



Predicting The Incidence Rate And Case Fatality Rate Of COVID-19 in Italy

by

Mike Weltevrede (ANR 756479)

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:
August 10, 2020

Abstract

TODO

Acknowledgements

TODO

Contents

1	Introduction	1
2	Problem description	2
3	Methodology	3
3.1	SIR model	3
3.2	Model 1: Within-Region Spread	5
3.3	Model 2: Weighted Within-Region Spread	7
3.4	Model 3: Within and Between-Region Spread	8
3.5	Model 4: Full Model (Weighted Within and Between-Region Spread) . .	9
3.6	Model selection	9
3.7	Modelling undocumented infections	10
4	Dataset	19
4.1	Geographical structure of Italy	19
4.2	Coronavirus data	19
4.3	Independent variables	22
5	Results	25
5.1	Model 1: Within-Region Spread	25
5.2	Model 2: Weighted Within-Region Spread	28
5.3	Model 3: Within and Between-Region Spread	29
5.4	Model 4: Full Model (Weighted Within and Between-Region Spread) . .	30
6	Conclusion	31
7	Future research	32
	Appendices	36
A	Abbreviations	36
B	Tables	37
B.1	Illustrative tables	37
B.2	Results from Model 1: Within-Region Spread	37
B.3	Results from Model 2: Weighted Within-Region Spread	40
B.4	Results from Model 3: Within and Between-Region Spread	40
B.5	Results from Model 4: Full Model (Weighted Within and Between-Region Spread)	43

C	Figures	43
C.1	Plots of β_{within} over time	43
D	Derivations	48
D.1	Calculation of population variables	48
D.2	Functional forms for modelling undocumented infections	48
D.2.1	Linear function	49
D.2.2	General quadratic function	49
D.2.3	Special case quadratic formula: downwards opening	52
D.2.4	Special case quadratic formula: upwards opening	53
D.2.5	Cubic function	54

1 Introduction

2 Problem description

In this section, we elaborate on the problem at hand, namely the epidemiological spread of SARS-CoV-2 and the disease it causes. By themselves, viruses were responsible for more deaths than all armed conflicts combined in the twentieth century (Adda, 2016). Since the beginning of 2020, the novel coronavirus SARS-CoV-2 (causing the viral disease COVID-19) has plagued the world. Starting from Wuhan, China, it has made its way to every single continent apart from Antarctica and (nearly) every country in the world. In response to SARS-CoV-2, governments have been implementing far-reaching measures to try and contain the virus, such as shutting down schools and restaurants, but also by locking down the entire country. On August 2, 2020, 17.8 million people were reported to have been infected with COVID-19, leading to 675 thousand consequent deaths. Only 12 sovereign member states of the United Nations reported no infections. For two of these countries, namely North Korea and Turkmenistan, these reports are suspected to be false.

Italy has been one of the most intensely struck countries by COVID-19. Until the end of March, it had the highest number of confirmed cases per 100,000 inhabitants. It was subsequently taken over by Spain. Italy remained the second most struck country until May 1, when the United States took over. On July 3, 2020, it had the ninth highest absolute number of confirmed cases, after the United States, Brazil, Russia, India, Peru, Chile, the United Kingdom, and Spain. Despite dropping in this ranking, Italy reported the second highest global death-to-cases ratio of 14.45% (34,818 deaths to 240,961 cases), only after the United Kingdom, which reports a death-to-cases ratio of 15.50% (43,995 deaths to 283,757 cases). The third highest death-to-cases ratio of 12.24% (29,189 deaths to 238,511 cases) was reported by Mexico. 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). Due to the extreme nature of the pandemic in Italy and the availability of enough data, this thesis chooses to focus on Italy. Specifically, we focus on modelling on the level of regions rather than on a nation-wide approach.

We are basing our models on specifications as used by Adda (2016). In the paper, Adda (2016) investigates the spread of several viral diseases in the past, namely for influenza, gastroenteritis, and chickenpox. The key additions made are, firstly, that a spatial spillover effect is considered and, secondly, that some sort of weighting on the parameters is allowed on the basis of region specific variables. These models have not previously been applied to SARS-CoV-2 and can show interesting insights compared to other models.

3 Methodology

In this section, we explain the methodology applied in this thesis. We discuss our models and the thought process behind them. In Section 3.1, we describe the most commonly used model in epidemiology: the SIR model. In Section 3.2, we present a model ignoring effects across regions and for which the transmission rate parameter is determined by the previous infectives. Subsequently, in Section 3.3, we extend this model by allowing the transmission rate parameter to be weighted by several other factors. After this, Section 3.4 presents a model that takes effects across regions into account for which the transmission rate parameter is determined by the previous infectives. The last model is presented in Section 3.5, which takes effects across regions into account as well as allowing for the transmission rate parameter to be weighted by other factors. In Section 3.6, we consider how to do model selection for these four models to determine the best set of regressors to use. Lastly, Section 3.7 describes how undocumented infections are modelled.

3.1 SIR model

Adda (2016) starts from the Standard Inflammatory Response (SIR) model, the most commonly used model in epidemiology (Anderson & May, 1992; Kermack & McKendrick, 1927). We will follow the notation by Keeling and Rohani (2011). That is, 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. Keeling and Rohani (2011) furthermore define X to be the number of susceptible individuals, Y to be the number of infectives, and Z to be the number of recovered individuals, so that $S = X/N$, $I = Y/N$, and $R = Z/N$, where N is the total population size. As such, at any point in time, we have that

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

$$X, Y, Z \in [0, N] \text{ and } X + Y + Z = N.$$

The SIR model is postulated in continuous time, i.e. the equations in (3.1), (3.2), and (3.3) 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 fraction of infectives) to which a flow is added or subtracted.

$$\frac{dS}{dt} = -\beta SI + wR \quad (3.1)$$

$$\frac{dI}{dt} = \beta SI - \gamma I \quad (3.2)$$

$$\frac{dR}{dt} = \beta I - wR \quad (3.3)$$

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. 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 βSI in equations (3.1) and (3.2), 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 decease. This relates to the term γI in equations (3.2) and (3.3).

Finally, we assume that there is a constant rate of change at which immune individuals lose their immunity. This is denoted by the term wR in equations (3.1) and (3.3). For instance, Adda (2016) mentions that w is set to 0 for chickenpox as individuals acquire a lifetime immunity while w 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 et al., 2020). This 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). Nonetheless, no definitive results have been shown. For simplicity's sake, we assume that lifelong immunity is achieved, or at least long enough to last through the time-scope of our analysis: we set $w = 0$.

One of the main measures resulting from the SIR model is the estimation of the basic reproduction number $R_0 := \beta/\gamma$. An epidemic is said to develop if $R_0 > 1$. This is clear because this is the case when $\beta > \gamma$, i.e. the spread of the virus exceeds the recovery rate: individuals become infected more quickly than they recover. 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, 2020), showing that this measure is also used communicated to citizens as a way of informing them whether the pandemic is tending to end.

3.2 Model 1: Within-Region Spread

Recall that the SIR model is postulated in continuous time. Adda (2016) discretizes the SIR model as in (3.8). Recall from (3.2) that $\frac{dI}{dt} = \beta SI - \gamma I$. As such, the discretized version (for a region r) is:

$$I_{r,t} - I_{r,t-1} = \beta S_{r,t-1} I_{r,t-1} - \gamma I_{r,t-1} \quad (3.4)$$

There are a few things to note. Firstly, if we want to estimate this equation's parameters, an error occurs. This is added to the model by an error term denoted by $\eta_{r,t}$:

$$I_{r,t} - I_{r,t-1} = \beta S_{r,t-1} I_{r,t-1} - \gamma I_{r,t-1} + \eta_{r,t} \quad (3.5)$$

Secondly, individuals that get infected do not immediately infect others because there is a so-called latent period, being the period between an infection and the moment that the infective is infectious. For COVID-19, the latent period is estimated to be approximately 2 days shorter than the incubation period (He et al., 2020). The incubation period is 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. The incubation period for COVID-19 is estimated to be above 2 and below 11.5 (Lauer et al., 2020), 12.5 (Q. Li et al., 2020), and 14 days (Linton et al., 2020). This is a large range, but this is not rare. For instance, the incubation period for chicken pox is estimated to be between 9 and 21 days (Papadopoulos, 2018). While the maximum incubation period is not agreed upon by Lauer et al. (2020) and Q. Li et al. (2020), their results on the median are similar. Lauer et al. (2020) report a median incubation period of 5.1 days (95% CI: 4.5 to 5.8 days), while Q. Li et al. (2020) report 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. Because the latent period is estimated to be shorter than the incubation period, there are infectives who are able to infect others before showing symptoms. We call these people pre-symptomatic, which is distinctive from asymptomatic people in the sense that asymptomatic people do not develop symptoms and pre-symptomatic people will develop symptoms but they develop a higher viral load just before said symptoms are apparent. On June 9, 2020, the World Health Organization (WHO) said that it is unclear whether asymptomatic people can actually spread the virus but that pre-symptomatic people may actually be able to infect others (Sutherland & Gretler, 2020). 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. Sutherland and Gretler (2020) moreover reiterate the WHO's

statement that studies have been done that show that asymptomatic people can spread the virus but that more research needs to be done to show how many of these infectious asymptomatic people exist. We discuss how we model pre-symptomatic individuals in Section 3.7. Adda (2016) models this transmission lag by making the lag on the right hand side of (3.5) dependent on the incubation period. This is denoted by the parameter τ :

$$I_{r,t} - I_{r,t-1} = \beta S_{r,t-\tau} I_{r,t-\tau} - \gamma I_{r,t-\tau} + \eta_{r,t} \quad (3.6)$$

For instance, Adda (2016) chooses τ equal to one week for acute diarrhea and flu-like illnesses as these have an incubation period of less than a week. Due to the results from Lauer et al. (2020), Q. Li et al. (2020), and Linton et al. (2020), indicating an incubation period of roughly five days, and the result from He et al. (2020) that the latent period is roughly two days shorter than the incubation period, we choose $\tau = 3$.

Thirdly, Adda (2016) adds regressors to the model as control variables, such as the region fixed effects, week effects and year effects in levels. Note that regressors can be added to the model to capture possible effects that would otherwise be included in the error, confounding the estimation of the transmission parameter β . Adda (2016) denotes this matrix of regressors by X , not to be confused with the notation by Keeling and Rohani (2011) for the number of infectives. For this reason, we instead denote the matrix of regressors as used by Adda (2016) by M . This leads to the following formulation:

$$I_{r,t} - I_{r,t-1} = \beta S_{r,t-\tau} I_{r,t-\tau} - \gamma I_{r,t-\tau} + \delta M_{r,t} + \eta_{r,t} \quad (3.7)$$

For our application, the data does not span multiple years. As such, we do not have year effects. Moreover, given that year effects are not available, week effects would capture a time trend. However, we do add a weekend effect. More information and reasoning is provided in Section 4.3.

Fourthly, there are two other key differences in the model by Adda (2016) that are not properly explained in the paper. First of all, Adda (2016) does not include the term $\gamma I_{r,t-\tau}$ in the model. Presumably, this is because Adda (2016) considers the number of new cases and, therefore, the number of recovered individuals do not impact that value. Second of all, Adda (2016) replaces the proportion of infectives $I_{r,t-\tau}$ by the number of new cases $X_{r,t-\tau} - X_{r,t-\tau-1}$ (where we follow the notation from Keeling and Rohani (2011)). Similarly, Adda (2016) puts the dependent variable to be the number of new cases instead of the number of infectives divided by the population (the incidence rate).

TODO: Why is this (not a problem) c.q. what do we do with this?

Defining $\Delta X_t := X_t - X_{t-1}$ and following the notation by Keeling and Rohani (2011), the within-region model as defined by Adda (2016), ignoring effects across regions, is given by:

$$\Delta X_{r,t} = \beta_{within} \Delta X_{r,t-\tau} S_{r,t-\tau} + \delta M_{r,t} + \eta_{r,t} \quad (3.8)$$

The model is estimated by ordinary least squares (OLS). The moment condition that needs to be satisfied due to the strict exogeneity assumption is

$$E[\eta_{r,t}(\beta_{within}\Delta X_{r,t-\tau}S_{r,t-\tau} + \delta M_{r,t})] = 0.$$

A general assumption that is made, is that the idiosyncratic error $\eta_{r,t}$ is uncorrelated with the regressors in the matrix $M_{r,t}$. That is, we assume that $E[\eta_{r,t} | M_{r,t}] = 0$. Now note that we need to only consider the relation between $\eta_{r,t}$ and $\Delta X_{r,t-\tau}S_{r,t-\tau}$. The reason why we assume that $E[\eta_{r,t} | \Delta X_{r,t-\tau}S_{r,t-\tau}] = 0$ is that, for a large enough lag τ , the error is not correlated with past data at that lag. That is, the people that are classified as infectives at time $t - \tau$ do not have an effect on the error that we make when considering the infectives at time t under a correct model specification. As such, we assume that the moment condition holds.

3.3 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 β 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 (3.8). The resulting model is a completely new addition. After defining the between-regions model in Section 3.4, we apply the same methodology to obtain the full weighted model in Section 3.5.

Let the tensor W contain K region-specific variables that may influence the transmission rate β . As such, we now allow for β_{within} to differ for these K variables. In Section 4, we elaborate on how these variables included in W are specifically defined and selected. For instance, we include the number of rail travellers, which changes over time. For this, we use data from the Google Mobility Report (Google LLC, 2020). We define M and η in the same way as in Section 3.2. Taking this into account, the weighted within-region model is defined as:

$$\Delta X_{r,t} = \Delta X_{r,t-\tau}S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k + \delta M_{r,t} + \eta_{r,t} \quad (3.9)$$

The moment condition that needs to be satisfied due to the strict exogeneity assumption is

$$E \left[\eta_{r,t} \left(\Delta X_{r,t-\tau}S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k + \delta M_{r,t} \right) \right] = 0.$$

The same reasoning as in Section 3.2 applies with regards to the assumption that $E[\eta_{r,t} | X_{r,t}] = 0$. On the assumption that $E[\eta_{r,t} | \Delta X_{r,t-\tau}S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k] = 0$,

we should realize that the only difference with the reasoning in Section 3.2 is the addition of the weighting matrix W . Since these are external factors, such as socioeconomic variables, we believe that these will be uncorrelated with the error. Combining this with the arguments from Section 3.2, we assume that the moment condition holds.

3.4 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. The following model is defined:

$$\Delta X_{r,t} = \beta_{within} \Delta X_{r,t-\tau} S_{r,t-\tau} + \beta_{between} S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} + \delta M_{r,t} + \eta_{r,t} \quad (3.10)$$

TODO: Why may this specification not be good

The moment condition that needs to be satisfied due to the strict exogeneity assumption is

$$E \left[\eta_{r,t} \left(\beta_{within} \Delta X_{r,t-\tau} S_{r,t-\tau} + \beta_{between} S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} + \delta M_{r,t} \right) \right] = 0.$$

In the same way as in Section 3.2, we can assume that $E[\eta_{r,t} \mid M_{r,t}] = 0$ and $E[\eta_{r,t} \mid \Delta X_{r,t-\tau} S_{r,t-\tau}] = 0$. Following the same reasoning as before, we assume that the number of infectives who come into contact with susceptibles in other regions at a certain time is not correlated with the error if the lag is large enough. As such, we assume that the moment condition holds.

In (3.10), the transmission parameter β is now allowed to be different within and between regions. Adda (2016) estimates (3.10) 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. For this reason, we only consider OLS for this model.

