



# Predicting The Incidence Rate And Case Fatality Rate Of The Novel Coronavirus SARS-CoV-2

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 (Tilburg University)  
dr. George Knox (Tilburg University)

Date:  
May 14, 2020



# Contents

<b>1</b>	<b>Acknowledgements</b>	<b>2</b>
<b>2</b>	<b>Management summary</b>	<b>2</b>
<b>3</b>	<b>Introduction</b>	<b>2</b>
<b>4</b>	<b>Problem description</b>	<b>2</b>
<b>5</b>	<b>Materials</b>	<b>3</b>
<b>6</b>	<b>Results</b>	<b>8</b>
<b>7</b>	<b>Conclusion</b>	<b>9</b>
	<b>References</b>	<b>10</b>
<b>A</b>	<b>Tables</b>	<b>11</b>

# 1 Acknowledgements

## 2 Introduction

### 3 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. We are basing our model on specifications as used by Adda (2016). In his paper, Adda investigates the spread of several viral diseases in the past, namely influenza, gastroenteritis, and chickenpox. He starts from the Standard Inflammatory Response (SIR) model

to model the incidence rate  $Inc_{r,t}$  for several viruses, being the percentage of the population in a region  $r$  who have the virus at a time  $t$ :

$$\begin{aligned} Inc_{r,t} = & Inc_{r,t-lag} S_{r,t-lag} \sum_{k=1}^K a_{within}^k W_{r,t-lag}^k \\ & + \sum_{c \neq r} Inc_{c,t-lag} S_{r,t-lag} \sum_{k=1}^{\tilde{K}} a_{between}^k \widetilde{W}_{r,c,t-lag}^k \\ & + X_{r,t} \delta + \eta_{r,t} \end{aligned} \quad (1)$$

Adda models the susceptible population as the total population who currently do not have the virus and who are not immune. That is, a certain proportion of immune people lose their immunity and become susceptible again. At this point, we will assume that all recovered patients achieve immunity. This assumption can be challenged because it is currently still unknown whether immunity is always achieved, especially among those who have had only light to medium symptoms. However, it is estimated that COVID-19 antibodies will remain in a patient's system for two to three years, based on what is known about other coronaviruses, but it is too early to know for certain (Leung, 2020). As such, we believe our assumption is generally valid.

That is, let  $S$  denote the fraction of individuals who are susceptible to contracting the disease,  $I$  the fraction of individuals who are infected, and  $R$  the fraction of individuals who have recovered but are still immune. Then:

TODO: update to our current setting

$$\begin{cases} \frac{dI(t)}{dt} = \alpha S(t)I(t) - \beta I(t) \\ \frac{dR(t)}{dt} = \beta I(t) - \lambda R(t) \\ \frac{dS(t)}{dt} = -\alpha S(t)I(t) + \lambda R(t) \end{cases}$$

Notice that Adda also models interaction between regions using the matrix  $\widetilde{W}_{r,c}$ . At first, we will neglect interactions between regions. The model becomes:

## 4 Dataset

The NUTS classification (Nomenclature of territorial units for statistics) is a hierarchical system for dividing up the economic territory of the EU and the UK (*NUTS - Nomenclature Of Territorial Units For Statistics - Background*, n.d.). In this thesis we focus our attention on Italy, the epicenter of coronavirus cases in Europe. Italy consists of 21 so-called *regioni* (regions), comparable to the 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 *North West*, *North East*, *Centre*, *South*, and *the Islands*. The third-level NUTS regions are 110 provinces, which are subregions of the *regioni*.

