

## On COVID-19 Modelling

**Extension of the article submitted to DMV, work in progress**

**Robert Schaback**

August 1, 2020 , additions at the end, corrections in red, except for renumbering references and equations.

**Abstract** This is an analysis of the COVID-19 pandemic by comparably simple mathematical and numerical methods. The final goal is to predict the peak of the epidemic outbreak per country with a reliable technique. The difference to other modelling approaches is to stay extremely close to the available data, using as few hypotheses and parameters as possible.

For the convenience of readers, the basic notions of modelling epidemics are collected first, focusing on the standard SIR model. Proofs of various properties of the model are included. But such models are not directly compatible with available data. Therefore a special variation of a SIR model is presented that directly works with the data provided by the Johns Hopkins University. It allows to monitor the registered part of the pandemic, but is unable to deal with the hidden part. To reconstruct data for the unregistered Infected, a second model uses current experimental values of the infection fatality rate and a data-driven estimation of a specific form of the recovery rate. All other ingredients are data-driven as well. This model allows predictions of infection peaks.

Various examples of predictions are provided for illustration. They show what countries have to face that are still expecting their infection peak. Running the model on earlier data shows how closely the predictions follow the transition from an uncontrolled outbreak to the mitigation situation by non-pharmaceutical interventions like contact restrictions.

**Keywords** Epidemiology, SIR model, ordinary differential equations

**Mathematics Subject Classification (2010)** 92D30 · 92D25 · 93C15 · 34A34

---

Prof. (em.) Dr. Robert Schaback  
Institut für Numerische und Angewandte Mathematik,  
Universität Göttingen, Lotzestraße 16-18, 37083 Göttingen  
<http://num.math.uni-goettingen.de/schaback>  
E-mail: [schaback@math.uni-goettingen.de](mailto:schaback@math.uni-goettingen.de)

## 1 Introduction and Overview

During an epidemic outbreak like COVID-19, everybody wants to know how hard the impact will be. In particular:

- What is the health risk for me, my family, our friends, the city, the country, and the world?
- Is the health system prepared properly?
- Should households fill up their reserves in time?

This is a situation that asks for mathematics, like in the old times when mathematicians were needed to predict floods or solstices. Such predictions should be based on data and arguments, and they should provide well-supported suggestions for what to do. To understand the process and to make predictions, it should be modelled, and the model should be computable. Then predictions will be possible, and reality will decide later whether the model and the predictions were useful. Many models are possible, and the approach presented here is just one of them. The specific goal is to stay as close as possible to the available data, but it turns out that the available data are not directly usable for the standard models that give the basic understanding. To this end, two extensions to the standard SIR model are developed that get closer to the available data and finally are able to make data-driven predictions.

The beginning is made in section 2 with an introduction to standard terms like *Basic Reproduction Number*, *Herd Immunity Threshold*, and *Doubling Time*, together with some critical remarks on their use in the media. These notions are based on the standard SIR model for epidemics that is treated in quite some detail, including proofs for most of the mathematical properties. Experts can skip over this completely. Readers interested in the predictions should jump right away to section 5. For simplicity, the presentation ignores all delay-related issues like *incubation period* and *serial interval*.

To bridge the gap between model and data, Section 3 describes the Johns Hopkins data source with its limitations and flaws, and then presents a variation of a SIR model that can be applied directly to the data. It allows to estimate basic parameters, including the Basic Reproduction Number. But since the Johns Hopkins data provide no information about the unregistered cases and the Susceptibles, the model cannot yield reliable predictions of peaks of epidemics.

Therefore section 4 combines the data-compatible model of section 3 with a SIR model dealing with the unknown Susceptibles and the unregistered Infectious. This needs extra parameters that must be extracted from the literature. The first is the *infection fatality rate*, as provided e.g. by an der Heiden/Buchholz [10], Streeck et al. [24], Verity et al. [25]. Section 4.3.1 pairs it with the *case fatality rate* and shows how the latter can be deduced from the Johns Hopkins data. Like in Bommer/Vollmer [1], their combination gives a detection rate for the confirmed cases.

Section 4.4 introduces the second additional parameter: a recovery rate that can be directly used in the model and estimated from the infection fatality rate and the observable case fatality and case death rates. However, this parameter is not needed for prediction, just for determination of the unknown variables from the known data as long as the latter are available.

Then section 5 combines all of this into a larger model that makes predictions under the assumption that there are no further changes to the parameters by political action. It estimates the parameters of a full SIR model from the available Johns-Hopkins data by the techniques of section 4, using two additional technical parameters: the number of days used backwards for estimation of constants, and the number of days in which recovery or death can be expected on average, for estimation of case fatality and recovery rates. This is where time delays enter, but not into the model, only into internal estimation procedures. After the data-driven estimation of these parameters, the prediction uses only the infection fatality rate. All other ingredients are derived from the Johns Hopkins data.

Results are presented in section 5. Given the large uncertainties in the Johns-Hopkins data, the predictions are rather plausible. However, reality will have the final word on this prediction model.

The paper closes with a summary and a list of open problems.

## 2 Classical SIR Modelling

This contains the basic notions for modelling epidemics, defined and explained in mathematical terms. In particular, there will be a rigid mathematical underpinning of what is precisely meant when media talk about

- *flattening the epidemic outbreak (mitigation)*,
- *basic reproduction number*,
- *Herd Immunity Threshold*, and
- *doubling time*,

pointing out certain abuses of these notions. This will not work without calculus, but things were kept as simple as possible. Readers from outside the mathematics community should take the opportunity to brush up their calculus knowledge. Experts should go over to section 3.

### 2.1 The Model

The simplest standard *SIR* model of epidemics, due to Kermack-McKendrick [15] in 1927 and easily retrievable from the Wikipedia [27], deals with three variables

Susceptible (*S*), Infectious (*I*), and Removed (*R*).

The Removed cannot infect anybody anymore, being either dead or immune. This is the viewpoint of bacteria or viruses. The difference between death and immunity of subjects is totally irrelevant for them: they cannot proliferate anymore in both cases. The SIR model cannot say anything about death rates of persons.

The Susceptible are not yet infected and not immune, while the Infectious can infect Susceptibles. Individuals move by infection from *S* to *I*, and by death or healing from *I* to *R*. The three classes *S*, *I*, and *R* are disjoint and add up to a fixed total population count  $N = S + I + R$ . All of these are ideally assumed to be smooth functions

of time  $t$ , and satisfy the differential equations

$$\begin{aligned}\dot{S} &= -\beta \frac{S}{N} I, \\ \dot{I} &= +\beta \frac{S}{N} I - \gamma I, \\ \dot{R} &= \gamma I.\end{aligned}\tag{1}$$

where the dot stands for the time derivative, and where  $\beta$  and  $\gamma$  are positive parameters. The product  $\frac{S}{N}I$  models the probability that an Infectious meets a Susceptible and is actually infected.

