Comparing the accuracy of several network-based COVID-19 prediction algorithms

Massimo A. Achterberg^{a,*}, Bastian Prasse^a, Long Ma^a, Stojan Trajanovski^b, Maksim Kitsak^a, Piet Van Mieghem^a

^a Faculty of Electrical Engineering, Mathematics and Computer Science, Delft University of Technology, P.O. Box 5031, 2600 GA Delft, The Netherlands ^b Microsoft Inc., 2 Kingdom St, London W2 6BD, United Kingdom

Abstract

Researchers from various scientific disciplines have attempted to forecast the spread of the Coronavirus Disease 2019 (COVID-19). The proposed epidemic prediction methods range from basic curve fitting methods and traffic interaction models to machine-learning approaches. If we combine all these approaches, we obtain the Network Inference-based Prediction Algorithm (NIPA). In this paper, we analyse a diverse set of COVID-19 forecast algorithms, including several modifications of NIPA. Among the diverse set of algorithms that we evaluated, original NIPA performs best on forecasting the spread of COVID-19 in Hubei, China and in the Netherlands. In particular, we show that network-based forecasting is superior to any other forecasting algorithm.

Keywords: Epidemiology, Network inference, Forecast accuracy, Bayesian methods, SIR model, Time series methods, Machine learning methods

1 1. Introduction

In December, 2019, the SARS-CoV-2 virus, which causes the Coronavirus Disease 2019 (COVID-19), emerged in the Chinese province Hubei. The number of COVID-19 infected cases mainly in China rose dramatically to almost 80,000 the end of February. From China, COVID-19 quickly spread throughout the

Preprint submitted to International Journal of Forecasting

^{*}Corresponding author Email address: m.a.achterberg@tudelft.nl (Massimo A. Achterberg)

whole world, with almost ten million cases at the end of June, 2020. Many
countries imposed a nation-wide lockdown to slow down the spread of COVID19. A reliable forecast of the pandemic outbreak is key for targeted disease
countermeasures and for the appropriate design of an exit strategy to lift the
lockdown.

Unfortunately, just as weather forecasts, the prediction of epidemic out-11 breaks is subject to fundamental limits [1]. One aspect is the limited avail-12 ability of data, because epidemic time series are relatively short and carrying 13 out medical tests on a large scale is challenging. Also, the final number of in-14 fected cases is highly sensitive to initial perturbations [2]. Nonetheless, many 15 methods have been developed and applied to forecast the spread of COVID-19. 16 Perhaps the simplest approach is based on fitting the number of infections to a 17 sigmoid curve, such as the logistic function [3, 4], Hill function [5] or Gompertz 18 function [6]. Using nonlinear regression, the parameters of the sigmoid curve 19 can be estimated. For the comparison of prediction algorithms in this work, we 20 evaluate the prediction based on the logistic function. The logistic function is of 21 particular interest, because the logistic function is the (approximate) solution 22 for the number of infected cases [7] in the Susceptible-Infected-Susceptible (SIS) 23 epidemic model and the number of removed cases in the Susceptible-Infected-24 Removed (SIR) epidemic model [2, 8]. 25

By fitting the number of infected cases to a sigmoid curve, we implicitly assume that the spread in a particular region is independent of other regions, which contrasts the strong interconnectedness of our modern world. Networkbased techniques take into account the interaction between different regions, which is due to the movement of people.

The interaction can be described by network G with N nodes. Each node iin the network G represents a particular region (country, province, municipality or city), and the link $a_{ij} \in \{0, 1\}$ represents the existence of an interaction from region j to region i, specified by a link weight β_{ij} denoting the infection probability from region j to region i. The self-infection probability within a region iis given by β_{ii} , which we expect to be dominant over the other infection probabilities, because the interaction within a region is stronger than the interaction with other regions. The $N \times N$ infection probability matrix B, with elements β_{ij} is, however, unknown and must be derived from past observations of the epidemic. We will address this issue in more detail in Section 2.

Throughout this work, we often use "the number of infected cases", which 41 we understand as "the number of cases reported by local authorities". The 42 asymptomatic individuals, who do not feel sick and even do not know that they 43 are infected and infectious, are not reported and can infect others unnoticed. To 44 gain understanding of the percentage of asymptomatic cases, a possibility is to 45 test the population at random with, for example, a blood test. For COVID-19, 46 the fraction of asymptomatic cases is estimated to be as large as 80% [9]. Since 47 the number of asymptomatic cases cannot be determined on a daily basis, we 48 confine ourselves to the number of reported cases in this work. 49

Many scientific disciplines investigate and forecast the spread of COVID-50 19. Statistical approaches are commonly based on Kalman filtering [10] or 51 consider Bayesian approaches [11]. Network-based approaches consider aero-52 plane networks, daily commute traffic or cell phone traffic [12]. Data scien-53 tists apply machine learning algorithms, like adaptive neuro-fuzzy inference sys-54 tem [13] or Long Short-Term Memory (LSTM) [14]. Mathematicians perform 55 parameter estimation on compartmental models like the Susceptible-Infected-56 Removed model (SIR) [14, 15] or the Susceptible-Exposed-Infected-Removed 57 (SEIR) model [16]. 58

Most epidemic models forecast the number of infected cases as a point forecast (generally: the mean of a distribution) rather than a complete distribution. All models in this work have been designed to provide point forecasts, but can be generalised to provide prediction intervals. We discuss this topic further in Section 2.

The focus of this work is the comparison of a diverse set of methods to forecast the spread of COVID-19, ranging from fitting closed-form epidemic curves and comprehensive machine-learning algorithms to network-based approaches. We focus on the spread of COVID-19, but we emphasise that all methods can be

applied to general epidemic outbreaks. We show that pure machine-learning and 68 network-agnostic algorithms or epidemiological models are inferior to algorithms 69 which combine multiple approaches that rely on the underlying network topol-70 ogy. In particular, the Network Inference-based Prediction Algorithm (NIPA) 71 is superior to any other algorithm that we have evaluated. In Section 2 we ex-72 plain eight forecast algorithms to predict the future number of COVID-19 cases. 73 Thereafter, we show their accuracy in two selected regions: Hubei (China) and 74 the Netherlands in Section 3 and discuss the strengths and weaknesses of each 75 algorithm. Finally, we conclude in Section 4. 76

77 2. Prediction algorithms

The spread of COVID-19 can be measured in terms of the daily number of 78 reported cases. We model the course of the epidemic by an SIR compartmental 79 model, where each individual is either Susceptible (healthy), Infected (can in-80 fected the susceptible) or Removed (recovered or died). We denote the (discrete) 81 time by k = 1, ..., n where n is the total number of observation days. The first 82 COVID-19 case was reported on day k = 1. Given that nearly all governments 83 report their epidemic data once a day, we take a time step of 1 day as a natural 84 choice and investigate the effect of the time step on the prediction accuracy 85 in Appendix G. The Susceptible-Infected-Removed (SIR) epidemic model with 86 time-varying spreading parameters is given by: 87

Definition 1 (SIR Epidemic Model [8, 17, 18]). The viral state $v_i[k] = (S_i[k], \mathcal{I}_i[k], \mathcal{R}_i[k])^T$ of every region *i* evolves in discrete time k = 1, 2, ..., n according to

$$\mathcal{I}_{i}[k+1] = (1-\delta_{i})\mathcal{I}_{i}[k] + (1-\mathcal{I}_{i}[k] - \mathcal{R}_{i}[k])\sum_{j=1}^{N}\beta_{ij}[k]\mathcal{I}_{j}[k]$$
(1)

$$\mathcal{R}_i[k+1] = \mathcal{R}_i[k] + \delta_i \mathcal{I}_i[k], \qquad (2)$$

and the fraction of susceptible individuals follows as

$$\mathcal{S}_i[k] = 1 - \mathcal{I}_i[k] - \mathcal{R}_i[k]. \tag{3}$$

Here, $\beta_{ij}[k] \ge 0$ denotes the infection probability from region j to region i at time k, and $\delta_i > 0$ denotes the curing probability of region i.

The spread of COVID-19 cannot be described exactly by the SIR equations 90 (1). The COVID-19 pandemic evolves in continuous time, whereas the SIR 91 model (1) evolves in discrete time, with a time step of 1 day. Additionally, the 92 SIR model (1) is unable to describe phenomena like personal social distanc-93 ing, nation-wide lockdowns and the availability of vaccinations. Each of these 94 model assumptions introduces model errors. Prior to the introduction of sev-95 eral forecasting algorithms, we explain how model errors can be used to obtain 96 prediction intervals for the forecasted number of infected cases. 97

As described in [19], we obtain the fraction of susceptible $S_i[k]$, infectious $\mathcal{I}_i[k]$ and removed $\mathcal{R}_i[k]$ individuals in every region *i* from the observed infections $y_i[k]$. We aim to find the best possible forecast $\hat{y}_i[k]$ for the cumulative number of infected cases $y_i[k]$ for every region *i* and time *k*. In this work, we discuss eight prediction methods.

