firstly, we do not have sufficient information on the effect of the weather on the virus. That is, SARS-CoV-2 has only been quite apparent since January 2020 and there has not been enough fluctuation over time in temperatures to show a necessary effect that can be disentangled from, for example, policies being effective in driving the virus back. Secondly, we challenge the social reasons entirely, although not quantitatively. In our view, bad weather conditions in themselves are not likely to be strong enough instruments for the viral spread. That is, even if they are indeed exogenous with respect to the error term and that they are correlated with the viral spread, we expect this to not be quite strong.

TODO: This is currently a claim and I have not looked at Adda's quantitative tests for these instruments.

For this reason, we only consider OLS for this model.

### 3.4 Model 4: Full Model

We now incorporate the between-region effects as well as the weighting of the transmission parameter. In addition to (3.2), we now also put weights on the between-region transmission parameter by some possibly influential variables. Let the tensor  $\tilde{W}$  contain  $\tilde{K}$  variables that now can influence the transmission rate  $\alpha_{within}$  between two regions  $r$  and  $c$ .

TODO: Possibly consider not following Adda's notation with the tildes and use something like V and L instead of  $\widetilde{W}$  and  $\widetilde{K}$ , respectively.

$$\begin{aligned}
I_{r,t} = & I_{r,t-lag} S_{r,t-lag} \sum_{k=1}^K \alpha_{within}^k W_{r,t-lag}^k \\
& + \sum_{c \in R \setminus r} I_{c,t-lag} S_{r,t-lag} \sum_{k=1}^{\widetilde{K}} \alpha_{between}^k \widetilde{W}_{r,c,t-lag}^k \\
& + X_{r,t} \delta + \eta_{r,t}
\end{aligned} \tag{3.4}$$

## 4 Dataset

In this section, we will outline the structure of the data that is used and how it was retrieved. Firstly, we discuss the structure of Italian regions and the reasons why we chose to use Italy as our region of interest. Subsequently, we will look at the data on COVID-19, such as the incidence rate. Here, we also discuss how we tackled possibly errors in the data, as well as missing values. Lastly, we consider the variables that are included in the weighted models in sections 3.2 and 3.4.

Italy has been one of the most intensely struck countries in the world. On June 6th, 2020, it had the seventh highest absolute number of cases, after the United States, Brazil, Russia, the United Kingdom, Spain, and India. Despite dropping in this positioning, Italy reports the highest death-to-cases ratio of 14.47%, followed closely by the United Kingdom, which reports a death-to-cases ratio of 14.18%. The sudden onset of the spread of SARS-CoV-2 put immense pressure on the Italian hospitals, especially in the northern regions such as Lombardy. This forced patients with coronavirus-caused pneumonia to be sent home as well as literal collapses of overworked healthcare workers (Horowitz, 2020).

The Italian ministry of Health Services (Ministero della Salute) has posted daily reports containing tables with a detailed numerical overview of new cases, active intensive care (IC) patients, and tests executed, all divided up between the second-level NUTS regions (also called NUTS 2 regions). The NUTS classification (Nomenclature of territorial units for statistics) is a hierarchical system for dividing up the economic territory of the European Union (EU) and the United Kingdom (Eurostat, 2020a) as used by Eurostat, the statistical office of the EU. Italy consists of 21 so-called *regioni* (regions), comparable to Dutch provinces. These constitute the second-level NUTS regions (also called NUTS 2 regions), where the region of *Trentino-Alto Adige* is split into two regions: *Provincia Autonoma di Bolzano/Bozen* and *Provincia Autonoma di Trento*. Italy's first-level NUTS regions are defined as groups of regions, namely *Nord-Ovest* (North West), *Nord-Est* (North East), *Centro* (Center), *Sud* (South), and *Isole* (Islands). The third-level NUTS regions are 107 provinces, which are subregions of the *regioni*. Ideally, we would want to have coronavirus data on the NUTS 3 regions since many policies are introduced at that level, such as a lockdown put into place on March 7th, 2020 until the strict national lockdown was instated. Unfortunately, the data was not reported at this granular level. As such, we chose to use the NUTS 2 regions.