3.5 Model 4: Full Model (Weighted Within and Between-Region Spread)

We now incorporate the between-region effects as well as the weighting of the transmission parameter. In addition to (3.9), we now also put weights on the between-region transmission parameter by some possibly influential variables. Let the tensor \widetilde{W} contain \tilde{K} variables that now can influence the transmission rate β_{within} between two regions r and c .

$$\begin{aligned}\Delta X_{r,t} = & \Delta X_{r,t-\tau} S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k \\ & + S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} \sum_{k=1}^{\tilde{K}} \beta_{between}^k \widetilde{W}_{r,c,t-\tau}^k \\ & + \delta M_{r,t} + \eta_{r,t}\end{aligned}\tag{3.11}$$

The moment condition that needs to be satisfied due to the strict exogeneity assumption is

$$E \left[\eta_{r,t} \left(\Delta X_{r,t-\tau} S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k + S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} \sum_{k=1}^{\tilde{K}} \beta_{between}^k \widetilde{W}_{r,c,t-\tau}^k + \delta M_{r,t} \right) \right] = 0.$$

Following the same reasoning as in Section 3.3, we note that the error is assumed to be uncorrelated with the weighting matrices. We explained in Section 3.2 why the error is assumed to be uncorrelated with the interaction between infectives and susceptibles. As such, we assume that the moment condition holds.

3.6 Model selection

For model selection, we use the Akaike Information Criterion or AIC (Akaike, 1974). The AIC for a particular model is defined as

$$AIC = -2 \log(ML) + 2k,\tag{3.12}$$

where ML denotes the maximum likelihood for the model and k denotes the number of parameters in the model. In contrast, one could also consider the Bayesian Information Criterion or BIC (Schwarz et al., 1978). Schwarz et al. (1978) developed it as an alternative to the Akaike Information Criterion. The BIC is defined as

$$BIC = -2 \log(ML) + k \log(n), \quad (3.13)$$

where n denotes the sample size. Both the AIC and BIC are used as the minimizer in the model selection. That is, the model that is picked by the model selection procedure is the one with the lowest AIC or BIC. When choosing between the two methods, one should realize that they have different properties, particularly related to consistency. The AIC tends to select a larger model than the BIC. Moreover, if the true model is included in the set of candidate models, and under some additional assumptions, the BIC will select the true model with probability one as n goes to infinity whereas the AIC is not consistent. On the other hand, if the true model is not in the set of candidate models, clearly no method can possibly select the true model. However, the AIC is efficient in the sense that it will asymptotically select the model that minimizes the mean prediction error while the BIC is not efficient (Vrieze, 2012). Proponents of using the AIC over the BIC argue that this shows that the AIC is to be preferred because it is virtually impossible for the true model to be constructed because *“all models are wrong”* (Box, 1976). That does not mean that reality cannot be modelled; some models can be useful despite not being perfectly true. Burnham and Anderson (2002) state that *“A model is a simplification or approximation of reality and hence will not reflect all of reality. [...] While a model can never be “truth,” a model might be ranked from very useful, to useful, to somewhat useful to, finally, essentially useless”*. Lastly, Vrieze (2012) shows by simulation that the BIC can fail in finite sample sizes even if the true model is in the candidate set. This is because the BIC has a higher maximum risk, defined as the mean squared error of estimating the true covariance matrix. Because we believe that, indeed, the true model generating the data will quite likely not be included in our candidate set, we use the AIC to perform model selection.

3.7 Modelling undocumented infections

A common concern with the spread of viruses, especially one so rapidly spreading as SARS-CoV-2, is that there is no possibility to test the entire population on whether they are infected because the testing capacity is simply not there. If this were possible, then all individuals who were tested to be positive could be isolated and the spread of the virus would be dampened tremendously. However, since this is not possible, there are likely many infectives in society who spread the virus but who are undocumented. In China, around 86% of the infections went undocumented (R. Li et al., 2020). R. Li et al. (2020) also estimate that these were also contagious, with around 55% of the contagiousness of documented infectives. This was investigated during the period from January 10 till

January 23, 2020, so considering a lack of major restrictions such as travel bans. R. Li et al. (2020) make the important note that these results are indeed highly dependent on the specific situation in the country of interest, for instance due to differences in testing, case definition, and reporting. Nonetheless, even if these numbers are lower for other cases, such as Italy under lockdown, this research shows that undocumented infections should be taken into account.

In this section, we aim to model the undocumented infections. Note that, by definition, there is no data on the amount of undocumented infections because, otherwise, these cases would indeed be documented. As such, some assumptions need to be made since we cannot apply *supervised learning* methods (being models where there is a data on a dependent variable to predict) to determine the number of undocumented infections. Firstly, we assume that the amount of undocumented individuals is decreasing as the testing capacity increases. Similarly, the amount of documented individuals increases in the testing capacity. The logic behind this is clear: as more people are tested, more infectives move from being undocumented to being documented. Secondly, as mentioned, R. Li et al. (2020) consider that there are no major restrictions. As we know, Italy has been under a strict national lockdown. This was imposed on March 10, 2020. The restrictions were relaxed around May 18, when businesses were allowed to reopen and citizens were allowed free movement within the region they live in, although they were still barred from travelling to other regions unless they had an essential motive (Severgnini, 2020). However, we do not take this into account in this thesis besides including the indicator variable for the lockdown, as described in Section 4.3. Future research could be done to include these restrictions more robustly.

At a point in time t , we denote the testing capacity by TC_t . In Section 4.2, we explain how a measure of the testing capacity is obtained. The total number of infected people at time t is denoted by X_t . This group can be subdivided into the documented infections DI_t and the undocumented infections UI_t such that $DI_t + UI_t = X_t$. Therefore, we can denote the documented and undocumented infections as proportions of the total number of infected people, at any point in time. As mentioned before, this proportion may change over time as the testing capacity increases, among others. This proportion is therefore defined as a function of the testing capacity over time:

$$f_t := f(TC_t), \tag{3.14}$$

such that

$$\begin{cases} DI_t &= f_t X_t \\ UI_t &= (1 - f_t) X_t. \end{cases}$$

Notice that the undocumented infectives can then be written as $UI_t = \frac{1-f_t}{f_t} DI_t$.

There are some properties that (3.14) should satisfy and some assumptions that we make. These are as follows:

- (A1) Since f_t is a proportion, we need to have $f_t \in [0, 1]$.
- (A2) If no one is tested, we assume that there are a certain minimum amount of documented infections, denoted by $f^{min} \in [0, 1]$. That is,

$$f(0) = f^{min}.$$

Note that at any point in time, it should hold that

$$\begin{aligned} DI_t + UI_t &< N_t \\ \iff DI_t + \frac{1 - f_t}{f_t} DI_t &< N_t \\ \iff \frac{1}{f_t} DI_t &< N_t \\ \iff f_t &> \frac{DI_t}{N_t}, \end{aligned}$$

so f^{min} should be chosen to be larger than $\frac{DI_t}{N_t}$. The fact that this is true should be clear. If f_t would be lower than the fraction of the population that is documented to be infective, then the total number of infectives in a population would exceed the total number of people living in that population, which is not possible.

- (A3) Denote the total population at time t as N_t . Then, if there is enough testing capacity such that the entire population can be tested, we assume that all infections will be documented. That is,

$$f(N_t) = 1.$$

This also assumes that the tests that are executed are perfect at determining whether someone actually is infected. However, it is common knowledge that such tests have a certain rate of false positives and negatives. In the case of COVID-19 specifically, positive screening tests are not followed-up (as is usually common practice to confirm a diagnosis) because of scarcity in testing resources and/or prioritization of allocating tests to the sickest patients (Frasier, 2020). Moreover, BMJ (2020) reports that serological tests for COVID-19 carry with them risks of bias and heterogeneity in their accuracy. Therefore, they state that these serological tests should only be used cautiously for clinical decision making and epidemiological surveillance. For this reason, one could choose to relax the assumption and assume $f(N_t) = f^{max}$ for some $f^{max} \in [0, 1]$ set to be a more reasonably perceived value.

- (A4) As mentioned earlier in this section, f_t needs to be monotonically increasing in TC_t . That is, the proportion of infectives that are documented is increasing in the testing capacity. Mathematically, this means that

$$f'(N_t) \geq 0.$$

We test several functional forms of the function f_t . Derivations are given in appendix D.

- **Linear form**

$$f_t = \frac{1 - f^{min}}{N_t} TC_t + f^{min}. \quad (3.15)$$

- **Quadratic form**

We specify three functional forms for a quadratic form. First of all, a general form. After this, we discuss two special cases.

- For the general quadratic form, we assume without loss of generality that $f\left(\frac{1}{2}N_t\right) = \gamma$ for some $\gamma \in \left[\frac{1}{4} + \frac{3}{4}f^{min}, \frac{3}{4} + \frac{1}{4}f^{min}\right]$. Then the formula becomes:

$$f_t = \frac{2 - 4\gamma + 2f^{min}}{N_t^2} TC_t^2 + \frac{4\gamma - 1 - 3f^{min}}{N_t} TC_t + f^{min}. \quad (3.16)$$

If $\gamma \in \left[\frac{1}{4} + \frac{3}{4}f^{min}, \frac{1}{2} + \frac{1}{2}f^{min}\right)$, the function is upwards opening. If $\gamma \in \left(\frac{1}{2} + \frac{1}{2}f^{min}, \frac{3}{4} + \frac{1}{4}f^{min}\right]$, the function is downwards opening. If $\gamma = \frac{1}{2} + \frac{1}{2}f^{min}$, then the formula simplifies to the linear specification. In appendix D.2.2, we explain why γ cannot be below $\frac{1}{4} + \frac{3}{4}f^{min}$ or above $\frac{3}{4} + \frac{1}{4}f^{min}$.

- We assume that the vertex (i.e. the extremum) is the point $(N_t, 1)$, i.e. the parabola is downwards opening. Note that any quadratic function can be rewritten to the so-called vertex form $f(x) = a(x - h)^2 + k$, where the vertex of the function is (h, k) . Choosing this special case means that there will be no unknown parameters needed to define the function because we know the location of the vertex and a known point $(0, f^{min})$ on the parabola. We can then derive that the formula becomes:

$$f_t = \frac{f^{min} - 1}{N_t^2} TC_t^2 - \frac{2(f^{min} - 1)}{N_t} TC_t + f^{min}. \quad (3.17)$$

Note that this is equivalent to (3.16) for $\gamma = \frac{3}{4} + \frac{1}{4}f^{min}$. Therefore, this is a boundary case for a downwards opening quadratic function.

- For the same reason as for the previous specification, we assume that the vertex is the point $(0, f^{min})$, i.e. the parabola is upwards opening. We can then derive that the formula becomes:

$$f_t = \frac{1 - f^{min}}{N_t^2} TC_t^2 + f^{min}. \quad (3.18)$$

Note that this is equivalent to (3.16) for $\gamma = \frac{1}{4} + \frac{3}{4}f^{min}$. Therefore, this is a boundary case for an upwards opening quadratic function.

- **Cubic form**

For the cubic form, we assume without loss of generality that $f(\frac{1}{4}N_t) = \gamma_1$ and $f(\frac{1}{2}N_t) = \gamma_2$ for some $\gamma_1, \gamma_2 \in (0, 1)$ such that $\gamma_1 < \gamma_2$. Then the formula becomes:

$$\begin{aligned} f(TC_t) = & \frac{8 + 64\gamma_1 - 48\gamma_2 - 24f^{min}}{3N_t^3}TC_t^3 + \frac{-2 - 32\gamma_1 + 20\gamma_2 + 14f^{min}}{N_t^2}TC_t^2 \\ & + \frac{1 + 32\gamma_1 - 12\gamma_2 - 21f^{min}}{3N_t}TC_t + f^{min}, \end{aligned} \quad (3.19)$$

No bounds on γ_1 and γ_2 have been set. Particularly, there are combinations of γ_1 and γ_2 for which the codomain of f_t on $TC_t \in [0, N_t]$ may not be the interval $[0, 1]$, violating assumption (A1), and for which the function is not monotonically increasing, violating assumption (A4). One could derive explicit conditions on possible combinations for γ_1 and γ_2 such that this is not the case but this is not done in this thesis.

These definitions can easily be generalised to be applicable to regions by considering the total population in a region $N_{r,t}$ instead of the total population N_t . Then, the function would be dependent on r as well:

$$f_{r,t} := f(TC_{r,t}). \quad (3.20)$$

such that

$$\begin{cases} DI_{r,t} &= f_{r,t}X_{r,t} \\ UI_{r,t} &= (1 - f_{r,t})X_{r,t}. \end{cases}$$

In Figure 3.1, we specify several functional forms for the specifications as mentioned above. Figure 3.1a shows four different functional forms for the quadratic functional forms while Figure 3.1b shows four different functional forms for the cubic specification.

Note that not all of the plots in Figure 3.1 are meant to be realistic portrayals. They simply show how the functions behave as the parameters change. Moreover, recall that there are combinations of γ_1 and γ_2 for the cubic representation for which assumptions (A1) and (A4) are violated. Figure 3.1b shows that $\gamma_1 = 0.35$ and $\gamma_2 = 0.7$ cause the function to exceed the maximum value allowed for $f_{r,t}$ of 1, despite decreasing so that $f(N_{r,t}) = 1$. A combination of $\gamma_1 = 0.65$ and $\gamma_2 = 0.7$ creates a non-monotonic functional form. As explained earlier in this section, this is not desirable. Henceforth, if we would use

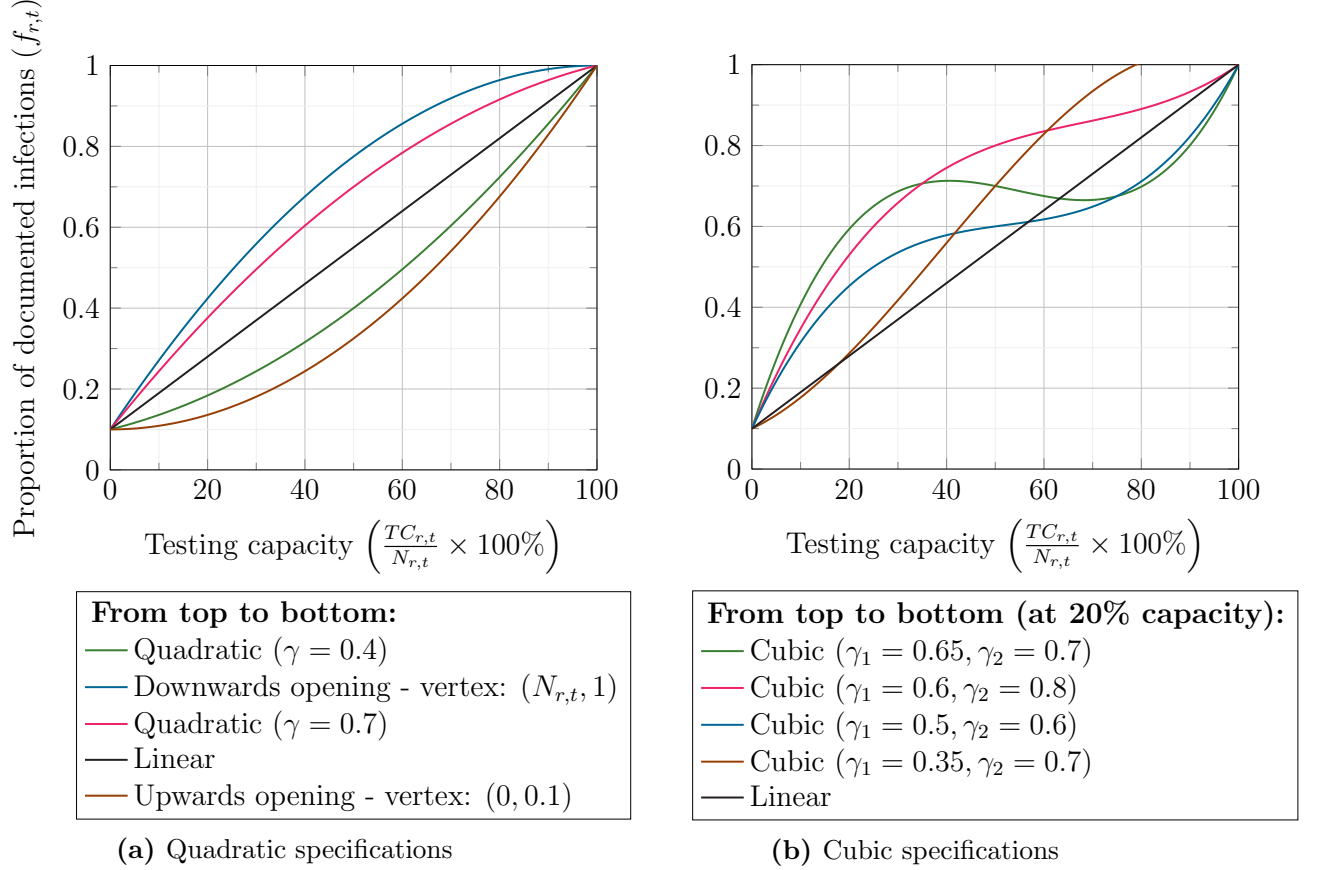


Figure 3.1. Functional forms for the proportion of documented infectives ($f^{\min} = 0.1$)

a cubic form, the values of γ_1 and γ_2 should first be tested by means of a plot, for instance.