¹⁰³ 2.1. Potential generalisation to prediction intervals

Before introducing the different prediction methods, we emphasise that this work focusses on *short-term* point forecasts. The long-term epidemic behaviour is very random, and providing forecast intervals is essential to give a complete picture of the long-term viral spread [20]. Extending the point forecast methods in this work to prediction intervals is outside the scope of this work. Nonetheless, we consider it valuable to conceptually discuss an extension of the SIR equations (1) to allow for the computation of prediction intervals. Any real epidemic does not follow the SIR model (1) exactly. Instead, the infection state $\mathcal{I}_i[k]$ evolves from time k to k + 1 as

$$\mathcal{I}_{i}[k+1] = (1-\delta_{i})\mathcal{I}_{i}[k] + (1-\mathcal{I}_{i}[k] - \mathcal{R}_{i}[k])\sum_{j=1}^{N}\beta_{ij}[k]\mathcal{I}_{j}[k] + w_{i}[k], \quad (4)$$

where $w_i[k]$ denotes the *model error* of region *i* at time *k*, see also Appendix A. The equations (4) can be used as a basis for prediction intervals with a Monte ¹⁰⁶ Carlo approach. Define the $N \times 1$ error vector as $w[k] = (w_1[k], ..., w_N[k])^T$ ¹⁰⁷ and the $N \times 1$ infection vector as $\mathcal{I}[k] = (\mathcal{I}_1[k], ..., \mathcal{I}_N[k])^T$ for all times k. ¹⁰⁸ Then, based on equation (4), the past observations $\mathcal{I}[1], ..., \mathcal{I}[n]$ and the errors ¹⁰⁹ w[1], ..., w[n - 1], point forecast algorithms provide an estimate of the viral ¹¹⁰ state $\mathcal{I}[k]$ at future times k > n.

Conceptually, a prediction interval for the future viral state $\mathcal{I}_i[k]$ can be obtained by two steps. First, obtain random samples from the distribution of the model errors w[1], ..., w[n - 1]. Second, for every sample of the errors w[1], ..., w[n - 1], obtain a point forecast of the future viral states $\mathcal{I}[k]$. The prediction intervals for the future viral state $\mathcal{I}[k]$ can be obtained from the ensemble of point forecasts.

The details of the outlined method for obtaining prediction intervals are beyond the scope of this paper. Two particular challenges are the determination of the distribution of the model errors w[k] and the implementation of a computationally efficient sampling method.

121 2.2. Sigmoid curves

The logistic function is a well-known example of an epidemiological sigmoid curve [3, 7]. We assume the cumulative number of infected cases $y_i[k]$ in region *i* at time *k* to follow a logistic function:

$$y_i[k] = \frac{y_{\infty,i}}{1 + e^{-K_i(k - t_{0,i})}},\tag{5}$$

where $y_{\infty,i}$ is the long-term fraction of infections, K_i is the logistic growth rate and $t_{0,i}$ is the inflection point, also known as the epidemic peak. The parameters $y_{\infty,i}$, K_i and $t_{0,i}$ are estimated for each region separately using a nonlinear curve fitting procedure, which is explained in Appendix F. Other sigmoid curves, like the Hill function and Gompertz function, are also discussed in Appendix F.

127 2.3. LSTM

Recurrent neural networks [21] (RNNs) have been used in various tasks related to sequences [22], time series analysis and forecasting, speech recognition

or natural language processing [23] and it has been demonstrated they achieve 130 state-of-the-art performance. Long Short-term Memory (LSTM) networks [24] 131 are specific types of RNNs that resolved the long-standing problem in the past 132 for long-term dependencies caused by the difference in input growth which in 133 turns leads to vanishing or exploding gradients in neural networks backpropa-134 gation. LSTM introduces additional input, output and optional forget gates as 135 interfaces with additional weights on the top of standard input data and hidden 136 weights in the standard RNN unit. There are several variations [25, 26] for 137 the LSTM networks, just to mention few: with or without forget gate and a 138 "peephole connection"; that perform better in one or another task [27]. For the 139 internal mechanism between the gates and the exact mathematical relations, we 140 refer to [28] or [26]. In this work, we utilize the most common one - an LSTM 141 with a forget gate. In the simulations, we use an LSTM with sequence and hid-142 den sizes both equal to four in a single LSTM layer (e.g. it is possible to stack 143 few LSTM layers which leads to more overfitting), a learning rate of 0.1 and 144 Adam optimizer [29], with mean square error loss in 2000 epochs of training. 145

146 2.4. NIPA

Network-based approaches take into account the interactions between dif-147 ferent regions. However, the contact network G is unknown (and consequently 148 also the infection probability matrix B) and must be inferred from the epidemic 149 outbreak. The Network Inference-based Prediction Algorithm (NIPA) was orig-150 inally proposed in [18], and we applied an adaption of NIPA to the spread of 151 COVID-19 in Hubei, China [19] and Italy [30]. NIPA consists of two steps. 152 First, the underlying infection matrix B is inferred from the epidemic outbreak. 153 Second, the infection matrix B and the estimated curing rates δ_i for every node i 154 are used to forecast the outbreak by iterating the Susceptible-Infected-Removed 155 (SIR) model on the estimated infection matrix B. Even though NIPA success-156 fully forecasted the spread of COVID-19 in the Chinese province Hubei, the 157 underlying infection matrix B cannot be inferred [31]. 158

¹⁵⁹ 2.5. NIPA on each region separately

As a benchmark model, we apply NIPA on each region separately, which we name *NIPA separate*. NIPA separate is a machine-learning method based on the SIR model, but does not consider the interaction between different regions.

163 2.6. NIPA static prior

The formulation of NIPA can be extended to include knowledge on the underlying contact network. We use a time-independent traffic network (with the corresponding traffic intensity matrix M) to obtain a prior for the infection probability matrix B as

$$B_{\text{prior}} = \text{diag}\left(c_1, \dots, c_N\right) M. \tag{6}$$

We explain our motivation for the prior infection matrix B_{prior} in Appendix B. The positive scalars $c_1, ..., c_N$ are unknown and are set by cross-validation. We assume that the true infection matrix B is normally distributed around the prior infection matrix B_{prior} . Based on the prior infection matrix B_{prior} and the observations of the COVID-19 spread, we obtain the Bayesian estimate $B_{\text{posterior}}$ by solving the optimisation problem

$$B_{\text{posterior}} = \underset{B}{\operatorname{argmax}} \operatorname{Pr}\left[B|y[1], ..., y[n]\right]$$
(7)
s.t.
$$\sum_{j=1}^{N} \beta_{ij} \leq 1, \quad i = 1, ..., N,$$

where y[k] is the observed $N \times 1$ infection vector $y[k] = (y_1[k], ..., y_N[k])^T$ at all times k = 1, ..., n. Using the estimated infection matrix $B_{\text{posterior}}$ and the estimated curing rates δ_i for every region i, we forecast the outbreak by iterating the SIR model. For details on NIPA static prior, we refer to Appendix C.

168 2.7. NIPA dynamic prior

Many countries have imposed some kind of lockdown, in which the free movement of people is significantly restricted. Thus, the true contact network G

Table 1: All algorithms discussed in this paper. *If the algorithm is based on a phenomenological epidemic process, like the SIR model. **If the algorithm is able to forecast small perturbations in the global trend. ***If the spread between different regions is considered.

Algorithm	${\rm Epidemiology}^*$	$Adaptive^{**}$	Network***
NIPA	\checkmark	\checkmark	\checkmark
NIPA separate	\checkmark	\checkmark	×
NIPA static prior	\checkmark	\checkmark	\checkmark
NIPA dynamic prior	\checkmark	\checkmark	\checkmark
Logistic function	\checkmark	×	×
Hill function	\checkmark	×	×
Gompertz function	\checkmark	×	×
LSTM	×	\checkmark	×

varies over time. We use a time-varying traffic matrix M[k] as an approximation for the prior infection matrix $B_{\text{prior}}[k]$, whose entries equal

$$B_{\text{prior}}[k] = \text{diag}\left(c_1, ..., c_N\right) M[k] \tag{8}$$

for all times k. The positive scalars $c_1, ..., c_N$ are unknown and are set by holdout validation. We propose a Bayesian approach called *NIPA dynamic prior* to estimate the true infection matrix B[k] from the time series of infected cases $y_i[k]$ and the prior infection matrix $B_{\text{prior}}[k]$. Using the estimated time-varying infection matrix $B_{\text{posterior}}[k]$ and the curing rates δ_i for each region *i*, we forecast the outbreak by iterating the SIR model. Appendix D explains the technical details of NIPA dynamic prior.

A challenge to NIPA dynamic prior is the unavailability of the contact network in the future. Hence, we assume the traffic matrix to remain constant after the last observation point n: $B_{\text{prior}}[n+k] = B_{\text{prior}}[n]$ for all k > 0.

¹⁷⁹ 3. Evaluation of the prediction performance

We evaluate the prediction accuracy of the methods discussed in Section 2 by forecasting the spread of COVID-19 in a selected number of regions. We set the maximal forecast horizon to six days, because of the difficulty of predictingepidemic outbreaks [2].

Each prediction algorithm produces a forecast $\hat{y}_i[k]$ for the cumulative number of infected cases $y_i[k]$ for every region *i* at time *k*. To quantify the prediction error at time *k*, we use the Symmetric Mean Absolute Percentage Error (sMAPE)

$$e_{\text{sMAPE}}[k] = \frac{1}{N} \sum_{i=1}^{N} \frac{|y_i[k] - \hat{y}_i[k]|}{(y_i[k] + \hat{y}_i[k])/2},$$
(9)

which is commonly used in forecasting [32]. Furthermore, we quantify the Percentage Error (PE)

$$e_{\text{PE},i}[k] = \frac{y_i[k] - \hat{y}_i[k]}{y_i[k]},$$
(10)

for every region i and time k to investigate over- and underestimations. We consider the spread of COVID-19 in two regions: The cities in Hubei, China and the provinces in the Netherlands. These regions cannot be regarded as full representatives of the spread of COVID-19, let alone general infectious diseases. Rather, these regions illustrate the strengths and weaknesses of our methods.

