

First-principles machine learning modelling of COVID-19

Luca Magri^{*1,2} and Nguyen Anh Khoa Doan^{3,4}

¹University of Cambridge, Department of Engineering, Cambridge CB2 1PZ, United Kingdom

²Institute for Advanced Study, Technical University of Munich, Garching 85748, Germany (visiting fellow)

³Department of Mechanical Engineering, Technical University of Munich, Garching 85747, Germany

⁴Institute for Advanced Study, Technical University of Munich, Garching 85748, Germany

Abstract

Background: The coronavirus disease 2019 (COVID-19) has changed the world since the World Health Organization declared its outbreak on 30th January 2020, recognizing the outbreak as a pandemic on 11th March 2020. As often said by politicians and scientific advisors, the objective is “to flatten the curve”, or “push the peak down”, or similar wording, of the virus spreading. Central to the official advice are mathematical models and data, which provide estimates on the evolution of the number of infected, recovered and deaths. The accuracy of the models is improved day by day by inferring the contact, recovery, and death rates from data (confirmed cases).

Methods: A data-driven model trained with *both* data *and* first principles is proposed. The model can quickly be re-trained any time that new data becomes available.

Data: John Hopkins University CSSE has been collecting global data from official organizations, such as the World Health Organization, Italy Ministry of Health, and others [1].

Results: The outputs of the analysis are the estimates of infected, recovered and deaths due to COVID-19, as well as the contact, recovery, death rates, basic reproduction number (R_0) and doubling times. The following case studies are analysed: United Kingdom, Italy, Germany, France, Spain, Belgium, USA, New York City, China, and the World. A summary of the results is shown in Table 2. A fast exponential growth in the absence of intervention is found for all cases.

Discussion: The method can be applied to more detailed epidemic models with virtually no conceptual modification.

Acknowledgements: L. Magri is advising the Scientific Pandemic Influenza Group on Modelling (SPI-M) through the Royal Society’s Rapid Assistance in Modelling the Pandemic (RAMP) initiative (<https://epcced.github.io/ramp/>).

Competing interests: The authors declare no competing interests.

1 Introduction

In December 2019, a cluster of unexplained pneumonia cases in Wuhan, the capital of Hubei province in the People’s Republic of China, resulted into a global pandemic by 11 March 2020, as declared by the

*Corresponding author: lm547@cam.ac.uk

World Health Organization [2]. The disease is caused by a single-stranded RNA coronavirus (Severe acute respiratory syndrome coronavirus 2, SARS-CoV-2) similar to the pathogen responsible for SARS (severe acute respiratory syndrome) and MERS (Middle East respiratory syndrome). The disease caused by this virus has been named COVID-19 (Coronavirus Disease 2019). On 19th April 2020, 1:00 BST, the World Health Organization [2] reported 2,203,927 confirmed cases, 148,749 confirmed deaths, and 213 countries / territories with cases [2].

To control the epidemic, aggressive measures have been implemented worldwide, for example, self-isolation of confirmed and suspected cases, contact tracing and tracking, and social distancing. According to the data, the most draconian measures have managed (or are managing) to suppress (or substantially mitigate) the epidemic. Examples are the localised lockdown of the Hubei region in China (23rd-24th January 2020) [3]; and the national lockdowns of Italy (9th March 2020) [4], Spain (14th March 2020) [5]; the United Kingdom (24th March 2020) [6], among others.

Scientific advice typically relies on estimates of the contact, recovery and death rates. This information is summarized in the basic reproduction number, R_0 , which is the average number of new infections generated by a single infected person within a susceptible population. Estimates of COVID-19 R_0 are variable due to the different methods, models and parameters employed, as well as the databases used [7]. As reported in [7], most official sources estimates R_0 to fall in the range 2 – 3. *Flattening the curve* or *keeping the peak down*, or similar wording, which have been extensively used by governments to level with a lay audience, can be achieved by either reducing the contact rate, β , or by increasing the recovery rate, γ [8]. The latter can be achieved with a vaccine or a cure, which is not presently available. Therefore, to *flatten the curve*, Governments are acting on minimising the contact rate [8]. The objective of this paper is threefold. First, a model that optimally combines data from official databases and first principles of an epidemic model is proposed. Second, the model is applied to provide quantitative estimates on the contact, recovery, death rates; the basic reproduction number, R_0 ; the doubling times; and the evolution of the number of infected, recovered, deaths, and susceptible. Ten cases are analysed: United Kingdom, Italy, Germany, France, Spain, Belgium, USA, New York City, China, and the World. Third, predictions of future dynamics are provided. Although the results are consistent with the first principles and working assumptions used, they are affected by uncertainty because of biases in the data, such as errors in reporting, changes in case definition and testing regime, [7], and modelling assumptions. However, as argued by [7], the fast growth rate and large numbers likely make small biases negligible; and multiplicative corrections, such as constant under-reporting, affect the observed trend only weakly. The paper is structured as follows. The method is presented in Sec. 2 and the results are shown in Sec. 3.

2 Methods and data

In first-principles machine learning modelling, we need first principles and data (*machine learning*) to generate a model. Section 2.1 introduces the first principles and working assumptions, Sec. 2.2 describes the data, and Sec. 2.3 formulates the problem as a constrained optimization problem. The proposed solution method is presented in Sec. 2.4.

2.1 Epidemic model: First principles

The COVID-19 infectious disease is an epidemic [9]. To model an epidemic, suitable groups (also known as compartments [10]) are defined to cover the entire population of a country. Because (i) the epidemic has a (relatively) short time scale, for which the new births can be neglected; (ii) the number of deaths is small as compared with the entire population; and (iii) travel restrictions are enforced, the population, N , is assumed to be constant. The population of a country is divided into mutually exclusive groups: susceptible (S), infected (I), deceased (D), and recovered (R) (Fig. 1). In this model, the deaths are due to COVID-19. Every group is assumed to have the same characteristics, i.e., the groups are homogeneous. Every susceptible person can contract the virus (the immune group is neglected). These working assumptions can be relaxed in more complex models [11, 9].

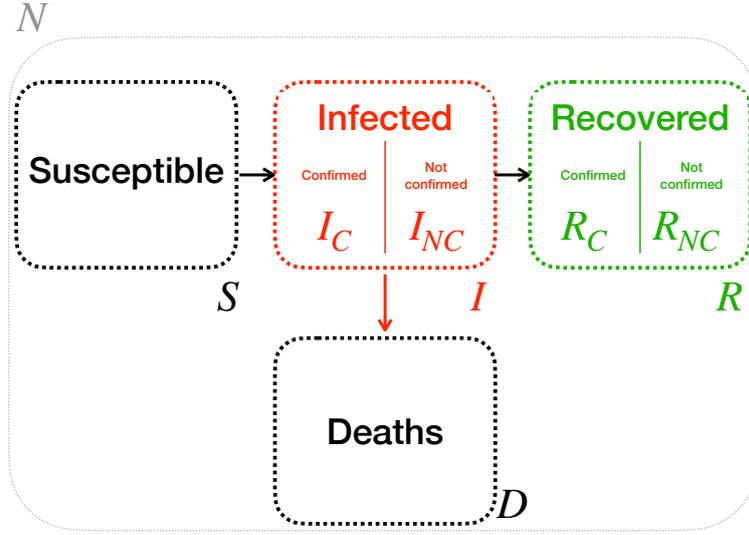


Figure 1: Population and groups in the SIRD-epidemic model.

Mathematically,

$$N = S + R + I + D. \quad (1)$$

Equation (1) is a continuity equation. The population N , which does not vary in time, is the sum of the groups S , I , R , D , which vary in time. This compartmental approach is known as the SIR-epidemic model with vital dynamics and constant population [10, 12]. The model will be called the SIRD-model for brevity. The working assumptions and first principles are mathematically expressed by four ordinary differential equations (ODEs) with time-varying parameters (non-autonomous dynamical system)

$$\dot{S} = -\beta \frac{I}{N} S, \quad (2)$$

$$\dot{I} = \beta \frac{I}{N} S - (\mu + \gamma) I \quad (3)$$

$$\dot{R} = \gamma I \quad (4)$$

$$\dot{D} = \mu I \quad (5)$$

subject to initial conditions S_0 , I_0 , R_0 and D_0 . The symbol $\dot{\cdot}$ denotes the time derivative, d/dt . In compact form

$$\dot{\mathbf{q}} = \mathbf{F}(\mathbf{q}; \boldsymbol{\alpha}) \quad (6)$$

$$\mathbf{q} = \mathbf{q}_0 \quad \text{at } t = 0 \quad (7)$$

where \mathbf{F} is the model (i.e., the SIRD equations), and

$$\mathbf{q} \equiv [S, I, R, D]^T, \quad (8)$$

$$\boldsymbol{\alpha} \equiv [\beta, \gamma, \mu]^T, \quad (9)$$

are the column vectors of the state and parameters, respectively. I/N is the probability to come into contact with an infected individual; β is the average number of contacts per person per unit of time weighed by the transmissibility (contact rate); γ is the average number of recovered people per unit of time (recovery rate); μ is the average number of deaths due to COVID-19 per unit of time (death

rate). These parameters are time dependent and depend on several variables, such as governmental policies (lockdown, school/university closures, social distancing, etc.), heterogeneity in the population (age, life style, herd immunity, hygiene standards, etc.), and properties of the epidemic (virus genome, spreading mechanisms, etc.). The SIRD parameters estimate the epidemic time scales: $1/\gamma$ is the average time to recover; $1/\beta$ is the average time between one contact (with an infected) and another; and $1/\mu$ is the average time to decrease (for those who do not recover). The basic reproduction ratio¹, $R_0 \equiv \beta/\gamma$, is the expected number of secondary infections from a single infection entering a population where all members are susceptible [9]. If $R_0 > 1$, the number of infected increases (Eq. (3)). If $R_0 < 1$, the disease does not grow on average. The total number of new cases per unit of time due to the contact of S susceptible people with infected people is $\beta I/N \cdot S$. This is the only nonlinear term of the equations. (Other nonlinearities are hidden in the time dependence of the parameters β , γ and μ .)

Equations (2)-(5) are interpreted as follows. The first equation is the rate of change of the susceptible group. The number of susceptible, S , changes faster in time if there are more infected people, I and more susceptible that can be infected, S . Clearly, the susceptible group is constant in time if the contact rate of the virus is zero, and/or if the number of infected is zero, and/or if the number of susceptible is zero. The second equation is the time derivative of the continuity equation. It expresses the fact that, in this epidemic model, the population N is assumed to be constant. The third equation is the rate of change of recovered people. The number of recovered is proportional to the number of infected, I , because a recovered person must have been infected. The fourth equation is the rate of change of the deceased group. The number of deaths is proportional to the number of infected, I , because a deceased individual must have been infected (in this model).

2.2 Data sources

The reliable data is about the number of confirmed infected, I_c , and confirmed deaths, D_c , which are arranged in a vector

$$\mathbf{q}_c \equiv [I_c, D_c]^T. \quad (10)$$

The data on the confirmed recovered, R_c , was discontinued because it was deemed inaccurate². The data used here is publicly available in the CSSEGISandData/COVID-19 GitHub repository³, which collects data from official sources and organizations.

2.3 Problem formulation

The calculation of the groups' dynamics and time-varying epidemic parameters is a constrained optimization problem:

Calculate

$$\mathbf{q}, \alpha \quad (11)$$

to minimize

$$E_d \equiv \lambda_1 \|I - I_c\|^2 + \lambda_2 \|D - D_c\|^2 \quad (12)$$

subject to

$$\text{an epidemic model.} \quad (13)$$

¹ $R_0 \equiv \beta/(\gamma + \mu)$. Because μ is sufficiently small, it will be neglected.

²https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data/csse_covid_19_time_series

³https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data/csse_covid_19_daily_reports

The epidemic model used in this paper is provided by Eqs. (6) and (7), however, more detailed models can be used. $\|\cdot\|$ is a norm, λ_1 and λ_2 are user-defined normalization factors. The loss function, E_d , measures the error between the candidate solution (I, D) and the data (I_c, D_c) . Among all the possible candidate solutions, only the solutions that fulfil the epidemic model (Eqs. (6) and (7)) will be accepted. The cumulative confirmed number of cases is the dataset used. This is a quantity to be preferred over the daily increase of confirmed cases because it is smoother, i.e., it is not significantly affected by random fluctuations, in contrast with the daily increase. The algorithm that solves this constrained optimization problem is presented in Sec. 2.4.

2.4 First-principles machine learning epidemic modelling

A data-driven model combined with first principles is proposed. This is referred to as *first-principles machine learning* for brevity. The data-driven algorithm is an optimal interpolator, while the epidemic model helps to obtain parameters that are consistent with the model. This synergistic combination helps to reduce the uncertainty in the predictions, which are as good as the employed epidemic model and the accuracy of the data.

The first-principles machine learning epidemic modelling is based on the combination of an ODE-solver, which time-advances the SIRD model in Eqs. (2)-(5) (first principles), and a feedforward neural network (machine learning), which performs the assimilation of data with the epidemic model to learn the parameters' vector $\alpha(t)$ (Fig. 2) and predict the state, $q(t)$. The Neural Network (NN) receives as an input the entire time series of total confirmed infected cases $\{I_c(t)\}_{t=0}^{N_t}$ and total confirmed deceased $\{D_c(t)\}_{t=0}^{N_t}$ up until the 17th of April 2020. The time $t = 0$ corresponds to the day when the first infection was recorded, and N_t is the number of days from $t = 0$ to the 17th of April 2020. From the time series, $\{I_c(t)\}_{t=0}^{N_t}$ and $\{D_c(t)\}_{t=0}^{N_t}$, the NN infers the time evolution of the parameters of the SIRD model, i.e. $\{\hat{\beta}(t)\}_{t=0}^{N_t}$, $\{\hat{\gamma}(t)\}_{t=0}^{N_t}$ and $\{\hat{\mu}(t)\}_{t=0}^{N_t}$, where $\hat{\cdot}$ denotes the quantity estimated by the neural network. Consistently with (9), the parameters $\hat{\beta}$, $\hat{\gamma}$, $\hat{\mu}$ are cast in the vector $\hat{\alpha} \equiv [\hat{\beta}, \hat{\gamma}, \hat{\mu}]^T$. Subsequently, $\{\hat{\alpha}\}$ is fed into the time-integration of the SIRD model with initial condition $q_0 = [N_0 - I_0 - D_0, I_0, 0, D_0]^T$ where N_0 is the population of the country analysed (Table 1), whereas I_0 and D_0 are the confirmed infections and deaths on the day of the first confirmed cases, respectively. Finally, the time-integration of the SIRD model provides the state \hat{q} .