Next, we argue which of these forms is most appropriate. As mentioned at the beginning of this section, we cannot estimate which form would fit the data best because there is, by definition, no data on the undocumented infections. As such, we argue which functional form to use by theoretical logic rather than a mathematical approach. Before that, there are two things to notice. Firstly, note that it is difficult to test 100% of the population without some rigorous metric, such as making it obligatory to get the test. Secondly, the shape of the functional form may differ depending on the basic reproduction number R_0 , as defined in Section 2. R_0 estimates how many people an infective will on average infect. If $R_0 > 1$, a person is estimated to infect more than one person and an epidemic is expected to develop. In this case, we expect that an increased testing capacity will have a larger immediate effect. We assume that a person who has been tested positive adheres to the common guidelines that they should self-quarantine. Consequently, this infective does not infect other people who would otherwise become undocumented

infectives. For the remainder of this argument, we assume that $R_0 > 1$. The reason for this is that there are a variety of methods to estimate R_0 and that we cannot reasonably make our own model easily dependent on these varying results. Future research could be conducted regarding a two-step approach.

The main question that we need to ask ourselves is whether the impact of a change in testing capacity is different relative to the initial testing capacity. That is, if the testing capacity is low and we increase it, does that have a larger effect on the proportion of documented infectives than when testing capacity is high and we increase it by the same amount?

We first argue why a downwards opening quadratic function fits the requirements well. Note that when a large proportion of the population has been tested, the pool of untested people, who are potentially infectious, is smaller. The probability that they, in isolation of other effects, are infected is lower. The argument for this is as in the previous paragraph: assuming that the people close to them who were tested positive (be that family, acquaintances, or those that they would perhaps run into at the supermarket) do indeed self-isolate, they would not have been able to be in contact with them and they have a lower chance to be infected. When a small number of people is tested and suddenly the testing capacity is increased, a larger pool of people who had symptoms and could previously not be tested, now have access to a test. The people who are now most likely to get tested positive have strong symptoms. As they are now tested positive, we assume they self-quarantine and cannot infect other people. Therefore, the functional form that fits this argument best is a downwards opening quadratic function.

One could also consider the cubic representation with $\gamma_1 = 0.6$ and $\gamma_2 = 0.8$, or some similar parameter values, as in Figure 3.1b. There, we see similar behaviour at the start of the graph where there is a sharp increase, after which it levels out. The difference is found when the last proportion of the population is tested, leading to a sudden sharp increase in the proportion of documented infections. An argument in favour of this specification is that it may be difficult to track down and convince the last proportion of the population to take a test who, at that point, may be infectious. For instance, these may simply be people who refuse to take such a test, whether those reasons are grounded or not. There may also be people who underestimate their symptoms or their importance and who, even though they are encouraged to get tested, believe that they do not need to be. For instance, they may feel that others need to get the test more. If these people are to be reached, a more rigorous program is needed and this may cause the sharp rise as a high testing capacity is reached.

Weighing these two specifications off, we believe that the former argument is more general and stable, where the second argument is quite specific and whose validity may

differ across countries. In general, of all possible fitting solutions, the one with the least number of assumptions needed is to be preferred. Therefore, we opt to use a downwards opening quadratic functional form over a cubic form.

Lastly, the question is what to choose for the parameter γ , if anything. Recall that (3.16) and (3.17) are equivalent when $\gamma = \frac{3}{4} + \frac{1}{4}f^{min}$, meaning that (3.17) is the most extreme case possible and that the slope cannot be constructed to be more steep. To be general, we choose (3.16) to be our functional form with an unknown parameter γ , denoted by $f_{r,t}(\gamma)$.

Now that we have chosen our functional form, we can adapt the models (3.8)-(3.11) to include these undocumented infections. Let us take the within-region spread model (3.8) as an example. Recall that this model was given as

$$I_{r,t} = \beta_{within} I_{r,t-\tau} S_{r,t-\tau} + X_{r,t} \delta + \eta_{r,t}.$$

Using that $I_{r,t} = \frac{DI_{r,t}}{f_{r,t}}$, this becomes

$$\frac{DI_{r,t}}{f_{r,t}(\gamma)} = \beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} + X_{r,t} \delta + \eta_{r,t}. \quad (3.21)$$

We can rewrite (3.21) as follows

$$\begin{aligned} DI_{r,t} &= f_{r,t}(\gamma) \left(\beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} + X_{r,t} \delta + \eta_{r,t} \right) \\ \iff DI_{r,t} &= \beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} f_{r,t}(\gamma) + X_{r,t} \delta f_{r,t}(\gamma) + \eta_{r,t} f_{r,t}(\gamma). \end{aligned}$$

The moment conditions that need to hold are:

$$\begin{aligned} E \left[\eta_{r,t} f_{r,t}(\gamma) \left(\beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} f_{r,t}(\gamma) + X_{r,t} \delta f_{r,t}(\gamma) \right) \right] &= 0 \\ \iff E \left[\eta_{r,t} f_{r,t}^2(\gamma) \left(\beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} + X_{r,t} \delta \right) \right] &= 0 \\ \iff f_{r,t}^2(\gamma) E \left[\eta_{r,t} \left(\beta_{within} \frac{DI_{r,t-\tau}}{f_{r,t-\tau}(\gamma)} S_{r,t-\tau} + X_{r,t} \delta \right) \right] &= 0. \end{aligned}$$

Since $f_{r,t}^2(\gamma)$ is simply a scaling function, regardless of the chosen parameter, it has no influence on the dependence between the error and the regressors. As such, it can be taken out of the expectation term. Subsequently, we can divide both sides of the equation by $f_{r,t}^2(\gamma)$ to obtain the original moment condition of (3.8):

$$E [\eta_{r,t} (\beta_{within} I_{r,t-\tau} S_{r,t-\tau} + X_{r,t} \delta)] = 0.$$

Therefore, the scaling of the infectives by using our functional form, has no additional impact on the moment conditions. A similar logic applies to the moment conditions for (3.9)-(3.11) so that their moment conditions also do not depend on $f_{r,t}(\gamma)$.

Lastly, we are interested in investigating the relationship between $TC_{r,t}$ and $f_{r,t}(\gamma)$ over time for all regions and to compare these. Because the population size differs over the regions, this is likely to impact the absolute number of tests executed. As such, instead of comparing $f_{r,t}(\gamma)$ to $TC_{r,t}$, we compare it to $TC_{r,t}/N_{r,t}$. The results are shown in Figure 3.2.

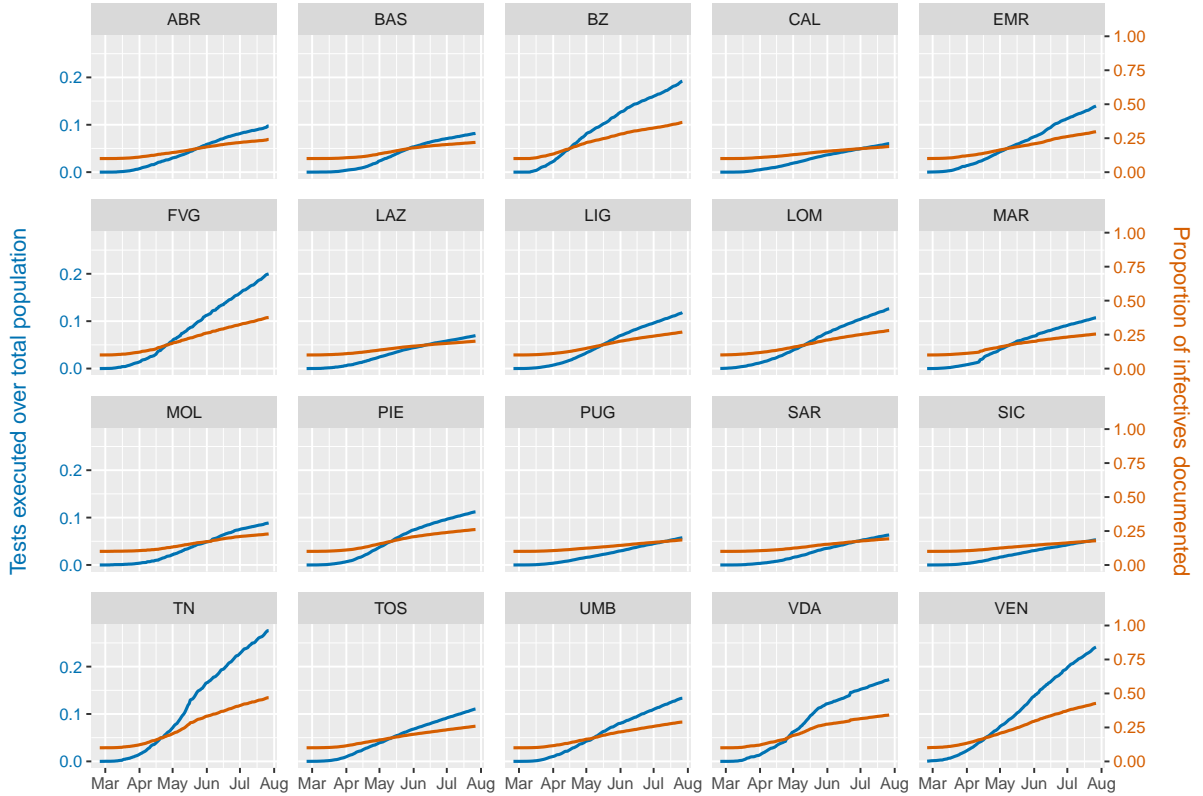


Figure 3.2. Total number of people tested over the total population ($TC_{r,t}/N_{r,t}$) versus proportion of infectives that are documented $f_{r,t}(\gamma)$

In this figure, we can see that the pattern of the relationship between the two variables is similar over time for different groups of regions. Importantly, we note that the proportion of infectives that go undocumented decreases over time. This is logical because the testing capacity increases over time and we have assumed a monotonic relationship.

4 Dataset

In this section, we outline the structure of the data that is used and how it was retrieved. Firstly, we discuss the structure of Italian regions in Section 4.1. Subsequently, we look at the data on COVID-19 such as the incidence rate in Section 4.2. Here, we also discuss how we tackled possibly errors in the data, as well as missing values. Lastly, Section 4.3 discusses the variables that are included in the weighted models in Sections 3.3 and 3.5.

4.1 Geographical structure of Italy

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* (Trento-South Tyrol) 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*, comparable to *Het Gooi*, *Twente*, or the *Achterhoek* in the Netherlands.

4.2 Coronavirus data

The *Presidenza del Consiglio dei Ministri - Dipartimento della Protezione Civile* (Presidency of the Council of Ministers - Department of Civil Protection), hereafter referred to as the Department of Civil Protection, has posted daily reports containing tables with a detailed numerical overview of new cases, active intensive care (IC) patients, tests executed, and more (Rosini, 2020). This data is all divided up between the NUTS 2 regions. 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 7, 2020 until the strict national lockdown was instated. Unfortunately, the data was not reported at this granular level. As such, we choose to use the NUTS 2 regions.

For $R = 21$ Italian regions, we retrieved the data on the coronavirus from February 25, 2020, until July 17, 2020, leading to a total amount of time observations of $T = 144$. The statistics that are of interest to us are:

- New amount of current positive cases (*nuovi_positivi*);
- Total amount of deaths (*deceduti*);
- Total amount of recoveries (*dimessi_guariti*);

TODO:
Up-
date
ac-
cord-
ingly
if this
changes,
also
 $T=x$

- Total amount of positive cases (*totale_casi*);
- Total amount of tests performed (*tamponi*);
- Total number of people tested (*casi_testati*).

The report also contains, for instance, the number of active ICU cases (*terapia_intensiva*) and the number of hospitalized people who showed symptoms (*ricoverati_con_sintomi*).¹ There are two notes to make. Firstly, the data source states that the new amount of current positive cases at time t is defined as the first difference of the total amount of positive cases: ($totale_casi_t - totale_casi_{t-1}$). However, this is not always the case. To illustrate, we consider the region of Abruzzo on June 16 till June 18. The daily number of positive tests equal 1, 0, and -1, respectively, while the number of new confirmed cases equal 2, 2, and 1, respectively. This is likely a small measurement or computational error. We take the first difference of the total amount of positive cases to define the number of confirmed cases. Secondly, the semantic difference between the total amount of tests performed (*tamponi*) and the total amount of people tested (*casi_testati*) is that the latter indicates the number of unique persons that were tested because 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 Department of Civil Protection. 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 positive for COVID-19 before or after their death, regardless of whether they died inside or outside the hospital, assuming that these deaths were reported. 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). In the UK, there is a radical difference between the

¹Official data descriptions of all variables can be found at <https://github.com/pcm-dpc/COVID-19/blob/master/dati-andamento-covid19-italia.md>

total number of deaths until June 28 with a positive test result (43,575 deaths), the total number of deaths until June 19 where COVID-19 is mentioned on the death certificate (53,858 deaths), and the total number of deaths until June 19 over and above the usual number at that time of the year (65,132 deaths) (BBC News, 2020). This shows that the UK reports deaths due to COVID-19 on the death certificates even for people who were not tested positive. Moreover, there are many excess deaths over the usual number that may or may not be due to COVID-19 that are now not counted in the official reports.