189 3.1. Hubei, China

We evaluate the prediction accuracy first on the Chinese province Hubei. In 190 December 2019, the first cases of COVID-19 were detected in Wuhan, the capital 191 of Hubei. The first case outside Wuhan was reported on January 21. From 192 January 24 onwards, the whole province Hubei was under lockdown, prohibiting 193 any non-urgent travels. On February 15, the local government in Hubei changed 194 the diagnosing policy, causing an erratic increase in the number of reported cases 195 on February 15. Therefore we restrict ourselves to the period from January 21 196 to February 14. The reported cases are provided by the Health Commission of 197 Hubei [33]. The majority of COVID-19 patients were reported in Wuhan, as 198 shown in Figure 1. We have removed Shennongjia from our analysis, because 199 of the small number of infections in that region. 200

For NIPA static prior, we require a traffic network describing the interactions between the cities in Hubei. The Chinese company Baidu provides an estimate



Figure 1: The left figure shows the geographical map of Hubei. The darker the city, the more infections per 100,000 inhabitants on February 14. The three cities with the most infections on February 14 are displayed on the right.

of the number of commuters between all cities in Hubei on a daily basis [34]. The static prior is set proportional to the traffic network on January 21, which corresponds to day k = 1.

Figure 2 shows the prediction accuracy over time for different forecast algorithms. The horizontal axis shows the date *d*. We have forecasted the disease several days ahead in time, using all available information from January 22 until *d*. For example, the right-most point in Figure 2a includes data from January 22 to February 13 to forecast the situation on February 14.

The sMAPE error in Figure 2 tends to decrease as time evolves, because a 211 growing amount of data is available. Furthermore, the total number of infected 212 cases quickly increases, whereas the daily infected cases increase at a lower 213 rate, indicating sub-exponential growth [2, 35]. Sub-exponential growth will 214 inevitably reduce the sMAPE error, because sMAPE is a relative error metric. 215 On the other hand, the prediction accuracy decreases rapidly if the forecast 216 horizon is enlarged. Especially the number of cases for five and six days ahead 217 in time around February 1 cannot be predicted accurately, which is illustrated 218 by Figure 2e and 2f, respectively. 219

The logistic function performs generally worse than the other algorithms, for which several reasons may exist. First, by fitting a logistic curve, we assume the number of cases to follow the SIR model closely [2, 8]. Hence, we do not

allow any individual or governmental responses to COVID-19, which typically 223 flattens the (logistic) curve. Second, the logistic function ignores the spread 224 between regions, which further deteriorates the prediction accuracy. Third, the 225 logistic function is symmetric around the epidemic peak at $k = t_0$; the increase 226 and decrease of the number of cases around the peak is equal. Most epidemic 227 outbreaks of COVID-19 show a rapid increase and a more gradual decrease of the 228 daily number of cases. A possible reason is that most lockdowns are enforced 229 immediately, whereas lockdown measures are lifted gradually. Occasionally, 230 the Hill function [5] and Gompertz function [6] are used to predict epidemic 231 outbreaks, because they allow asymmetry around the epidemic peak. In this 232 work, we focus on the logistic function because of its relation to the solution of 233 the SIR and SIS model, and discuss the Hill function and the Gompertz function 234 in Appendix F. 235

The performance of LSTM is moderately good, but LSTM fails to find an accurate forecast around January 31. Since the time series is the shortest at the left part of Figure 2, less data is available to train LSTM. Pure machine-learning algorithms are known to yield a lower prediction accuracy than other methods if the time series is short [36].

The prediction accuracy of all NIPA methods in Figure 2 is similar, although 241 NIPA static prior is considerably worse around February 4 for the prediction 242 of three or more days ahead in time. A possible reason is as follows. The im-243 pact of the nation-wide lockdown on January 24 is captured incorrectly by the 244 static prior, whereas the original NIPA method has more freedom to adjust its 245 contact network accordingly and NIPA dynamic prior receives a more tailored 246 prior to the current situation. Another reason is that the prior network (dy-247 namic or static) may deviate significantly from the true infection matrix. Under 248 ideal circumstances, namely that the epidemic outbreak exactly follows the SIR 249 model, we show that NIPA static prior outperforms NIPA in Appendix E. 250

Figure 2 also shows that the negligence of the network interaction by NIPA separate decreases the prediction accuracy compared to NIPA. Hence, a networkbased approach appears beneficial for forecasting. We summarise the results in





Figure 2: The prediction accuracy for the situation in Hubei, China. The subplots show the prediction accuracy for a forecast horizon of (a) 1 day, (b) 2 days, (c) 3 days, (d) 4 days, (e) 5 days and (f) 6 days for the prediction algorithms from Section 2.

Another interesting topic is *forecast bias*: The tendency to systematically overestimate or underestimate the true number of infected cases. Using the Percentage Error (PE) we estimate the bias for all prediction algorithms for region i at time k. The surface error plots in Figure 3 show the Percentage Error as a function of time for a 4-days ahead prediction. The logistic function and LSTM show the largest deviation around the mean, especially around February

1, which is in agreement with Figure 2. Furthermore, Figure 3 illustrates that 261 the logistic function and LSTM systematically underestimate the true number 262 of cases. On the other hand, NIPA static prior appears to overestimate the 263 true number of cases. A possible reason is the following. The static network is 264 taken to be proportional to the traffic flow before the lockdown measures. When 265 the lockdown is introduced, the static prior remains constant, so the algorithm 266 overestimates the true result. After some time, the newly collected data shows 267 evidence that the prior is not very accurate, so NIPA static prior ignores the 268 prior and uses the data instead, which improves the forecast accuracy again. 260

270 3.2. The Netherlands

As a second case study, we regard the spread of COVID-19 in the Nether-271 lands. The first case was diagnosed on February 27, who had visited Italy the 272 week before. After February 27, the number of cases grew rapidly, as depicted 273 in Figure 4. The epidemic peak was observed at the end of March, and the daily 274 number of cases has dropped ever since. We consider the spread of COVID-19 on 275 a provincial level, for which data is available from the Dutch National Institute 276 for Public Health and the Environment, called RIVM [37]. The Netherlands is 277 subdivided into twelve provinces, for which the RIVM reports the daily number 278 of new infections. Since the number of infected cases increased more gradually 279 in the Netherlands than in Hubei, China, the total epidemic period is longer 280 and more data points are available. A more gradual increase in the number of 281 cases should be beneficial for the prediction accuracy. 282

For NIPA static prior, we require a traffic network as an approximation for 283 the interaction between the provinces. Statistics Netherlands (Centraal Bureau 284 voor de Statistiek) reports the number of people m_{ij} working in province i and 285 living in province j, averaged over one year [38]. We use the Google Mobil-286 ity Data "Workplaces" to estimate the time-varying traffic network for each 287 province in the Netherlands [39]. Google reports the percentage decrease of 288 traffic $p_i[k]$ on day k in province i compared to an ordinary day between Jan-289 uary 3 and February 6, 2020. During the lockdown, we expect $p_i[k] < 1$ because 290



Figure 3: The surface error plots for a 4-days forecast horizon versus time. The subfigures show (a) NIPA, (b) NIPA separate, (c) NIPA static prior, (d) NIPA dynamic prior, (e) Logistic function and (f) LSTM.

of the lockdown measures. Then we construct the time-dependent traffic matrix as follows: $m_{ij}[k] = m_{ij} \cdot p_i[k]$.

The prediction accuracy for the Netherlands is outlined in Figure 5. Before April 1, the situation in the Netherlands is similar to Hubei, where the NIPA methods perform better, but there exist large deviations in the prediction accuracy. After April 1, the accuracy of the NIPA methods is nearly identical. In other words, the influence of the initial static/dynamic network on the prediction is small. The main reason is that the NIPA algorithms are trained on a



Figure 4: The left figure shows the geographical map of the Netherlands. The darker the province, the more infections per 100,000 inhabitants on May 19. The four provinces with the most infections on May 19 are displayed on the right.

growing amount of infection data as time advances. Among the best performing methods over the whole period are traditional NIPA and NIPA separate,
whereas the logistic function and LSTM show the worst performance.

The prediction accuracy of NIPA separate and NIPA are comparable, except at the left-hand side of Figure 5. A possible reason is that the spread of the coronavirus is at the beginning mainly dominated by interprovincial interactions. After the imposing of the lockdown at the end of March, the interaction between the provinces is lowered significantly, so the spreading mainly takes place within each province.

308 4. Conclusion

We have compared the prediction accuracy of eight algorithms to forecast the spread of COVID-19. We summarise the results in Table 2. The error in Table 2 is obtained by averaging over all sMAPE forecast errors for forecast horizons between one and six days. Fitting a sigmoid curve, like the logistic function, performs the worst of all methods. The main reasons for the low prediction accuracy are the imposed symmetry around the epidemic peak and the negligence of the interaction between regions. Other sigmoid curves, such as the Hill



Figure 5: The prediction accuracy for the situation in the Netherlands. The subplots show the prediction accuracy (a) 1 day ahead, (b) 2 days ahead, (c) 3 days ahead, (d) 4 days ahead, (e) 5 days ahead and (f) 6 days ahead.

³¹⁶ function and the Gompertz function, perform slightly better than the logistic ³¹⁷ function, but perform worse than most other algorithms. The machine-learning ³¹⁸ algorithm Long Short-Term Memory (LSTM) is not based on any phenomeno-³¹⁹ logical epidemic processes nor considers provincial interactions. Table 2 shows ³²⁰ that the prediction accuracy of LSTM is comparable to the Hill and Gompertz ³²¹ functions.



of machine learning, phenomenological epidemiology (SIR model) and considers 323 the interaction between different regions. Table 2 illustrates that the prediction 324 accuracy of NIPA is considerably better than any other algorithm. Applying 325 NIPA for each region separately (NIPA separate) yields a forecast error which is 326 comparable to LSTM. We conclude that a network-based approach is beneficial 327 for an accurate forecast. We have also shown that choosing a time-varying or 328 static prior close to the true contact network may improve the forecast accuracy 329 of NIPA. Surprisingly, the inclusion of a time-varying or static prior in NIPA 330 on real infection data is not beneficial for the forecast accuracy for the consid-331 ered regions. Among several reasons, the chosen prior might be an inaccurate 332 estimate of the true contact network. 333

In a practical setting, such as the current COVID-19 pandemic, policymakers might prefer to anticipate to a worst-case scenario. In that case, an asymmetric error metric that penalises underestimations more significantly than overestimations may be more suitable.