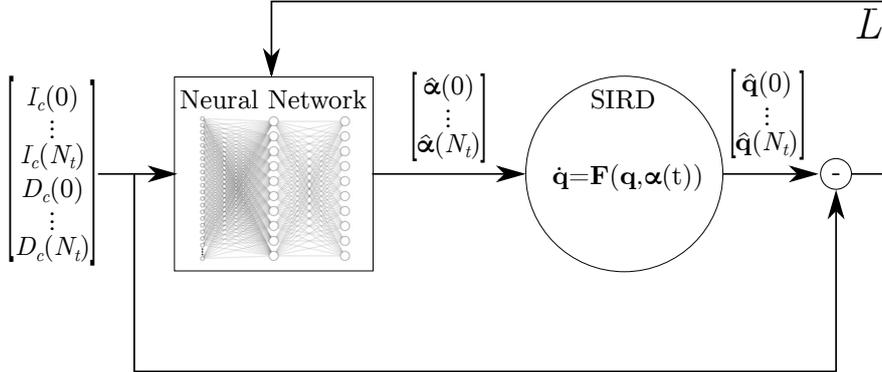


Figure 2: First-principles machine learning architecture for epidemic modelling. The graph of the neural network is pictorial.

Algorithmically, the architecture is trained as follows:

1. **First guess on the parameters.**
From the dataset $\{I_c(t)\}_{t=0}^{N_t}$ and $\{D_c(t)\}_{t=0}^{N_t}$, a set of constant parameters $\alpha_0 \equiv [\beta_0, \gamma_0, \mu_0]^T$ is obtained by nonlinear regression of the data, I_c and D_c , during the initial exponential growth only. This time window is $[0, t = \text{Regr}]$ (Table 1).

2. Initialization of the neural network.

The neural network is pre-trained to output the set of constant parameters, α_0 . This set of parameters ensures that the initial state of the neural network is consistent with the initial exponential growth, which makes the time integration of the SIRD model robust. Unless otherwise specified, the neural network consists of 1 layer with 8 neurons^a. The time evolution of the SIRD parameters is obtained by nonlinear combination of the neurons with a sigmoid activation.

3. Training of the neural network.

The entire architecture, which consists of the neural network and the SIRD time-integrator, is optimized by a gradient-based optimizer (L-BFGS-B optimizer [13]) to minimize the loss function

$$\begin{aligned}
 L = & \underbrace{\sum_{t=0}^{N_t} \left((\log(I_c(t)) - \log(\hat{I}(t)))^2 + (\log(D_c(t)) - \log(\hat{D}(t)))^2 \right)}_{E_{d1}} + \\
 & \underbrace{0.01 \frac{\log(\max(I_c))}{\max(I_c)} \sum_{t=0}^{N_t} \left((I_c(t) - \hat{I}(t))^2 + (D_c(t) - \hat{D}(t))^2 \right)}_{E_{d2}} + \\
 & \underbrace{100 \frac{\log(\max(I_c))}{\max(\alpha_0)} \sum_{t=0}^{N_t-1} \left((\hat{\beta}(t) - \hat{\beta}(t+1))^2 + (\hat{\gamma}(t) - \hat{\gamma}(t+1))^2 + 100(\hat{\mu}(t) - \hat{\mu}(t+1))^2 \right)}_{E_r} + \\
 & \underbrace{100 \frac{\log(\max(I_c))}{\max(\alpha_0)} \left((\hat{\beta}(0) - \beta_0)^2 + (\hat{\gamma}(0) - \gamma_0)^2 + 100(\hat{\mu}(0) - \mu_0)^2 \right)}_{E_0} \quad (14)
 \end{aligned}$$

The loss function is composed of four terms, which can be interpreted as follows:

- E_{d1} is the error in a log-scale between the prediction and the available data (infected and deaths). This removes noisy fluctuations from the solution.
- E_{d2} is the error in a linear scale between the prediction and the available data (infected and deaths).
- E_r is a regularization term, which prevents large discontinuities in the time-variation of the SIRD parameters from occurring, making the evolution smoother. The regularization factor before the sum in E_r is an empirical scaling factor to ensure that the orders of magnitude of E_r and E_d are comparable. The factor 100 before the terms with $\hat{\mu}$ ensures that the parameters of the SIRD model have a comparable order of magnitude.
- E_0 constrains the initial values of $\hat{\alpha}$ to be close to the first guess obtained at step 1. This ensures that, in the early stage of the epidemic, the growth is largely exponential with parameters that are nearly constant. A typical convergence of the optimizer is shown in Fig. 3.

^aArchitectures with 4 to 64 neurons provide the same accuracy (result not shown).

3 Results

	NY	Italy	Germany	UK	Spain	USA	France	China	Belgium	World
N_0 [Mil]	5.8	60.36	83.02	66.56	46.94	327.2	66.99	1386	11.46	7777.06
Regr	20	40	60	60	50	70	60	15	55	20

Table 1: Country populations, N_0 (from Census databases on Google) and time (days) of initial exponential growth, Regr.

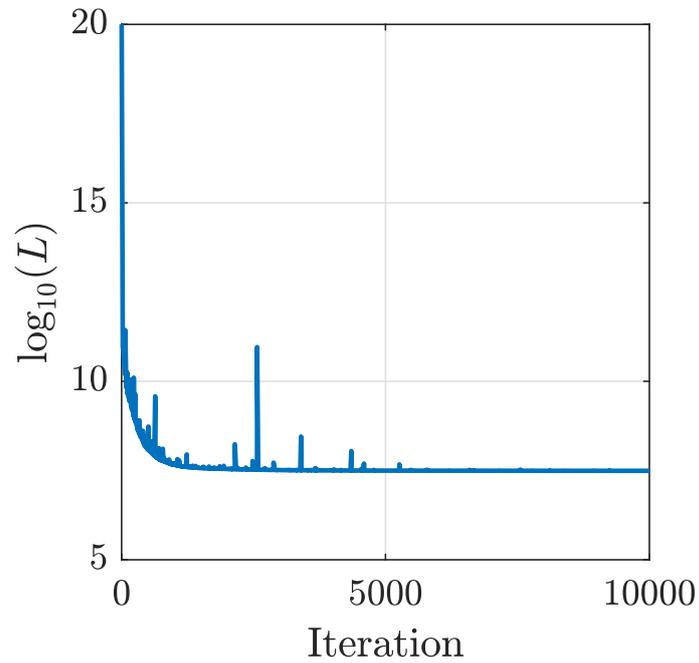


Figure 3: Typical evolution of the loss function during the o-eps-converted-to.pdf optimization process (step 3). World data.

3.1 United Kingdom

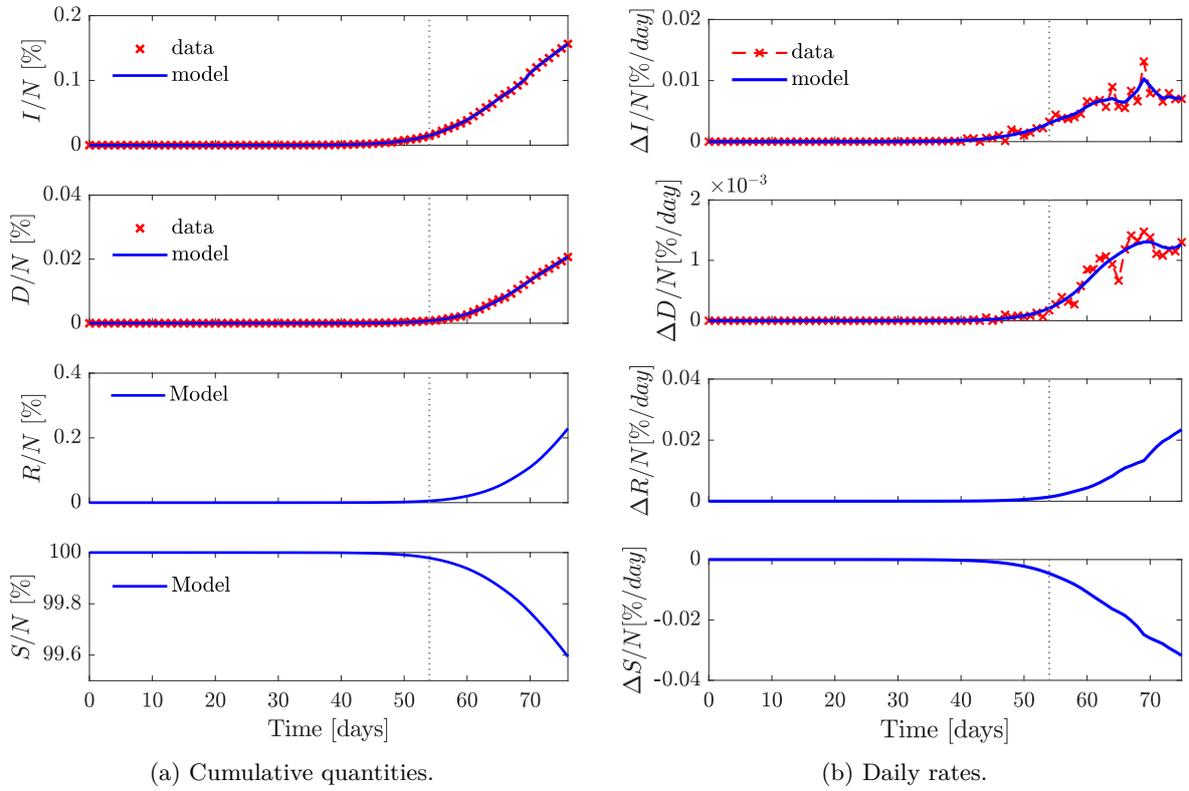
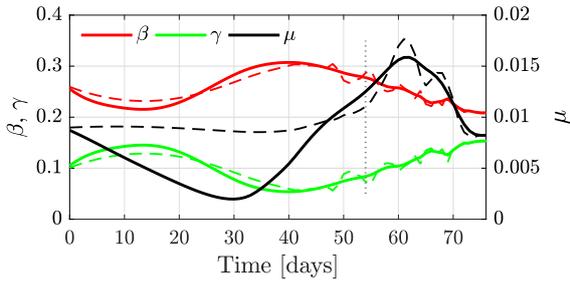
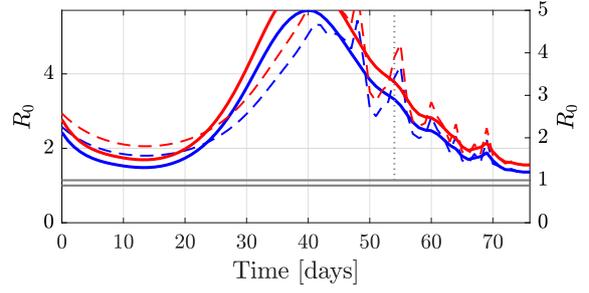


Figure 4: United Kingdom (day 0 = 31st January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

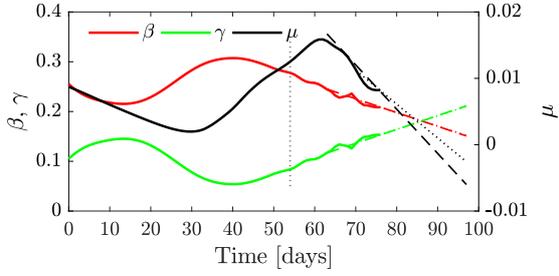


(a) Time-varying contact rate (β), recovery rate (γ), and death rate (μ).

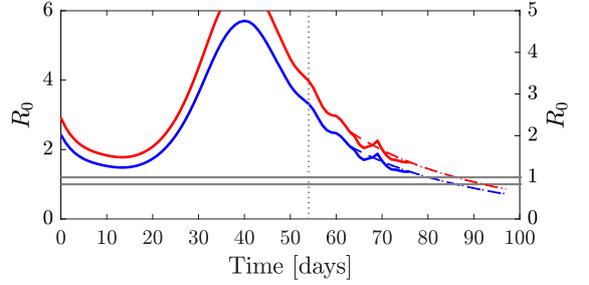


(b) Basic reproduction number. The blue (red) curve corresponds to the left (right) vertical axis.

Figure 5: United Kingdom (day 0 = 31st January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.



(a) Extrapolated trends of the time-varying contact rate (β), recovery rate (γ), and death rate (μ) with average slope over the last seven days (dotted lines) and fourteen days (dashed lines).



(b) Extrapolated trend of the basic reproduction number with average slope over the last seven days (dotted lines) and fourteen days (dashed lines).