Data was gathered from various sources. The specific information on the coronavirus in Italian regions was retrieved from the Ministero della Salute (the Italian ministry of health services), who publish daily reports under a title similar to *Covid-19, i casi in Italia 17 aprile ore 18*, where *17 aprile* would be updated to the relevant date (Salute, n.d.). 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*)

It should be noted that the death statistics for Italy do not include the total amount of coronavirus victims who died outside hospitals, including dozens who died in different nursing homes across the country. Therefore, the official death statistics are considered an underestimate (Stancati & Sylvers, 2020). However, we did not model the amount of deaths or the death rate, so this does not impact our analysis. Nonetheless, we do take into account an underestimation of other information. For instance, not all people infected with COVID-19 are tested. These are not only symptomatic patients but also asymptomatic Italians. Moreover, it is unclear how the government collects this information. If

regions or provinces submit this information to the government each day, there may be provinces who fail to submit their data for a certain day. Despite this, we assume that this official information is accurate and representative of the region itself.

Two issues that we want to address are missing data and the correction of data. In the official publications that we use, data that was wrongly published on a day  $t - 1$  is corrected by subtracting the error from 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 five 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. 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  $t+1$ .

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 1. As such, missing data is simply imputed with a value of 0.

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

Table 1: Number of confirmed cases around a day  $t$  with missing data

Regressors were obtained from Eurostat, which is the statistical office of the European Union. Statistical data, broken down to the three NUTS levels, are published on their website (*Eurostat Regions Database*, n.d.). The data can be freely filtered according to time period, geolocation (being the NUTS regions), and other aspects depending on the data, such as sex, age, or the unit of measure. The specification of the regressors we used can be found in Table 2.

Continue here - fill in the descriptions

Regressor	Description	Unit of measure
air_passengers_arrived	x	Number
air_passengers_departed	x	Number
tourist_arrivals	x	Number
broadband_access	x	Percentage of population
death_rate_diabetes	x	Number per 100,000 inhabitants
death_rate_influenza	x	Number per 100,000 inhabitants
death_rate_chd	x	Number per 100,000 inhabitants
death_rate_cancer	x	Number per 100,000 inhabitants
death_rate_pneumonia	x	Number per 100,000 inhabitants
available_beds	x	Number
maritime_passengers_disembarked	x	Number
maritime_passengers_embarked	x	Number
risk_of_poverty_or_social_exclusion	x	Percentage of population
weekend	x	Binary indicator
weekNumber	x	Number

Table 2: Specification of regressors

We need to make sure that there is no large correlation between regressors. Specifically, we concur that there are people who often have multiple diseases at the same time.

	Diabetes	Respiratory	Hypertension	Cancer	CHD	Pneumonia	TB
Diabetes		0.14	0.07	0.15	-0.23	0.36	0.20
Respiratory	0.14		0.07	0.71	-0.45	0.69	-0.09
Hypertension	0.07	0.07		0.11	0.19	0.02	-0.09
Cancer	0.15	0.71	0.11		-0.02	0.64	0.18
CHD	-0.23	-0.45	0.19	-0.02		-0.40	0.13
Pneumonia	0.36	0.69	0.02	0.64	-0.40		-0.02
TB	0.20	-0.09	-0.09	0.18	0.13	-0.02	

**Note, the following is old and is for the specification of Adda. The specification of  $W$  will likely be in  $X$  for the other models.**

The spatial weighting matrix  $W_r$  has the following structure:

$$W_r = \begin{bmatrix} V_r & C_r \end{bmatrix},$$

where  $V_r$  consists of  $K_V$  time-varying regressors and  $C_r$  consists of  $K_C$  time-constant regressors, so  $V_r \in \mathbb{R}^{T \times K_V}$  and  $C_r \in \mathbb{R}^{T \times K_C}$ . Taking an example:

$$W_r = \begin{bmatrix} V_r^{\text{schools closed}} & V_r^{\text{lockdown started}} & C_r^{\text{hospital beds}} & C_r^{\text{internet access}} \end{bmatrix}.$$

We note that the descriptive data (like demographics and economic data) that we use is assumed to be time-constant during the coronacrisis (due to lack of data). The time-varying information that we use consists binary indicators for whether certain policies (such as closing down schools or instigating a lockdown) were implemented. As such,  $W_r$  mostly contains time-constant information.