Table 2: All algorithms discussed in this paper. The Netherlands is abbreviated as NL. *As input, each algorithm requires the population size N_i of each region i and a time series of the infected cases $y_i[k]$ in each region i at every time k.

Algorithm	Additional input [*]	Error (Hubei)	Error (NL)	Bias
NIPA	-	0.122	0.0381	
NIPA separate	-	0.129	0.0487	
NIPA static prior	static network	0.135	0.0384	over
NIPA dynamic prior	dynamic network	0.129	0.0429	
Logistic function	-	0.186	0.0735	under
Hill function	-	0.142	0.0531	
Gompertz function	-	0.141	0.0528	
LSTM	-	0.160	0.0570	under

338 Acknowledgements

- ³³⁹ LM is supported by the China scholarship council.
- ³⁴⁰ This work has been supported by the Universiteitsfonds Delft in the program
- ³⁴¹ TU Delft COVID-19 Response Fund.

342 References

- [1] K. R. Moran, G. Fairchild, N. Generous, K. Hickmann, D. Osthus, R. Pried horsky, J. Hyman, S. Y. Del Valle, Epidemic Forecasting is Messier
- ³⁴⁵ Than Weather Forecasting: The Role of Human Behavior and Internet
- ³⁴⁶ Data Streams in Epidemic Forecast, The Journal of Infectious Diseases
- ³⁴⁷ 214 (suppl4) (2016) S404-S408. doi:10.1093/infdis/jiw375.

348 URL https://doi.org/10.1093/infdis/jiw375

- B. Prasse, M. A. Achterberg, P. Van Mieghem, Fundamen tal Limits of Predicting Epidemic Outbreaks, retrieved from
 https://www.nas.ewi.tudelft.nl/people/Piet/papers/TUD2020410_
- $_{\tt 352}$ prediction_limits_epidemic_outbreaks.pdf (2020).
- 353 URL https://www.nas.ewi.tudelft.nl/people/Piet/papers/ 354 TUD2020410_prediction_limits_epidemic_outbreaks.pdf
- ³⁵⁵ [3] P. F. Verhulst, Recherches mathématiques sur la loi d'accroissement de la
 ³⁵⁶ population (1845) 1–45.
- 357 URL http://gdz.sub.uni-goettingen.de/dms/load/img/?PPN= 358 PPN129323640_0018
- [4] K. Roosa, Y. Lee, R. Luo, A. Kirpich, R. Rothenberg, J. Hyman, P. Yan,
 G. Chowell, Short-term Forecasts of the COVID-19 Epidemic in Guangdong
 and Zhejiang, China: February 13–23, 2020, J. Clin. Med. 9 (2020) 596.
 doi:10.3390/jcm9020596.
- 363 URL https://www.mdpi.com/2077-0383/9/2/596

364	[5]	A. Hill, Proceedings of the physiological society: January
365		22, 1910, The Journal of Physiology 40 (suppl) (1910) i–vii.
366		doi:10.1113/jphysiol.1910.sp001386.
367		URL https://physoc.onlinelibrary.wiley.com/doi/abs/10.1113/
368		jphysiol.1910.sp001386
369	[6]	B. Gompertz, On the nature of the function expressive of the law of human
370		mortality, and on a new mode of determining the value of life contingencies,
371		Philosophical Transactions of the Royal Society of London 115 (1825) 513– $$
372		583.
373		URL http://www.jstor.org/stable/107756
374	[7]	P. Van Mieghem, Approximate formula and bounds for the time-varying
375		susceptible-infected-susceptible prevalence in networks, Phys. Rev. E 93
376		(2016) 052312. doi:10.1103/PhysRevE.93.052312.
377		URL https://link.aps.org/doi/10.1103/PhysRevE.93.052312
378	[8]	W. O. Kermack, A. G. McKendrick, A contribution to the mathematical
379		theory of epidemics, Proc. R. Soc. Lond. A 115 (1927) 700–721. doi:
380		10.1098/rspa.1927.0118.
381	[9]	M. Day, Covid-19: four fifths of cases are asymptomatic, China figures
382		indicate, BMJ 369 (2020). doi:10.1136/bmj.m1375.
383		URL https://www.bmj.com/content/369/bmj.m1375
384	[10]	Q. Yang, C. Yi, A. Vajdi, L. W. Cohnstaedt, H. Wu, X. Guo, C. M. Scoglio,
385		Short-term forecasts and long-term mitigation evaluations for the COVID-
386		19 epidemic in Hubei Province, China, med Rxiv (2020). doi:10.1101/ $\$
387		2020.03.27.20045625.
388		URL https://www.medrxiv.org/content/early/2020/03/30/2020.03.
389		27.20045625
390	[11]	L. Lorch, W. Trouleau, S. Tsirtsis, A. Szanto, B. Schölkopf, M. Gomez-
391		Rodriguez, A Spatiotemporal Epidemic Model to Quantify the Effects of
392		Contact Tracing, Testing, and Containment (2020). arXiv:2004.07641.

- ³⁹³ [12] S. Y. Chang, E. Pierson, P. W. Koh, J. Gerardin, B. Redbird, D. Grusky,
- J. Leskovec, Mobility network modeling explains higher SARS-CoV-2 infec-
- tion rates among disadvantaged groups and informs reopening strategies,
 medRxiv (2020). doi:10.1101/2020.06.15.20131979.
- ³⁹⁷ URL https://www.medrxiv.org/content/early/2020/06/17/2020.06.
 ³⁹⁸ 15.20131979
- [13] M. Al-qaness, A. Ewees, H. Fan, M. Abd El Aziz, Optimization Method
 for Forecasting Confirmed Cases of COVID-19 in China, J. Clin. Med. 9
 (2020) 674. doi:10.3390/jcm9030674.
- 402 URL https://www.mdpi.com/2077-0383/9/3/674
- [14] Z. Yang, Z. Zeng, K. Wang, S. Wong, W. Liang, M. Zanin, P. Liu, X. Cao,
 Z. Gao, Z. Mai, J. Liang, X. Liu, S. Li, Y. Li, F. Ye, W. Guan, Y. Yang,
 F. Li, S. Luo, Y. Xie, B. Liu, Z. Wang, S. Zhang, Y. Wang, N. Zhong, J. He,
 Modified SEIR and AI prediction of the epidemics trend of COVID-19 in
- 407 China under public health interventions, Journal of Thoracic Disease 12 (3)
 408 (2020).
- 409 URL http://jtd.amegroups.com/article/view/36385
- [15] A. Kergassner, C. Burkhardt, D. Lippold, S. Nistler, M. Kergassner,
 P. Steinmann, D. Budday, S. Budday, Meso-scale modeling of COVID19 spatio-temporal outbreak dynamics in Germany, medRxiv (2020). doi:
 10.1101/2020.06.10.20126771.
- 414 URL https://www.medrxiv.org/content/early/2020/06/17/2020.06.
 415 10.20126771
- [16] S. He, Y. Peng, K. Sun, SEIR modeling of the COVID-19 and its dynamics,
 Nonlinear Dynamics (2020). doi:10.1007/s11071-020-05743-y.
- [17] M. Youssef, C. Scoglio, An individual-based approach to SIR epidemics in
 contact networks, Journal of Theoretical Biology 283 (1) (2011) 136 144.
 doi:10.1016/j.jtbi.2011.05.029.

421 URL http://www.sciencedirect.com/science/article/pii/

422 S0022519311002815

[18] B. Prasse, P. Van Mieghem, Network Reconstruction and Prediction of
Epidemic Outbreaks for General Group-Based Compartmental Epidemic
Models, IEEE Transactions on Network Science and Engineering, to appear
(2020).

427 URL https://ieeexplore.ieee.org/document/9069319

- [19] B. Prasse, M. A. Achterberg, L. Ma, P. Van Mieghem, Network-inferencebased prediction of the COVID-19 epidemic outbreak in the Chinese
 province Hubei, Applied Network Science (35) (2020). doi:10.1007/
 s41109-020-00274-2.
- 432 URL https://link.springer.com/epdf/10.1007/s41109-020-00274-2
- ⁴³³ [20] P. Cirillo, N. N. Taleb, Tail risk of contagious diseases, Nat. Phys. 16 (2020)
 ⁴³⁴ 606-613. doi:10.1038/s41567-020-0921-x.
- 435 URL https://www.nature.com/articles/s41567-020-0921-x
- 436 [21] J. L. Elman, Finding Structure in Time, Cognitive Science 14 (2) (1990)
 437 179-211.
- ⁴³⁸ [22] I. Goodfellow, Y. Bengio, A. Courville, Deep Learning, MIT Press, 2016.
- ⁴³⁹ [23] T. Young, D. Hazarika, S. Poria, E. Cambria, Recent Trends in Deep Learn⁴⁴⁰ ing Based Natural Language Processing [Review Article], IEEE Computa⁴⁴¹ tional Intelligence Magazine 13 (2018) 55–75.
- [24] S. Hochreiter, J. Schmidhuber, Long Short-Term Memory, Neural Computation 9 (8) (1997) 1735–1780. doi:10.1162/neco.1997.9.8.1735.
- ⁴⁴⁴ [25] F. A. Gers, J. Schmidhuber, LSTM recurrent networks learn simple
 ⁴⁴⁵ context-free and context-sensitive languages, IEEE Transactions on Neural
 ⁴⁴⁶ Networks 12 6 (2001) 1333–40.