TODO: Add table

### 4.1 Coronavirus data

As mentioned, the specific information on the coronavirus in Italian regions was retrieved from the Ministero della Salute, who publish daily reports under a title similar to *Covid-19, i casi in Italia: 10 giugno ore 18*, where *10 giugno* (June 10th) would be updated to the relevant date (Ministero della Salute, 2020b). These daily reports were posted with tables containing the following information per region:

- Hospitalized with symptoms (*Ricoverati con sintomi*)
- Active intensive care patients (*Terapia intensiva*)
- Home isolated active cases (*Isolamento domiciliare*)
- Total number of active cases (*Totale attualmente positivi*)
- Dismissed/recovered (*Dimessi/guariti*)
- Deceased (*Deceduti*)
- Total confirmed cases (*Casi totali*)

- Increase in total confirmed cases - compared to the previous day (*Incremento casi totali - rispetto al giorno precedente*)
- Total amount of tests executed (*Tamponi*)
- Total amount of persons tested (*Casi testati*)
- Increase in total amount of tests executed (*Incremento tamponi*)

The difference between the total amount of tests executed and the total amount of persons tested is that the latter indicates the number of unique persons that were tested. That is, individuals could have been tested more than once. Do note that *tamponi* is a good indication of the *testing capacity* as the number of tests that Italy is able to execute. Henceforth, when the term *testing capacity* is used, this refers to *tamponi*, unless indicated otherwise.

It should be noted that there is a measurement error in the number of infectives, as is the case in any other country. This is because there is no possibility that every citizen can be tested for COVID-19. For that reason, the actual number of infectives is higher than the official count as reported in the tables of the Ministero della Salute.

TODO: Talk about the functional form modelling of the unreported cases later

With respect to the reported death statistics, there is a distinction between Italy and some other European countries. Namely, the Italian numbers include deaths of all patients who were tested for COVID-19 before or after their death, regardless of whether they died inside or outside the hospital. In contrast, other countries may only count deaths in hospitals. French death counts, for instance, only have included deaths at hospitals and clinics caring for patients, excluding people who die at home or in care homes, although the French president Emmanuel Macron did announce that these centers would be tracked from the first week of April onward (Sevillano, 2020). Moreover, Italian data makes no distinction between people who died because of COVID-19 or simply had the disease but who died from other causes (also referred to as comorbidities). Patients who had pre-existing conditions actually make up around 96% of the total death count in Italy (Istituto Superiore di Sanità, 2020). In some other countries, such as Germany, a distinction between these two groups is actually made (Caccia, 2020).

We also make the note that it is unclear how the Ministero della Salute collects its information. If regions or provinces submit this information to the government each day, there may be areas that fail to submit their data for a certain day. Despite this, we assume that this official information is accurate and representative of the region for which it has been reported. If this is not the case, the numbers in the report on the next day will compensate for the error on the day before. In the official publications that we use, data that was wrongly published on a day  $t - 1$  is corrected by subtracting the error from or adding the error to the cases from day  $t$ . As such, if the error is larger than the number of new cases, the reported amount of new cases is negative. It happened fifteen times that a number was reported to be negative. Since negative numbers are not logical in the context of confirmed cases or deaths, we correct this by subtracting the error from the day before and set the previously negative number to 0. In the case that the error on day  $t$  is larger than the number on  $t - 1$ , for instance if a value of -10 is reported on day  $t$  while the value for day  $t - 1$  is less than 10, we propagate the error to multiple lags until this issue no longer occurs. For non-negative corrected numbers, we do not have a way to detect which these are and we cannot reasonably assume how this number should be split up among day  $t$  and  $t + 1$ . As such, these are left as is. One note that should be made is a highly negative value of -229 reported for the region of Campania on June 12, 2020, whereas the number of new cases in the week before that date only ranges from 0 to 5. We assume that this corrects for all errors in the past, not just those near to June 12. Propagating this error backwards as described before would lead to zero new cases per day for Campania from May 13, 2020 onward.

Q: Because the report above is highly unrealistic and we have no way of knowing how this correction is distributed across the days, should we delete the region of Campania entirely?

TODO: Update accordingly if this changes

Regarding missing data, there are only three cases, namely for Abruzzo on March 10, Puglia on March 16, and Campania on March 18. Given that faulty data is also corrected as described before, we assume that the cases missing on day  $t$  are added to those of day  $t + 1$ . This is confirmed by higher values compared to the expected trend, as seen in Table 4.1. As such, missing data is simply imputed with a value of 0.