Figure 6: United Kingdom (day 0 = 31st January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

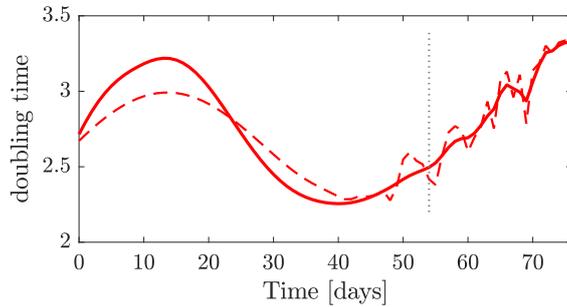


Figure 7: United Kingdom (day 0 = 31st January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

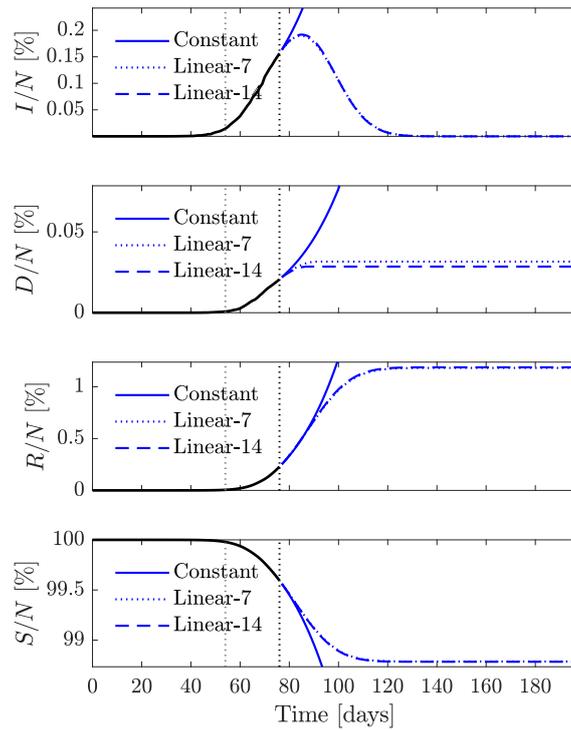


Figure 8: United Kingdom (day 0 = 31st January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 4a.

3.2 Italy

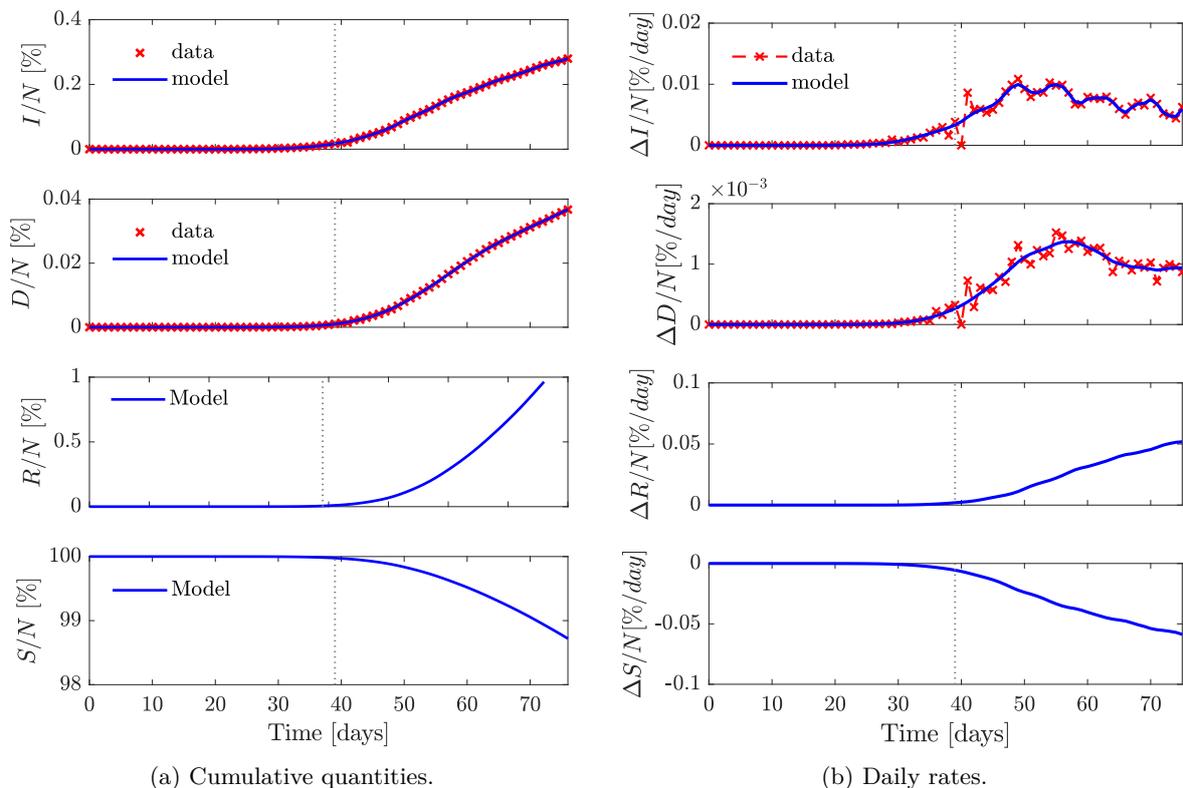


Figure 9: Italy (day 0 = 31st January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

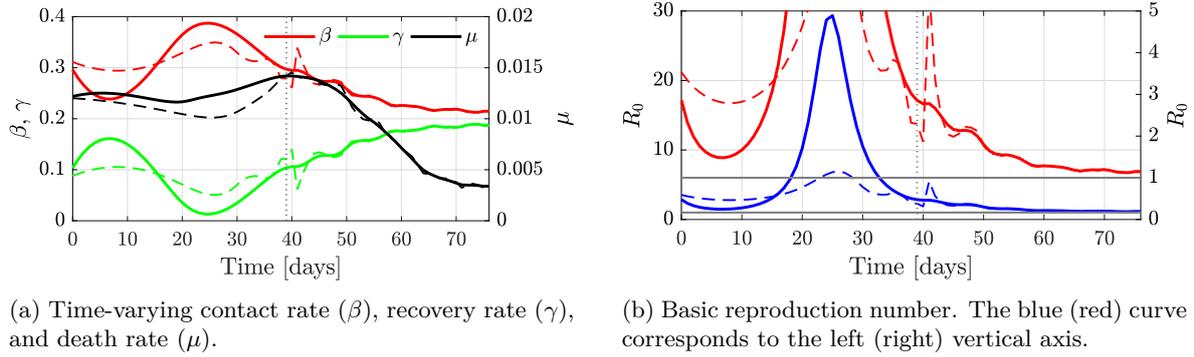


Figure 10: Italy (day 0 = 31st January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

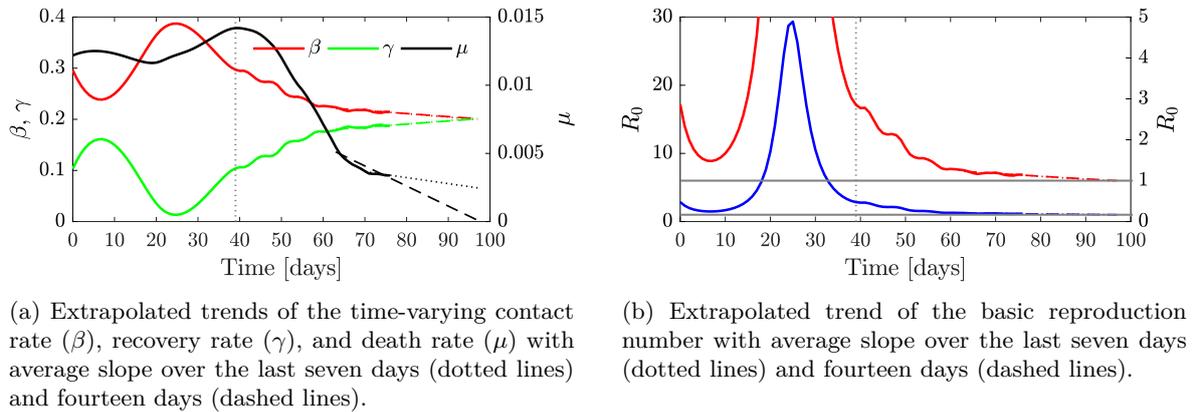


Figure 11: Italy (day 0 = 31st January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

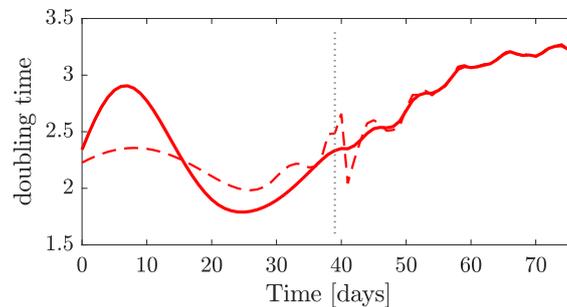


Figure 12: Italy (day 0 = 31st January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

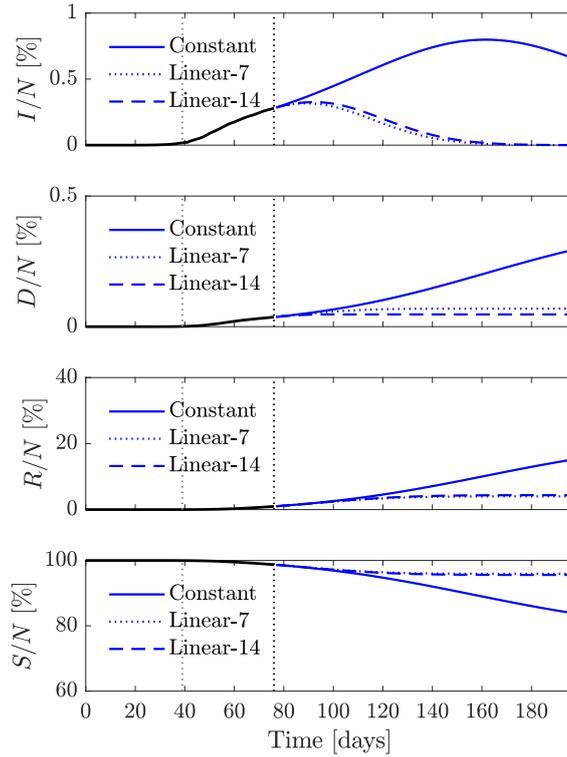


Figure 13: Italy (day 0 = 31st January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 9a.

3.3 Germany

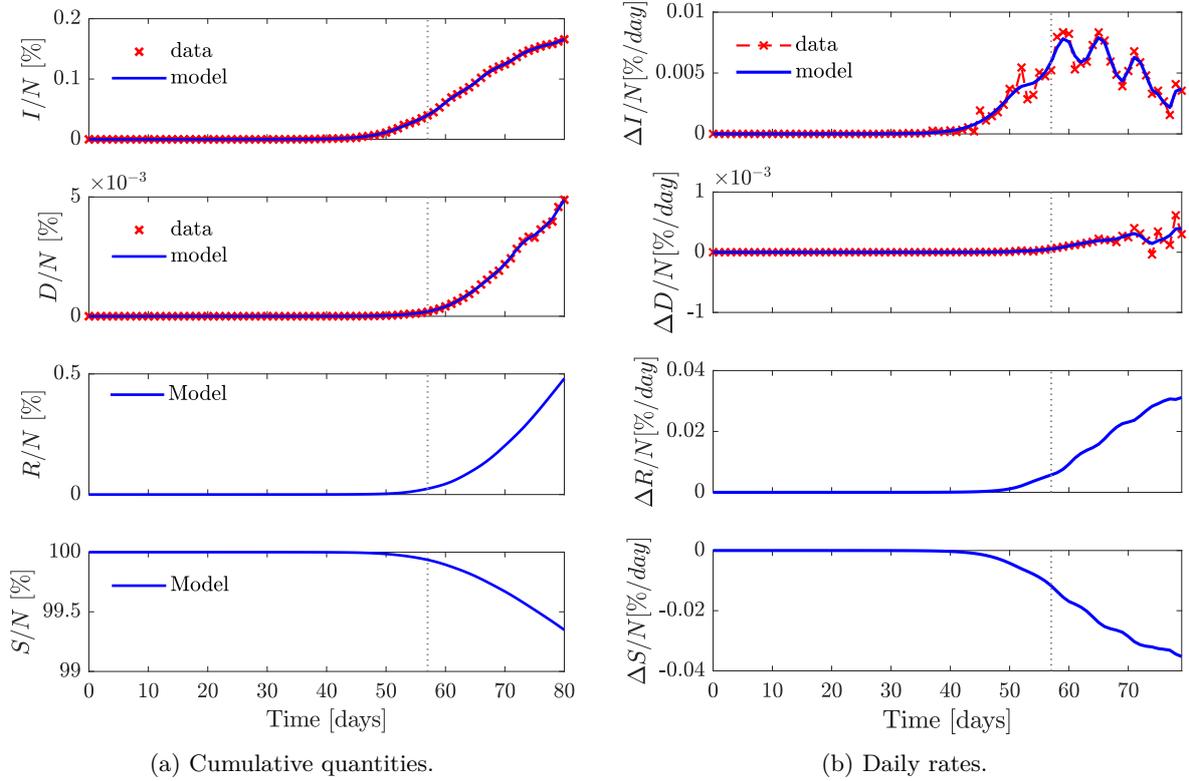


Figure 14: Germany (day 0 = 27th January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