Managing an SIR epidemic means *modifying* the constants  $\beta$  and  $\gamma$ . This is why one should see the parameters as control variables, and we shall treat them even as time series from section 3 on.

Note further that the Removed of the SIR model are not the Recovered of the Johns Hopkins data that we treat later, and the SIR model does not account for the Confirmed counted there. Similarly, there is no direct relation to the data published by the Robert Koch Institute. It is a major problem to match models with the available data, and we shall explain the latter to some detail in section 3. The inventors Kendrick and McKermack fitted their model already in 1927 [15] to data from the plague in Bombay 1905-1906.

## 2.2 Other Models

In many publications concerning COVID-19 (e.g. an der Heiden/Buchholz [10], Dandekar/Barbasthatis [2], De Brouwer et al. [3], Friston et al. [7], Khailaie et al. [16], Kucharski et al. [17], Maier/Brockmann [18]), the SIR model is extended by Exposed  $E$  that are infected, but not (yet) infectious. This introduces an additional parameter and would require dealing with a latency delay properly. We avoid this complication to keep the model as simple as possible. Note that there are extensions of SIR models with 14 to 21 parameters, e.g. Friston et al. [7], Giordano et al. [8], Khailaie et al. [16]. Fitting model parameters in the above papers is partially done numerically and partially by Bayesian approaches using Markov chain sampling of prior distributions. Here, we avoid fitting and time delays as far as possible.

Conceptually different are the agent-based model that is used by Ferguson et al. [6] for parameter estimation, and the approach of Mohring et al. [19] working consistently with time delays.

## 2.3 Simple Properties of the SIR Model

Since  $\dot{N} = \dot{S} + \dot{I} + \dot{R} = 0$  holds in (1), the equation  $N = S + I + R$  is kept valid at all times. The term  $\beta \frac{S}{N}I$  moves Susceptibles to Infectious, while  $\gamma I$  moves Infectious to Removed. Thus  $\beta$  represents an *infection rate* while the *removal rate*  $\gamma$  accounts for either healing or fatality after infection, i.e. immunity. Political decisions about

reducing contact probabilities will affect  $\beta$ , while  $\gamma$  resembles the balance between the medical aggressivity of the infection and the quality of the health care system.

As long as the Infectious  $I$  are positive, the Susceptibles  $S$  are decreasing, while the Removed  $R$  are increasing. Excluding the trivial case of zero Infectious from now on, the Removed and the Susceptible will be strictly monotonic. Therefore we can use them to re-parameterise the model at certain places.

The SIR model is not really dependent on the total population  $N$ . Moreover, if we scale time by  $\tau := t \cdot \gamma$  and go over to *relative* quantities

$$\begin{aligned} s(\tau) &:= \frac{S(\tau/\gamma)}{N}, \\ r(\tau) &:= \frac{R(\tau/\gamma)}{N}, \\ i(\tau) &:= \frac{I(\tau/\gamma)}{N}, \end{aligned}$$

we get the new system

$$\begin{aligned} s'(\tau) &= \frac{ds}{d\tau} = -\frac{\beta}{\gamma}s(\tau)i(\tau) = -R_0s(\tau)i(\tau) \\ i'(\tau) &= \frac{di}{d\tau} = \left(\frac{\beta}{\gamma}s(\tau) - 1\right)i(\tau) = (R_0s(\tau) - 1)i(\tau) \\ r'(\tau) &= \frac{dr}{d\tau} = i(\tau) \end{aligned} \quad (2)$$

only containing the *Basic Reproduction Number*

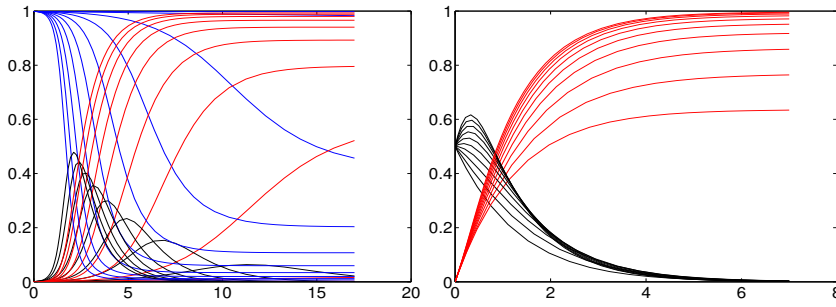
$$R_0 := \frac{\beta}{\gamma} \quad (3)$$

that will turn out to be of central importance. Both  $\beta$  and  $\gamma$  vary under a change of time scale in (1), but the basic reproduction number is invariant. Physically,  $\beta$  and  $\gamma$  have the dimension  $time^{-1}$ , but  $R_0 = \beta/\gamma$  and the new “time” parameter  $\tau$  in (2) are dimensionless. Another interpretation of (2) is that after a time scale one can assume  $\gamma = 1$  and  $R_0 = \beta$ . We call  $\tau$  the *unit removal parameter*, because its unit can be seen as the average time needed to get removed, i.e. either dead or immune. We use a prime to denote derivatives with respect to  $\tau$ . But in all later sections that make real-world interpretations, we have to use real time, and then we shall go back to (1).

A standard mathematical trick is to divide the first equation by the third to get

$$\begin{aligned} \frac{ds}{dr} &= -R_0s, \\ s(r) &= s(r(0)) \exp(-R_0(r - r(0))). \end{aligned} \quad (4)$$

We shall use (4) in section 2.11 to study the long-term behaviour of solutions. The introduction of (4) is a typical pitfall for mathematics: it is a nice theoretical simplification, but it obscures the most interesting practical aspect, in this case the fraction  $i$  of infectious persons in the population. The same holds for the simplification by setting  $d\tau = \gamma \frac{t}{N} dt$  that is ignored here, leaving it to interested readers.



**Fig. 1** Some typical SIR system solutions, relative to the total population. See the explanation in section 2.4. The peaked curves for the Infectious are “flattened” for small  $R_0$ .

## 2.4 Examples

Figure 1 shows a series of test runs of a SIR model. Recall that the relative Recovered  $r$  are increasing from zero, and the relative Susceptibles  $s$  are decreasing down from one. The relative Infectious  $i$  are in between and can possibly show a sharp peak that everybody tries to avoid. We shall deal with the mathematics of the peak in sections 2.8, 2.13, and 2.14, while the rest of the paper focuses on data-driven predictions of peaks. The Infectious are usually not covered by the media who tend to focus on the cumulative number of confirmed cases, containing the Removed.