TODO: Update accordingly if this changes

Table 4.1: Number of confirmed cases around a day  $t$  with missing data

	Abruzzo	Puglia	Campania
Day $t - 1$	8	64	60
Day $t + 1$	46	110	192
Day $t + 2$	5	43	97

## 4.2 Independent variables

Independent variables, or regressors, were obtained from Eurostat, which is the statistical office of the European Union (Eurostat, 2020b). Statistical data, broken down to the three NUTS levels, are published on their website. The data can be freely filtered according to year, geolocation (being the NUTS regions), and other aspects depending on the data, such as sex, age, or the unit of measure. Unfortunately, this data is not available daily and is often not up-to-date. That is, sometimes data is available up to 2016. For each variable, we kept the most recent data and assumed that this would be representative for the present. In Table 4.2 we mention per variable in what year the most recent observations were.

We distinguish three sets of regressors, as mentioned in section 3. Firstly, we have a set of control variables included in the tensor  $X_{r,t}$  which are not assumed to have a (large) effect on the transmission parameter  $\alpha$ . Secondly, the tensor  $W_{r,t}$  consists of variables that are assumed to affect the transmission within regions. Lastly, the matrix  $\widetilde{W}_{c,r,t}$  contains variables that are assumed to affect the transmission between regions. The specification of these regressors can be found in Table 4.2.

TODO: Fix this and look up the actual maximum year per variable

TODO: Insert  $\widetilde{W}$  variables and possibly move around variables to  $X$

Table 4.2: Specification of regressors

Matrix	Variable	Year	Description
$X_{r,t}$	weekend	n/a	Binary indicator denoting if the day is on the weekend (Saturday or Sunday)
	weekNumber	n/a	The calendar week number
$W_{r,t}$	airPassengersArrived	2018	Number of air passengers arrived
	touristArrivals	2018	Number of tourist arrivals
	broadbandAccess	2019	Percentage of population that has access to broadband internet
	deathRateDiabetes	2016	Number of deaths from diabetes per 100,000 inhabitants
	deathRateInfluenza	2016	Number of deaths from influenza per 100,000 inhabitants
	deathRateChd	2016	Number of deaths from coronary heart disease per 100,000 inhabitants
	deathRateCancer	2016	Number of deaths from cancer per 100,000 inhabitants

Table 4.2 continues on next page

Table 4.2 continued from previous page

Matrix	Variable	Year	Description
	deathRatePneumonia	2016	Number of deaths from pneumonia per 100,000 inhabitants
	availableBeds	2018	Number of hospital beds
	riskOfPovertyOrSocialExclusion	2018	Percentage of population at risk of poverty or social exclusion
$\widetilde{W}_{c,r,t}$			

One of the most important aspects in interpreting the results of a regression analysis is that interpretations are made under the *ceteris paribus* assumption. That is, we look at the effect of a change in one variable while holding all other variables constant. Because of this, there should be no large correlation between our independent variables. If there would be a large correlation between some regressors, then it is not possible to consider a change in one variable without causing a change in some other variable(s). Specifically for our case, we concur that there are people who often have multiple diseases at the same time and that there is likely a large correlation between the various death rates. To investigate this, we consider the correlation matrix in Figure 4.1. As described before, these variables are unfortunately not varying over time but they do vary over the regions. Because we are using the region-wise correlation, do note that a small sample size of  $R = 21$  is used. Therefore, the numbers should be taken with a grain of salt.

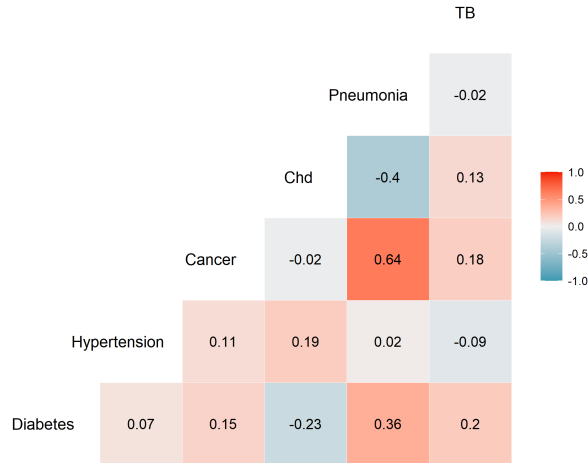


Figure 4.1: Correlation matrix of the discharge rates for various comorbidities of COVID-19

Figure 4.1 shows us that the largest correlation is 0.64 and occurs between the discharge rates of pneumonia and cancer. We also see a relatively high correlation of -0.4 between the discharge rates of pneumonia and CHD. For this reason, we remove the discharge rate of pneumonia from the model.