We also make the note that it is unclear how the Department of Civil Protection 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 or do so inaccurately. For instance, different regions may adhere to different principles when deciding whether a death is classified as being due to COVID-19. 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 or, otherwise, the error will be assumed to be consistently applied to the data received from that region. 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 the number of confirmed cases or the number of deaths was reported to be negative. 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 close to June 12. Propagating this error backwards as described before would lead to zero new cases per day for Campania from May 13 until June 12. Since we have no reason to know how this error is distributed, we remove the region of Campania from our dataset. Another solution could be to distribute the error according to the daily number of cases relative to the total amount of cases until June 12.

TODO: There are no more missing values because we use a new data source. However, there are some cases where the data is zero but this is illogical, for instance BAS_tested on July 5 and 12. Update this section accordingly.

Regarding missing data, there are only three cases, namely for Abruzzo on March 10,

TODO:
Up-
date
ac-
cord-
ingly
if this
changes

TODO:
Up-
date
ac-
cord-
ingly
if this
changes

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. Because we have no way of knowing how these values are distributed over the two days, missing data is simply imputed with a value of 0. One could, on the other hand, assume a certain spread, such as fifty-fifty.

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

4.3 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 as described in Section 4.1, 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 only available on an annual basis and is often not up-to-date. That is, sometimes data is available only up to 2016. For each variable, we keep the most recent data and assume 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 β . 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: Insert \widetilde{W} variables; if still applicable

Table 4.2. Specification of regressors

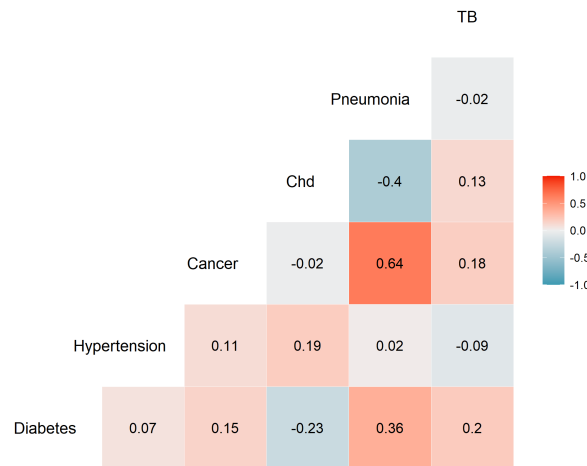
Matrix	Variable	Year	Description
$X_{r,t}$	weekend	n/a	Binary indicator denoting if the day is on the weekend (Saturday or Sunday)
$W_{r,t}$	touristArrivals	2018	Number of tourist arrivals
	deathRateDiabetes	2016	Number of deaths from diabetes per 100,000 inhabitants

Table 4.2 continues on next page

Table 4.2 continued from previous page

Matrix	Variable	Year	Description
$\widetilde{W}_{c,r,t}$	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
	deathRatePneumonia	2016	Number of deaths from pneumonia per 100,000 inhabitants
	availableBeds	2018	Number of hospital beds
	populationDensity	2018	Amount of people per square kilometer

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.



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

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 dis-

charge rates of pneumonia and cancer. We also see a relatively high correlation of -0.4 between the discharge rates of pneumonia and coronary heart disease. For this reason, we remove the discharge rate of pneumonia from the model.

5 Results

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

5.1 Model 1: Within-Region Spread

In this section, we present the results for the within-region spread model. Recall that this was given in equation (3.8) as:

$$I_{r,t} = \beta_{within} I_{r,t-\tau} S_{r,t-\tau} + X_{r,t} \delta + \eta_{r,t}$$



Firstly, we present the results where the data is pooled to a national level. Subsequently, results are presented for the models per region to which model selection is applied with the Akaike Information Criterion (AIC). For both result sets, we present the results from the regular model as well as modelling the undocumented infections with a quadratic form with $\gamma = 0.7$ and $f^{min} = 0.1$ as in (3.16).

Naively, one could consider constructing a model for the entire nation of Italy. Even though this does not take into account regional differences, as described in Section 2, it may achieve good results if regions are sufficiently similar. The results from estimating (3.8) with OLS for the national data are given in Table 5.1.

TODO: Update dates in table caption

Table 5.1. Estimates from Model 1: Within-region spread on a national level. Data spans February 25 till August 9, 2020. Undocumented infections are modelled using the quadratic specification with $\gamma = 0.7$ and $f^{min} = 0.1$.

	Regular model				Modelling undocumented infections			
	Estimate	Std. Error	t value	p value	Estimate	Std. Error	t value	p value
Intercept	3.88	96.20	0.04	0.97	43.56	948.18	0.05	0.96
Weekend	515.81	136.82	3.77	0.0002***	5044.71	1349.91	3.74	0.0003***
β_{within}	0.92	0.04	25.93	0.0000***	0.92	0.04	26.03	0.0000***

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

Table 5.1 shows an estimate for β_{within} of 0.92 that is statistically highly significant at a 1% significance level, whether undocumented infections are modelled or not. Note that it is quite small compared to the other estimates. This is because this represents the estimated effect of only a unit change in $I_{t-\tau} S_{t-\tau}$. Because Italy has many inhabitants, this means that a unit change is relatively small. A similar effect will be observed for the regional models.

Of course, this model does not take into account effects specific to regions. In Table B.2 in Appendix B, we present the results from running the model on each region separately with the same model specification for each region. It is clear that the same model might not be suitable for all regions. That is, we should apply model selection to the individual models. To execute model selection, we use the AIC and we make sure that the term for β_{within} remains in the model. The models also retain an intercept. As such, model selection is solely performed on whether the weekend dummy should be included. In Table 5.2, we present the results. The results from using the BIC versus the AIC are presented in Table B.3

TODO: Replace p -values with standard errors, add stars and parentheses

TODO: Update dates in table caption

Table 5.2. Estimates from Model 1: Within-region spread per region with model selection by AIC. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are modelled using the quadratic specification with $\gamma = 0.7$ and $f^{min} = 0.1$.

Region	Regular model			Modelling undocumented infections		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
National	0.92	3.88	515.81	0.92	43.56	5044.71
	0.00	0.97	0.00	0.00	0.96	0.00
ABR	0.62	6.76	8.16	0.62	67.36	80.01
	0.00	0.03	0.10	0.00	0.03	0.10
BAS	0.62	0.71	1.44	0.62	7.03	14.42
	0.00	0.14	0.08	0.00	0.15	0.08
BZ	0.69	6.16		0.69	59.72	
	0.00	0.01		0.00	0.01	
CAL	0.56	2.72	4.37	0.56	26.99	43.42
	0.00	0.06	0.07	0.00	0.06	0.07
EMR	0.89	6.79	58.19	0.89	66.36	569.70
	0.00	0.63	0.01	0.00	0.63	0.01
FVG	0.67	8.11		0.67	79.80	
	0.00	0.00		0.00	0.00	
LAZ	0.87	5.28	11.08	0.87	52.05	110.11
	0.00	0.17	0.05	0.00	0.18	0.05
LIG	0.80	10.41	14.50	0.80	102.59	141.85
	0.00	0.08	0.09	0.00	0.08	0.09
LOM	0.83	62.77	189.03	0.83	614.62	1844.96
	0.00	0.21	0.01	0.00	0.22	0.01
MAR	0.85	2.49	17.05	0.85	24.44	169.85
	0.00	0.54	0.01	0.00	0.55	0.01
MOL	0.34	1.46	2.32	0.35	14.41	22.91
	0.00	0.02	0.03	0.00	0.02	0.03
PIE	0.87	12.83	64.98	0.87	126.77	636.20

Table 5.2 continues on next page

Table 5.2 continued from previous page

Region	Regular model			Modelling undocumented infections (Q)		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
PUG	0.00	0.44	0.01	0.00	0.43	0.01
	0.80	6.59		0.80	65.48	
SAR	0.00	0.02		0.00	0.02	
	0.62	3.81		0.62	37.86	
SIC	0.00	0.01		0.00	0.01	
	0.78	5.17		0.78	51.38	
TN	0.00	0.02		0.00	0.02	
	0.42	15.91	15.63	0.43	152.57	153.48
TOS	0.00	0.01	0.09	0.00	0.01	0.09
	0.85	3.10	28.52	0.85	30.58	281.93
UMB	0.00	0.61	0.00	0.00	0.61	0.00
	0.76	1.50	3.89	0.76	14.77	38.22
VDA	0.00	0.32	0.13	0.00	0.32	0.13
	0.47	2.72	6.62	0.47	26.65	65.45
VEN	0.00	0.08	0.01	0.00	0.08	0.01
	0.85	12.08	33.54	0.85	118.04	324.08
	0.00	0.31	0.06	0.00	0.31	0.07

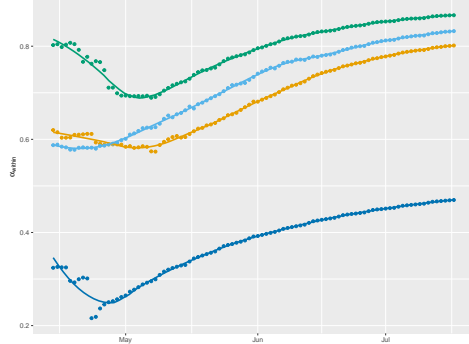
Significance levels: * = 0.05, ** = 0.01, *** = 0.001

Indeed, we see that the AIC gives a varying model selection per region. As mentioned, all models retain the intercept and the term $I_{t-\tau}S_{t-\tau}$ in the model. In 16 out of 21 cases, the entire model is selected. In the other 5 cases, the weekend dummy is excluded. Table 5.2 also shows that the estimate for β_{within} varies vastly over the regions, indicating that a pooled national model is not suitable. If we do not model undocumented infections, it ranges from 0.34 for Molise till 0.92 for the national model. Excluding the national model, the highest estimate of β_{within} is found for the region of Emilia-Romagna (EMR) with 0.89. The estimated parameters for the other variables vary a bit more, but this is likely due to the differences in the population count, thereby affecting the magnitude of the dependent variable $I_{r,t}$. If undocumented infections are modelled, the estimates for β_{within} do not change much compared to when only the documented infections are used. The other estimated parameters are higher, roughly by a factor 10.

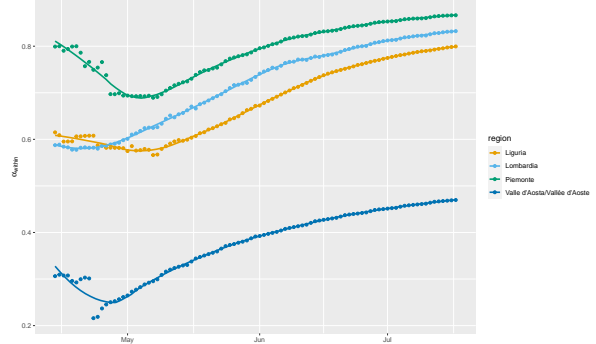
We are also interested in looking at the estimate of β_{within} over time. That is, if we keep adding data, do we see an interesting effect in its progression? We use at least 50 data points. In Figure 5.1 we present plots for the regions in the *Nord-Ovest* (North-West) NUTS 1 region. Plots for the other NUTS 1 regions can be found in Appendix C.1. Each point in the graphs in Figure 5.1 is the estimate of β_{within} when only data before that date was used.

TODO: Add national plot and text

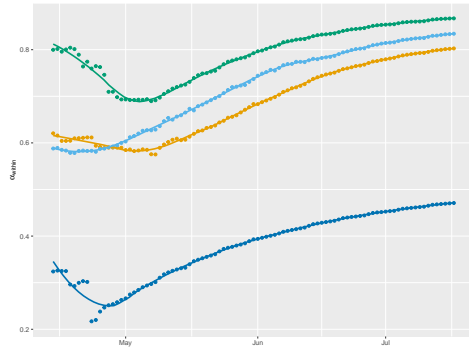
Update
ac-
cord-
ingly



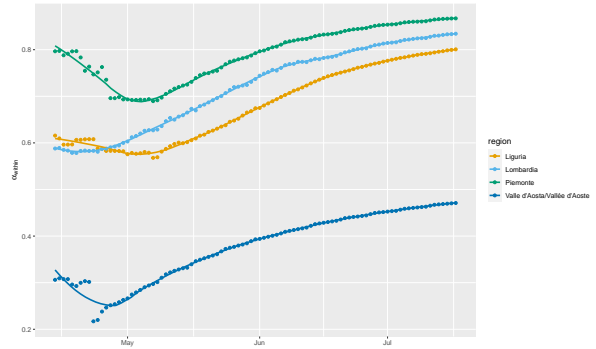
(a) Without model selection



(b) With model selection by AIC



(c) Without model selection; with modelling undocumented infections



(d) With model selection by AIC; with modelling undocumented infections

Figure 5.1. Progression of β_{within} over time for the *Nord-Ovest* (North-West) NUTS 1 region

In these figures, we see that the value for β_{within} is increasing and levels out as more data is added. Unfortunately, this is not a positive progression. Recall that β_{within} represents the transmission rate. Of course, we would like to see that the transmission rate decreases over time.

TODO: Explain more and possibly consider using a rolling window.

5.2 Model 2: Weighted Within-Region Spread

In this section, we present the results for the weighted within-region spread model. Recall that this was given in equation (3.9) as:

$$I_{r,t} = I_{r,t-\tau} S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k + X_{r,t} \delta + \eta_{r,t}$$

TODO: Results still need to be inserted.

5.3 Model 3: Within and Between-Region Spread

In this section, we present the results for the within and between-region spread model. Recall that this was given in equation (3.10) as:

$$I_{r,t} = \beta_{within} I_{r,t-\tau} S_{r,t-\tau} + \beta_{between} S_{r,t-\tau} \sum_{c \in R \setminus r} I_{c,t-\tau} + X_{r,t} \delta + \eta_{r,t}$$

Notice that it does not make sense to consider a national model. Because we do not consider countries outside of Italy, the set $R \setminus r$ is empty if we consider r to be the entire country of Italy. This would mean that the national model for the within and between-region spread model is equivalent to the national model for the within-region spread model. As such, in this section, we only consider the model applied to the regions. In Table B.4 in Appendix B, we present the results from running the model on each region separately without applying model selection. The results from applying model selection with AIC are presented in Table 5.3.

TODO: Replace p -values with standard errors, add stars and parentheses

TODO: Update dates in table caption

Table 5.3. Estimates from Model 3: Within and between-region spread per region with model selection by AIC. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are modelled using the quadratic specification with $\gamma = 0.7$ and $f^{min} = 0.1$.