- ⁴⁴⁷ [26] F. A. Gers, J. Schmidhuber, F. Cummins, Learning to Forget: Continual
 ⁴⁴⁸ Prediction with LSTM, Neural Computation 12 (10) (2000) 2451–2471.
 ⁴⁴⁹ doi:10.1162/089976600300015015.
- ⁴⁵⁰ [27] R. Jozefowicz, W. Zaremba, I. Sutskever, An Empirical Exploration of
 ⁴⁵¹ Recurrent Network Architectures, in: F. Bach, D. Blei (Eds.), In Proc. of
 ⁴⁵² ICML (32nd International Conference on Machine Learning), Vol. 37 of
 ⁴⁵³ PMLR, Lille, France, 2015, pp. 2342–2350.
- ⁴⁵⁴ [28] Y. Yu, X. Si, C. Hu, J. Zhang, A Review of Recurrent Neural Networks:
 ⁴⁵⁵ LSTM Cells and Network Architectures, Neural Computation 31 (7) (2019)
 ⁴⁵⁶ 1235–1270, pMID: 31113301. doi:10.1162/neco_a_01199.
- ⁴⁵⁷ [29] D. P. Kingma, J. Ba, Adam: A Method for Stochastic Optimization, In
 ⁴⁵⁸ Proc of ICLR (International Conference for Learning Representations), San
 ⁴⁵⁹ Diego, 2015 abs/1412.6980 (2015).
- [30] C. Pizzuti, A. Socievole, B. Prasse, P. Van Mieghem, Network-based Pre diction of COVID-19 Epidemic Spreading in ItalySubmitted (2020).
- [31] B. Prasse, Ρ. Van Mieghem, Predicting Dynamics on Net-462 works Hardly Depends on the Topology, available on ArXiv: 463 https://arxiv.org/abs/2005.14575 (2020). arXiv:2005.14575. 464
- [32] R. J. Hyndman, A. B. Koehler, Another look at measures of forecast
 accuracy, International Journal of Forecasting 22 (4) (2006) 679–688.
 doi:10.1016/j.ijforecast.2006.03.001.
- 468 URL http://www.sciencedirect.com/science/article/pii/
 469 S0169207006000239
- 470 [33] News from the Health Commission of Hubei, retrieved on February 16,
- 471 2020 from http://wjw.hubei.gov.cn/fbjd/dtyw (2020).
- 472 URL http://wjw.hubei.gov.cn/fbjd/dtyw

- 473 [34] Baidu Migration website, retrieved on February 16, 2020 from https://
- 474 qianxi.baidu.com/2020/ (2020).
- 475 URL https://qianxi.baidu.com/2020/
- 476 [35] B. F. Maier, D. Brockmann, Effective containment explains subexponential
- 477 growth in recent confirmed COVID-19 cases in China, Science 368 (6492)
- 478 (2020) 742-746. doi:10.1126/science.abb4557.
- 479 URL https://science.sciencemag.org/content/368/6492/742
- [36] S. Makridakis, E. Spiliotis, V. Assimakopoulos, The M4 Competition: 100,000 time series and 61 forecasting methods, International Journal of Forecasting 36 (1) (2020) 54 74, m4 Competition.
 doi:10.1016/j.ijforecast.2019.04.014.
- 484 URL http://www.sciencedirect.com/science/article/pii/
 485 S0169207019301128
- ⁴⁸⁶ [37] RIVM, Actuele informatie over het nieuwe coronavirus (COVID⁴⁸⁷ 19), retrieved on May 25, 2020 from https://www.rivm.nl/
 ⁴⁸⁸ coronavirus-covid-19/actueel (2020).
- 489 URL https://www.rivm.nl/coronavirus-covid-19/actueel
- ⁴⁹⁰ [38] CBS, Banen van werknemers naar woon- en werkregio, retrieved on May
- 29, 2020 from https://opendata.cbs.nl/statline/#/CBS/nl/dataset/
 83628NED/table?ts=1583844319444 (2018).
- 493 URL https://opendata.cbs.nl/statline/#/CBS/nl/dataset/
 494 83628NED/table?ts=1583844319444
- ⁴⁹⁵ [39] G. LLC, Google COVID-19 Community Mobility Reports, retrieved on May
- 496 25, 2020 from https://www.google.com/covid19/mobility/ (2020).
- 497 URL https://www.google.com/covid19/mobility/
- 498 [40] P. E. Paré, J. Liu, C. L. Beck, B. E. Kirwan, T. Başar, Analysis, Estimation,
- and Validation of Discrete-Time Epidemic Processes, IEEE Transactions on
 Control Systems Technology 28 (1) (2020) 79–93.

- [41] S. Boyd, L. Vandenberghe, Convex Optimization, Cambridge University
 Press, 2004. doi:10.1017/CB09780511804441.
- [42] R. Tibshirani, Regression Shrinkage and Selection via the Lasso, Journal
 of the Royal Statistical Society. Series B (Methodological) 58 (1) (1996)
 267–288.
- ⁵⁰⁶ URL http://www.jstor.org/stable/2346178
- ⁵⁰⁷ [43] P. Van den Driessche, J. Watmough, Reproduction numbers and
 ⁵⁰⁸ sub-threshold endemic equilibria for compartmental models of dis⁵⁰⁹ ease transmission, Mathematical Biosciences 180 (1) (2002) 29 48.
 ⁵¹⁰ doi:10.1016/S0025-5564(02)00108-6.
- 511URLhttp://www.sciencedirect.com/science/article/pii/512S0025556402001086
- [44] M. Kiskowski, G. Chowell, Modeling household and community transmission of Ebola virus disease: Epidemic growth, spatial dynamics and insights
 for epidemic control, Virulence 7 (2) (2016) 163–173, pMID: 26399855.
 doi:10.1080/21505594.2015.1076613.
- ⁵¹⁷ URL https://doi.org/10.1080/21505594.2015.1076613
- [45] C. P. Winsor, The Gompertz Curve as a Growth Curve, Proceedings of the
 National Academy of Sciences 18 (1) (1932) 1–8. doi:10.1073/pnas.18.
- ⁵²⁰ 1.1.
- URL https://www.pnas.org/content/18/1/1

522 Appendix A. SIR Epidemic Model

The SIR epidemic model is defined in Definition 1. The COVID-19 pandemic does not exactly follow the SIR epidemic model. Instead, at every time k, the fraction of COVID-19 infections in region i obeys

$$\mathcal{I}_{i}[k+1] = (1-\delta_{i})\mathcal{I}_{i}[k] + \mathcal{S}_{i}[k] \sum_{j=1}^{N} \beta_{ij}[k]\mathcal{I}_{j}[k] + w_{i}[k].$$
(A.1)

Here, $w_i[k]$ denotes the *model error* of region *i* at time *k*. Under Assumption 2, the model errors $w_i[k]$ are identically distributed at all times *k* and for every region *i*:

Assumption 2. The model error $w_i[k]$ is normally distributed as

$$w_i[k] \sim \mathcal{N}\left(0, \sigma_w^2\right).$$
 (A.2)

Furthermore, the model errors $w_i[k]$, $w_j[\tilde{k}]$ are stochastically independent for all times $k \neq \tilde{k}$ and regions $i \neq j$.

Assumption 3. For every node *i*, the curing probabilities satisfy $\delta_i \leq 1$, and, at every time $k \in \mathbb{N}$, the infection probabilities $\beta_{ij}[k]$ satisfy

$$\sum_{j=1}^{N} \beta_{ij}[k] \le 1. \tag{A.3}$$

⁵²⁸ Under Assumption 3, the fractions $S_i[k]$, $\mathcal{I}_i[k]$ and $\mathcal{R}_i[k]$ remain in [0, 1] at ⁵²⁹ every time k as stated by Lemma 4, which is inspired by [40, Lemma 1] and has ⁵³⁰ been proved for time-invariant infection probabilities β_{ij} in [19].

Lemma 4 ([19]). Suppose that $\mathcal{I}_i[1] \ge 0$, $\mathcal{R}_i[1] \ge 0$ and $\mathcal{I}_i[1] + \mathcal{R}_i[1] \le 1$ for every node *i*. Then, under Assumption 3, it holds that $\mathcal{I}_i[k] \ge 0$, $\mathcal{R}_i[k] \ge 0$ and $\mathcal{I}_i[k] + \mathcal{R}_i[k] \le 1$ at every time $k \in \mathbb{N}$ for every node *i*.

Proof. We prove Lemma 4 by induction. Suppose that at time k for every node i it holds that

$$\mathcal{I}_i[k] \ge 0 \tag{A.4}$$

and

$$\mathcal{R}_i[k] \ge 0 \tag{A.5}$$

and

$$\mathcal{I}_i[k] + \mathcal{R}_i[k] \le 1. \tag{A.6}$$

Under Assumption 3 it holds that $0 \leq \delta_i \leq 1$ and $\beta_{ij} \geq 0$. Thus, we obtain from the SIR governing equations (1) and (A.6) that both $\mathcal{I}_i[k+1]$ and $\mathcal{R}_i[k+1]$ equal a sum of positive addends, which implies that

$$\mathcal{I}_i[k+1] \ge 0 \tag{A.7}$$

and

$$\mathcal{R}_i[k+1] \ge 0. \tag{A.8}$$

Furthermore, we obtain for every node i that

$$\mathcal{I}_{i}[k+1] + \mathcal{R}_{i}[k+1] = \mathcal{I}_{i}[k] + \mathcal{R}_{i}[k] + (1 - \mathcal{I}_{i}[k] - \mathcal{R}_{i}[k]) \sum_{j=1}^{N} \beta_{ij}[k] \mathcal{I}_{j}[k].$$
(A.9)

From (A.4), (A.5) and (A.6), we obtain that $\mathcal{I}_i[k] + \mathcal{R}_i[k] \in [0, 1]$. Since (A.5) and (A.6) imply that $\mathcal{I}_i[k] \leq 1$, it holds that

$$\sum_{j=1}^{N} \beta_{ij}[k] \mathcal{I}_j[k] \le 1 \tag{A.10}$$

⁵³⁴ under Assumption 3. Thus, $\mathcal{I}_i[k+1] + \mathcal{R}_i[k+1] \leq 1$, since the right side of ⁵³⁵ (A.9) is a convex combination of 1 and $\sum_{j=1}^N \beta_{ij}[k]\mathcal{I}_j[k] \in [0,1]$.