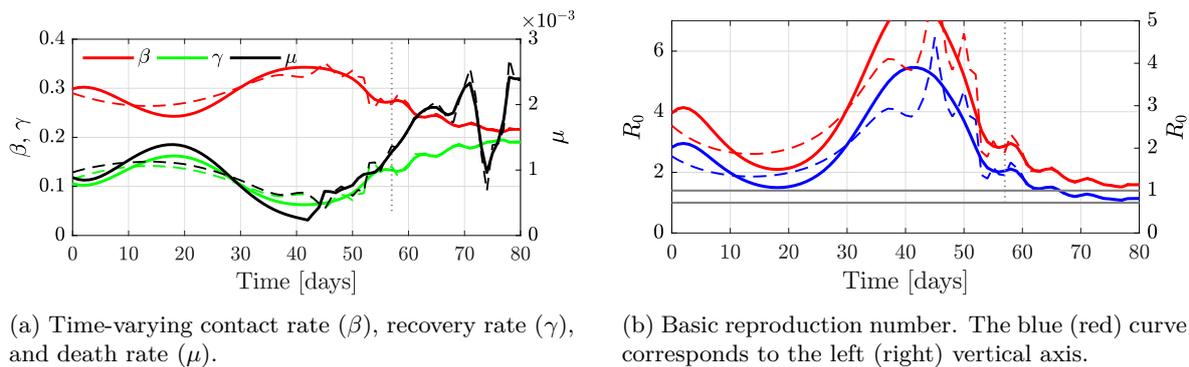


Figure 15: Germany (day 0 = 27th January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

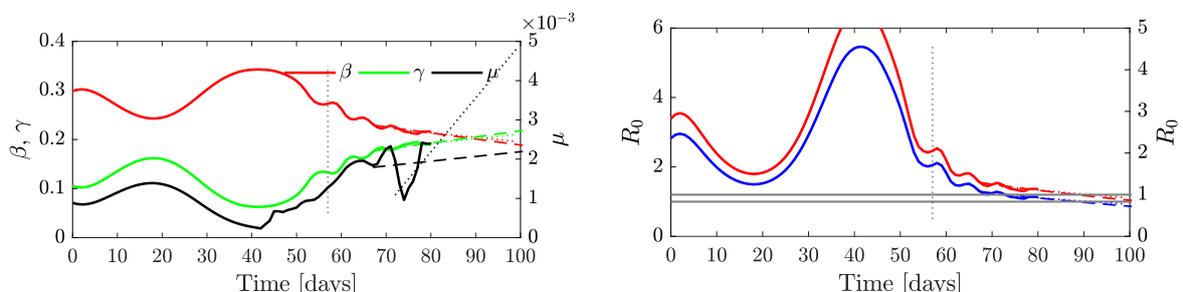


Figure 16: Germany (day 0 = 27th January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

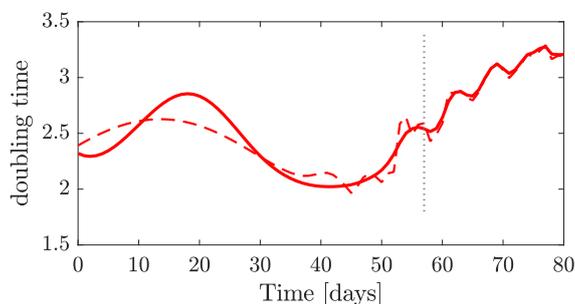


Figure 17: Germany (day 0 = 27th January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

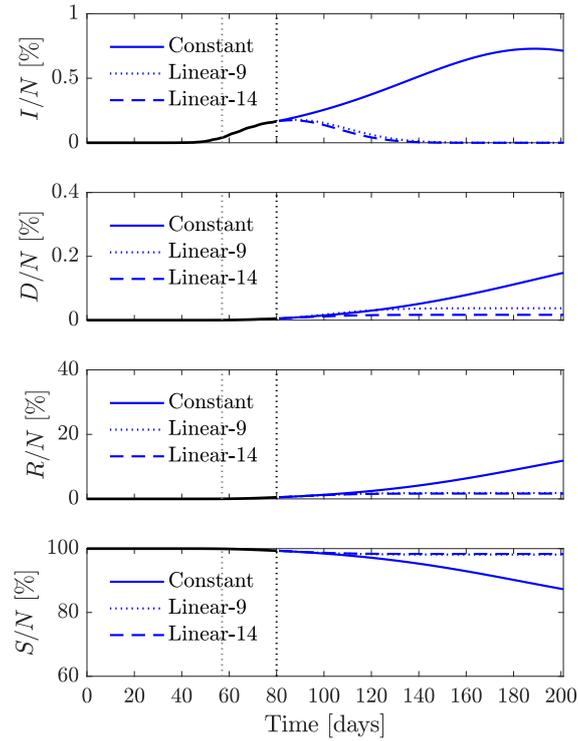


Figure 18: Germany (day 0 = 27th January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last nine days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 14a.

3.4 France

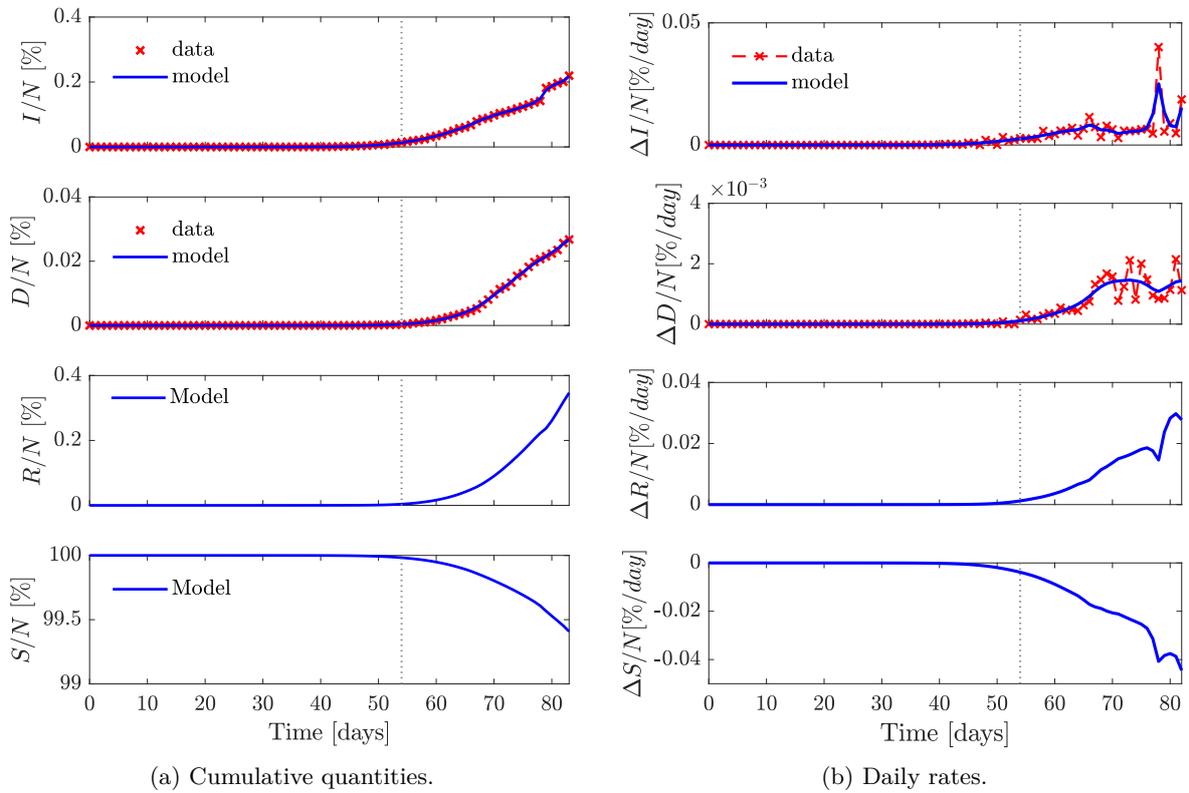


Figure 19: France (day 0 = 24th January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

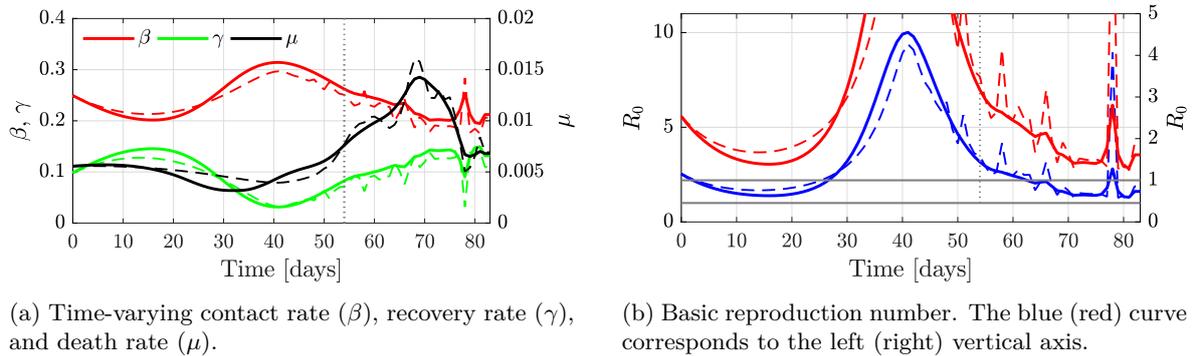


Figure 20: France (day 0 = 24th January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

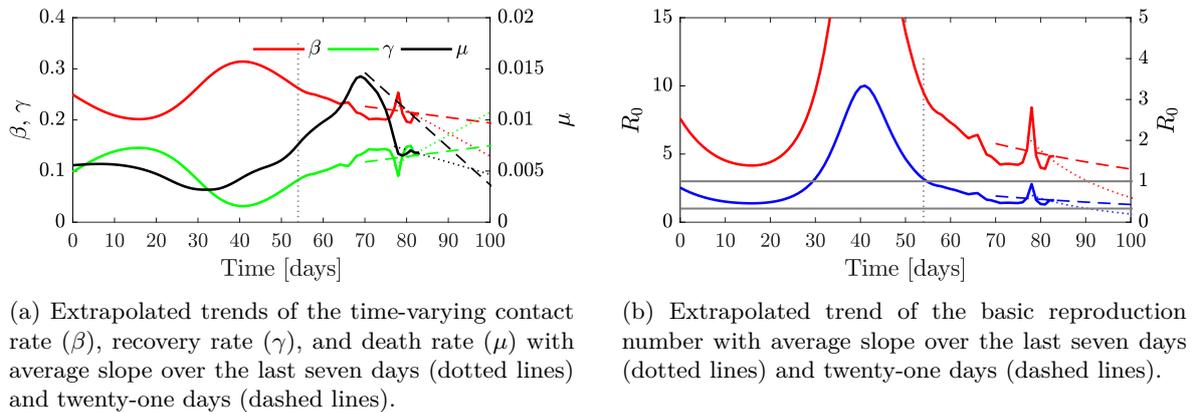


Figure 21: France (day 0 = 24th January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

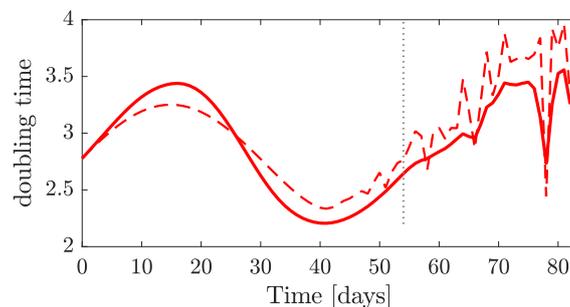


Figure 22: France (day 0 = 24th January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

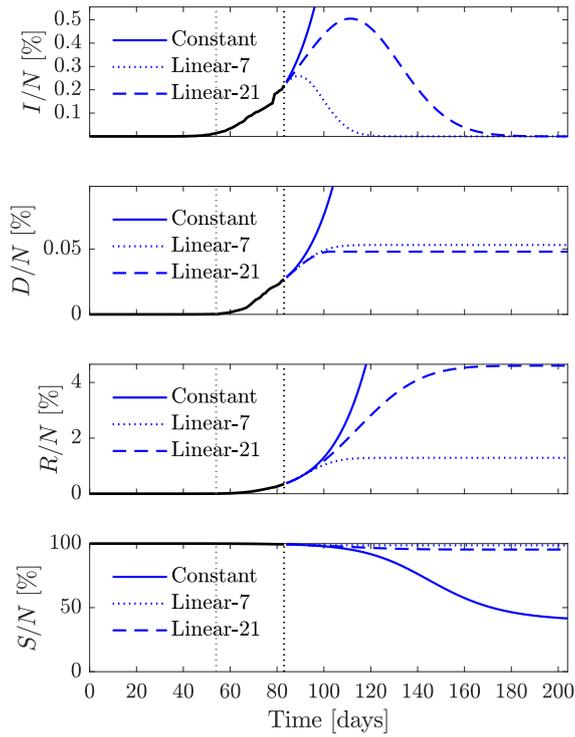


Figure 23: France (day 0 = 24th January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), twenty-one days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 19a.