Region	Regular model				Modelling undocumented infections			
	β_{within}	$\beta_{between}$	Intercept	Weekend	β_{within}	$\beta_{between}$	Intercept	Weekend
ABR	-0.35	0.02	-7.32	9.05	-0.35	0.02	-72.06	88.40
	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.01
BAS	0.05	0.00	-1.46	1.45	0.05	0.00	-14.53	14.36
	0.56	0.00	0.00	0.03	0.59	0.00	0.00	0.03
BZ	0.13	0.01	-4.56	6.34	0.13	0.01	-44.16	60.60
	0.19	0.00	0.08	0.08	0.19	0.00	0.08	0.09
CAL	-0.02	0.01	-3.14	4.54	-0.02	0.01	-31.36	44.73
	0.85	0.00	0.03	0.02	0.81	0.00	0.03	0.02
EMR	0.59	0.05	-5.33	62.13	0.58	0.05	-52.98	608.06
	0.00	0.03	0.72	0.00	0.00	0.03	0.72	0.00
FVG	-0.20	0.02	-5.03	5.29	-0.20	0.02	-49.49	51.87
	0.05	0.00	0.05	0.15	0.06	0.00	0.05	0.15
LAZ	-0.18	0.04	0.17	16.75	-0.18	0.04	2.35	164.98
	0.02	0.00	0.95	0.00	0.02	0.00	0.92	0.00
LIG	0.07	0.04	-0.25	23.72	0.06	0.04	-2.23	232.51
	0.37	0.00	0.95	0.00	0.42	0.00	0.96	0.00
LOM	0.79	0.03	64.59	189.06	0.79	0.02	630.09	1845.45
	0.00	0.64	0.20	0.01	0.00	0.68	0.21	0.01

Table 5.3 continues on next page

Table 5.3 continued from previous page

Region	Regular model				Modelling undocumented infections			
	β_{within}	$\beta_{between}$	Intercept	Weekend	β_{within}	$\beta_{between}$	Intercept	Weekend
MAR	0.90	-0.00	3.39	16.86	0.90	-0.00	33.39	168.14
	0.00	0.65	0.46	0.01	0.00	0.65	0.46	0.01
MOL	0.16	0.00	-0.17	2.40	0.17	0.00	-1.70	23.58
	0.06	0.00	0.81	0.02	0.05	0.00	0.80	0.02
PIE	0.45	0.08	-18.95	71.91	0.44	0.08	-184.59	703.31
	0.00	0.00	0.17	0.00	0.00	0.00	0.18	0.00
PUG	-0.03	0.02	-7.52	7.78	-0.03	0.02	-74.24	76.09
	0.75	0.00	0.00	0.03	0.74	0.00	0.00	0.03
SAR	0.10	0.01	-2.04		0.09	0.01	-20.69	
	0.34	0.00	0.18		0.36	0.00	0.17	
SIC	-0.03	0.02	-10.78	8.32	-0.04	0.02	-107.82	81.99
	0.75	0.00	0.00	0.01	0.70	0.00	0.00	0.01
TN	-0.01	0.02	-2.24	14.68	-0.01	0.02	-23.04	142.66
	0.90	0.00	0.69	0.07	0.89	0.00	0.67	0.06
TOS	0.16	0.04	-16.23	29.90	0.16	0.04	-159.91	294.53
	0.11	0.00	0.01	0.00	0.12	0.00	0.01	0.00
UMB	0.63	0.00	-0.27	4.07	0.63	0.00	-2.74	39.97
	0.00	0.13	0.88	0.11	0.00	0.13	0.88	0.11
VDA	-0.03	0.01	-3.69	7.58	-0.03	0.01	-36.51	74.60
	0.77	0.00	0.03	0.00	0.76	0.00	0.02	0.00
VEN	-0.07	0.10	-24.51	31.20	-0.06	0.10	-235.63	301.93
	0.58	0.00	0.03	0.04	0.62	0.00	0.03	0.04

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

TODO: Add interpretation.

5.4 Model 4: Full Model (Weighted Within and Between-Region Spread)

In this section, we present the results for the full weighted within and between-region spread model. Recall that this was given in equation (3.11) as:

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

TODO: Results still need to be inserted.

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.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE transactions on automatic control*, 19(6), 716–723.
- Anderson, R. M., & May, R. M. (1992). *Infectious diseases of humans: Dynamics and control*. Oxford University Press.
- BBC News. (2020). *Death rate ‘back to normal’ in UK*. Retrieved July 1, 2020, from <https://www.bbc.com/news/health-53233066/>
- BMJ. (2020). *Diagnostic accuracy of serological tests for COVID-19: Systematic review and meta-analysis*. Retrieved July 13, 2020, from <https://www.bmj.com/content/370/bmj.m2516/>
- Box, G. E. (1976). Science and statistics. *Journal of the American Statistical Association*, 71(356), 791–799.
- Burnham, K. P., & Anderson, D. R. (2002). A practical information-theoretic approach. *Model selection and multimodel inference*, 2nd ed. Springer, New York, 2.
- Caccia, F. (2020). *Coronavirus, “il conteggio dei morti varia da paese a paese. la Germania esclude chi ha altre patologie”*. Retrieved June 11, 2020, 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
- Eurostat. (2020a). *Eurostat regional data background*. Retrieved June 11, 2020, from <https://ec.europa.eu/eurostat/web/regions/background>
- Eurostat. (2020b). *Eurostat regional statistics database*. Retrieved June 11, 2020, from <https://ec.europa.eu/eurostat/web/regions/data/database>
- Frasier, S. L. (2020). *Coronavirus antibody tests have a mathematical pitfall*. Retrieved June 19, 2020, from <https://www.scientificamerican.com/article/coronavirus-antibody-tests-have-a-mathematical-pitfall/>
- Google LLC. (2020). *Google COVID-19 community mobility reports*. <https://www.google.com/covid19/mobility/>
- He, X., Lau, E. H., Wu, P., Deng, X., Wang, J., Hao, X., Lau, Y. C., Wong, J. Y., Guan, Y., Tan, X., Et al. (2020). Temporal dynamics in viral shedding and transmissibility of COVID-19. *Nature medicine*, 26(5), 672–675.
- Horowitz, J. (2020). *Italy’s health care system groans under coronavirus — a warning to the world*. Retrieved June 11, 2020, from <https://www.nytimes.com/2020/03/12/world/europe/12italy-coronavirus-health-care.html>
- Istituto Superiore di Sanità. (2020). *Caratteristiche dei pazienti deceduti positivi all’infezione da SARS-CoV-2 in Italia*. Retrieved June 11, 2020, 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., Azman, A. S., Reich, N. G., & 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.
- Leung, H. (2020). *What we know about coronavirus immunity and reinfection*. Retrieved June 9, 2020, from <https://time.com/5810454/coronavirus-immunity-reinfection/>
- Li, Q., Guan, X., Wu, P., Wang, X., Zhou, L., Tong, Y., Ren, R., Leung, K. S., Lau, E. H., Wong, J. Y., Et al. (2020). Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *New England Journal of Medicine*.
- Li, R., Pei, S., Chen, B., Song, Y., Zhang, T., Yang, W., & Shaman, J. (2020). Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (SARS-CoV-2). *Science*, 368(6490), 489–493.
- Linton, N. M., Kobayashi, T., Yang, Y., Hayashi, K., Akhmetzhanov, A. R., Jung, S.-m., Yuan, B., Kinoshita, R., & 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. (2020). *Coronavirus: Contagion rate R_0 below 1. prudence needed in phase two says ISS*. Retrieved June 11, 2020, from http://www.salute.gov.it/portale/news/p3_2_1_1_1.jsp?lingua=italiano&menu=notizie&p=dalministero&id=4717
- Papadopoulos, A. (2018). *Chickenpox: Practice essentials, background, pathophysiology*. Retrieved June 22, 2020, from <https://emedicine.medscape.com/article/1131785-overview/>
- Rosini, U. (2020). *COVID-19*. Retrieved July 4, 2020, from <https://github.com/pcm-dpc/COVID-19/tree/master/legacy/dati-regioni>
- Schwarz, G. Et al. (1978). Estimating the dimension of a model. *The annals of statistics*, 6(2), 461–464.
- Severgnini, C. (2020). *Discorso di Conte in conferenza stampa, le riaperture dal 18 maggio: “corriamo un rischio calcolato”*. Retrieved June 18, 2020, from corriere.it/politica/20_maggio.16/discorso-conte-conferenza-stampa-oggi-decreto-18-maggio-1e810142-9785-11ea-ba09-20ae073bed63.shtml
- Sevillano, E. (2020). *Tracking the coronavirus: Why does each country count deaths differently?* Retrieved June 11, 2020, from <https://english.elpais.com/society/2020->

03-30/tracking-the-coronavirus-why-does-each-country-count-deaths-differently.html

Sutherland, J., & Gretler, C. (2020). *WHO now says role of silent virus spreaders remains unclear*. Retrieved June 18, 2020, from <https://www.bloomberg.com/news/articles/2020-06-09/who-says-symptomless-spread-is-rare-in-jolt-to-virus-efforts>

Vrieze, S. I. (2012). Model selection and psychological theory: A discussion of the differences between the Akaike information criterion (AIC) and the Bayesian information criterion (BIC). *Psychological methods*, 17(2), 228.

Worldometer. (2020). *Italy population*. Retrieved August 3, 2020, from <https://www.worldometers.info/world-population/italy-population/>

Appendices

A Abbreviations

The tables in this appendix present commonly used abbreviations in this thesis, including the regional abbreviations.

Table A.1. Abbreviations for the Italian regions.

Abbreviation	Italian name	English name
ABR	Abruzzo	Abruzzo
BAS	Basilicata	Basilicata
BZ	Alto Adige or Provincia Autonoma di Bolzano/Bozen	South Tyrol or Province of Bolzano
CAL	Calabria	Calabria
CAM	Campania	Campania
EMR	Emilia-Romagna	Emilia-Romagna
FVG	Friuli Venezia Giulia	Friuli Venezia Giulia
LAZ	Lazio	Lazio
LIG	Liguria	Liguria
LOM	Lombardia	Lombardy
MAR	Marche	Marche
MOL	Molise	Molise
PIE	Piemonte	Piedmont
PUG	Puglia	Apuglia
SAR	Sardegna	Sardinia
SIC	Sicilia	Sicily
TN	Trentino or Provincia Autonoma di Trento	Trentino or Province of Trento
TOS	Toscana	Tuscany
UMB	Umbria	Umbria
VDA	Valle d'Aosta/Vallée d'Aoste	Aosta Valley
VEN	Veneto	Veneto

Table A.2. Commonly used abbreviations in this thesis.

Abbreviation	Full name	Description (if applicable)
SARS-CoV-2	Severe Acute Respiratory Syndrome Coronavirus 2	
COVID-19	Coronavirus Disease 2019	
SIR model	Standard Inflammatory Response model	
OLS	Ordinary Least Squares	
AIC	Akaike Information Criterion	
BIC	Bayesian Information Criterion	

Table A.2 continues on next page

Table A.2 continued from previous page

Abbreviation	Full name	Description (if applicable)
NUTS	Nomenclature des Unités Territoriales Statistiques	Nomenclature of territorial units for statistics, a geocode standard for referencing the subdivisions of countries for statistical purposes.

B Tables

B.1 Illustrative tables

These tables are for quick illustration purposes in the thesis meetings and will be removed / adapted in due time.

TODO: Remove this in due time

Table B.1. Results for Model 1, 2 and 3 for the region of Lombardy

	<i>Models 1,2, and 3 for Lombardy</i>		
	(1)	(2)	(3)
Weekend	29.893 (58.887)	32.870 (58.974)	27.793 (57.570)
α_{within}	0.845*** (0.043)	0.908*** (0.077)	0.588*** (0.098)
$\alpha_{within}^{railTravelers}$		-178.391 (182.373)	
$\alpha_{between}$			0.163*** (0.056)
Intercept	85.690** (40.350)	123.299** (55.738)	90.637** (39.480)
Observations	164	164	164
R ²	0.710	0.711	0.724
Adjusted R ²	0.706	0.706	0.719
AIC	2,385.266	2,386.288	2,378.795
Residual Std. Error	343.111 (df = 161)	343.157 (df = 160)	335.406 (df = 160)
<i>Note:</i>		*p<0.1; **p<0.05; ***p<0.01	

B.2 Results from Model 1: Within-Region Spread

In Section 5.1, we presented the results from the within-region spread model (3.8):

$$\Delta X_{r,t} = \beta_{within} \Delta X_{r,t-\tau} S_{r,t-\tau} + \delta M_{r,t} + \eta_{r,t}.$$

TODO: Replace p -values with standard errors, add stars and parentheses

TODO: Update dates in table caption

Table B.2. Estimates from Model 1: Within-region spread per region without model selection. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are modelled using the quadratic specification with $\gamma = 0.7$ and $f^{min} = 0.1$.

Region	Regular model			Modelling undocumented infections		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
National	0.92	3.88	515.81	0.92	43.56	5,044.71
	0.00	0.97	0.00	0.00	0.96	0.00
ABR	0.62	6.76	8.16	0.62	67.36	80.01
	0.00	0.03	0.10	0.00	0.03	0.10
BAS	0.62	0.71	1.44	0.62	7.03	14.42
	0.00	0.14	0.08	0.00	0.15	0.08
BZ	0.69	5.19	3.36	0.69	50.68	31.59
	0.00	0.04	0.42	0.00	0.04	0.43
CAL	0.56	2.72	4.37	0.56	26.99	43.42
	0.00	0.06	0.07	0.00	0.06	0.07
EMR	0.89	6.79	58.19	0.89	66.36	569.70
	0.00	0.63	0.01	0.00	0.63	0.01
FVG	0.67	6.33	5.66	0.67	62.32	55.85
	0.00	0.04	0.23	0.00	0.04	0.23
LAZ	0.87	5.28	11.08	0.87	52.05	110.11
	0.00	0.17	0.05	0.00	0.18	0.05
LIG	0.80	10.41	14.50	0.80	102.59	141.85
	0.00	0.08	0.09	0.00	0.08	0.09
LOM	0.83	62.77	189.03	0.83	614.62	1,844.96
	0.00	0.21	0.01	0.00	0.22	0.01
MAR	0.85	2.49	17.05	0.85	24.44	169.85
	0.00	0.54	0.01	0.00	0.55	0.01
MOL	0.34	1.46	2.32	0.35	14.41	22.91
	0.00	0.02	0.03	0.00	0.02	0.03
PIE	0.87	12.83	64.98	0.87	126.77	636.20
	0.00	0.44	0.01	0.00	0.43	0.01
PUG	0.80	5.34	4.20	0.80	53.11	41.55
	0.00	0.09	0.39	0.00	0.09	0.39
SAR	0.62	3.04	2.58	0.62	30.18	25.68
	0.00	0.05	0.32	0.00	0.05	0.32
SIC	0.78	3.86	4.51	0.78	38.40	44.67
	0.00	0.13	0.29	0.00	0.13	0.29
TN	0.42	15.91	15.63	0.43	152.57	153.48
	0.00	0.01	0.09	0.00	0.01	0.09
TOS	0.85	3.10	28.52	0.85	30.58	281.93
	0.00	0.61	0.00	0.00	0.61	0.00

Table B.2 continues on next page

Table B.2 continued from previous page

Region	Regular model			Modelling undocumented infections (Q)		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
UMB	0.76	1.50	3.89	0.76	14.77	38.22
	0.00	0.32	0.13	0.00	0.32	0.13
VDA	0.47	2.72	6.62	0.47	26.65	65.45
	0.00	0.08	0.01	0.00	0.08	0.01
VEN	0.85	12.08	33.54	0.85	118.04	324.08
	0.00	0.31	0.06	0.00	0.31	0.07

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

Table B.3. Estimates from Model 1: Within-region spread per region with model selection by AIC versus BIC. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are not modelled.