TODO: Cite a source on low sample size w.r.t. correlations

## 5 Results

In this section, we present the results from the models as presented in section 3.

Table 5.1: Estimates from Model 1: Within-Region Spread

	Estimate	Std. Error	<i>t</i> value	<i>p</i> value	
Intercept	$6.381 \times 10^{-3}$	$1.330 \times 10^{-3}$	4.800	$1.67 \times 10^{-6}$	***
Weekend	$4.274 \times 10^{-6}$	$1.893 \times 10^{-6}$	2.258	0.024045	*
Week number	$-5.162 \times 10^{-8}$	$1.542 \times 10^{-7}$	-0.335	0.737882	
Median age	$-1.346 \times 10^{-4}$	$2.811 \times 10^{-5}$	-4.789	$1.76 \times 10^{-6}$	***
Basilicata	$-5.528 \times 10^{-5}$	$9.657 \times 10^{-6}$	-5.725	$1.15 \times 10^{-8}$	***
Bolzano/Bozen	$-4.251 \times 10^{-4}$	$9.281 \times 10^{-5}$	-4.581	$4.84 \times 10^{-6}$	***
Calabria	$-3.132 \times 10^{-4}$	$5.934 \times 10^{-5}$	-5.278	$1.41 \times 10^{-7}$	***
Campania	$-5.419 \times 10^{-4}$	$1.097 \times 10^{-4}$	-4.939	$8.31 \times 10^{-7}$	***
Emilia-Romagna	$5.420 \times 10^{-5}$	$7.430 \times 10^{-6}$	7.295	$3.87 \times 10^{-13}$	***
Friuli-Venezia Giulia	$2.602 \times 10^{-4}$	$5.373 \times 10^{-5}$	4.842	$1.35 \times 10^{-6}$	***
Lazio	$-1.495 \times 10^{-4}$	$2.851 \times 10^{-5}$	-5.243	$1.69 \times 10^{-7}$	***
Liguria	$5.166 \times 10^{-4}$	$1.014 \times 10^{-4}$	5.094	$3.74 \times 10^{-7}$	***
Lombardia	$-5.591 \times 10^{-5}$	$2.290 \times 10^{-5}$	-2.442	0.014683	*
Marche	$8.592 \times 10^{-5}$	$1.750 \times 10^{-5}$	4.910	$9.64 \times 10^{-7}$	***
Molise	$8.479 \times 10^{-5}$	$2.030 \times 10^{-5}$	4.176	$3.05 \times 10^{-5}$	***
Piemonte	$2.241 \times 10^{-4}$	$3.973 \times 10^{-5}$	5.640	$1.87 \times 10^{-8}$	***
Puglia	$-2.351 \times 10^{-4}$	$4.527 \times 10^{-5}$	-5.194	$2.20 \times 10^{-7}$	***
Sardegna	$1.069 \times 10^{-4}$	$2.842 \times 10^{-5}$	3.760	$1.73 \times 10^{-4}$	***
Sicilia	$-3.418 \times 10^{-4}$	$6.762 \times 10^{-5}$	-5.055	$4.58 \times 10^{-7}$	***
Trento	$-1.176 \times 10^{-4}$	$3.399 \times 10^{-5}$	-3.459	$5.50 \times 10^{-4}$	***
Toscana	$1.388 \times 10^{-4}$	$3.124 \times 10^{-5}$	4.444	$9.18 \times 10^{-6}$	***
Umbria	$1.301 \times 10^{-4}$	$2.861 \times 10^{-5}$	4.546	$5.70 \times 10^{-6}$	***
Valle d'Aosta	$1.490 \times 10^{-4}$	$2.035 \times 10^{-5}$	7.321	$3.21 \times 10^{-13}$	***
Veneto	NA	NA	NA	NA	
$\alpha_{within}$	0.2915	0.01819	16.026	$< 2 \times 10^{-16}$	***

Significance levels: \* = 0.05, \*\* = 0.01, \*\*\* = 0.001

Notice that median age and the regional variables are perfectly collinear. TODO: Remove median age, unfortunately

The estimate for  $\alpha_{within}$  is 0.2915 and is statistically highly significant.