In both plots we set  $r(0) = 0$ ,  $\gamma = 1$ , and let  $R_0 = \beta$  vary from 0.1 to 5. The difference between the figures lies in the initial value  $i(0)$ . Left, due to a realistically small  $i(0) = 0.001$ , one cannot see the decaying peak-less cases of  $i$  near startup for  $R_0 < 1$ , while the right-hand plot has  $i(0) = 1/2$  and shows them. Decreasing  $R_0 \searrow 1$  flattens the peaks of the Infectious  $i$ , and there is no peak for  $R_0 \leq 1$ . Furthermore, one can observe that  $i$  always decays to zero, while  $s$  and  $r$  tend to fixed positive levels in the long run. The final level of  $r$  is particularly interesting because part of it is the total death toll. It decreases when  $R_0$  decreases. We shall prove all of this later. When countries change parameters by administrative actions like a shutdown, they jump to a more flat  $i$  curve, e.g. at an intersection point.

From the system, one can also infer that  $r$  has an inflection point where  $i$  has its maximum, since  $r'' = i'$ . If only  $r$  would be observable, one could locate the peak of  $i$  via the inflection point of  $r$ . Finally, note that small initial values  $i(0)$  of  $i$  delay the peak considerably, no matter how large  $R_0$  is. We shall prove this in section 2.14.

## 2.5 Interpretation of the Basic Reproduction Number $R_0$

Media often say that  $R_0$  gives the number of persons an average Infectious infects while being infectious. This is a rather mystical statement that needs underpinning. In the SIR system (1) the quantity

$$\frac{1}{\gamma} = \frac{I}{\dot{R}}$$

is a value that has the physical dimension of time. It describes the ratio between current Infectious and current newly Removed, and thus can be seen as the average time needed for an Infectious to get Removed, i.e. the average time that an Infectious can infect others. This is why we called the dimensionless  $\tau = t \cdot \gamma$  the *unit removal parameter* in section 2.3. Correspondingly,

$$\dot{I} + \gamma I = \dot{I} + \dot{R} = \beta \frac{S}{N} I$$

are the newly Infected, and therefore

$$\frac{1}{\beta} \frac{N}{S} = \frac{I}{\dot{I} + \dot{R}}$$

can be seen as the time it needs for an average Infectious to generate a new Infectious. The ratio  $R_t := \frac{\beta}{\gamma} \frac{S(t)}{N}$  then gives how many new Infectious can be generated by an Infectious while being infectious. This is the time-dependent *Reproduction Number*, but it is only close to  $R_0$  if  $S(t) \approx N$ , i.e. at the start of an outbreak. A correct statement is that  $R_0$  is the average number of infections an Infectious generates while being infectious, but within an unlimited supply of Susceptibles.

To let less new Infectious be generated, administrative actions try to change the parameters of the epidemic towards small  $R_0$ . We shall see that this is correct from a mathematical viewpoint as well, and we shall study the influence of  $R_0$  to quite some detail.

The above interpretation of  $R_0$  shows two major ways to make  $R_0$  small: reducing the number of possibly infective contacts, and reducing the time an Infectious has to infect others. The second works by putting all infectious persons into strict quarantine, while first can be done by reducing contacts of all persons, even the Susceptibles, and reducing the infection probability for each contact, e.g. by wearing masks.

SIR-based models of the COVID-19 pandemics estimate  $R_0$  between 2 and 6 during an uncontrolled outbreak (see e.g. the Robert Koch-Institute [21], De Brouwer et al. [3], Dehning et al. [5], and Maier/Brockmann [18]), while *non-pharmaceutical interventions* (NPI) bring  $R_0$  below 1. We shall see examples in 3.3.2 and 5.2.

The use of the Basic Reproduction Number  $R_0$  in the media suggests that large  $R_0$  are generally serious, because each Infectious infects several people. This is only true at the beginning of an outbreak, because then there are enough Susceptibles. But it will turn out in section 2.8 that the Infectious will always finally go to zero, whatever the Basic Reproduction Number is. See Figure 1 as well.

## 2.6 Conditions for Outbreaks

The first interesting question in a beginning epidemic is:

Will there be a serious outbreak, or will the infection disappear quickly?

Therefore we first look at the initial conditions for the model. Since everything is invariant under an additive time *shift*, we can start at time 0, and since time *scales* are irrelevant to the problem at startup, we can use the simplified system (2).

The relative Infectious  $i$  in (2) do not increase right from the start if  $\dot{I}(0) \leq 0$ , i.e.

$$s(0) \leq \frac{1}{R_0}, \quad (5)$$

and then they decrease further since the Susceptibles  $s$  must decrease and

$$\frac{i(\tau)'}{i(\tau)} = (\log i(\tau))' = R_0 s(\tau) - 1 < R_0 s(0) - 1 \leq 0. \quad (6)$$

There is no outbreak, and this must occur for all initial conditions if  $R_0 \leq 1$ . But if  $R_0 > 1$ , the outbreak depends on the initial condition (5). Altogether, outbreaks are fully characterised by

$$1 > s(0) > \frac{1}{R_0}. \quad (7)$$

## 2.7 Herd Immunity Threshold

In connection with an outbreak, the *Herd Immunity Threshold*

$$HIT = 1 - \frac{1}{R_0}$$

is often mentioned. The background question is:

If an uninfected population is threatened by an infection with Basic Reproduction Number  $R_0$ , what is the number of immune persons needed to prevent an outbreak right from the start?

In the idealised situation  $i(0) = 0$  and  $s(0) + r(0) = 1$ ,

$$r(0) = 1 - \frac{1}{R_0} = HIT$$

follows from (5) and (7) as the threshold between outbreak and decay for the relative Removed. This does not refer to a whole epidemic scenario. It is to be checked *before* anything happens, and useless within a developing epidemic, whatever the media say.

## 2.8 The Peak

In the outbreak case (7), the main questions are:

- When will the Infectious reach their maximum?
- How large will the maximal value be?



More generally, we ask for a time  $t_I$  or a unit removal parameter  $\tau_I = \gamma t_I$  where the Infectious  $i$  are positive and do not change. Then we have

$$0 = \frac{di}{d\tau}(\tau_I) = (R_0 s(\tau_I) - 1)i(\tau_I), \quad (8)$$

and the monotonicity of  $s$  implies uniqueness of  $\tau_I$  and

$$s(\tau_I) = \frac{1}{R_0}. \quad (9)$$

If  $i$  would increase without reaching a maximum in finite time, the first equation of (2) would imply that  $s$  goes exponentially to zero, but then there is a  $\tau_I$  with (9), and (8) follows. Summarising, this proves that whenever there is an outbreak by (7), there is a unique maximum of the relative Infectious  $i$  that we call the *peak* from now on. Behind the peak, or apart from any outbreak situation, the Infectious must go exponentially to zero due to (6), because the Susceptibles continue to decrease, no matter how large  $R_0$  is.