3.5 Spain

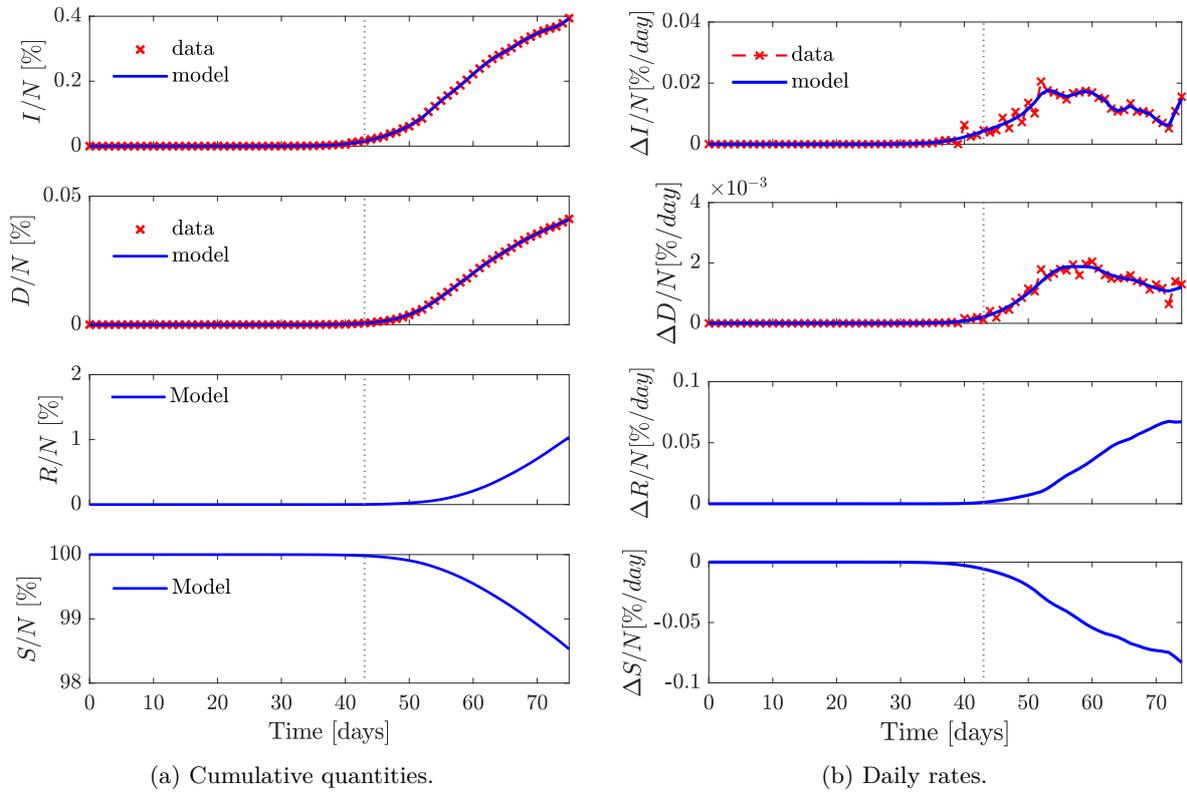


Figure 24: Spain (day 0 = 1st February 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

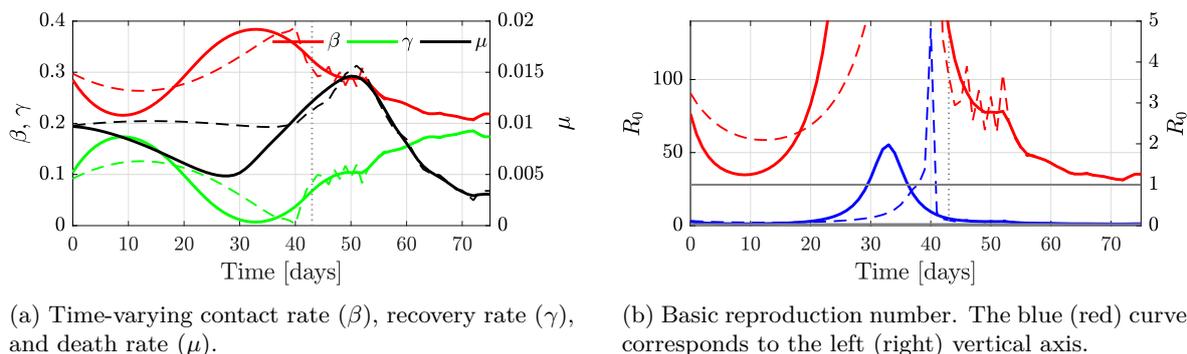


Figure 25: Spain (day 0 = 1st February 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

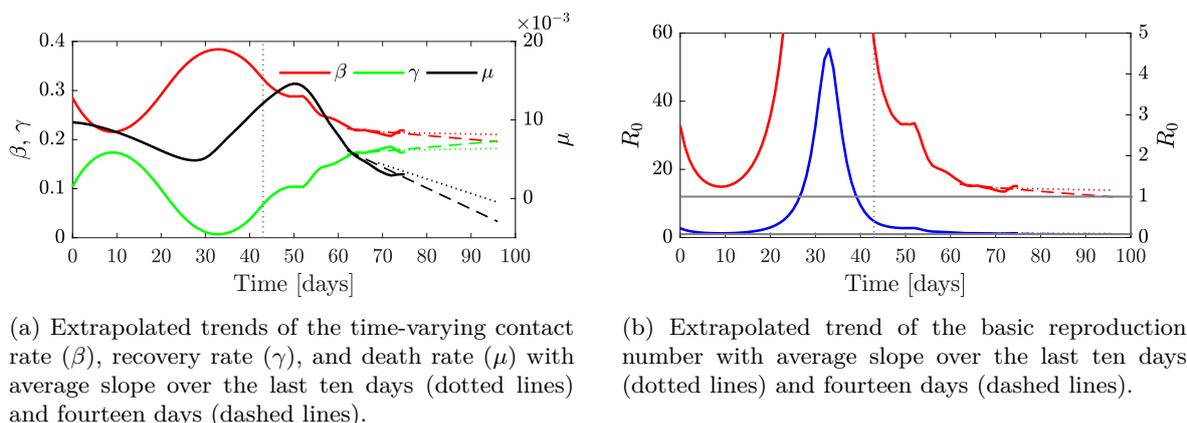


Figure 26: Spain (day 0 = 1st February 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

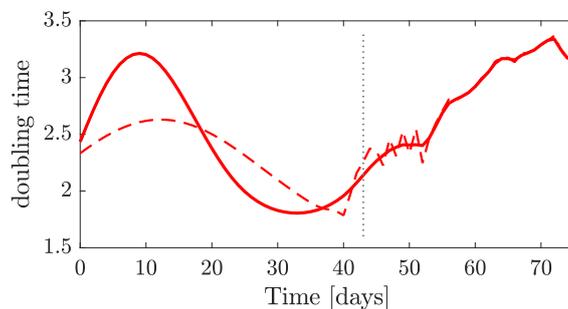


Figure 27: Spain (day 0 = 1st February 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

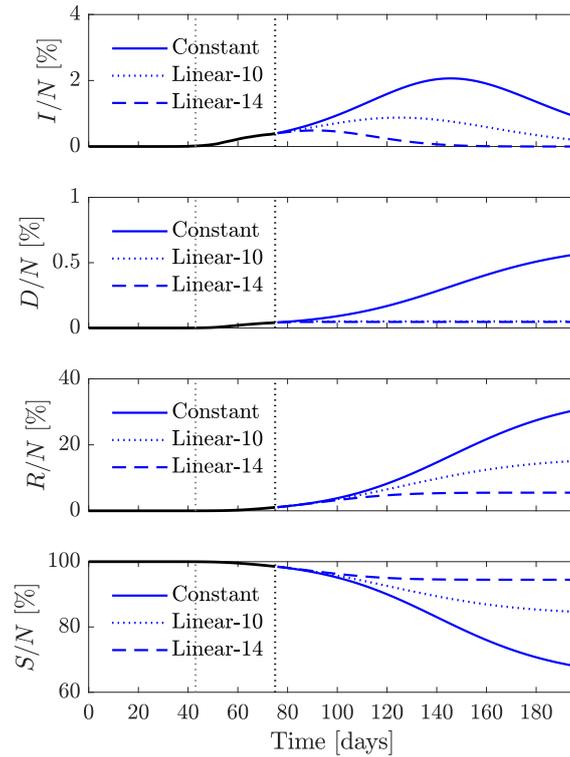


Figure 28: Spain (day 0 = 1st February 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last ten days (dotted lines), two days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 24a.

3.6 Belgium

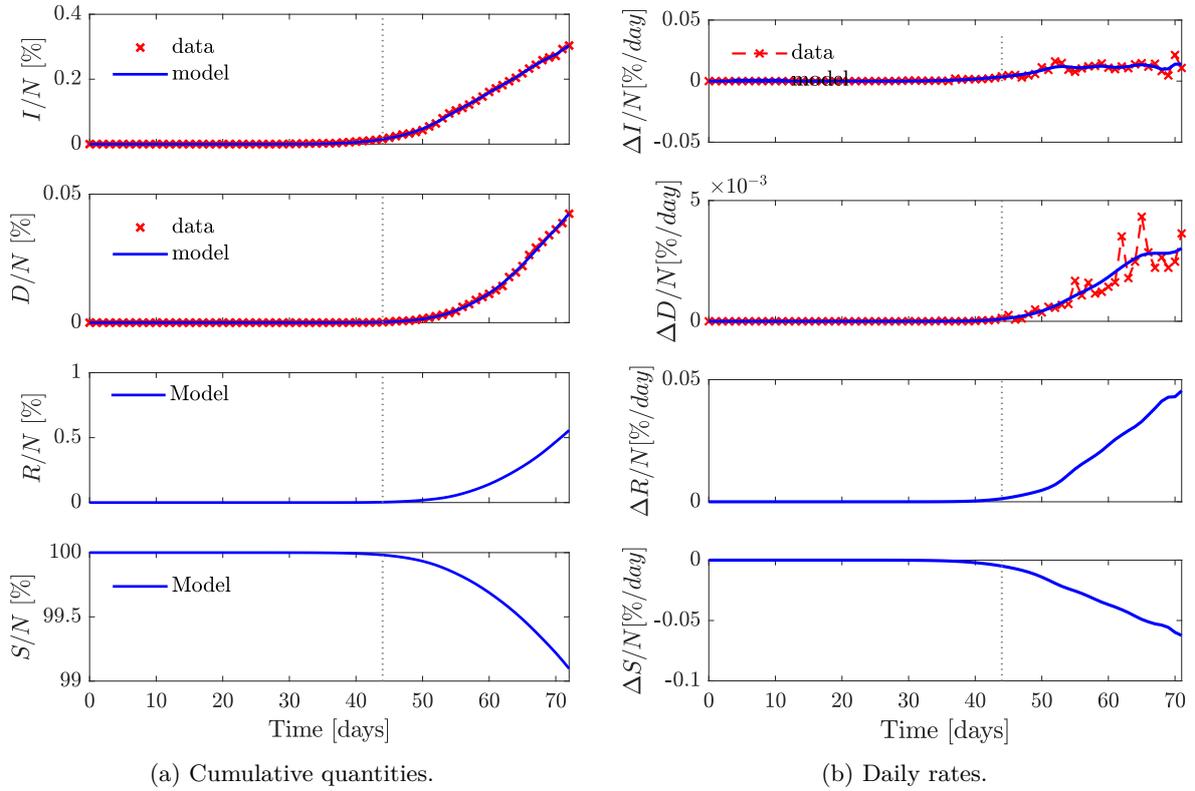


Figure 29: Belgium (day 0 = 4th February 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

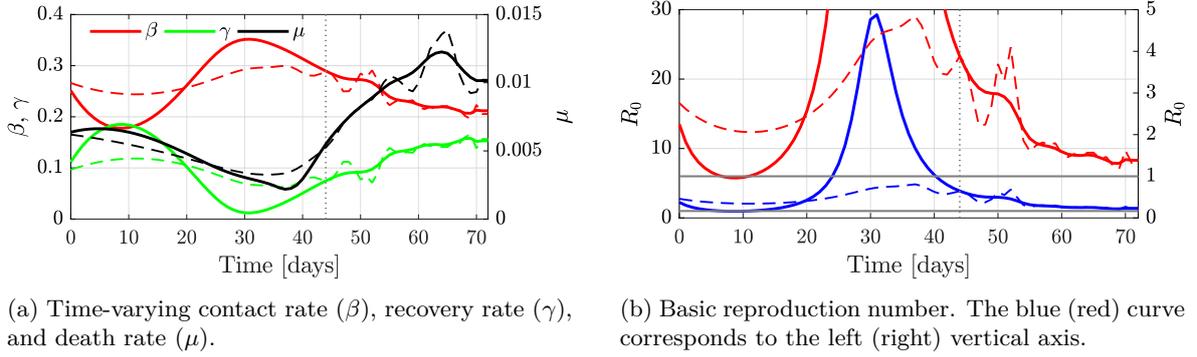


Figure 30: Belgium (day 0 = 4th February 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

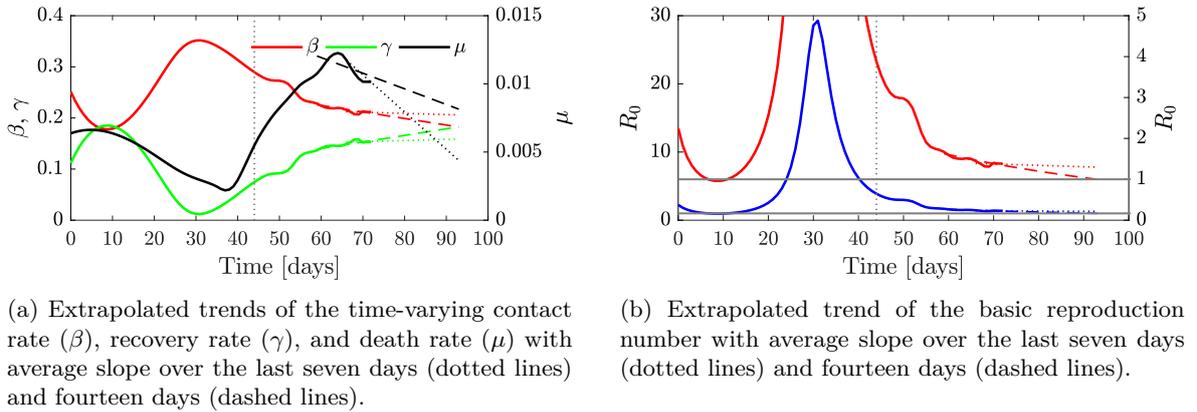


Figure 31: Belgium (day 0 = 4th February 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

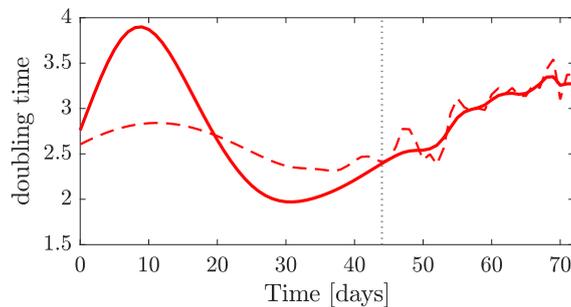


Figure 32: Belgium (day 0 = 4th February 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

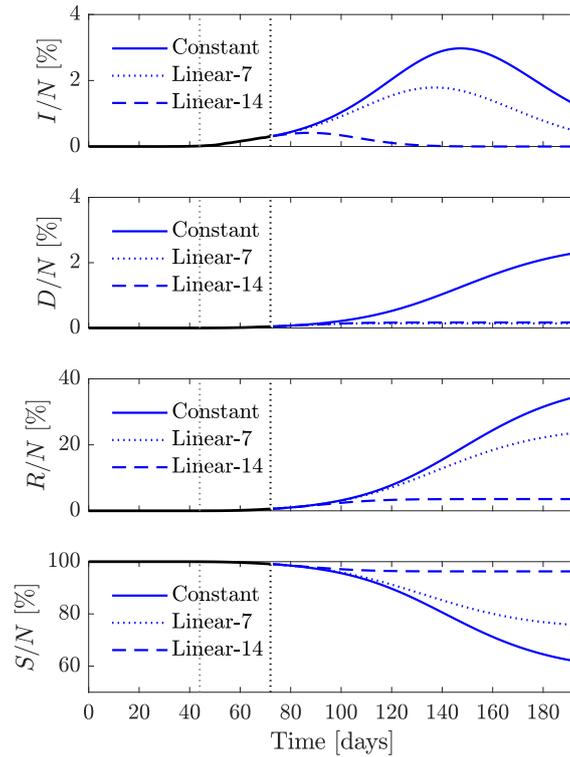


Figure 33: Belgium (day 0 = 4th February 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 29a.