⁵³⁶ Appendix B. Motivation for the static and dynamic prior

We intend to give a short motivation for (6). Suppose that each individual has on average $\langle d \rangle$ contacts (here $\langle \cdot \rangle$ denotes the average) in the population. If a person is infected and its neighbours are healthy, the person can infect any of its neighbours independently with probability p. Hence, the total number of infections follows a Bernoulli distribution

$$\Pr[m] = {\binom{\langle d \rangle}{m}} p^m (1-p)^{\langle d \rangle - m}. \tag{B.1}$$

In case $\langle d \rangle$ is large and $\lambda \equiv p \langle d \rangle$ is small, we can approximate (B.1) by a Poisson distribution

$$\Pr[m] = e^{-\lambda} \frac{\lambda^m}{m!}.$$
(B.2)

If there are N visiting, infected individuals, which may all infect the population independently, the resulting distribution is the sum of independent, identically distributed Poisson distributions, which is again a Poisson distribution with $\langle m \rangle = N \lambda$.

We denote the number of people living in region j and travelling for work to region i by m_{ij} . Each individual has $\langle d \rangle$ contacts and can infect each individual with probability p. Then region j has on average $m_{ij} \langle d \rangle p$ new infections, provided that no two individuals who visit the same region j have contact to the same people. In particular, the fraction of new infections that region i gets from region j is given by

$$\beta_{ij} = \frac{m_{ij} \langle d \rangle p}{N_i}.$$
 (B.3)

If we define $c_i = \frac{\langle d \rangle p}{N_i}$, we obtain equation (6).

542 Appendix C. Details on NIPA static prior

We assume that the infection matrix B is normally distributed around the prior B_{prior} , whose elements equal $b_{\text{prior},ij} = c_i m_{ij}$:

Assumption 5. Every non-diagonal element β_{ij} , where $i \neq j$, of the matrix B is normally distributed as

$$\Pr\left[\beta_{ij}\right] = \begin{cases} \alpha_i \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left(-\frac{1}{2\sigma_i^2} \left(\beta_{ij} - c_i m_{ij}\right)^2\right) & \text{if } 0 \le \beta_{ij} \le 1, \\ 0 & \text{otherwise.} \end{cases}$$
(C.1)

Here c_i denotes the proportionality constant, and the constant α_i is set such that

$$\int_{\mathbb{R}} \Pr\left[\beta_{ij}\right] d\beta_{ij} = 1.$$
 (C.2)

The normal distribution (C.1) is cut off for values outside of [0, 1], since the infection probability β_{ij} cannot be outside the interval [0, 1]. The standard deviation σ_i is a measure for the accuracy of the prior distribution (C.1). Both the proportionality constant c_i and the standard deviation σ_i are unknown. Assumption 5 implies that the diagonal elements β_{ii} of the matrix B are uniformly distributed in the interval [0, 1].

We obtain the estimate $B_{\text{posterior}}$ of the contact network by a Bayesian (or maximum a posteriori) approach. Given the observed $N \times 1$ infection vector $\mathcal{I}[k] = (\mathcal{I}_1[k], ..., \mathcal{I}_N[k])^T$ at all times k = 1, ..., n, we pose the optimisation problem

$$B_{\text{posterior}} = \underset{B}{\operatorname{argmax}} \operatorname{Pr} \left[B \big| \mathcal{I}[1], ..., \mathcal{I}[n] \right]$$
(C.3)
s.t.
$$\sum_{j=1}^{N} \beta_{ij} \leq 1, \quad i = 1, ..., N.$$

With the constraint in (C.3), we ensure that the predictions of the infections satisfy $0 \leq \mathcal{I}_i[k] \leq 1$, see Lemma 4 in Appendix A. We define the $(n-1) \times 1$ vector V_i and the $(n-1) \times N$ matrix F_i as [19]

$$V_{i} = \begin{pmatrix} \mathcal{I}_{i}[2] - (1 - \delta_{i})\mathcal{I}_{i}[1] \\ \vdots \\ \mathcal{I}_{i}[n] - (1 - \delta_{i})\mathcal{I}_{i}[n - 1] \end{pmatrix}$$
(C.4)

and

$$F_{i} = \begin{pmatrix} \mathcal{S}_{i}[1]\mathcal{I}_{1}[1] & \dots & \mathcal{S}_{i}[1]\mathcal{I}_{N}[1] \\ \vdots & \ddots & \vdots \\ \mathcal{S}_{i}[n-1]\mathcal{I}_{1}[n-1] & \dots & \mathcal{S}_{i}[n-1]\mathcal{I}_{N}[n-1] \end{pmatrix}.$$
 (C.5)

We obtain the Bayesian estimate $B_{\text{posterior}}$ by solving a constrained linear leastsquares problem. Proposition 6 is an adaptation of the Bayesian interpretation in [31].

Proposition 6. Under Assumptions 2 and 5, the Bayesian estimation problem

(C.3) is equivalent to solving the optimisation problem

$$\min_{\substack{\beta_{i1},\ldots,\beta_{iN}}} \left\| V_i - F_i \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{iN} \end{pmatrix} \right\|_2^2 + \rho_i \sum_{j=1, j \neq i}^N (\beta_{ij} - c_i m_{ij})^2$$
s.t. $0 \le \beta_{ij} \le 1, \quad j = 1, ..., N,$

$$\sum_{j=1}^N \beta_{ij} \le 1,$$
(C.6)

for every region *i*, where the penalisation parameter equals $\rho_i = \sigma_w^2 / \sigma_i^2$.

Proof. The objective function of the optimisation problem (C.3) is equivalent to

$$\hat{B} = \underset{B}{\operatorname{argmax}} \log \left(\Pr\left[B\right] \right) + \sum_{k=2}^{n} \log \left(\Pr\left[\mathcal{I}[k] \middle| \mathcal{I}[k-1], B\right] \right).$$
(C.7)

In the following, we rewrite the two terms in (C.7). First, with (C.1), it holds that

$$\log\left(\Pr\left[B\right]\right) = \begin{cases} \sum_{i=1}^{N} \sum_{j=1}^{N} \log\left(\alpha_{i}\right) - \log\left(\sqrt{2\pi}\sigma_{i}\right) - \frac{1}{2\sigma_{i}^{2}} \left(\beta_{ij} - c_{i}m_{ij}\right)^{2} & \text{if } 0 \leq \beta_{ij} \leq 1 \ \forall i, j, \\ -\infty & \text{otherwise.} \end{cases}$$

$$(C.8)$$

Neither the term $\log(\alpha_i)$ nor the term $\log(\sqrt{2\pi\sigma_i})$ depend on the matrix B. Furthermore, the prior $\log(\Pr[B])$ is finite only if $0 \leq \beta_{ij} \leq 1$ for all regions i, j. Thus, the optimisation problem (C.7) is equivalent to

$$\hat{B} = \underset{B}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{j=1}^{N} -\frac{1}{2\sigma_i^2} \left(\beta_{ij} - c_i m_{ij}\right)^2 + \sum_{k=2}^{n} \log\left(\Pr\left[\mathcal{I}[k] \middle| \mathcal{I}[k-1], B\right]\right)$$

s.t. $0 \le \beta_{ij} \le 1, \quad i = 1, ..., N, \quad j = 1, ..., N.$ (C.9)

Second, since the model errors $w_i[k]$ are stochastically independent for different

regions i, we can rewrite the second term in the objective of (C.9) as

$$\log\left(\Pr\left[\mathcal{I}[k]\big|\mathcal{I}[k-1],B\right]\right) = \sum_{i=1}^{N} \log\left(\Pr\left[\mathcal{I}_{i}[k]\big|\mathcal{I}[k-1],B\right]\right)$$
(C.10)

$$= \sum_{i=1}^{N} \log \left(\Pr\left[w_i[k] = \Delta_i[k] \right] \right), \qquad (C.11)$$

where the second equality follows from (A.1) and by defining

$$\Delta_i[k] = \mathcal{I}_i[k] - (1 - \delta_i) \,\mathcal{I}_i[k - 1] + \mathcal{S}_i[k - 1] \sum_{j=1}^N \beta_{ij} \mathcal{I}_j[k - 1].$$
(C.12)

Under Assumption 2, the model error $w_i[k]$ follows the normal distribution. Thus, it holds that

$$\log\left(\Pr\left[w_i[k] = \Delta_i[k]\right]\right) = -\log\left(\sqrt{2\pi}\sigma_w\right) - \frac{1}{2\sigma_w^2}\Delta_i^2[k].$$
(C.13)

The term $\log(\sqrt{2\pi}\sigma_w)$ is independent of the matrix *B*. Thus, it follows from (C.10) and (C.13) that the second term in the objective of (C.9) can be replaced by

$$\sum_{i=1}^{N} \sum_{k=2}^{n} \frac{1}{2\sigma_{w}^{2}} \Delta_{i}^{2}[k] = \sum_{i=1}^{N} \frac{1}{2\sigma_{w}^{2}} \left\| V_{i} - F_{i} \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{iN} \end{pmatrix} \right\|_{2}^{2}, \quad (C.14)$$

where the equality follows from the definition of the vector V_i and the matrix F_i in (C.4) and (C.5), respectively. Hence, the optimisation problem (C.9) becomes

$$\hat{B} = \underset{B}{\operatorname{argmin}} \sum_{i=1}^{N} \frac{1}{2\sigma_{w}^{2}} \left\| V_{i} - F_{i} \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{iN} \end{pmatrix} \right\|_{2}^{2} + \sum_{i=1}^{N} \frac{1}{2\sigma_{i}^{2}} \sum_{j=1}^{N} (\beta_{ij} - c_{i}m_{ij})^{2}$$

s.t. $0 \le \beta_{ij} \le 1, \quad i = 1, ..., N, \quad j = 1, ..., N.$ (C.15)

The problem (C.15) can be optimised independently for every region *i*. Thus, we obtain, after multiplication with $2\sigma_w^2$, the equivalent optimisation problems

for every region i as

$$\min_{\substack{\beta_{i1},\ldots,\beta_{iN}}} \left\| V_i - F_i \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{iN} \end{pmatrix} \right\|_2^2 + \frac{\sigma_w^2}{\sigma_i^2} \sum_{j=1}^N \left(\beta_{ij} - c_i m_{ij} \right)^2 \quad (C.16)$$
s.t. $0 \le \beta_{ij} \le 1, \quad j = 1, \dots, N.$

⁵⁵⁵ By identifying $\rho_i = \sigma_w^2 / \sigma_i^2$, we obtain that (C.16) with the constraint $\sum_{j=1}^N \beta_{ij} \leq$ ⁵⁵⁶ 1 is equivalent to the constrained linear least-squares problem (C.6).