Determining the peak is theoretically difficult, and in practice it requires good estimates for  $\beta$  and  $\gamma$ . Mathematical results on the peak will be in sections 2.13 and 2.14, while data-driven predictions follow in section 5.2

In real life it is highly important to avoid the peak situation, and this can only be done by administrative measures that change  $\beta$  and  $\gamma$  in (1) to the situation  $\beta < \gamma$ . This is what management of epidemics is all about, provided that an epidemic follows the SIR model. We shall see how countries perform.

In the peak situation of (8) and (9), the fraction

$$1 - \frac{1}{R_0} = 1 - s(\tau_I) = r(\tau_I) + i(\tau_I) \geq i(\tau_I) \quad (10)$$

of the relative Non-Susceptible at the peak is exactly the Herd Immunity Threshold. Thus it is correct to say that if the Immune of a population are below the Herd Immunity Threshold at startup, and if the Basic Reproduction Number is larger than one, the sum of the Immune and the Infectious will rise up to the Herd Immunity Threshold and then the Infectious will decay. This is often stated imprecisely in the media.

## 2.9 Analysing the Outbreak

When an outbreak starts, almost everybody is susceptible, i.e.  $s(0) \approx 1$ , and then

$$i' = R_0 s - 1 \approx R_0 - 1$$

models an exponential outbreak with exponent  $R_0 - 1 > 0$  in unit removal parametrisation, with a solution

$$i(\tau) \approx i(0) \exp((R_0 - 1)\tau).$$

If this is done in real time  $t$  and discrete time steps  $\Delta t$ , the system (1) yields

$$\frac{I(t + \Delta t)}{I(t)} \approx \exp((\beta - \gamma)\Delta t).$$

The severity of the outbreak in real time is not controlled by  $R_0 = \beta/\gamma$ , but rather by  $\beta - \gamma$ . Publishing single values  $I(t)$  does not give any information about  $\beta - \gamma$ . Better is the ratio of two subsequent values

$$\frac{I(t_2)}{I(t_1)} \approx \exp((\beta - \gamma)(t_2 - t_1)), \quad (11)$$

and if this gets smaller over time, the outbreak gets less dramatic because  $\beta - \gamma$  gets smaller. But (11) is by no means a correct way to estimate  $R_0$ .

Therefore, really useful information about an outbreak must concern  $I$ , but should not consist of single values. Increments in percent are much better, because their logarithm is proportional to  $\beta - \gamma$ . However, it needs increments of increments to see whether administrative actions are successful by changing  $\beta - \gamma$ . This is what the media rarely provided during the outbreak. On the positive side, the severity of a future outbreak in unit removal **parameterisation** is described correctly by estimates of  $R_0 > 1$ , if these have a solid mathematical and experimental basis. All changes of  $R_0$  should be carefully monitored.

## 2.10 Doubling Time

Another information used by media during an outbreak is the *doubling time*, i.e. how many days it takes until daily values double. It is  $n\Delta t$  with the number  $n$  from

$$2 = \frac{I(t + n\Delta t)}{I(t)} \approx \exp((\beta - \gamma)n\Delta t) = (\exp((\beta - \gamma)\Delta t))^n$$

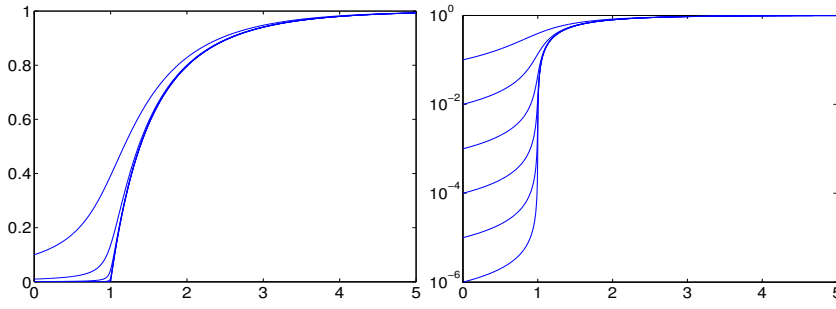
or

$$n = \frac{\log 2}{(\beta - \gamma)\Delta \tau},$$

i.e. it is inversely proportional to  $\beta - \gamma$ . If political action doubles the doubling time, it halves  $\beta - \gamma$ . If politicians do this repeatedly, they never reach  $\beta < \gamma$ , and they never escape an exponential outbreak if they do this any finite number of times. Extending the doubling time will never prevent a peak, it only postpones it and hopefully flattens it. When presenting a doubling time, media should always point out that this makes only sense during an exponential outbreak. And it is not related to the basic reproduction number  $R_0 = \beta/\gamma$ , but to the difference  $\beta - \gamma$ .

## 2.11 Long-term Behaviour

Aside from the peak, it is interesting to know the portions of the population that get either permanently removed (by death or immunity) or never come into contact with the infection. This concerns the long-term behaviour of the Removed and the Susceptibles. Figure 1 demonstrates how  $r$  and  $s$  level out under all circumstances shown, but is this always true, and what is the final ratio? And if one has additional information on the percentage of casualties within the Removed, what is the total death toll in the long run?



**Fig. 2** The asymptotic level  $r_\infty$  of the relative Removed as a function of  $R_0$  for  $s(0) = 0.9, 0.99, 0.999$  etc. as curves from the top. Right: logarithmic scale.

Going back to (4), we get

$$s(r) = s(0) \exp(-R_0 r) \quad (12)$$

when assuming  $r(0) = 0$  at startup. Since  $r$  is increasing, it has a limit  $0 < r_\infty \leq 1$  for  $\tau \rightarrow \infty$ , and in this limit

$$s_\infty = s(0) \exp(-R_0 r_\infty)$$

holds, together with the condition  $r_\infty + s_\infty = 1$ , because there are no more Infectious. The transcendental equation

$$s(0) \exp(-R_0 r_\infty) = 1 - r_\infty \quad (13)$$

has a unique solution in  $(0, 1)$  dependent on  $s(0) < 1$  and  $R_0$ . Therefore the Infectious always go to zero, but Susceptibles always remain. Then a new infection can always arise as soon as an infected person enters the sanitised population. The outbreak risk is dependent on the portion  $s_\infty = 1 - r_\infty$  of the Susceptibles by (5). This illustrates the importance of vaccination, e.g. against measles or influenza.