3.7 USA

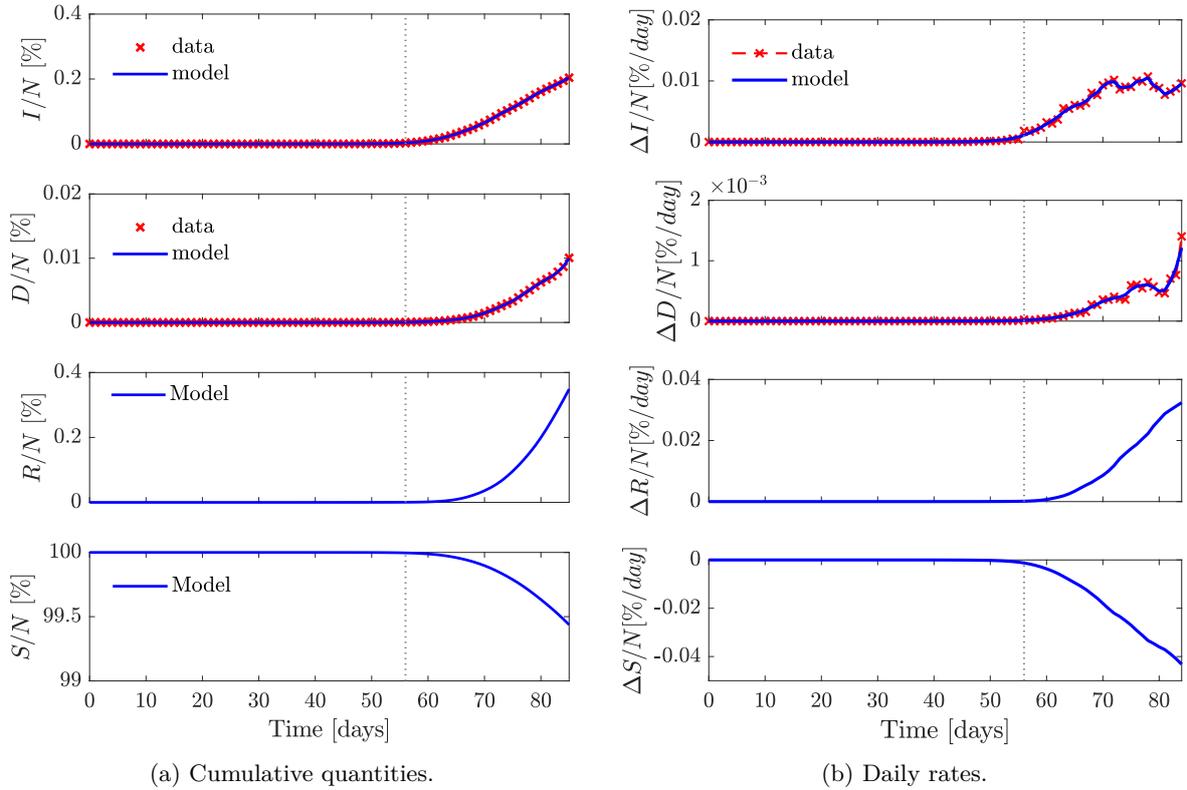
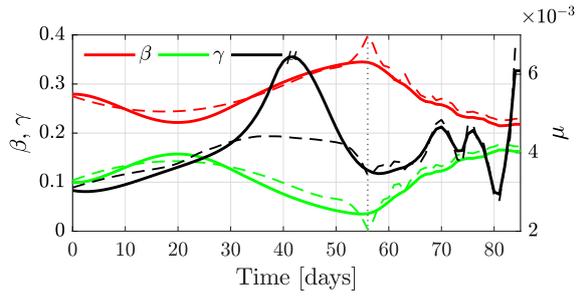
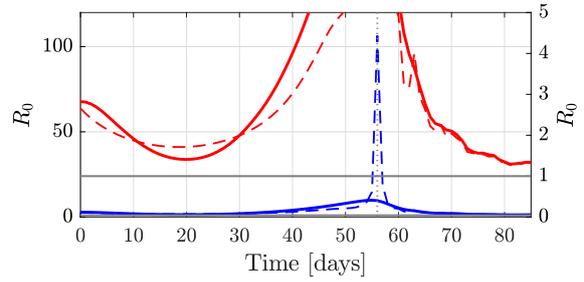


Figure 34: USA (day 0 = 22nd January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

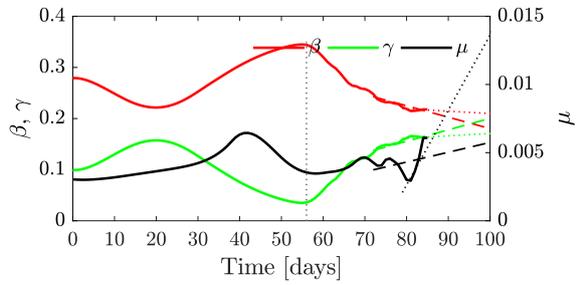


(a) Time-varying contact rate (β), recovery rate (γ), and death rate (μ).

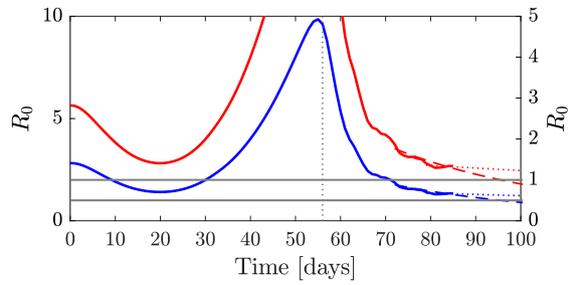


(b) Basic reproduction number. The blue (red) curve corresponds to the left (right) vertical axis.

Figure 35: USA (day 0 = 22nd January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.



(a) Extrapolated trends of the time-varying contact rate (β), recovery rate (γ), and death rate (μ) with average slope over the last seven days (dotted lines) and fourteen days (dashed lines). The positive slope of the death rate is a consequence of an anomaly in the data on confirmed deaths. The cause of the anomaly is not known to the authors.



(b) Extrapolated trend of the basic reproduction number with average slope over the last seven days (dotted lines) and fourteen days (dashed lines).

Figure 36: USA (day 0 = 22nd January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

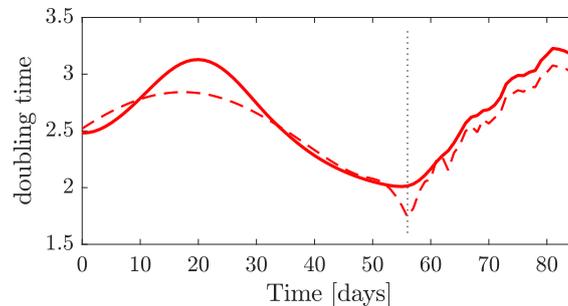


Figure 37: USA (day 0 = 22nd January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

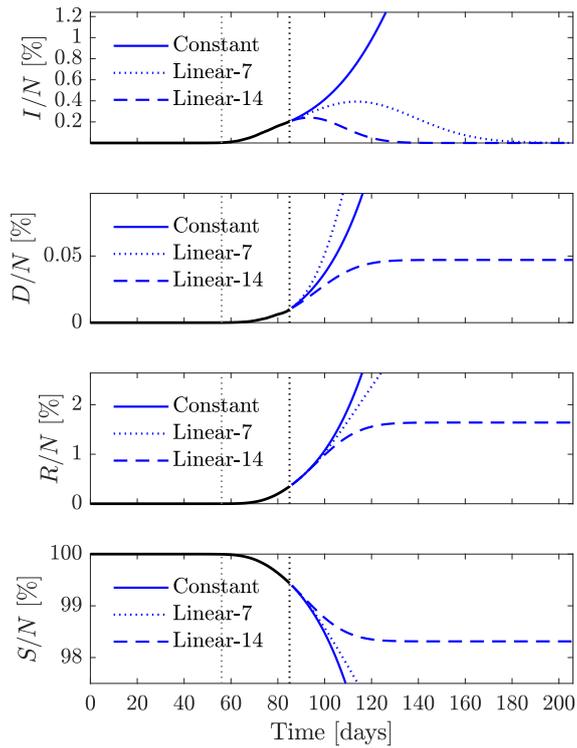


Figure 38: USA (day 0 = 22nd January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 34a.

3.8 New York City

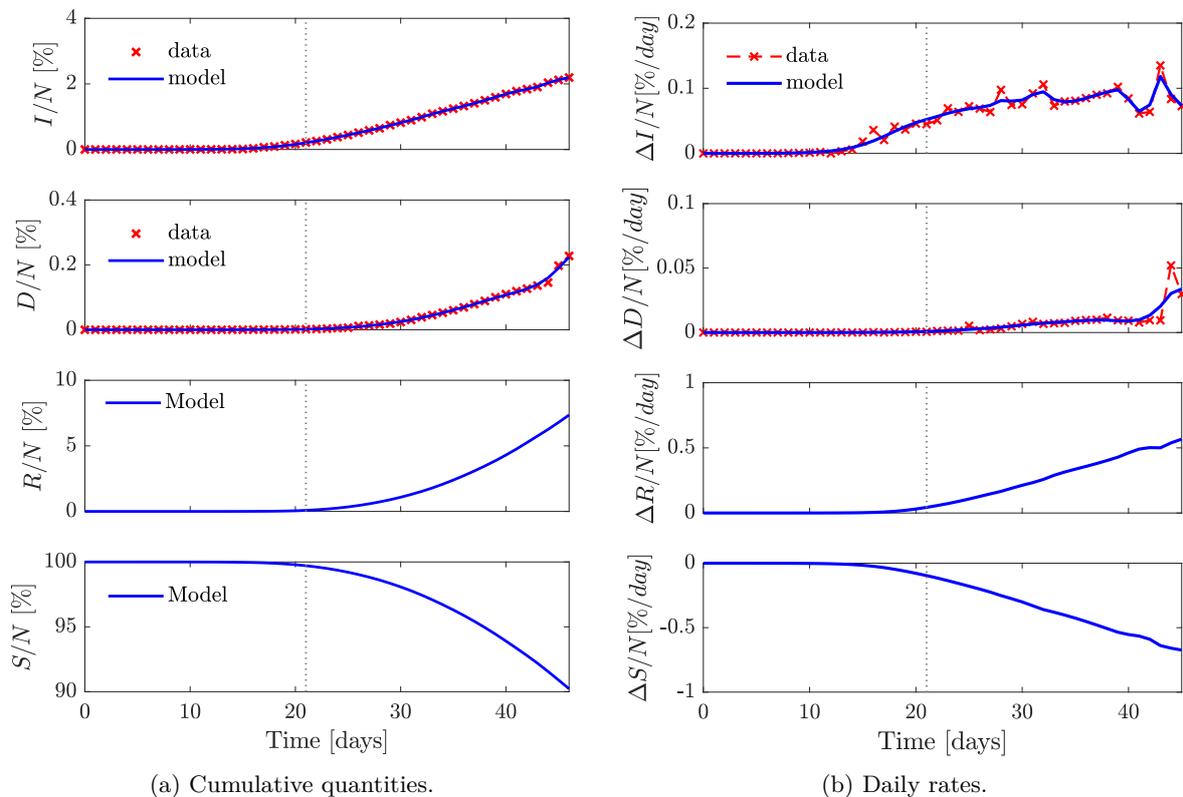


Figure 39: New York City (day 0 = 2nd March 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