The first term in the objective of (C.6) measures the fit to the observed epidemic data. The second term measures the deviation of the infection rates β_{ij} from the prior (C.1). The scalar parameter ρ_i balances the two terms: if the prior (C.1) is very accurate or the model errors $w_i[k]$ are large, then ρ_i should be large. The optimal value of the parameter ρ_i equals to the ratio of the unknown variances σ_w^2 and σ_i^2 of the model errors $w_i[k]$ and the prior (C.1), respectively. The optimisation problem (C.6) is convex and can be solved efficiently [41]. To obtain the solution to (C.6) numerically, we make use of the Matlab command lsqlin. We stress the similarity of the optimisation problem (C.6) to the *least absolute shrinkage and selection operator* (LASSO) of Tibshirani [42], which is the basis of NIPA without prior [19]. Instead of the second least-squares term in the objective of (C.6), LASSO considers the ℓ_1 -norm penalisation term

$$\rho_i \sum_{j=1, j \neq i}^N |\beta_{ij}|. \tag{C.17}$$

⁵⁵⁷ In fact, NIPA without prior can also be interpreted as a Bayesian estimation ⁵⁵⁸ approach [31].

559 Appendix C.1. Pseudocode

To solve the optimisation problem (C.6) for the infection rates $\beta_{i1}, ..., \beta_{iN}$, we must specify three unknown variables. First, the curing rate δ_i of region *i*, which determines the fractions $S_i[k]$ and $\mathcal{R}_i[k]$ of susceptible and recovered individuals, respectively [19]. Second, we must specify the parameter ρ_i . Third, also the proportionality constant c_i of the prior (C.1) is not known. We perform cross-validation to set the three unknown variables δ_i , ρ_i , c_i .

NIPA static prior is similar to NIPA without prior, except for two alterations. 566 First, we solve the constrained linear least-squares problem (C.6) instead of 567 LASSO. Second, additionally to the parameter ρ_i and the curing rate δ_i , for 568 Bayesian NIPA there is one more unknown variable, namely the proportionality 569 constant c_i , which is a parameter of the prior distribution (C.1). To determine 570 the constant c_i , we consider 50 logarithmically equidistant candidate values in 571 the set $\Psi = \{c_{\min}, ..., c_{\max}\}$. The minimal and the maximal values are set to 572 $c_{\min} = 0.01$ and $c_{\max} = 100$, respectively. We set the value of c_i by cross-573 validation. To obtain the epidemic outbreak prediction of Bayesian NIPA, we 574 execute [19, Algorithm 1], where [19, Algorithm 2] is replaced by Algorithm 1 575 stated below. 576

577 Appendix D. Details on NIPA dynamic prior

We assume that the time-varying infection rates $\beta_{ij}[k]$ are proportional to the known population flow $m_{ij}[k]$. More precisely, we assume that the infection rates $\beta_{ij}[k]$ for all regions i, j, when $i \neq j$, equal

$$\beta_{ij}[k] = c_i m_{ij}[k] \tag{D.1}$$

for some unknown proportionality constant $c_i > 0$. Furthermore, we assume the self-infection probabilities β_{ii} do *not* change over time k. With (D.1), the SIR model in Definition 1 yields that

$$\mathcal{I}_{i}[k+1] = (1-\delta_{i})\mathcal{I}_{i}[k] + \beta_{ii}\mathcal{S}_{i}[k]\mathcal{I}_{i}[k] + c_{i}\mathcal{S}_{i}[k] \sum_{j=1, j\neq i}^{N} m_{ij}[k]\mathcal{I}_{j}[k] + w_{i}[k].$$
(D.2)

578 Appendix D.1. Maximum-Likelihood Estimation

To predict the infectious state $\mathcal{I}_i[k]$ with (D.2), we must estimate the constants c_i , the self-infection probabilities β_{ii} and the curing rates δ_i . We define

Algorithm 1 NIPA static prior

- 1: **Input:** curing probability δ_i ; viral state $v_i[k]$ for k = 1, ..., n; infection state vector $\mathcal{I}[k]$ for k = 1, ..., n
- Output: infection probability estimates β_{i1}(δ_i), ..., β_{iN}(δ_i); mean squared error MSE(δ_i)
- 3: Compute V_i and F_i
- 4: $\rho_{\max,i} \leftarrow 2 \| F_i^T V_i \|_{\infty}$
- 5: $\rho_{\min,i} \leftarrow 10^{-4} \rho_{\max,i}$
- 6: $\Theta_i \leftarrow 100$ logarithmically equidistant values from $\rho_{\min,i}$ to $\rho_{\max,i}$
- 7: $\Psi \leftarrow 50$ logarithmically equidistant values from $c_{\min} = 0.01$ to $c_{\max} = 100$
- 8: for $\rho_i \in \Theta_i$ do
- 9: for $c_i \in \Psi$ do
- 10: estimate $MSE(\delta_i, \rho_i, c_i)$ by 3-fold cross-validation on F_i, V_i and solving (C.6) on the respective training set
- 11: **end for**
- 12: end for
- 13: $(\rho_{\text{opt},i}, c_{\text{opt},i}) \leftarrow \underset{\rho_i \in \Theta_i, c_i \in \Psi}{\operatorname{argmin}} \operatorname{MSE}(\delta_i, \rho_i, c_i)$
- 14: $(\beta_{i1}(\delta_i), ..., \beta_{iN}(\delta_i)) \leftarrow$ the solution to (C.6) on the whole data set F_i, V_i for

 $\rho_i = \rho_{\text{opt},i} \text{ and } c_i = c_{\text{opt},i}$

15: $MSE(\delta_i) \leftarrow MSE(\delta_i, \rho_{opt,i}, c_{opt,i})$

the $N \times 1$ vectors $c = (c_1, ..., c_N)^T$ and $b = (\beta_{11}, ..., \beta_{NN})^T$. We pose the estimation problem in a maximum-likelihood sense as

$$\max_{\substack{c,b} \\ \text{s.t.}} \quad \Pr\left[\mathcal{I}[1], ..., \mathcal{I}[n] \middle| c, b\right]$$

$$\text{s.t.} \quad c_i \ge 0, \quad i = 1, ..., N,$$

$$\beta_{ii} \ge 0, \quad i = 1, ..., N,$$

$$\beta_{ii} + c_i \sum_{j=1, j \ne i}^{N} m_{ij}[k] \le 1 \quad i = 1, ..., N, k = 1, ..., n.$$

$$(D.3)$$

The last constraint in (D.3) ensures that the predictions of the infections satisfy $\mathcal{I}_i[k] \leq 1$, see Lemma 4. From the maximum likelihood problem (D.3) we derive, for every region *i*, the LASSO optimisation problem as

$$\min_{c_{i},\beta_{ii}} \sum_{k=1}^{n-1} \left(\mathcal{I}_{i}[k+1] - (1-\delta_{i})\mathcal{I}_{i}[k] - \beta_{ii}\mathcal{S}_{i}[k]\mathcal{I}_{i}[k] - c_{i}\mathcal{S}_{i}[k] \sum_{j=1, j\neq i}^{N} m_{ij}[k]\mathcal{I}_{j}[k] \right)^{2} + \rho_{i}(\beta_{ii} + c_{i})$$
s.t. $c_{i} \geq 0$,
 $\beta_{ii} \geq 0$,
 $\beta_{ii} + c_{i} \sum_{j=1, j\neq i}^{N} m_{ij}[k] \leq 1, \quad k = 1, ..., n.$
(D.4)

Here, we denote the regularisation parameter by $\rho_i \geq 0$, which aims to avoid overfitting. The greater the parameter ρ_i , the smaller the estimates of the coefficients β_{ii}, c_i . If the regularisation parameter $\rho_i = 0$, then solving the LASSO (D.4) for every node *i* is equivalent to solving the maximum-likelihood problem (D.3). (The equivalence of the optimisation problem (D.3) and the LASSO (D.4) can be derived analogously to Proposition 6.)

To solve the optimisation problem (D.4) for the constants c_i and β_{ii} , we must specify two unknown variables. First, the curing rate δ_i of region *i*, which determines the fractions $S_i[k]$ and $\mathcal{R}_i[k]$ of susceptible and recovered individuals, respectively [19]. Second, we must specify the parameter ρ_i . We perform holdout cross-validation to set the unknown variables δ_i and ρ_i : The training set follows from the first 80% of the observations, and the validation set equals the

- ⁵⁹¹ last 20% of the observations. In pseudocode, NIPA dynamic prior is given by
- ⁵⁹² Algorithm 2.