To see how  $r_\infty$  and  $s_\infty = 1 - r_\infty$  behave as functions of  $R_0$  and  $s(0)$ , we solve the equation (13) by the Lambert  $W$  function to get

$$r_\infty = 1 + \frac{1}{R_0} W(-s(0)R_0 \exp(-R_0)) \quad (14)$$

with a surprising behaviour. See Figure 2 for illustration. Left, the curves for unrealistically small initial values  $s(0) = 0.9, 0.99$  and  $0.999$  for Susceptibles can still be distinguished from the more interesting curves below that coincide for all  $s(0)$  closer to one and have a sharp turn at  $R_0 = 1$ . The logarithmic plot to the right shows that for  $R_0 < 1$  the curves separate, and that it pays off significantly to have  $R_0 < 1$  for  $s(0)$  close to one.

This has some serious implications, if the model is correct for an epidemic situation. When politicians try to “flatten the curve” by bringing  $R_0$  below 1 at some early time when the Susceptibles are still abundant, the asymptotic rate  $r_\infty$  of Removed will be *dramatically* smaller than for any other situation, because one stays left of the

sharp turn in Figure 2. This is particularly important if the rate of fatalities within the Removed is high.

Large values of  $R_0$  lead to large relative values of Removed to Susceptible in the limit. The consequence is that systems with large  $R_0$  have a dramatic outbreak and lead to a large portion of Removed. This is good news if the rate of fatalities within the Removed is low, but very bad news otherwise. When pressing  $R_0$  below one, the risk of re-infection rises due to the larger portion of Susceptibles, but the deaths contained in the Removed are kept low.

The decay situation (5) implies that  $s_\infty \leq 1/R_0$  holds, and consequently

$$r_\infty = 1 - s_\infty \geq 1 - \frac{1}{R_0} = HIT.$$

Therefore the final rate of the Removed is not smaller than the Herd Immunity Threshold. This is good news for possible re-infections, but only if the death rate among the Removed is small enough.

## 2.12 Asymptotic Exponential Decay

If we go back to (6) for a unit removal parameter  $\tau_D$  where  $i$  decreases, in an outbreak or not, we have  $R_0 s_\infty \leq R_0 s(\tau_D) < 1$  and then

$$i(\tau_D) \exp((R_0 s_\infty - 1)(\tau - \tau_D)) \leq i(\tau) \leq i(\tau_D) \exp((R_0 s(\tau_D) - 1)(\tau - \tau_D))$$

for all  $\tau \geq \tau_D$ . Therefore the exponential decay in unit removal parametrisation is not ruled by  $R_0 - 1$  as in the outbreak case with  $R_0 > 1$ , but rather by  $R_0 s_\infty - 1$ . This also holds for large  $R_0$  because  $s_\infty$  counteracts. The bell shapes of the peaked  $i$  curves are not symmetric with respect to the peak. Inserting (14), the relative Infectious always decay asymptotically exponentially like

$$\exp((R_0 s_\infty - 1)\tau) = \exp((W(-s(0)R_0 \exp(-R_0)) - 1)\tau) \text{ for } \tau \rightarrow \infty$$

with the Lambert  $W$  function. By MAPLE, the slowest decay arises for  $R_0 = 1$ .

## 2.13 Maximal Infectious at the Peak

At the peak of the Infectious  $i$  at  $\tau_I$  in an outbreak (7) with  $r(0) = 0$  we know

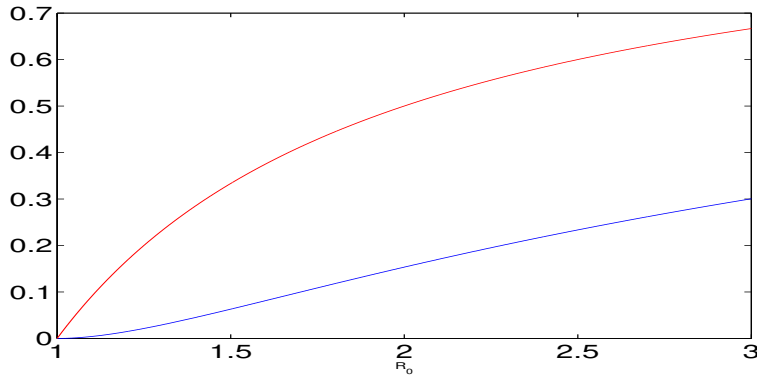
$$s(\tau_I) = \frac{1}{R_0} = s(r(\tau_I)) = s(0) \exp(-R_0 r(\tau_I))$$

from (9) and (4), and get the Removed at the peak as

$$r(\tau_I) = \frac{1}{R_0} \log(s(0)R_0). \quad (15)$$

Then the exact value of the Infectious  $i$  at the peak is

$$i(\tau_I) = 1 - s(\tau_I) - r(\tau_I) = 1 - \frac{1}{R_0} - \frac{1}{R_0} \log(s(0)R_0), \quad (16)$$



**Fig. 3** The effect of  $R_0$  on the peak value  $i(\tau_i)$  of Infectious.

improving (10). Note that the log is positive due to the outbreak condition (7). It is remarkable that the *value* of  $i$  at the peak does not depend on initial conditions, while the next section proves that the *position* of the peak does.

For standard infections that have starting values  $s(0) = S(0)/N$  very close to one, the maximal ratio of Infectious is

$$i(\tau_i) \approx 1 - \frac{1}{R_0} - \frac{1}{R_0} \log(R_0).$$

Figure 3 shows the behaviour of this function, as the lower curve. A value of  $R_0 = 4$  leads to a maximum of more than 40% of the population infectious at a single time. If 5% need hospital care, a country needs hospital beds for 2% of the population around peak time. This disaster calls for mitigation by lowering  $R_0$ .

The upper curve leaves the log term out, i.e. it marks the rate (9) of the Susceptibles at the peak, and by (10) the difference is the rate  $r(\tau_i)$  of the Recovered at the peak. It also marks the extreme case in (7) with  $R_0 s(0) = 1$ , i.e. having the smallest possible initial value of  $s(0)$  for a given  $R_0$  to generate an outbreak. Therefore all  $s(0)$ -dependent possibilities vary between the two curves.