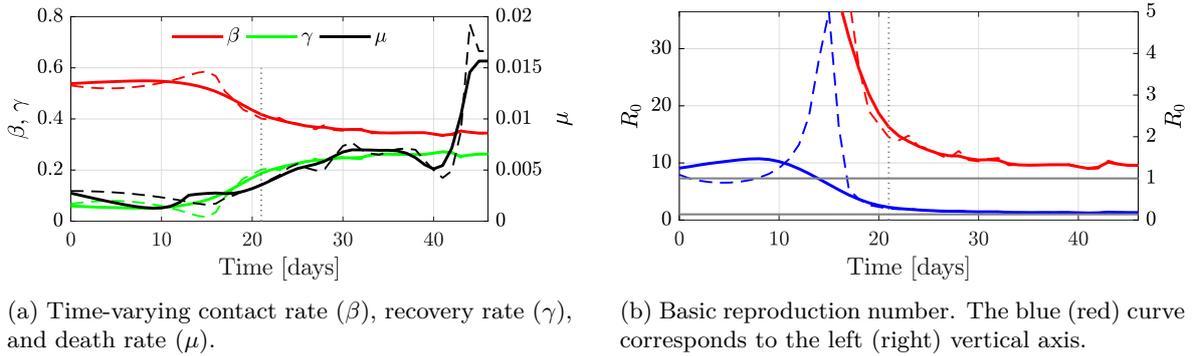


Figure 40: New York City (day 0 = 2nd March 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

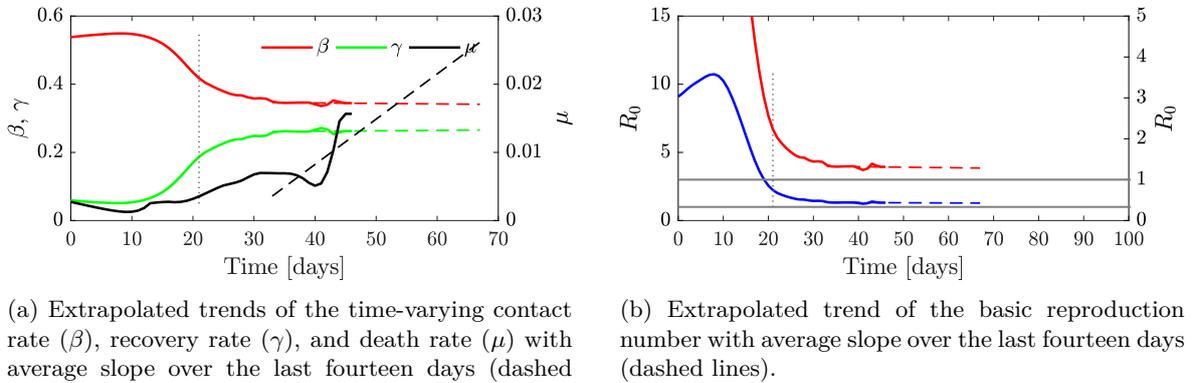


Figure 41: New York City (day 0 = 2nd March 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

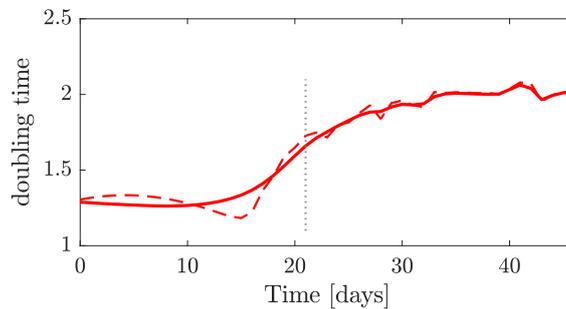


Figure 42: New York City (day 0 = 2nd March 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

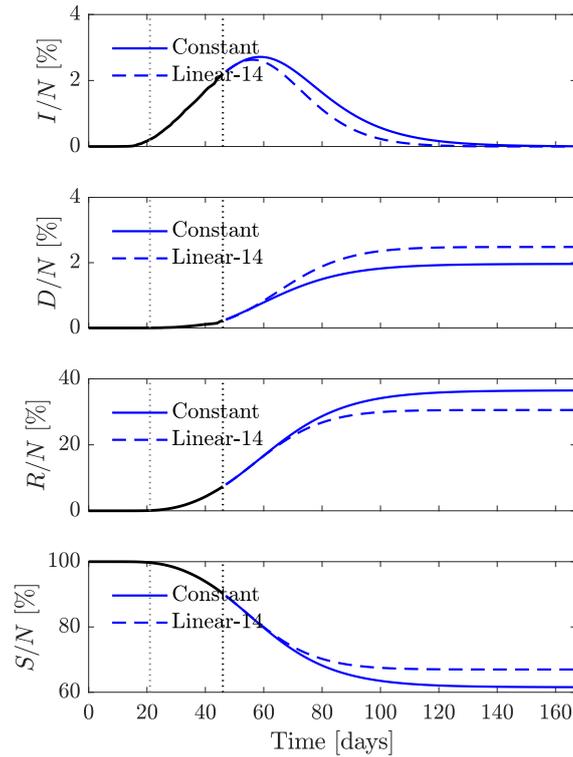


Figure 43: New York City (day 0 = 2nd March 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 39a.

3.9 China

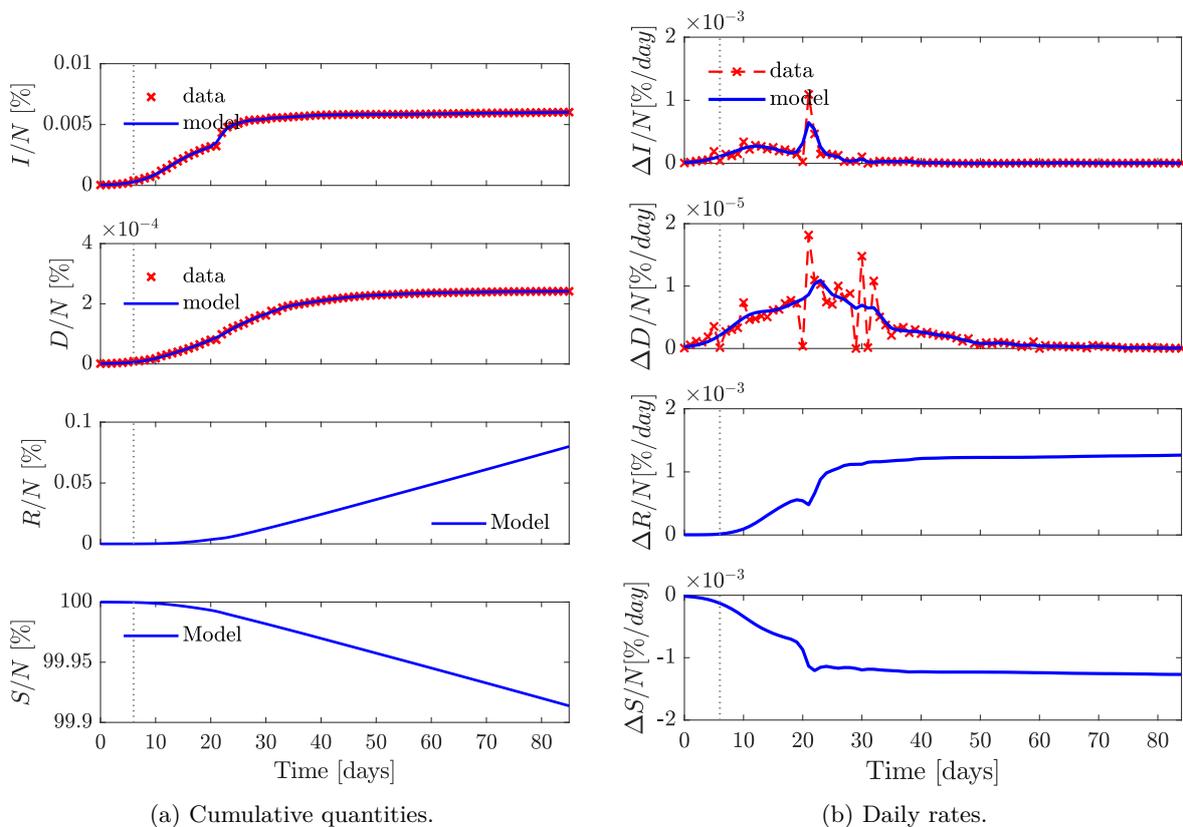


Figure 44: China (day 0 = 22nd January 2020): First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible. The vertical dotted lines indicate the day of lockdown.

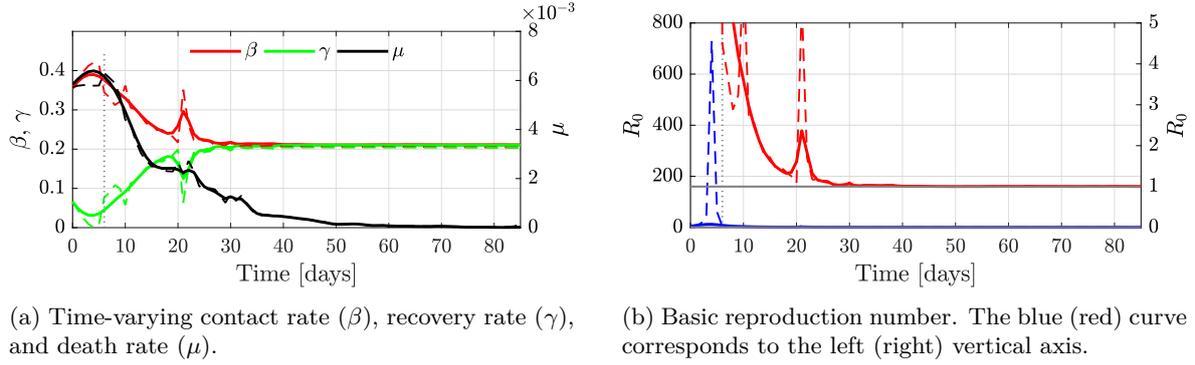


Figure 45: China (day 0 = 22nd January 2020): SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14). The vertical dotted lines indicate the day of lockdown.

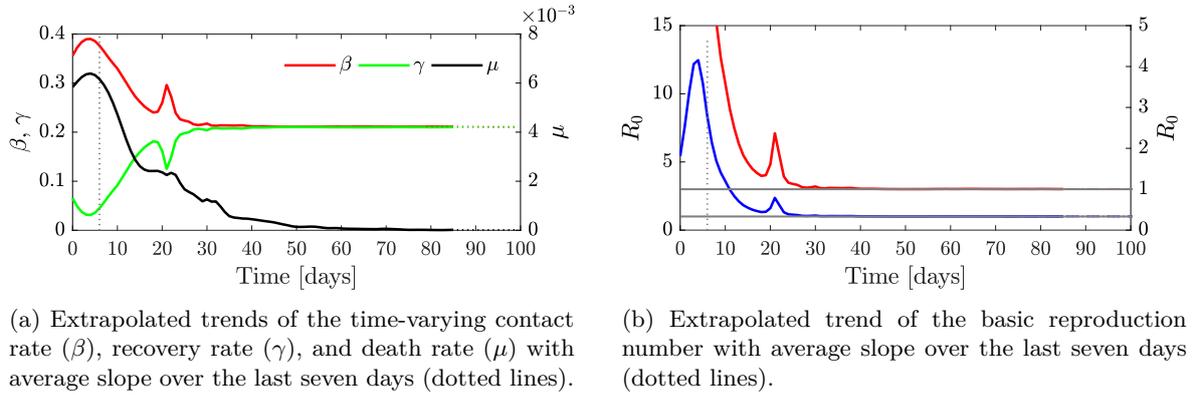


Figure 46: China (day 0 = 22nd January 2020): Extrapolated trends of the SIRD parameters. The vertical dotted lines indicate the day of lockdown.

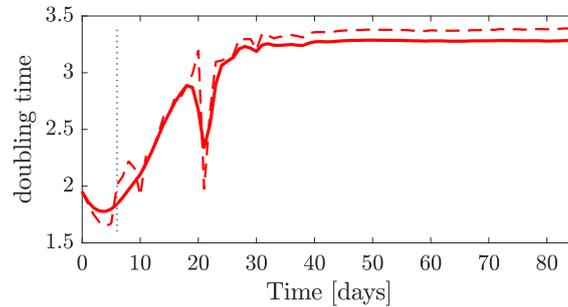


Figure 47: China (day 0 = 22nd January 2020): Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

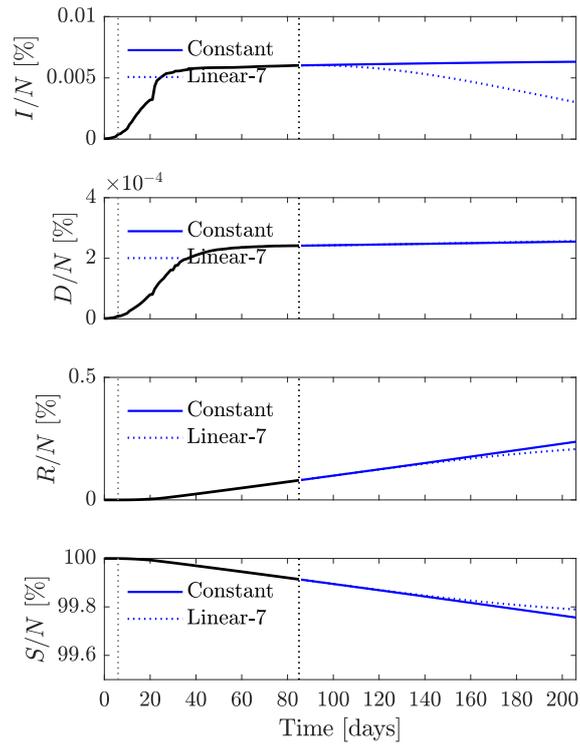


Figure 48: China (day 0 = 22nd January 2020): From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines) and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The left vertical dotted line represents the day of lockdown, the right vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 44a.