We will use the following specifications for the weights and regressors:

- $W_{r,t-lag}$  contains  $K := K_V + K_C$  region-specific variables that potentially influence the transmission rate of SARS-CoV-2 within a region  $r$ . We split these in several categories:

#### **Economic**

- The amount of freight being transported by plane from and to the region (not available interregionally).
- The amount of freight being transported by ship from and to the region (not available interregionally).
- The amount of arrivals at tourist accommodations.
- The GDP at current market prices per inhabitant.
- The disposable income per inhabitant.
- The amount of journeys made for transport of freight by road by loading and unloading region.

#### **Demographics, social, etcetera**

- The area size.
- The median age and median age squared.
- The population number.
- The percentage of people at risk of poverty or social exclusion.
- The percentage of people with broadband access.
- The percentage of people who used internet to contact the public authorities in the last year.



- The percentage of people that attained a certain education level.

### **Medical**

- The average length-of-stay in a hospital.
- The crude death rate for several different diseases.
- The number of health personnel (doctors and nurses).
- The number of hospital beds.

### **Travelling**

- The number of passengers travelling by plane from and to the region (not available interregionally).
- The number of passengers travelling by ship from and to the region (not available interregionally).
- The length of railroads, motorways, navigable rivers, etcetera.
- $X_{r,t}$  contains certain fixed effects to control for, such as a binary indicator whether the day was on a weekend.

When we will also consider interactions between regions, we will define  $\widetilde{W}_{r,t-lag}$  to contain  $\tilde{K}$  variables that potentially influence the transmission rate of SARS-CoV-2 across regions:

- Amount of passengers that travelled from region  $c$  to region  $r$  via railroad.
- Amount of freight that travelled from region  $c$  to region  $r$  via railroad.
- A binary indicator indicating whether the regions border each other.
- The distance between the largest (most populous) cities in the regions.
- The population ratios.
- The log regional GDP ratios.