Region	Model selection with AIC			Model selection with BIC		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
National	0.92	3.88	515.81	0.92	3.88	515.81
	0.00	0.97	0.00	0.00	0.97	0.00
ABR	0.62	6.76	8.16	0.61	9.28	
	0.00	0.03	0.10	0.00	0.00	
BAS	0.62	0.71	1.44	0.61	1.15	
	0.00	0.14	0.08	0.00	0.01	
BZ	0.69	6.16		0.69	6.16	
	0.00	0.01		0.00	0.01	
CAL	0.56	2.72	4.37	0.55	4.05	
	0.00	0.06	0.07	0.00	0.00	
EMR	0.89	6.79	58.19	0.89	6.79	58.19
	0.00	0.63	0.01	0.00	0.63	0.01
FVG	0.67	8.11		0.67	8.11	
	0.00	0.00		0.00	0.00	
LAZ	0.87	5.28	11.08	0.87	8.64	
	0.00	0.17	0.05	0.00	0.02	
LIG	0.80	10.41	14.50	0.80	14.63	
	0.00	0.08	0.09	0.00	0.01	
LOM	0.83	62.77	189.03	0.83	62.77	189.03
	0.00	0.21	0.01	0.00	0.21	0.01
MAR	0.85	2.49	17.05	0.85	2.49	17.05
	0.00	0.54	0.01	0.00	0.54	0.01
MOL	0.34	1.46	2.32	0.34	1.46	2.32
	0.00	0.02	0.03	0.00	0.02	0.03
PIE	0.87	12.83	64.98	0.87	12.83	64.98
	0.00	0.44	0.01	0.00	0.44	0.01
PUG	0.80	6.59		0.80	6.59	
	0.00	0.02		0.00	0.02	

Table B.3 continues on next page

Table B.3 continued from previous page

Region	Model selection with AIC			Model selection with BIC		
	β_{within}	Intercept	Weekend	β_{within}	Intercept	Weekend
SAR	0.62	3.81		0.62	3.81	
	0.00	0.01		0.00	0.01	
SIC	0.78	5.17		0.78	5.17	
	0.00	0.02		0.00	0.02	
TN	0.42	15.91	15.63	0.41	20.95	
	0.00	0.01	0.09	0.00	0.00	
TOS	0.85	3.10	28.52	0.85	3.10	28.52
	0.00	0.61	0.00	0.00	0.61	0.00
UMB	0.76	1.50	3.89	0.75	2.65	
	0.00	0.32	0.13	0.00	0.04	
VDA	0.47	2.72	6.62	0.47	2.72	6.62
	0.00	0.08	0.01	0.00	0.08	0.01
VEN	0.85	12.08	33.54	0.84	22.63	
	0.00	0.31	0.06	0.00	0.03	

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

B.3 Results from Model 2: Weighted Within-Region Spread

In Section 5.2, we presented the results from the within and between-region spread model (3.9):

$$\Delta X_{r,t} = \Delta X_{r,t-\tau} S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k + \delta M_{r,t} + \eta_{r,t}$$

TODO: Results still need to be inserted.

B.4 Results from Model 3: Within and Between-Region Spread

In Section 5.3, we presented the results from the within and between-region spread model (3.10):

$$\Delta X_{r,t} = \beta_{within} \Delta X_{r,t-\tau} S_{r,t-\tau} + \beta_{between} S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} + \delta M_{r,t} + \eta_{r,t}$$

This appendix contains additional tables with results for this model.

TODO: Replace p -values with standard errors, add stars and parentheses

TODO: Update 0.00 with scientific notation for $\beta_{between}$

TODO: Update dates in table caption

Table B.4. Estimates from Model 3: Within and between-region spread per region without model selection. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are modelled using the quadratic specification with $\gamma = 0.7$ and $f^{min} = 0.1$.

Region	Regular model				Modelling undocumented infections			
	β_{within}	$\beta_{between}$	Intercept	Weekend	β_{within}	$\beta_{between}$	Intercept	Weekend
ABR	-0.35 0.00	0.02 0.00	-7.32 0.00	9.05 0.01	-0.35 0.00	0.02 0.00	-72.06 0.00	88.40 0.01
BAS	0.05 0.56	0.00 0.00	-1.46 0.00	1.45 0.03	0.05 0.59	0.00 0.00	-14.53 0.00	14.36 0.03
BZ	0.13 0.19	0.01 0.00	-4.56 0.08	6.34 0.08	0.13 0.19	0.01 0.00	-44.16 0.08	60.60 0.09
CAL	-0.02 0.85	0.01 0.00	-3.14 0.03	4.54 0.02	-0.02 0.81	0.01 0.00	-31.36 0.03	44.73 0.02
EMR	0.59 0.00	0.05 0.03	-5.33 0.72	62.13 0.00	0.58 0.00	0.05 0.03	-52.98 0.72	608.06 0.00
FVG	-0.20 0.05	0.02 0.00	-5.03 0.05	5.29 0.15	-0.20 0.06	0.02 0.00	-49.49 0.05	51.87 0.15
LAZ	-0.18 0.02	0.04 0.00	0.17 0.95	16.75 0.00	-0.18 0.02	0.04 0.00	2.35 0.92	164.98 0.00
LIG	0.07 0.37	0.04 0.00	-0.25 0.95	23.72 0.00	0.06 0.42	0.04 0.00	-2.23 0.96	232.51 0.00
LOM	0.79 0.00	0.03 0.64	64.59 0.20	189.06 0.01	0.79 0.00	0.02 0.68	630.09 0.21	1,845.45 0.01
MAR	0.90 0.00	-0.00 0.65	3.39 0.46	16.86 0.01	0.90 0.00	-0.00 0.65	33.39 0.46	168.14 0.01
MOL	0.16 0.06	0.00 0.00	-0.17 0.81	2.40 0.02	0.17 0.05	0.00 0.00	-1.70 0.80	23.58 0.02
PIE	0.45 0.00	0.08 0.00	-18.95 0.17	71.91 0.00	0.44 0.00	0.08 0.00	-184.59 0.18	703.31 0.00
PUG	-0.03 0.75	0.02 0.00	-7.52 0.00	7.78 0.03	-0.03 0.74	0.02 0.00	-74.24 0.00	76.09 0.03
SAR	0.10 0.33	0.01 0.00	-3.05 0.07	3.20 0.17	0.09 0.35	0.01 0.00	-30.64 0.07	31.56 0.17
SIC	-0.03 0.75	0.02 0.00	-10.78 0.00	8.32 0.01	-0.04 0.70	0.02 0.00	-107.82 0.00	81.99 0.01
TN	-0.01 0.90	0.02 0.00	-2.24 0.69	14.68 0.07	-0.01 0.89	0.02 0.00	-23.04 0.67	142.66 0.06
TOS	0.16 0.11	0.04 0.00	-16.23 0.01	29.90 0.00	0.16 0.12	0.04 0.00	-159.91 0.01	294.53 0.00
UMB	0.63 0.00	0.00 0.13	-0.27 0.88	4.07 0.11	0.63 0.00	0.00 0.13	-2.74 0.88	39.97 0.11
VDA	-0.03 0.77	0.01 0.00	-3.69 0.03	7.58 0.00	-0.03 0.76	0.01 0.00	-36.51 0.02	74.60 0.00
VEN	-0.07 0.58	0.10 0.00	-24.51 0.03	31.20 0.04	-0.06 0.62	0.10 0.00	-235.63 0.03	301.93 0.04

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

TODO: Replace p -values with standard errors, add stars and parentheses

TODO: Update 0.00 with scientific notation for $\beta_{between}$

TODO: Update dates in table caption

Table B.5. Estimates from Model 3: Within and between-region spread per region with model selection by AIC versus BIC. Estimates are given with p -values in parentheses. Data spans February 25 till August 9, 2020. Undocumented infections are not modelled.

Region	Model selection with AIC				Model selection with BIC			
	β_{within}	$\beta_{between}$	Intercept	Weekend	β_{within}	$\beta_{between}$	Intercept	Weekend
ABR	-0.35	0.02	-7.32	9.05	-0.35	0.02	-7.32	9.05
	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.01
BAS	0.05	0.00	-1.46	1.45	0.04	0.00	-1.01	
	0.56	0.00	0.00	0.03	0.66	0.00	0.02	
BZ	0.13	0.01	-4.56	6.34	0.15	0.01	-2.46	
	0.19	0.00	0.08	0.08	0.15	0.00	0.29	
CAL	-0.02	0.01	-3.14	4.54	-0.02	0.01	-3.14	4.54
	0.85	0.00	0.03	0.02	0.85	0.00	0.03	0.02
EMR	0.59	0.05	-5.33	62.13	0.59	0.05	-5.33	62.13
	0.00	0.03	0.72	0.00	0.00	0.03	0.72	0.00
FVG	-0.20	0.02	-5.03	5.29	-0.21	0.02	-3.39	
	0.05	0.00	0.05	0.15	0.04	0.00	0.15	
LAZ	-0.18	0.04	0.17	16.75	-0.18	0.04	0.17	16.75
	0.02	0.00	0.95	0.00	0.02	0.00	0.95	0.00
LIG	0.07	0.04	-0.25	23.72	0.07	0.04	-0.25	23.72
	0.37	0.00	0.95	0.00	0.37	0.00	0.95	0.00
LOM	0.79	0.03	64.59	189.06	0.79	0.03	64.59	189.06
	0.00	0.64	0.20	0.01	0.00	0.64	0.20	0.01
MAR	0.90	-0.00	3.39	16.86	0.90	-0.00	3.39	16.86
	0.00	0.65	0.46	0.01	0.00	0.65	0.46	0.01
MOL	0.16	0.00	-0.17	2.40	0.16	0.00	-0.17	2.40
	0.06	0.00	0.81	0.02	0.06	0.00	0.81	0.02
PIE	0.45	0.08	-18.95	71.91	0.45	0.08	-18.95	71.91
	0.00	0.00	0.17	0.00	0.00	0.00	0.17	0.00
PUG	-0.03	0.02	-7.52	7.78	-0.02	0.02	-4.99	
	0.75	0.00	0.00	0.03	0.85	0.00	0.03	
SAR	0.10	0.01	-2.04		0.10	0.01	-2.04	
	0.34	0.00	0.18		0.34	0.00	0.18	
SIC	-0.03	0.02	-10.78	8.32	-0.03	0.02	-10.78	8.32
	0.75	0.00	0.00	0.01	0.75	0.00	0.00	0.01
TN	-0.01	0.02	-2.24	14.68	-0.03	0.02	2.41	
	0.90	0.00	0.69	0.07	0.73	0.00	0.63	
TOS	0.16	0.04	-16.23	29.90	0.16	0.04	-16.23	29.90
	0.11	0.00	0.01	0.00	0.11	0.00	0.01	0.00
UMB	0.63	0.00	-0.27	4.07	0.63	0.00	1.02	
	0.00	0.13	0.88	0.11	0.00	0.15	0.55	

Table B.5 continues on next page

Table B.5 continued from previous page

Region	Model selection with AIC				Model selection with BIC			
	β_{within}	$\beta_{between}$	Intercept	Weekend	β_{within}	$\beta_{between}$	Intercept	Weekend
VDA	-0.03	0.01	-3.69	7.58	-0.03	0.01	-3.69	7.58
	0.77	0.00	0.03	0.00	0.77	0.00	0.03	0.00
VEN	-0.07	0.10	-24.51	31.20	-0.08	0.11	-14.90	
	0.58	0.00	0.03	0.04	0.52	0.00	0.14	

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

B.5 Results from Model 4: Full Model (Weighted Within and Between-Region Spread)

In Section 5.4, we presented the results from the within and between-region spread model (3.11):

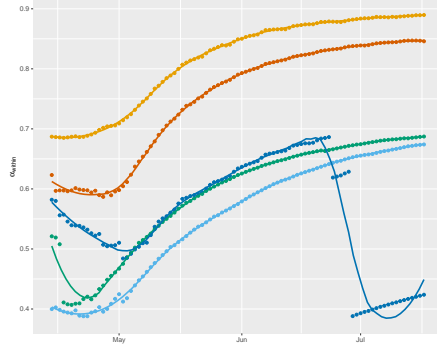
$$\begin{aligned}
\Delta X_{r,t} = & \Delta X_{r,t-\tau} S_{r,t-\tau} \sum_{k=1}^K \beta_{within}^k W_{r,t-\tau}^k \\
& + S_{r,t-\tau} \sum_{c \in R \setminus r} \Delta X_{c,t-\tau} \sum_{k=1}^{\tilde{K}} \beta_{between}^k \widetilde{W}_{r,c,t-\tau}^k \\
& + \delta M_{r,t} + \eta_{r,t}
\end{aligned}$$

TODO: Results still need to be inserted.

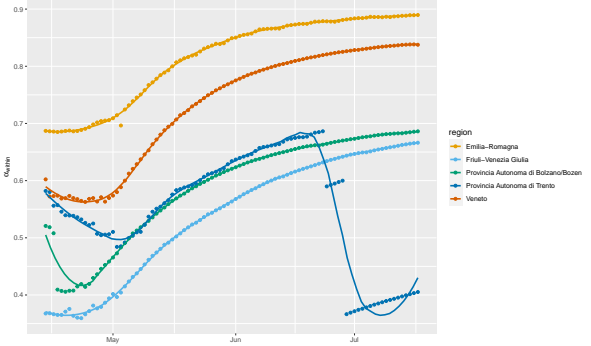
C Figures

C.1 Plots of β_{within} over time

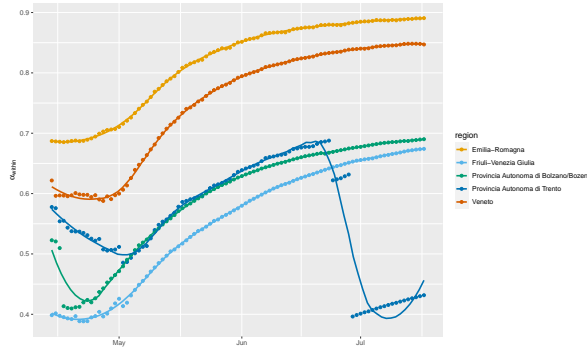
In Section 5.1 we presented the plots of β_{within} over time for the *Nord-Ovest* (North-West) NUTS 1 region for Model 1: within-region spread. In this section, we present the plots for the other NUTS 1 regions.