Table 5.2: Estimates from Model 3: Within and Between-Region Spread

	Estimate	Std. Error	<i>t</i> value	<i>p</i> value	
Intercept	$1.198 \times 10^{-3}$	$1.161 \times 10^{-3}$	1.032	0.302247	
Weekend	$1.282 \times 10^{-5}$	$1.664 \times 10^{-6}$	7.704	$1.84 \times 10^{-14}$	***
Week Number	$-6.939 \times 10^{-7}$	$1.391 \times 10^{-7}$	-4.989	$6.47 \times 10^{-7}$	***
Median Age	$-2.517 \times 10^{-5}$	$2.455 \times 10^{-5}$	-1.025	0.305426	
Basilicata	$-2.932 \times 10^{-5}$	$8.396 \times 10^{-6}$	-3.492	0.000487	***
Bolzano/Bozen	$-7.160 \times 10^{-5}$	$8.103 \times 10^{-5}$	-0.884	0.376982	

Table 5.2 continues on next page



Table 5.2 continued from previous page

	Estimate	Std. Error	<i>t</i> value	<i>p</i> value	
Calabria	$-6.970 \times 10^{-5}$	$5.190 \times 10^{-5}$	-1.343	0.179453	
Campania	$-1.182 \times 10^{-4}$	$9.582 \times 10^{-5}$	-1.234	0.217405	
Emilia-Romagna	$2.843 \times 10^{-5}$	$6.476 \times 10^{-6}$	4.390	$1.18 \times 10^{-5}$	***
Friuli-Venezia Giulia	$4.054 \times 10^{-5}$	$4.697 \times 10^{-5}$	0.863	0.388193	
Lazio	$-3.988 \times 10^{-5}$	$2.491 \times 10^{-5}$	-1.601	0.109556	
Liguria	$1.130 \times 10^{-4}$	$8.862 \times 10^{-5}$	1.275	0.202357	
Lombardia	$2.442 \times 10^{-5}$	$1.996 \times 10^{-5}$	1.223	0.221253	
Marche	$2.666 \times 10^{-5}$	$1.524 \times 10^{-5}$	1.750	0.080277	
Molise	$1.726 \times 10^{-6}$	$1.774 \times 10^{-5}$	0.097	0.922502	
Piemonte	$6.400 \times 10^{-5}$	$3.473 \times 10^{-5}$	1.843	0.065454	
Puglia	$-5.557 \times 10^{-5}$	$3.957 \times 10^{-5}$	-1.404	0.160308	
Sardegna	$9.863 \times 10^{-6}$	$2.475 \times 10^{-5}$	0.398	0.690297	
Sicilia	$-8.086 \times 10^{-5}$	$5.907 \times 10^{-5}$	-1.369	0.171134	
Trento	$7.956 \times 10^{-6}$	$2.966 \times 10^{-5}$	0.268	0.788506	
Toscana	$2.605 \times 10^{-5}$	$2.723 \times 10^{-5}$	0.957	0.338886	
Umbria	$9.453 \times 10^{-6}$	$2.502 \times 10^{-5}$	0.378	0.705624	
Valle d'Aosta	$6.558 \times 10^{-5}$	$1.780 \times 10^{-5}$	3.684	0.000234	***
Veneto	NA	NA	NA	NA	
$\alpha_{within}$	-0.1104	$1.972 \times 10^{-2}$	-5.598	$2.39 \times 10^{-8}$	***
$\alpha_{between}$	0.04845	$1.467 \times 10^{-3}$	33.036	$< 2 \times 10^{-16}$	***

Significance levels: \* = 0.05, \*\* = 0.01, \*\*\* = 0.001

The estimate for  $\alpha_{within}$  is -0.1104 and is statistically highly significant. Comparing this to the results from model 1, the sign has flipped. We also notice that the sign has flipped for the region of Trento. The estimate for  $\alpha_{between}$  is 0.04845, which is also statistically highly significant.