#### 2.14 Localising the Peak

Knowing now how large the peak is, we want to find out where it is. We write the unit removal parameter  $\tau$  as a function of  $r$  by  $\frac{d\tau}{dr} = \left(\frac{dr}{d\tau}\right)^{-1} = \frac{1}{i}$  and integrate from  $r = 0 = r(0)$  to  $r = r(\tau_i)$  to get the peak position

$$\tau_i = \int_0^{r(\tau_i)} \frac{1}{i(r)} dr = \int_0^{\log(s(0)R_0)/R_0} \frac{1}{1 - r - s(0)\exp(-R_0 r)} dr$$

as a nasty function of  $s(0)$  and  $R_0$ , using (2), (12), and  $1 = i(r) + s(r) + r$ . To prove that the peak moves towards zero for both limits  $R_0 \nearrow \infty$  and  $R_0 \searrow 1$ , we first observe

that  $i \geq i(0)$  holds left of the peak. Then we use (15) to get

$$\tau_I \leq \frac{r(\tau_i)}{i(0)} = \frac{1}{i(0)R_0} \log(s(0)R_0) \leq \frac{1}{i(0)R_0} \log(R_0) \leq \frac{1}{e \cdot i(0)} \approx \frac{0.37}{i(0)} \quad (17)$$

by inserting the maximum of  $\log(R_0)/R_0$  at  $e$ . The upper bound gets large when  $i(0)$  gets small, a realistic case by Figures 1 and 8. This calls for a lower bound.

For fixed  $s(0)$  and  $i(0)$  there will be a maximal peak position for a rather specific  $R_0$ . A MAPLE-based analysis shows that  $R_0(s(0)) = -W(-s(0)/e)^{-1}$  with Lambert's  $W$  function yields

$$\tau_I \geq \frac{0.3(1-i(0))}{\sqrt{i(0)}}.$$

Therefore the peak can indeed move arbitrarily far out for small  $i(0)$  and large  $s(0) = 1 - i(0)$ . There is not much leeway for smaller  $R_0$  to bring the peak position to zero for large  $s(0)$ , namely  $\frac{1}{s(0)} < R_0 < R_0(s(0))$ . Both bounds for  $R_0$  tend to one for  $s(0) \rightarrow 1$ .

The practical consequence is that keeping  $R_0 > 1$  close to one by mitigation is no good idea, because the peak can move far into the future for realistically small  $i(0)$ , delaying the epidemic in an intolerable way. Countries should go for  $R_0$  considerably smaller than one.

## 2.15 Turnaround Time

In a peak situation (7) one can consider the *turnaround parameter*  $\tau_T$  at which the Infectious  $i$  come back to their starting value  $i(0)$  behind the peak. At that point the population has accumulated more Removed, dead or immune. We calculate the integral

$$\int_0^\infty i(\tau) d\tau = \int_0^\infty r'(\tau) d\tau = r_\infty - r(0).$$

The rectangle of length  $\tau_T$  and height  $i(0)$  fits under the  $i$  curve, and therefore

$$i(0)\tau_T \leq r_\infty - r(0) \leq r_\infty \leq 1,$$

proving that the real turnaround *time*  $t_T = \tau_T/\gamma$  has a fixed bound  $t_T \leq r_\infty/(i(0)\gamma)$ . From Figure 2 one can see that making  $R_0$  smaller will decrease the bound via  $r_\infty$ .

## 2.16 Estimating and Varying Parameters

If real-time data for the SIR model (1) were fully available, one could solve for

$$\begin{aligned} \gamma &= \frac{\dot{R}}{I}, & b &:= \beta \frac{S}{N} = \frac{\dot{I} + \gamma I}{I} = \frac{\dot{I} + \dot{R}}{I}, \\ \beta &= \frac{N}{N-I-R} \cdot \frac{\dot{I} + \dot{R}}{I}, & R_0 &= \frac{N}{N-I-R} \cdot \frac{\dot{I} + \dot{R}}{\dot{R}} = -\frac{N \dot{S}}{S \dot{R}} = -\frac{1}{s} \frac{ds}{dr}, \end{aligned} \quad (18)$$

and we shall use this in section 3.3. The validity of a SIR model can be tested by checking whether the right-hand sides for  $\beta$ ,  $\gamma$  and  $R_0$  are roughly constant. If data

are sampled locally, e.g. before or after a peak, the above technique should determine the parameters for the global epidemic and be useful for either prediction or backward testing.

However, in pandemics like COVID-19, the parameters  $\beta$  and  $\gamma$  change over time by administrative action. This means that they should be considered as functions in the above equations, and then their changes may be used for conclusions about the influence of such actions. From this viewpoint, one can go back to the SIR model and consider  $\beta$  and  $\gamma$  as control functions that just describe the relation between the variables.

But the main argument against using (18) is that the data are rarely available. This is the concern of the next section.

### 3 Using Available Data

Now we confront the modelling of the previous section with available data. This is crucial for manoeuvring countries through the epidemics (Sentker [23])<sup>1</sup>. From now on we have to work in real time and go back to (1) instead of all mathematical simplifications.

#### 3.1 Johns Hopkins Data

We work with the COVID-19 data from the Johns Hopkins University at GitHub [9]. They are the only source that provides comparable data on a worldwide scale, namely