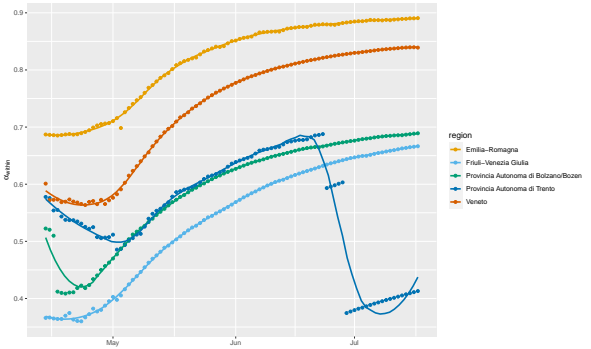
(a) Without model selection



(b) With model selection by AIC

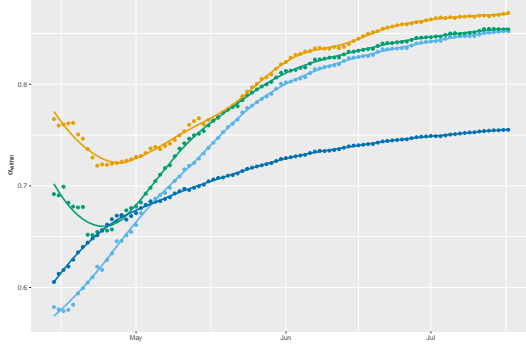


(c) Without model selection; with modelling undocumented infections

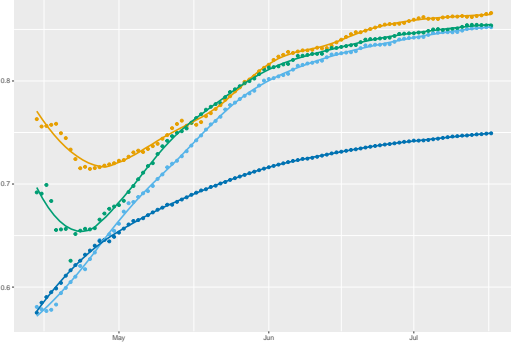


(d) With model selection by AIC; with modelling undocumented infections

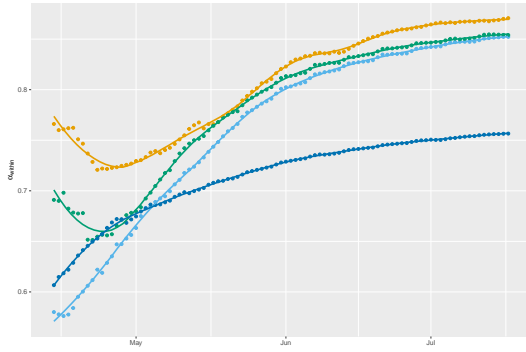
Figure C.1. Progression of β_{within} over time for the *Nord-Est* (North-East) NUTS 1 region



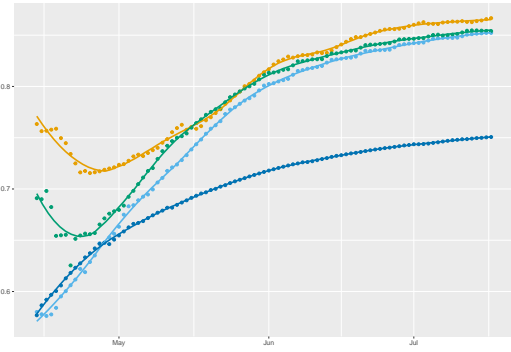
(a) Without model selection



(b) With model selection by AIC

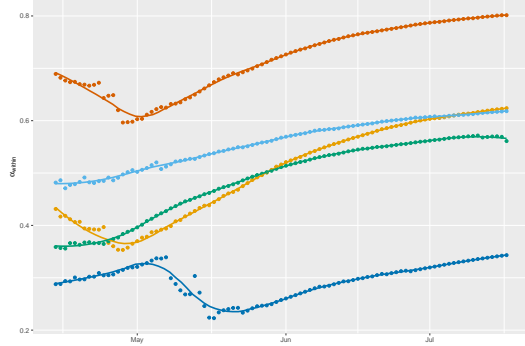


(c) Without model selection; with modelling undocumented infections

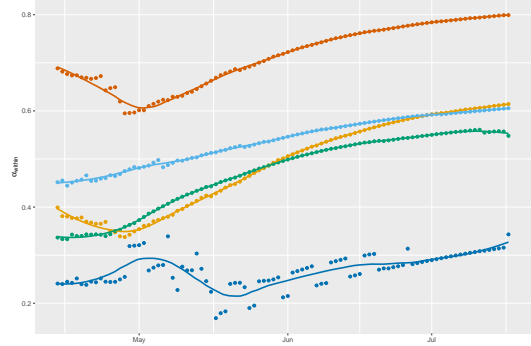


(d) With model selection by AIC; with modelling undocumented infections

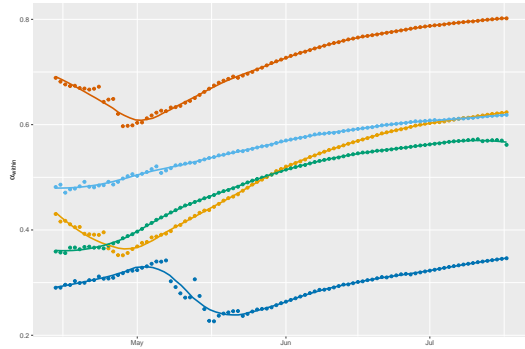
Figure C.2. Progression of β_{within} over time for the *Centro (IT)* (Centre) NUTS 1 region



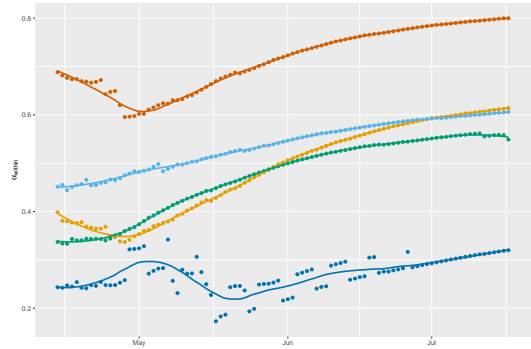
(a) Without model selection



(b) With model selection by AIC

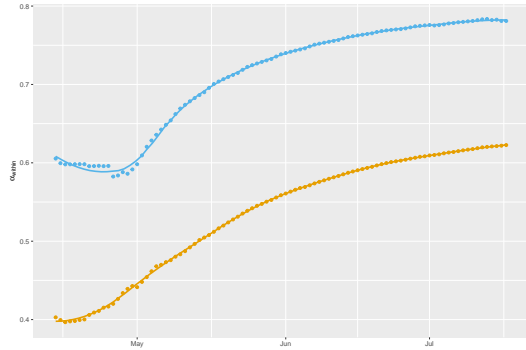


(c) Without model selection; with modelling undocumented infections

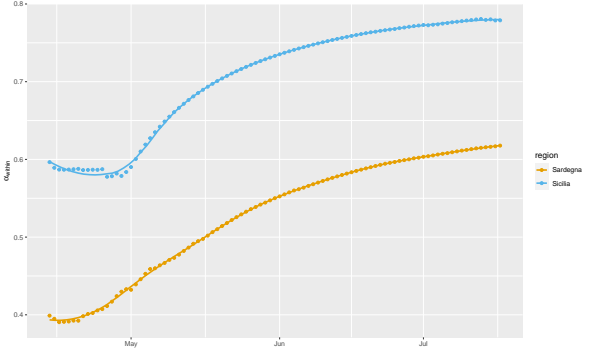


(d) With model selection by AIC; with modelling undocumented infections

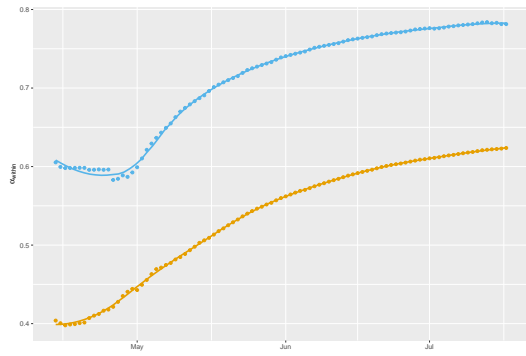
Figure C.3. Progression of β_{within} over time for the *Sud* (South) NUTS 1 region



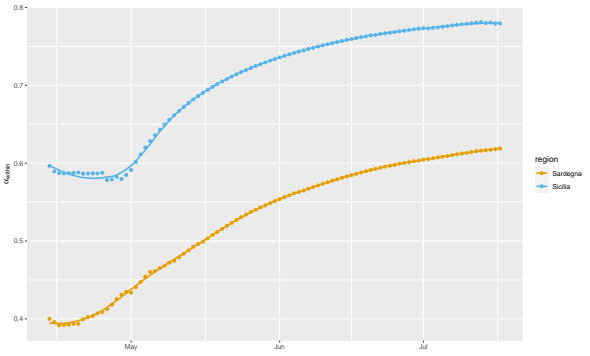
(a) Without model selection



(b) With model selection by AIC



(c) Without model selection; with modelling undocumented infections



(d) With model selection by AIC; with modelling undocumented infections

Figure C.4. Progression of β_{within} over time for the *Isole* (Islands) NUTS 1 region

D Derivations

D.1 Calculation of population variables

In this appendix, we will explain how the susceptible population and total population are calculated. Unfortunately, we do not have data on the total population per day. For this reason, we retrieve the latest population numbers per region from Eurostat (2020b), which are from January 1, 2019, and the yearly population growth rates for 2019 and 2020 from Worldometer (2020). For 2019, growth rate was equal to -0.13% and for 2020, excluding the deaths due to the pandemic, it was estimated to be equal to -0.15%. We only have the population growth rates available for the whole of Italy, not per region, unfortunately. As such, we assume that the growth rates are uniformly applicable to all regions. Of course, this is likely to introduce a small error since these growth rates differ over the regions. We assume that this error is negligible.

We denote the population of region r at time t by $N_{r,t}$. We denote the yearly population growth rates for 2019 and 2020 by g_{2019} and g_{2020} , respectively. Lastly, recall that the data for the pandemic starts at February 25, 2020. This is the 54th day of 2020, a leap year. As such, the population of region r on February 25, 2020 is calculated as:

$$N_{r,2020-02-25} = (1 + g_{2019})(1 + g_{2020})^{\frac{54}{366}} N_{r,2019-01-01} - D_{r,2020-02-25} \quad (\text{D.1})$$

where $D_{r,t}$ denotes the number of deaths in region r at time t .

Recall that the data reported at time t is reported with respect to the last 24 hours. As such, the susceptible population at time t can be calculated with the data at that same time. The susceptible population of region r at time t , denoted by $X_{r,t}$, is therefore calculated as follows:

$$X_{r,t} = N_{r,t} - Y_{r,t} - Z_{r,t} \quad (\text{D.2})$$

where $Y_{r,t}$ denotes the number of infectives and $Z_{r,t}$ denotes the number of removed individuals. Recall that Z is made up by adding the recovered individuals $R_{r,t}$ and the deceased individuals $D_{r,t}$. Because we use the calculation of $N_{r,t}$ as in the previous paragraph, the error discussed propagates into the calculation of $X_{r,t}$. However, as before, we assume that this error is negligible.

D.2 Functional forms for modelling undocumented infections

In this appendix, we give the derivations for the functional forms for modelling undocumented infections as discussed in Section 3.7.

D.2.1 Linear function

For modelling the undocumented infections, we want to construct a formula for a linear function that obeys the following assumptions:

- (I) $f(TC_t) = aTC_t + b$ for some $a, b \in \mathbb{R}$,
- (II) $f(0) = f^{min}$ for some $f^{min} \in [0, 1]$,
- (III) $f(N_t) = 1$

From assumption (II), we obtain that $b = f^{min}$. From assumption (III), we can then derive the value of a . The equation that we need to solve is:

$$aN_t + f^{min} = 1.$$

This is readily solved as $a = \frac{1-f^{min}}{N_t}$. As such, we have derived that

$$f(TC_t) = \frac{1 - f^{min}}{N_t} TC_t + f^{min}.$$

D.2.2 General quadratic function

For modelling the undocumented infections, we want to construct a general formula for a quadratic function that obeys the following assumptions:

- (I) $f(TC_t) = aTC_t^2 + bTC_t + c$ for some $a, b, c \in \mathbb{R}$,
- (II) $f(0) = f^{min}$ for some $f^{min} \in [0, 1]$,
- (III) $f(N_t) = 1$,
- (IV) $f(\beta N_t) = \gamma$ for some $\beta, \gamma \in (0, 1)$,
- (V) The vertex of the parabola should be to the right of N_t in the case of a downwards opening parabola and to the left of the origin in the case of an upwards opening parabola.

From assumption (II), we obtain that $c = f^{min}$. From assumptions (III) and (IV), we can then derive the values of a and b in terms of β , γ and N_t . The set of equations that we need to solve are:

$$\begin{cases} aN_t^2 + bN_t + f^{min} &= 1 \text{ (from assumption (III))} \\ a\beta^2 N_t^2 + b\beta N_t + f^{min} &= \gamma \text{ (from assumption (IV))} \end{cases} \quad (\text{D.3})$$

To solve (D.3), we can apply row reduction as follows:

$$\begin{aligned}
\left(\begin{array}{cc|c} N_t^2 & N_t & 1 - f^{min} \\ \beta^2 N_t^2 & \beta N_t & \gamma - f^{min} \end{array} \right) &\xrightarrow{r_2 - \beta^2 r_1} \left(\begin{array}{cc|c} N_t^2 & N_t & 1 - f^{min} \\ 0 & \beta(1 - \beta)N_t & \gamma - f^{min} - \beta^2 + \beta^2 f^{min} \end{array} \right) \\
&\xrightarrow{r_2 \div \beta(1 - \beta)} \left(\begin{array}{cc|c} N_t^2 & N_t & 1 - f^{min} \\ 0 & N_t & \frac{\gamma - f^{min} - \beta^2 + \beta^2 f^{min}}{\beta(1 - \beta)} \end{array} \right) \\
&\xrightarrow{r_1 - r_2} \left(\begin{array}{cc|c} N_t^2 & 0 & \frac{\beta - \gamma + (1 - \beta)f^{min}}{\beta(1 - \beta)} \\ 0 & N_t & \frac{\gamma - f^{min} - \beta^2 + \beta^2 f^{min}}{\beta(1 - \beta)} \end{array} \right) \\
&\xrightarrow{r_1 \div N_t^2} \left(\begin{array}{cc|c} 1 & 0 & \frac{\beta - \gamma + (1 - \beta)f^{min}}{\beta(1 - \beta)N_t^2} \\ 0 & 1 & \frac{\gamma - f^{min} - \beta^2 + \beta^2 f^{min}}{\beta(1 - \beta)N_t} \end{array} \right) \\
&\xrightarrow{r_2 \div N_t} \left(\begin{array}{cc|c} 1 & 0 & \frac{\beta - \gamma + (1 - \beta)f^{min}}{\beta(1 - \beta)N_t^2} \\ 0 & 1 & \frac{\gamma - f^{min} - \beta^2 + \beta^2 f^{min}}{\beta(1 - \beta)N_t} \end{array} \right)
\end{aligned}$$

As such, we have derived that

$$\begin{cases} a &= \frac{\beta - \gamma + (1 - \beta)f^{min}}{\beta(1 - \beta)N_t^2} \\ b &= \frac{\gamma - f^{min} - \beta^2 + \beta^2 f^{min}}{\beta(1 - \beta)N_t} \\ c &= f^{min}. \end{cases} \quad (D.4)$$

Firstly, note that this function is an upwards opening parabola if $a > 0$ and a downwards opening parabola if $a < 0$. For instance, we have that:

$$\begin{aligned}
a &> 0 \\
&\iff \frac{\beta - \gamma + (1 - \beta)f^{min}}{\beta(1 - \beta)N_t^2} > 0 \\
&\iff \beta - \gamma + (1 - \beta)f^{min} > 0 \\
&\iff \gamma < \beta + (1 - \beta)f^{min}
\end{aligned}$$

where we use that $\beta(1 - \beta)N_t^2 > 0$. Similarly, we have that $a < 0$ if $\gamma > \beta + (1 - \beta)f^{min}$.