## 5 Results

```

call:
lm(formula = fm, data = df_long)

Residuals:
    Min       1Q   Median       3Q      Max
-2.238e-04 -1.822e-05 -1.290e-06  1.333e-05  5.579e-04

coefficients:
(Intercept)
weekend1
weekNumber
factor(Code)BAS
factor(Code)BZ
factor(Code)CAL
factor(Code)CAM
factor(Code)EMR
factor(Code)FVG
factor(Code)LAZ
factor(Code)LIG
factor(Code)LOM
factor(Code)MAR
factor(Code)MOL
factor(Code)PIE
factor(Code)PUG
factor(Code)SAR
factor(Code)STC
factor(Code)TN
factor(Code)TOS
factor(Code)UMB
factor(Code)VDA
factor(Code)VEN
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(arrivals, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(broadbandAccess, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateDiabetes, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateRespiratory, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateHypertension, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateCancer, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateChd, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRatePneumonia, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(dischargeRateTB, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(availableBeds, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(maritimePassengersDisembarked, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(riskOfPovertyOrSocialExclusion, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(railTravelers, 1)
dplyr::lag(incidenceRate, 1):dplyr::lag(susceptibleRate, 1):dplyr::lag(medianAge, 1)

Estimate Std. Error t value Pr(>|t|)
6.739e-06 5.836e-06 1.155 0.248343
3.605e-06 1.861e-06 1.937 0.052873 .
1.109e-06 2.140e-07 5.180 2.43e-07 ***
-2.968e-05 6.966e-06 -4.261 2.13e-05 ***
2.100e-05 7.070e-06 2.970 0.003010 **
-3.069e-05 7.176e-06 -4.277 1.98e-05 ***
-3.533e-05 6.719e-06 -5.258 1.60e-07 ***
9.301e-06 7.312e-06 1.272 0.203537
-2.682e-05 7.634e-06 -3.514 0.000451 ***
-1.239e-05 6.798e-06 -1.852 0.064150 .
-1.884e-05 7.354e-06 -2.562 0.010488 *
1.630e-05 7.829e-06 2.082 0.037493 *
-1.862e-05 7.898e-06 -2.358 0.018458 *
-2.356e-05 7.436e-06 -3.169 0.001552 **
1.728e-05 6.833e-06 2.528 0.011541 *
-1.889e-05 7.680e-06 -2.460 0.013963 *
-3.666e-05 6.705e-06 -5.467 5.11e-08 ***
-3.029e-05 6.655e-06 -4.552 5.62e-06 ***
1.564e-05 6.889e-06 2.271 0.023258 *
-7.054e-06 7.411e-06 -0.952 0.341293
-3.252e-05 7.247e-06 -4.487 7.60e-06 ***
3.513e-05 7.017e-06 5.007 5.99e-07 ***
-2.298e-07 7.077e-06 -0.032 0.974096
7.041e+00 1.520e+00 4.631 3.85e-06 ***
3.931e+01 4.700e+00 8.363 < 2e-16 ***
-2.785e-01 2.367e-02 -11.765 < 2e-16 ***
-3.128e+01 4.806e+00 -6.508 9.46e-11 ***
-1.659e+02 1.448e+01 -11.457 < 2e-16 ***
3.623e+01 5.398e+00 6.711 2.47e-11 ***
-4.597e+01 1.191e+01 -3.860 0.000117 ***
-6.828e+00 2.563e+00 -2.664 0.007783 **
1.838e+02 1.403e+01 13.104 < 2e-16 ***
3.297e+01 4.853e+00 6.795 1.40e-11 ***
-1.650e+01 2.946e+00 -5.601 2.41e-08 ***
1.790e+00 1.768e+00 1.012 0.311512
1.440e-01 1.070e-02 13.459 < 2e-16 ***
3.885e+00 1.767e+00 2.199 0.028020 *
4.081e-01 4.461e-02 9.147 < 2e-16 ***

signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.92e-05 on 2124 degrees of freedom
(1 observation deleted due to missingness)
Multiple R-squared:  0.5414,    Adjusted R-squared:  0.5334
F-statistic: 67.76 on 37 and 2124 DF,  p-value: < 2.2e-16

```

Figure 1: Output Least-Squares Dummy Variables

## 6 Conclusion

## 7 Future research

## 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.
- Anselin, L. (2013). *Spatial econometrics: methods and models* (Vol. 4). Springer Science & Business Media.
- Baltagi, B. H., Song, S. H., & Koh, W. (2003). Testing panel data regression models with spatial error correlation. *Journal of econometrics*, 117(1), 123–150.
- Eurostat regions database*. (n.d.). Eurostat. Retrieved from <https://ec.europa.eu/eurostat/web/regions/data/database>
- Leung, H. (2020, Apr). *What we know about coronavirus immunity and reinfection*. Time Magazine. Retrieved from <https://time.com/5810454/coronavirus-immunity-reinfection/>
- Millo, G., & Piras, G. (2012). splm: Spatial panel data models in r. *Journal of Statistical Software*, 47(1). Retrieved from <https://www.jstatsoft.org/article/view/v047i01> doi: 10.18637/jss.v047.i01
- Nuts - nomenclature of territorial units for statistics - background*. (n.d.). Eurostat. Retrieved from <https://ec.europa.eu/eurostat/web/nuts/background>
- Salute, M. d. (n.d.). *Nuovo coronavirus*. Retrieved from <http://www.salute.gov.it/nuovocoronavirus>
- Stancati, M., & Sylvers, E. (2020, Apr). *Italy's coronavirus death toll is far higher than reported*. Dow Jones & Company. Retrieved from <https://www.wsj.com/articles/italys-coronavirus-death-toll-is-far-higher-than-reported-11585767179>

## A Tables