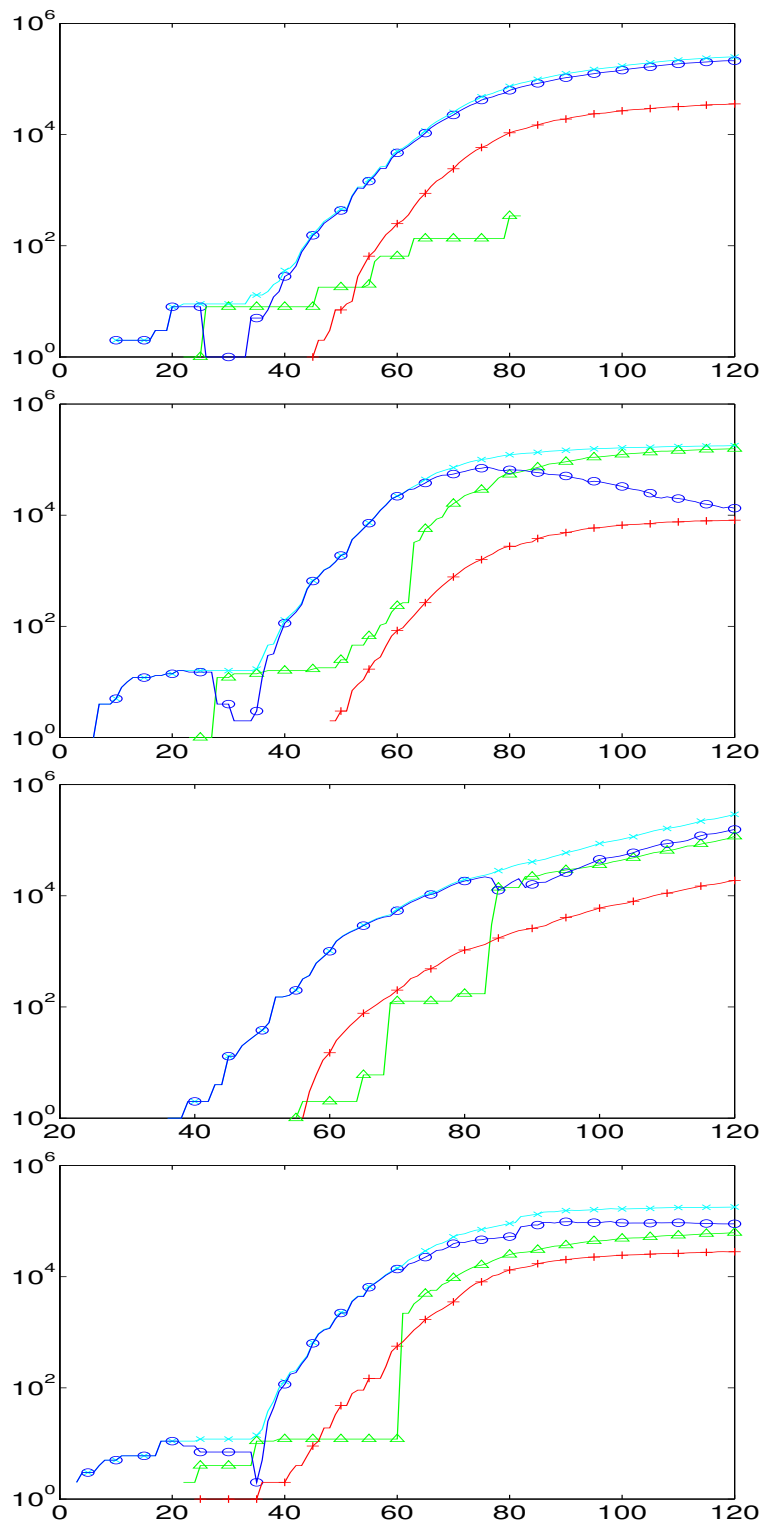
1. Confirmed ( $C$ ) or *cumulative infected*
2. Dead ( $D$ ), and
3. Recovered ( $R$ ), i.e. alive and immune,

as cumulative integer valued time series for days from Jan. 22nd, 2020. All these values are absolute numbers, not relative to a total population. Note that the unconfirmed cases and the Susceptibles are not accessible at all, while the Confirmed contain the Dead and the Recovered of earlier days.

The media, in particular German TV, present COVID-19 data in a rather debatable way. When mentioning Johns Hopkins data, they provide  $C$ ,  $D$ , and  $R$  separately without stating the most important figures, namely  $I = C - D - R$ , their change, and the change of their change. When mentioning data of the Infectious from the Robert Koch institute alongside, they do not say precisely that these are non-cumulative and should be compared to the  $I = C - R - D$  data of the Johns Hopkins University. And, in most cases during the outbreak, they did not mention the change of the change. Quite like all other media.

We take the data as presented, but there are many well-known flaws. In particular, the values for specific days are partly belonging to previous days, due to delays in the chains of data transmission in different countries. This is why, at some points, we

<sup>1</sup> Original text in German, April 16th: *Schnelle Modelle, die dem Abgleich mit der Wirklichkeit standhalten, sind eine wichtige Voraussetzung, das Land politisch durch die Seuche zu steuern.*



**Fig. 4** Raw Johns Hopkins data in logarithmic presentation up to day 120, from top: UK, Germany, Brazil, and France. Markers  $X$  for Confirmed,  $O$  for Infectious,  $\wedge$  for Recovered,  $+$  for Deaths, not on all data points.



shall apply some conservative smoothing of the data. Finally, there are inconsistencies that possibly need data changes. In particular, there are countries like Germany who deliver data of Recovered in a very questionable way. The law in Germany did not enforce authorities to collect data of Recovered, and the United Kingdom did not report numbers of Dead and Recovered from places outside the National Health System, e.g. from Senior's retirement homes. Both strategies have changed somewhat in the meantime, as of early May, but the data still keep these flaws. See Figure 4 for examples.

We might assume that the Dead plus the Recovered of the Johns Hopkins data are the Removed of the SIR model, and that the Infectious  $I = C - R - D$  of the Johns Hopkins data are the Infectious of the SIR model. But this is not strictly valid, because the Johns Hopkins data concern only registered cases.

On the other hand, one can take the radical viewpoint that facts are not interesting if they do not show up in the Johns Hopkins data. Except for the United Kingdom, the important figures concern COVID-19 casualties that are actually registered as such, others do not count, and serious cases needing hospitalisation or leading to death should not go unregistered. If they do in certain countries, using such data will not be of any help, unless other data sources are available.

An important point for what follows is that the data come as daily values. To make this compatible with differential equations, we shall replace derivatives by differences.

### 3.2 Examples

To get a first impression about the Johns Hopkins data, Figure 4 shows raw data up to day 120, May 21st. For better visibility, not all data points have markers. Here, and in all plots to follow, the  $x$  axis has the days after Jan. 22nd, 2020. It might be helpful to remember that day 100 is May 1st. The  $y$  axis is logarithmic, because then linearly increasing or decreasing parts in the figures correspond to exponentially increasing or decreasing numbers in the real data.

Many presentations in the media are non-logarithmic, and then all exponential outbreaks look similar. The most interesting data are the Infectious  $I = C - R - D$  marked by  $O$  that show a peak or not, and the cumulative casualties  $D$  marked by  $+$ . The data for other countries tell similar stories and are suppressed.

One can see in Figure 4 that Germany has passed the peak of the Infectious, while France is roughly at the peak and the United States and Brazil are still in an exponential outbreak. The early figures, below day 40, are rather useless, but then an exponential outbreak is visible in all cases. This outbreak changes its slope due to political actions, and we shall analyse this later. See Dehning et al. [5] for a detailed early analysis of slope changes.

There are strange anomalies in the Recovered ( $\wedge$  marker). France seems not to have delivered any data between days 40 and 58, Germany changed the data delivery policy between days 62 and 63, and the UK data for the Recovered are a mess. We shall avoid using data on the Recovered as much as possible.

It should be noted that the available medical results on the COVID-19 disease often state that Confirmed will die or survive after a more or less fixed number of days. This would imply that the curves marked + for the Dead and the curves marked  $\wedge$  for the Recovered should roughly follow the curves marked  $X$  for the Confirmed with a fixed but measurable delay. This is partially observable, but much less accurately for the Recovered.