Now note that our function is continuous. As such, we assume without loss of generality that $\beta = \frac{1}{2}$ and do the following derivations to deduce the values of γ for which assumption (V) holds. That is, we want to find the values of γ for which

$$f'(TC_t) = 0 \iff \begin{cases} TC_t \geq N_t \text{ for } \gamma > \frac{1}{2} + \frac{1}{2}f^{min} \\ TC_t \leq 0 \text{ for } \gamma < \frac{1}{2} + \frac{1}{2}f^{min}. \end{cases}$$

Firstly, assuming $\beta = \frac{1}{2}$, the expressions for a and b as in (D.4) reduce to:

$$\begin{cases} a &= \frac{\frac{1}{2} - \gamma + \frac{1}{2}f^{min}}{\frac{1}{4}N_t^2} \\ &= \frac{2 - 4\gamma + 2f^{min}}{N_t^2} \\ b &= \frac{\gamma - f^{min} - (\frac{1}{2})^2 + (\frac{1}{2})^2 f^{min}}{\frac{1}{4}N_t} \\ &= \frac{4\gamma - 1 - 3f^{min}}{N_t}. \end{cases} \quad (D.5)$$

We now need to derive the values of γ such that assumption (V) holds. That is:

$$\begin{aligned}
& f'(TC_t) = 0 \\
& \iff \frac{\partial aTC_t^2 + bTC_t + c}{\partial TC_t} = 0 \\
& \iff 2aTC_t + b = 0 \\
& \iff TC_t = -\frac{b}{2a}.
\end{aligned}$$

Using (D.5), we can fill out a and b to obtain:

$$TC_t = \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}}N_t.$$

Let $\gamma > \frac{1}{2} + \frac{1}{2}f^{min}$. Then, we need to derive γ such that

$$\begin{aligned}
& \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}}N_t \geq N_t \\
& \iff \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}} \geq 1.
\end{aligned}$$

Note that this is only the case if two conditions are satisfied:

$$\begin{cases} \text{sign}(1 - 4\gamma + 3f^{min}) &= \text{sign}(4 - 8\gamma + 4f^{min}) \end{cases} \quad (\text{D.6a})$$

$$\begin{cases} |1 - 4\gamma + 3f^{min}| &\geq |4 - 8\gamma + 4f^{min}| \end{cases} \quad (\text{D.6b})$$

Note that our assumption that $\gamma > \frac{1}{2} + \frac{1}{2}f^{min}$ is equivalent to $\gamma > \frac{2+2f^{min}}{4}$ which, in turn, is equivalent to $4 - 8\gamma + 4f^{min} < 0$. As such, (D.6a) tells us that both the numerator and denominator of the fraction are negative. Therefore, to satisfy (D.6a), we need that

$$\begin{aligned}
& 1 - 4\gamma + 3f^{min} < 0 \\
& \iff \gamma > \frac{1 + 3f^{min}}{4}
\end{aligned}$$

Since we assumed that $\gamma > 2 + 2f^{min}$, this is always satisfied because $f^{min} \in [0, 1]$ so that $1 + 3f^{min} < 2 + 2f^{min} < \gamma$. That brings us to the second condition (D.6b). Because we know that both parts of the fractions are negative, we can now solve for γ as follows:

$$\begin{aligned}
& \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}}N_t \geq N_t \\
& \iff 1 - 4\gamma + 3f^{min} \leq 4 - 8\gamma + 4f^{min} \\
& \iff \gamma \leq \frac{3 + f^{min}}{4} = \frac{3}{4} + \frac{1}{4}f^{min}.
\end{aligned}$$

Let $\gamma < \frac{1}{2} + \frac{1}{2}f^{min}$. Then, we need to derive γ such that

$$\begin{aligned} \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}}N_t &\leq 0 \\ \iff \frac{1 - 4\gamma + 3f^{min}}{4 - 8\gamma + 4f^{min}} &\leq 0. \end{aligned}$$

Note that this is only the case if one of the following two conditions is satisfied:

$$\begin{cases} 1 - 4\gamma + 3f^{min} \leq 0 & \text{and } 4 - 8\gamma + 4f^{min} > 0 \\ 1 - 4\gamma + 3f^{min} \geq 0 & \text{and } 4 - 8\gamma + 4f^{min} < 0 \end{cases} \quad \begin{matrix} \text{(D.7a)} \\ \text{(D.7b)} \end{matrix}$$

As before, note that our assumption that $\gamma > \frac{1}{2} + \frac{1}{2}f^{min}$ is equivalent to $4 - 8\gamma + 4f^{min} > 0$. As such, we know that the only condition that can be satisfied is (D.7a). Therefore, we need that

$$\begin{aligned} 1 - 4\gamma + 3f^{min} &\leq 0 \\ \gamma &\geq \frac{1 + 3f^{min}}{4} = \frac{1}{4} + \frac{3}{4}f^{min}. \end{aligned}$$

As such, we should have that $\gamma \in [\frac{1}{4} + \frac{3}{4}f^{min}, \frac{3}{4} + \frac{1}{4}f^{min}]$. When $\gamma \in [\frac{1}{4} + \frac{3}{4}f^{min}, \frac{1}{2} + \frac{1}{2}f^{min})$, the parabola we receive is upwards opening. On the other hand, when $\gamma \in (\frac{1}{2}, \frac{3}{4} + \frac{1}{4}f^{min}]$, the parabola we receive is downwards opening. When $\gamma = \frac{1}{2} + \frac{1}{2}f^{min}$, the function we receive is linear, since $a = \frac{2 - 4\gamma + 2f^{min}}{N_t^2} = 0$.

Conclusively, we have derived that

$$f(TC_t) = \frac{2 - 4\gamma + 2f^{min}}{N_t^2}TC_t^2 + \frac{4\gamma - 1 - 3f^{min}}{N_t}TC_t + f^{min},$$

under the assumption that $\beta = \frac{1}{2}$, with $\gamma \in [\frac{1}{4} + \frac{3}{4}f^{min}, \frac{3}{4} + \frac{1}{4}f^{min}]$.

D.2.3 Special case quadratic formula: downwards opening

For modelling the undocumented infections, we want to construct a formula for a downwards opening quadratic function that obeys the following assumptions:

- (I) $f(x) = ax^2 + bx + c$ for some $a, b, c \in \mathbb{R}$,
- (II) $f(0) = f^{min}$ for some $f^{min} \in [0, 1]$,
- (III) $f(N_t) = 1$,
- (IV) $f'(N_t) = 0$, i.e. the vertex of the parabola is found at $TC_t = N_t$.

Consider that any quadratic formula can be written as $f(TC_t) = a(TC_t - h)^2 + k$, which is called the vertex form, where the vertex (i.e. the extremum) of the function is (h, k) . By assumptions (III) and (IV), $h = N_t$ and $k = 1$. Therefore,

$$f(TC_t) = a(TC_t - N_t)^2 + 1.$$

Using assumption (II), we can solve this equation for a :

$$\begin{aligned} a(0 - N_t)^2 + 1 &= f^{min} \\ \iff aN_t^2 &= f^{min} - 1 \\ \iff a &= \frac{f^{min} - 1}{N_t^2} \end{aligned}$$

Therefore, the formula becomes:

$$\begin{aligned} f(TC_t) &= \frac{f^{min} - 1}{N_t^2} (TC_t - N_t)^2 + 1 \\ &= \frac{f^{min} - 1}{N_t^2} (TC_t^2 + N_t^2 - 2N_t TC_t) + 1 \\ &= \frac{(f^{min} - 1)(TC_t^2 + N_t^2 - 2N_t TC_t) + N_t^2}{N_t^2} \\ &= \frac{f^{min} - 1}{N_t^2} TC_t^2 - \frac{2(f^{min} - 1)}{N_t} TC_t + f^{min}. \end{aligned}$$

D.2.4 Special case quadratic formula: upwards opening

For modelling the undocumented infections, we want to construct a formula for an upwards opening quadratic function that obeys the following assumptions:

- (I) $f(x) = ax^2 + bx + c$ for some $a, b, c \in \mathbb{R}$,
- (II) $f(0) = f^{min}$ for some $f^{min} \in [0, 1]$,
- (III) $f(N_t) = 1$,
- (IV) $f'(0) = 0$, i.e. the vertex of the parabola is found at $TC_t = 0$.

Just as in appendix D.2.4, we use the vertex form $f(TC_t) = a(TC_t - h)^2 + k$. By assumptions (III) and (IV), $h = 0$ and $k = f^{min}$. Therefore,

$$f(TC_t) = a(TC_t - 0)^2 + f^{min} = aTC_t^2 + f^{min}.$$

Using assumption (II), we can solve this equation for a :

$$\begin{aligned} aN_t^2 + f^{min} &= 1 \\ \iff a &= \frac{1 - f^{min}}{N_t^2} \end{aligned}$$

Therefore, the formula becomes:

$$f(TC_t) = \frac{1 - f^{min}}{N_t^2} TC_t^2 + f^{min},$$

which is already in the form as in assumption (I).

D.2.5 Cubic function

For modelling the undocumented infections, we want to construct a general formula for a cubic function that obeys the following assumptions:

- (I) $f(x) = ax^3 + bx^2 + cx + d$ for some $a, b, c, d \in \mathbb{R}$,
- (II) $f(0) = f^{min}$ for some $f^{min} \in [0, 1]$,
- (III) $f(N_t) = 1$,
- (IV) $f(\beta_1 N_t) = \gamma_1$ and $f(\beta_2 N_t) = \gamma_2$ for some $\beta_1, \beta_2, \gamma_1, \gamma_2 \in [0, 1]$ and $\beta_1 < \beta_2, \gamma_1 < \gamma_2$.

From assumption (II), we obtain that $d = f^{min}$. From assumptions (III) and (IV), we can then derive the values of a , b , and c in terms of the β s, γ s, and N_t . The set of equations that we need to solve are:

$$\begin{cases} aN_t^3 + bN_t^2 + cN_t + f^{min} &= 1 \text{ (from assumption (III))} \\ a\beta_1^3 N_t^3 + b\beta_1^2 N_t^2 + c\beta_1 N_t + f^{min} &= \gamma_1 \text{ (from assumption (IV))} \\ a\beta_2^3 N_t^3 + b\beta_2^2 N_t^2 + c\beta_2 N_t + f^{min} &= \gamma_2 \text{ (from assumption (IV))} \end{cases} \quad (\text{D.8})$$

In appendix D.2.2, we first solved these equations and then assumed a value for β afterwards, without loss of generality. In this case, the equations would become immensely populated if we were to keep the derivation general. As such, we first assume without loss of generality that $\beta_1 = \frac{1}{4}$ and $\beta_2 = \frac{1}{2}$. To solve (D.8), we can then apply row reduction as follows:

$$\begin{aligned}
\left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ \beta_1^3 N_t^3 & \beta_1^2 N_t^2 & \beta_1 N_t & \gamma_1 - f^{min} \\ \beta_2^3 N_t^3 & \beta_2^2 N_t^2 & \beta_2 N_t & \gamma_2 - f^{min} \end{array} \right) &= \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ \frac{1}{64} N_t^3 & \frac{1}{16} N_t^2 & \frac{1}{4} N_t & \gamma_1 - f^{min} \\ \frac{1}{8} N_t^3 & \frac{1}{4} N_t^2 & \frac{1}{2} N_t & \gamma_2 - f^{min} \end{array} \right) \\
&\xrightarrow{r_2 \times 64} \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ N_t^3 & 4N_t^2 & 16N_t & 64\gamma_1 - 64f^{min} \\ N_t^3 & 2N_t^2 & 4N_t & 16\gamma_2 - 64f^{min} \end{array} \right) \\
&\xrightarrow{r_3 \times 8} \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ 0 & 3N_t^2 & 15N_t & -1 + 64\gamma_1 - 63f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 - 7f^{min} \end{array} \right) \\
&\xrightarrow{r_2 - r_1} \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ 0 & 3N_t^2 & 15N_t & -1 + 64\gamma_1 - 63f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 - 7f^{min} \end{array} \right) \\
&\xrightarrow{r_3 - r_1} \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 - 7f^{min} \\ 0 & 3N_t^2 & 15N_t & -1 + 64\gamma_1 - 63f^{min} \end{array} \right) \\
&\xrightarrow{r_2 \leftrightarrow r_3} \left(\begin{array}{ccc|c} N_t^3 & N_t^2 & N_t & 1 - f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 - 7f^{min} \\ 0 & 3N_t^2 & 15N_t & -1 + 64\gamma_1 - 63f^{min} \end{array} \right) \\
&\xrightarrow{r_1 - r_2} \left(\begin{array}{ccc|c} N_t^3 & 0 & -2N_t & 2 - 8\gamma_2 + 6f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 \\ 0 & 0 & 6N_t & 2 + 64\gamma_1 - 24\gamma_2 - 42f^{min} \end{array} \right) \\
&\xrightarrow{r_3 - 3r_2} \left(\begin{array}{ccc|c} N_t^3 & 0 & -2N_t & 2 - 8\gamma_2 + 6f^{min} \\ 0 & N_t^2 & 3N_t & -1 + 8\gamma_2 \\ 0 & 0 & 6N_t & 2 + 64\gamma_1 - 24\gamma_2 - 42f^{min} \end{array} \right) \\
&\xrightarrow{r_1 + \frac{1}{3}r_3} \left(\begin{array}{ccc|c} N_t^3 & 0 & 0 & \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3} \\ 0 & N_t^2 & 0 & -2 - 32\gamma_1 + 20\gamma_2 + 14f^{min} \\ 0 & 0 & 6N_t & 2 + 64\gamma_1 - 24\gamma_2 - 42f^{min} \end{array} \right) \\
&\xrightarrow{r_2 - \frac{1}{2}r_3} \left(\begin{array}{ccc|c} N_t^3 & 0 & 0 & \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3} \\ 0 & N_t^2 & 0 & -2 - 32\gamma_1 + 20\gamma_2 + 14f^{min} \\ 0 & 0 & 6N_t & 2 + 64\gamma_1 - 24\gamma_2 - 42f^{min} \end{array} \right) \\
&\xrightarrow{r_1 \div N_t^3} \left(\begin{array}{ccc|c} 1 & 0 & 0 & \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3N_t^3} \\ 0 & 1 & 0 & \frac{-2-32\gamma_1+20\gamma_2+14f^{min}}{N_t^2} \\ 0 & 0 & 1 & \frac{2+64\gamma_1-24\gamma_2-42f^{min}}{6N_t} \end{array} \right) \\
&\xrightarrow{r_2 \div N_t^2} \left(\begin{array}{ccc|c} 1 & 0 & 0 & \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3N_t^3} \\ 0 & 1 & 0 & \frac{-2-32\gamma_1+20\gamma_2+14f^{min}}{N_t^2} \\ 0 & 0 & 1 & \frac{2+64\gamma_1-24\gamma_2-42f^{min}}{6N_t} \end{array} \right) \\
&\xrightarrow{r_3 \div 6N_t} \left(\begin{array}{ccc|c} 1 & 0 & 0 & \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3N_t^3} \\ 0 & 1 & 0 & \frac{-2-32\gamma_1+20\gamma_2+14f^{min}}{N_t^2} \\ 0 & 0 & 1 & \frac{2+64\gamma_1-24\gamma_2-42f^{min}}{6N_t} \end{array} \right)
\end{aligned}$$

Conclusively, we have derived that

$$\begin{cases} a = \frac{8+64\gamma_1-48\gamma_2-24f^{min}}{3N_t^3} \\ b = \frac{-2-32\gamma_1+20\gamma_2+14f^{min}}{N_t^2} \\ c = \frac{2+64\gamma_1-24\gamma_2-42f^{min}}{6N_t} = \frac{1+32\gamma_1-12\gamma_2-21f^{min}}{3N_t} \\ d = f^{min} \end{cases} \quad (D.9)$$

so that

$$\begin{aligned}
f(TC_t) &= \frac{8 + 64\gamma_1 - 48\gamma_2 - 24f^{min}}{3N_t^3} TC_t^3 + \frac{-2 - 32\gamma_1 + 20\gamma_2 + 14f^{min}}{N_t^2} TC_t^2 \\
&\quad + \frac{1 + 32\gamma_1 - 12\gamma_2 - 21f^{min}}{3N_t} TC_t + f^{min},
\end{aligned}$$

under the assumption that $\beta_1 = \frac{1}{4}$ and $\beta_2 = \frac{1}{2}$.