Algorithm 2 NIPA dynamic prior

- 1: **Input:** curing probability δ_i ; viral state $v_i[k]$ for k = 1, ..., n; infection state vector $\mathcal{I}[k]$ for k = 1, ..., n
- 2: **Output:** infection probability estimates $\beta_{i1}(\delta_i), ..., \beta_{iN}(\delta_i)$; mean squared error $MSE(\delta_i)$
- 3: Compute V_i and F_i
- 4: $\rho_{\max,i} \leftarrow 2 \|F_i^T V_i\|_{\infty}$
- 5: $\rho_{\min,i} \leftarrow 10^{-4} \rho_{\max,i}$
- 6: $\Theta_i \leftarrow 100$ logarithmically equidistant values from $\rho_{\min,i}$ to $\rho_{\max,i}$
- 7: for $\rho_i \in \Theta_i$ do
- 8: estimate $MSE(\delta_i, \rho_i)$ by hold-out cross-validation on F_i, V_i and solving (D.4) on the respective training set
- 9: end for
- 10: $\rho_{\text{opt},i} \leftarrow \operatorname*{argmin}_{\rho_i \in \Theta_i} \mathrm{MSE}\left(\delta_i, \rho_i\right)$
- 11: $(\beta_{i1}(\delta_i), ..., \beta_{iN}(\delta_i)) \leftarrow$ the solution to (D.4) on the whole data set F_i, V_i for $\rho_i = \rho_{\text{opt},i}$
- 12: $MSE(\delta_i) \leftarrow MSE(\delta_i, \rho_{opt,i})$

⁵⁹³ Appendix E. NIPA static prior under perfect conditions

The original NIPA method is known to provide accurate predictions when the epidemic perfectly follows the SIR model [19, Supplementary Material 1]. Here, we intend to show that NIPA static prior performs even better if the prior matrix is close to the real infection matrix.

Suppose we generate data from an SIR epidemic as in Definition 1. We use a network with N = 10 nodes with an equal curing rate δ for each node: $\delta_i = 0.2$ for all *i*. We set the curing rate δ_i in the NIPA algorithms equal to the exact curing rates $\delta_i = 0.2$, such that both NIPA and NIPA static prior will always estimate the curing rates correctly. We consider infection probabilities β_{ij} which are uniformly distributed in the interval (0, 1). The effective reproduction number R_0 can be computed as [43]

$$R_0 = \text{maximum eigenvalue of } \left(B \cdot \text{diag}\left(\frac{1}{\delta_1}, ..., \frac{1}{\delta_N}\right) \right).$$
 (E.1)

We normalise B element-wise such that the basic reproduction number R_0 equals 2.0. Furthermore, we set the population size N_i for each region i equal to a uniformly distributed number in the interval $[10^5, 10^6]$ and start with initially $y_1[1] = 100$ infected cases in node 1 and all other nodes are healthy. Most importantly, we set the prior infection matrix B_{prior} equal to the exact infection matrix B, multiplied by some noise

$$B_{\text{prior},ij} = \beta_{ij} w_{ij}.$$
(E.2)

Here, w_{ij} is uniformly distributed in the interval [1,2]. The other parameters are the same as in the main article.

The result in Figure E.6 is clear: NIPA static prior is able to capture the dynamics much better than NIPA. Hence, we conclude that NIPA static prior in combination with a good prior yields a better prediction accuracy than the original NIPA method.



Figure E.6: The prediction for (a) NIPA and (b) NIPA static prior with generated SIR data based on Definition 1 on a 10-node network.

⁶⁰⁴ Appendix F. Sigmoid curves

In epidemiology, sigmoid curves are commonly used to forecast the future number of infected cases. The logistic function was developed by Verhulst in 1845 to explain the growth of the population in a specific region [3]. The logistic function is the most often used sigmoid curve in epidemiology, because the logistic function also follows as the (approximate) solution of the SIS and SIR model [2]. The logistic function assumes the cumulative number of infected cases $y_i[k]$ in region *i* and time *k* to follow

$$y_i[k] = \frac{y_{\infty,i}}{1 + e^{-K_i(k - t_{0,i})}}.$$
 (F.1)

where $y_{\infty,i}$ is the long-term fraction of infections, K_i is the logistic growth rate and $t_{0,i}$ is the inflection point, which is also known as the epidemic peak.

607

608

609

The Hill function was introduced in 1910 to describe the binding of molecules on surfaces [5]. Later, it was successfully applied to describe the spread of epidemics [44]. The Hill function assumes the cumulative number of infected cases $y_i[k]$ in region *i* at time *k* to follow

$$y_i[k] = \frac{y_{\infty,i}}{1 + \left(\frac{K_i}{k - t_{0,i}}\right)^{n_i}},$$
 (F.2)

where $y_{\infty,i}$ is the long-term fraction of infections, K_i is the Hill growth rate, n_i is the Hill coefficient and $t_{0,i}$ is the inflection point, also known as the epidemic peak.

The Gompertz function was introduced in 1825 to describe human mortality in a general population [6]. Later the Gompertz function was also used to describe the spread of epidemics [45]. The Gompertz function assumes the cumulative number of infected cases $y_i[k]$ in region *i* at time *k* to follow

$$y_i[k] = y_{\infty,i} e^{-c_i e^{-a_i k}},$$
 (F.3)

where $y_{\infty,i}$ is the long-term fraction of infections, c_i is a displacement factor (comparable to the inflection point) and a_i is the Gompertz growth rate. We describe the curve-fitting procedure here for the logistic function, but the parameters for any curve can be estimated analogously. Suppose that we have a time series of the cumulative number of reported cases $y_{\text{rep},i}[k]$ for k = 1, ..., n and for every region *i*. Then we minimise the Mean Square Error for each region separately;

$$(\hat{y}_{\infty,i}, \hat{K}_{i}, \hat{t}_{0,i}) = \min_{(y_{\infty,i}, K_{i}, t_{0,i})} \sum_{k=1}^{n} \left(y_{\text{rep},i}[k] - \frac{y_{\infty,i}}{1 + e^{-K_{i}(k-t_{0,i})}} \right)^{2},$$

s.t. $0 \le y_{\infty,i} \le N_{i},$
 $K_{i} \ge 0,$
 $t_{0,i} \ge 0,$ (F.4)

where N_i is the population of region *i*. We evaluate the nonlinear minimisation problem (F.4) by the command **GlobalSearch** in Matlab. As initial conditions, we provide $y_{\infty,i}^0 = y(t_{obs}), K_i = 1, t_{0,i} = t_{obs}$. The parameters $(y_{\infty,i}, K_i, n_i, t_{0,i})$ for the Hill function and $(y_{\infty,i}, c_i, a_i)$ for the Gompertz function can be estimated analogously.

Appendix G. The influence of the time step on the prediction accuracy

In the discrete-time SIR model (1), we use the time step $\Delta t = 1$ day. By ap-623 proximating a continuous-time process (the COVID-19 pandemic) by a discrete-624 time process (SIR model) we make a model error. We investigate the influence 625 of the time step on the prediction accuracy, by comparing the NIPA prediction 626 accuracy for various time steps, ranging from $\Delta t = 0.5$ days to $\Delta t = 3$ days. 627 Since the number of infected cases is (generally) reported once a day, the data 628 for the time step $\Delta t = 0.5$ days is obtained by linearly interpolating the number 629 of cumulative cases $y_i[k]$. For a time step $\Delta t = 1$ day and $\Delta t = 0.5$ days, we 630 smooth the raw data before calling the NIPA algorithm [19]. 631

For the time steps $\Delta t = 2$ days and $\Delta t = 3$ days, there are two possible methods. Method (A) assumes that the cumulative number of cases $y_i[k]$ is reported every two (or three) days, and is unreported on the intermediate days. Then we smooth the remaining data, whereafter the NIPA algorithm is used. In fact, we have omitted the data on the intermediate days. In contrast, method (B) first smooths all raw data. Thereafter, we only use the cumulative number of cases $y_i[k]$ every two or three days for a time step of two or three days, respectively. The main difference is that method (A) completely neglects the data on intermediate days, whereas method (B) first applies a smoother and then neglects the intermediate data.

Figure G.7 and G.8 show an exemplary situation from the Netherlands for 642 three initial dates. The configuration for the time step $\Delta t = 1$ day and $\Delta t = 0.5$ 643 days is equal in both figures. In the beginning of the COVID-19 outbreak, as 644 shown in Figure G.7a for method (A) and Figure G.8a for method (B), the 645 prediction accuracy is similar for all time steps. The small amount of available 646 data and the rapidly increasing number of cases hampers accurate forecasting. 647 As the epidemic evolves, method (A) and method (B) start to deviate. By 648 omitting data as in method (A), the sMAPE error in Figure G.7 increases for 649 two and three days quicker than for smaller time steps. Hence, removing data 650 causes the prediction accuracy to decrease. On the other hand, method (B) in 651 Figure G.8 shows similar behaviour for all time steps. We conclude that if the 652 amount of data is unchanged, the choice of the time step has limited effect on 653 the prediction accuracy. 654

40



Figure G.7: (Method A: First remove, then smooth) The NIPA prediction accuracy for the situation in the Netherlands for varying time steps Δt . The subplots show the forecast for (a) March 18, (b) April 5 and (c) April 23. For the time step $\Delta t = 2$ days or $\Delta t = 3$ days, the data is first removed and then smoothed.



Figure G.8: (Method B: First smooth, then remove) The NIPA prediction accuracy for the situation in the Netherlands for varying time steps Δt . The subplots show the forecast for (a) March 18, (b) April 5 and (c) April 23. For the time step $\Delta t = 2$ days or $\Delta t = 3$ days, the data is first smoothed and then removed.