3.10 World

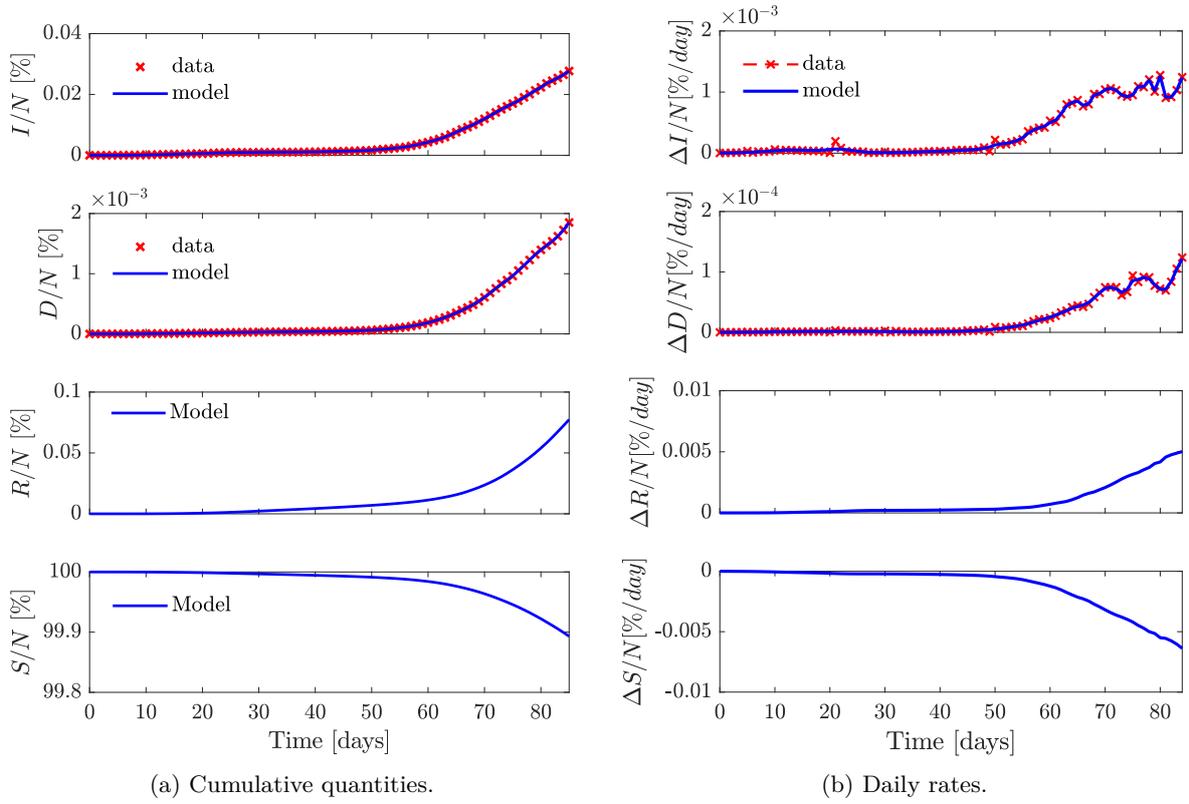


Figure 49: World (day 0 = 22nd January 2020: First and second rows: Validation of first-principles machine learning epidemic modelling. Third and fourth rows: Inference of recovered and susceptible.

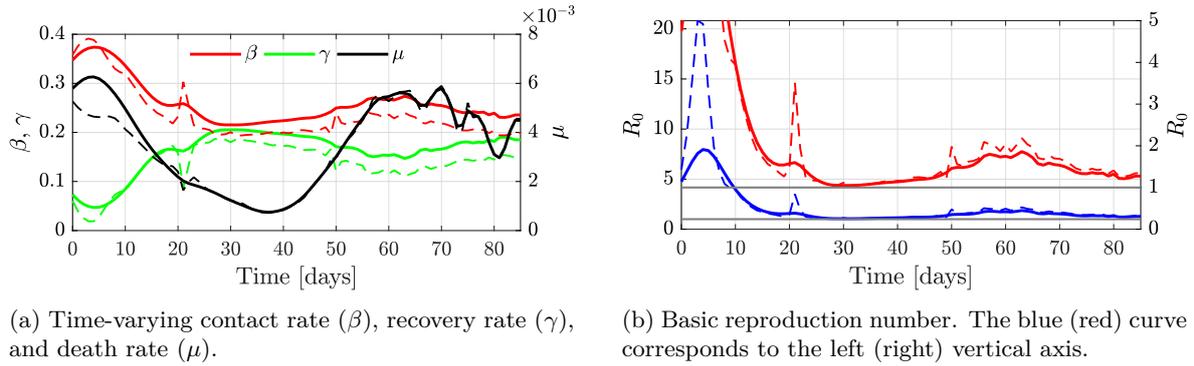


Figure 50: World (day 0 = 22nd January 2020: SIRD parameters. Neural network trained with the log (solid line) and without the log (dashed line) in the loss function (14).

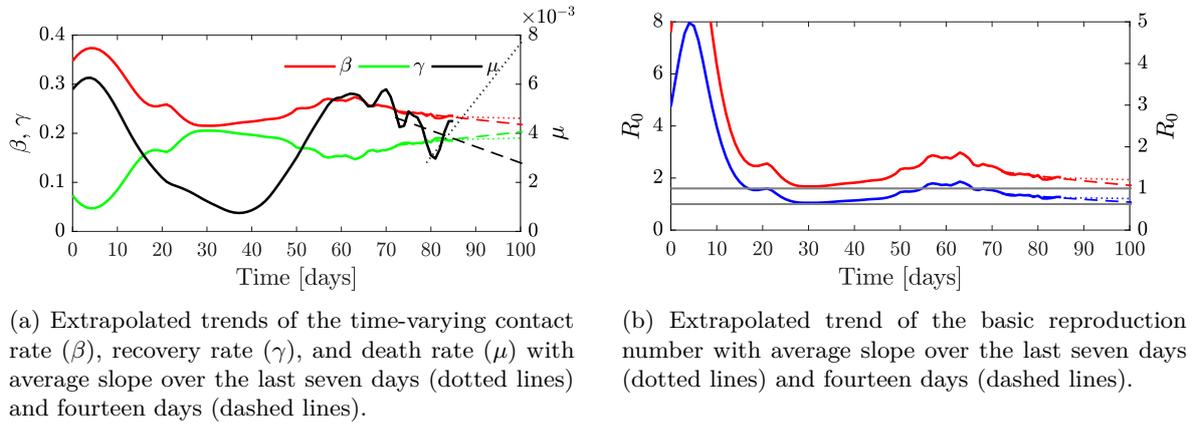


Figure 51: World (day 0 = 22nd January 2020: Extrapolated trends of the SIRD parameters.

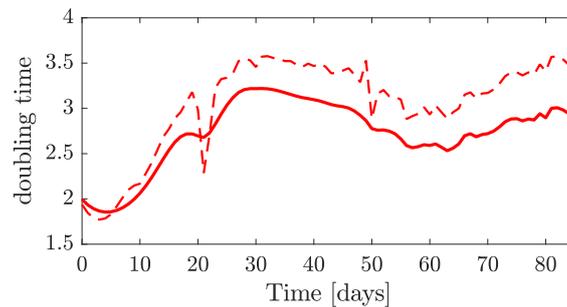


Figure 52: World (day 0 = 22nd January 2020: Doubling time with the log (solid line) and without the log (dashed line) in the loss function (14). The doubling time is calculated as $t = \log(2)/\beta(t)$. (To take into account the time derivative of $\beta(t)$, semi-parametric methods, e.g. [7], can be used.)

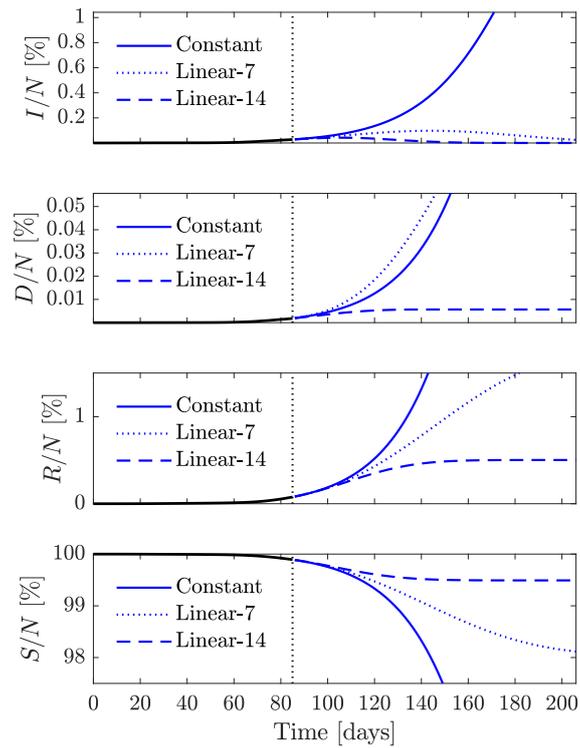


Figure 53: World (day 0 = 22nd January 2020: From the top: Blue lines indicate the extrapolated trends of the percentage of infected, recovered, deaths, and susceptible. Estimates with average slope over the last seven days (dotted lines), fourteen days (dashed lines), and with values of the parameters assumed to be constant and equal to the last day (solid lines). Black lines: The vertical line is the last day of the training data set, hence, the starting day for extrapolation. The black solid lines are taken from Fig. 49a.

Country	I_{max} [%]		D_{max} [%]		$R_0 = 1$
United Kingdom	2.8415	18/07/2020	2.1331	10/02/2021	-
	0.18959	26/04/2020	0.031601	02/05/2020	25/04/2020
	0.19187	26/04/2020	0.028537	29/04/2020	25/04/2020
Italy	0.79865	11/07/2020	0.41861	10/02/2021	-
	0.31652	30/04/2020	0.068895	02/07/2020	04/05/2020
	0.32581	03/05/2020	0.046895	08/05/2020	06/05/2020
Germany	0.72874	04/08/2020	0.26889	10/02/2021	-
	0.17774	25/04/2020	0.037317	01/11/2020	26/04/2020
	0.17436	23/04/2020	0.016698	08/10/2020	23/04/2020
France	6.9024	24/06/2020	2.9703	10/02/2021	-
	0.25882	23/04/2020	0.053117	12/06/2020	22/04/2020
	0.50498	15/05/2020	0.047987	07/05/2020	16/05/2020
Spain	2.064	27/06/2020	0.64766	10/02/2021	-
	0.87499	06/06/2020	0.050432	02/05/2020	21/07/2020
	0.48603	03/05/2020	0.046109	26/04/2020	05/05/2020
Belgium	2.9745	01/07/2020	2.6482	10/02/2021	-
	1.7878	21/06/2020	0.13973	24/05/2020	19/08/2020
	0.41385	02/05/2020	0.16903	17/07/2020	06/05/2020
USA	2.8442	08/07/2020	1.5018	10/02/2021	-
	0.3925	16/05/2020	0.54629	22/12/2020	12/06/2020
	0.24041	25/04/2020	0.04707	27/08/2020	26/04/2020
New York City	2.7161	30/04/2020	1.9625	10/02/2021	-
	-	-	-	-	-
	2.6298	27/04/2020	2.4807	16/12/2020	21/02/2021
China	0.0063773	06/11/2020	0.00027578	10/02/2021	-
	0.0060284	23/04/2020	0.0002657	10/02/2021	23/04/2020
	0.0089155	10/02/2021	0.0002416	22/04/2020	-
World	2.053	20/08/2020	0.86118	10/02/2021	-
	0.097413	15/06/2020	0.16569	10/02/2021	29/07/2020
	0.042022	08/05/2020	0.005671	13/06/2020	08/05/2020

Table 2: First column: Countries analysed. Second and third columns: Estimated maximum percentages of infected (I_{max}) and deaths (D_{max}) with dates. Fourth column: Estimate date on which the basic reproduction number becomes unity. For each country, the first / second / third row reports the estimate based on the extrapolation with constant parameters / linear parameters with average slope over a short window / linear parameters with average slope over a long window. The three different extrapolations provide an estimate of the range where the actual value lies. The results are consistent with the first principles and working assumptions of the SIRD model and the data (Sec. 2).

References

- [1] Novel Coronavirus (COVID-19) Cases, provided by John Hopkins University CSSE, <https://github.com/CSSEGISandData/COVID-19>, 2020.
- [2] <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>.
- [3] <https://www.who.int/docs/default-source/coronaviruse/who-china-joint-mission-on-covid-19-final-report.pdf>.
- [4] <https://www.gazzettaufficiale.it/eli/id/2020/03/09/20A01558/sg>.
- [5] https://administracion.gob.es/pag_Home/atencionCiudadana/Estado-de-alarma-crisis-sanitaria.html#.Xn3xj0dKjIU.
- [6] <https://www.bbc.co.uk/news/uk-52014472>.
- [7] Lorenzo Pellis, Francesca Scarabel, Helena B. Stage, Christopher E. Overton, Lauren H. K. Chappell, Katrina A. Lythgoe, Elizabeth Fearon, Emma Bennett, Jacob Curran-Sebastian, Rajenki Das, Martyn Fyles, Hugo Lewkowicz, Xiaoxi Pang, Bindu Vekaria, Luke Webb, Thomas House, and Ian Hall. Challenges in control of Covid-19: short doubling time and long delay to effect of interventions. 2020.
- [8] Neil M Ferguson et al. Impact of non-pharmaceutical interventions (NPIs) to reduce COVID19 mortality and healthcare demand. Technical Report arch, 2020.
- [9] Nicholas C. Grassly and Christophe Fraser. Mathematical models of infectious disease transmission. *Nature Reviews Microbiology*, 6(6):477–487, 2008.
- [10] W. Kermack and A. McKendrick. Contributions to the mathematical theory of epidemics. *Bulletin of Mathematical Biology*, 53(1-2):33–55, 1991.
- [11] Herbert W Hethcote. The mathematics of infectious diseases. *SIAM Review*, 42(4):599–653, 2000.
- [12] https://en.wikipedia.org/wiki/Compartmental_models_in_epidemiology.
- [13] Will Tribbey. *Numerical Recipes*, volume 35. Cambridge University Press, 3rd edition, 2010.