## 6 Conclusion

## **7 Future research**

TODO: Fix the papers with many authors

## References

- Adda, J. (2016). Economic activity and the spread of viral diseases: Evidence from high frequency data. *The Quarterly Journal of Economics*, 131(2), 891–941.
- Anderson, R. M., & May, R. M. (1992). *Infectious diseases of humans: dynamics and control*. Oxford University Press.
- Caccia, F. (2020, March). *Coronavirus, “il conteggio dei morti varia da paese a paese. La Germania esclude chi ha altre patologie”*. Retrieved 2020-06-11, from [https://www.corriere.it/cronache/20\\_marzo\\_22/coronavirus-il-conteggio-morti-varia-paese-paese-germania-esclude-chi-ha-altre-patologie-6a452e6a-6c19-11ea-8403-94d97cb6fb9f\\_preview.shtml](https://www.corriere.it/cronache/20_marzo_22/coronavirus-il-conteggio-morti-varia-paese-paese-germania-esclude-chi-ha-altre-patologie-6a452e6a-6c19-11ea-8403-94d97cb6fb9f_preview.shtml)
- Eurostat. (2020a). *Eurostat regional data background*. Retrieved 2020-06-11, from <https://ec.europa.eu/eurostat/web/regions/background>
- Eurostat. (2020b). *Eurostat regional statistics database*. Retrieved 2020-06-11, from <https://ec.europa.eu/eurostat/web/regions/data/database>
- He, X., Lau, E. H., Wu, P., Deng, X., Wang, J., Hao, X., ... others (2020). Temporal dynamics in viral shedding and transmissibility of COVID-19. *Nature medicine*, 26(5), 672–675.
- Horowitz, J. (2020, March). *Italy’s Health Care System Groans Under Coronavirus — a Warning to the World*. Retrieved 2020-06-11, from <https://www.nytimes.com/2020/03/12/world/europe/12italy-coronavirus-health-care.html>
- Istituto Superiore di Sanità. (2020, June). *Caratteristiche dei pazienti deceduti positivi all’infezione da SARS-CoV-2 in Italia*. Retrieved 2020-06-11, from <https://www.epicentro.iss.it/coronavirus/sars-cov-2-decessi-italia>
- Keeling, M. J., & Rohani, P. (2011). *Modeling infectious diseases in humans and animals*. Princeton University Press.
- Kermack, W. O., & McKendrick, A. G. (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London. Series A, Containing papers of a mathematical and physical character*, 115(772), 700–721.
- Kirkcaldy, R. D., King, B. A., & Brooks, J. T. (2020). COVID-19 and Postinfection Immunity: Limited Evidence, Many Remaining Questions. *JAMA*.
- Lauer, S. A., Grantz, K. H., Bi, Q., Jones, F. K., Zheng, Q., Meredith, H. R., ... Lessler, J. (2020). The incubation period of coronavirus disease 2019 (COVID-19) from publicly reported confirmed cases: estimation and application. *Annals of internal medicine*, 172(9), 577–582.
- Li, Q., Guan, X., Wu, P., Wang, X., Zhou, L., Tong, Y., ... others (2020). Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *New England Journal of Medicine*.
- Linton, N. M., Kobayashi, T., Yang, Y., Hayashi, K., Akhmetzhanov, A. R., Jung, S.-m., ... Nishiura, H. (2020). Incubation period and other epidemiological characteristics of 2019 novel coronavirus infections with right truncation: a statistical analysis of publicly available case data. *Journal of clinical medicine*, 9(2), 538.
- Ministero della Salute. (2020a, May). *Coronavirus: Contagion rate R0 below 1. Prudence needed in phase two says ISS*. Retrieved 2020-06-11, from [http://www.salute.gov.it/portale/news/p3\\_2\\_1\\_1\\_1.jsp?lingua=italiano&menu=notizie&p=dalministero&id=4717](http://www.salute.gov.it/portale/news/p3_2_1_1_1.jsp?lingua=italiano&menu=notizie&p=dalministero&id=4717)
- Ministero della Salute. (2020b, June). *Covid-19, i casi in Italia: 10 giugno ore 18*. Retrieved 2020-06-11, from <http://www.salute.gov.it/portale/nuovocoronavirus/dettaglioNotizieNuovoCoronavirus.jsp?lingua=italiano&menu=notizie&p=dalministero&id=4877>

Sevillano, E. (2020, March). *Tracking the coronavirus: why does each country count deaths differently?*  
Retrieved 2020-06-11, from <https://english.elpais.com/society/2020-03-30/tracking-the-coronavirus-why-does-each-country-count-deaths-differently.html>

# Appendices

## A Tables