### 3.3 The Johns Hopkins Data Model

We now define a model that works exclusively with the Johns Hopkins data, but comes close to a SIR model, without being able to use  $S$ . Since the SIR model does not distinguish between recoveries and deaths, we set in obvious notation

$$R_{SIR} \Leftrightarrow D_{JH} + R_{JH}$$

and let the Infectious be comparable, i.e.

$$I_{SIR} \Leftrightarrow I_{JH} := C_{JH} - D_{JH} - R_{JH}$$

which implies

$$(I + R)_{SIR} \Leftrightarrow C_{JH},$$

and we completely omit the Susceptibles. From now on, we shall drop the subscript  $JH$  when we use the Johns Hopkins data, but we shall use  $SIR$  when we go back to the SIR model.

Now we take (18) of section 2.16 and insert differences:

$$\begin{aligned} \gamma &= \frac{\dot{R}_{SIR}}{I_{SIR}} \\ &\approx \frac{(D+R)_{n+1} - (D+R)_n}{I_n} =: \gamma_n \\ b &:= \beta \frac{S_{SIR}}{N} = \frac{\dot{I}_{SIR} + \gamma I_{SIR}}{I_{SIR}} = \frac{\dot{I}_{SIR} + \dot{R}_{SIR}}{I_{SIR}}, \\ &\approx \frac{C_{n+1} - C_n}{I_n} =: b_n, \end{aligned}$$

defining time series  $\gamma_n$  and  $b_n$  that model  $\gamma$  and  $b = \beta \cdot S_{SIR}/N$  without knowing  $S_{SIR}$ . This is equivalent to the model

$$\begin{aligned} C_{n+1} - C_n &= b_n I_n, \\ I_{n+1} - I_n &= b_n I_n - \gamma_n I_n = (b_n - \gamma_n) I_n, \\ (R+D)_{n+1} - (R+D)_n &= \gamma_n I_n \end{aligned} \tag{19}$$

that maintains  $C = I + R + D$ , and we may call it a *Johns Hopkins Data Model*. It is very close to a SIR model if the time series  $b_n$  is not considered to be constant, but just an approximation of  $\beta \cdot S_{SIR}/N$ .

### 3.3.1 Estimating $R$

By brute force, one can take

$$r_n = \frac{b_n}{\gamma_n} = \frac{C_{n+1} - C_n}{R_{n+1} + D_{n+1} - R_n - D_n} \quad (20)$$

as a data-driven substitute for

$$\frac{\beta}{\gamma} \frac{S_{SIR}}{N} = R_0 \frac{S_{SIR}}{N}.$$

Then there is a rather simple observation:

If  $r_n$  is smaller than one, the Infectious decrease.

It follows using (20) via

$$\begin{aligned} I_{n+1} - I_n &= C_{n+1} - C_n - (R_{n+1} - R_n + D_{n+1} - D_n) \\ &= (r_n - 1)(R_{n+1} + D_{n+1} - R_n - D_n), \end{aligned}$$

but this is visible in the data anyway and not of much help.

Since  $r_n$  models  $R_0 \frac{S_{SIR}}{N}$ , it always underestimates  $R_0$ . This underestimation gets dramatic when it must be assumed that  $S_{SIR}$  gets seriously smaller than  $N$ .

At this point, it is not intended to forecast the epidemics. The focus is on extracting parameters from the Johns Hopkins data that relate to a background SIR-type model.

### 3.3.2 Example

Figure 5 shows  $R_0 \frac{S_{SIR}}{N}$  estimates via  $r_n$  for the last four weeks before day 120, i.e. March 21st. Except for the United States and Brazil, all countries were more or less successful in pressing  $r_n$  below one. In all cases,  $S_{SIR}/N$  is too close to one to have any influence. The variation in  $r_n$  is not due to the decrease in  $S_{SIR}/N$ , but should rather be attributed to political action. As mentioned above, the estimates for  $R_0$  by  $r_n$  are always optimistic.

For the figure, the raw Johns Hopkins data were smoothed by a double action of a  $1/4, 1/2, 1/4$  filter on the logarithms of the data. This smoother keeps constants and linear sections of the logarithm invariant, i.e. it does not change local exponential behaviour. This smoothing was not applied to Figure 4. It was by far not strong enough to eliminate the apparent 7-day oscillations that are **frequent** in the Johns Hopkins data, see Figure 5, Data from the Robert Koch Institute in Germany have even stronger 7-day variations.

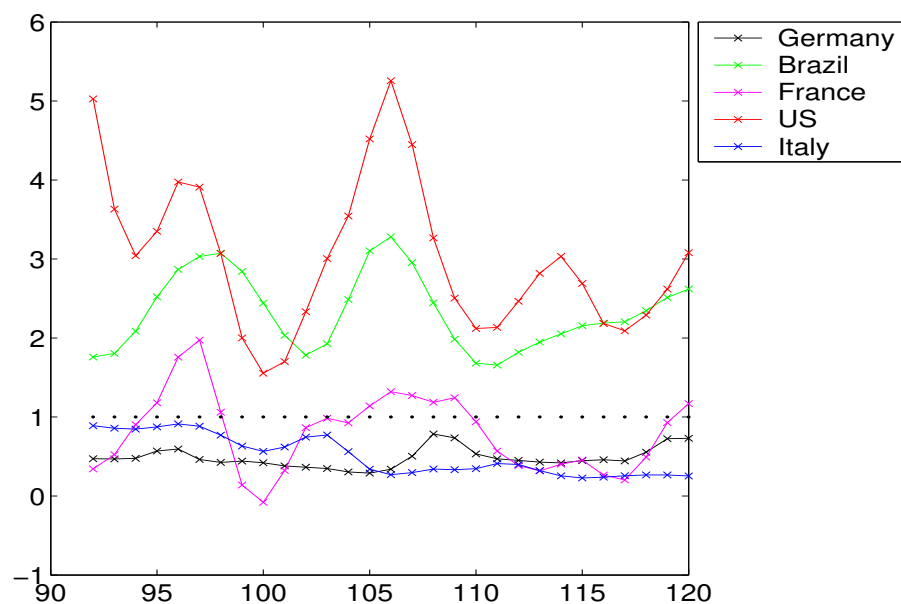


Fig. 5 Estimates of  $R_0$  via the time series  $r_n$  up to day 120

### 3.3.3 Properties of the Model

As long as  $r_n$  is roughly constant, the above approach will always model an exponential outbreak or decay, but never a peak, because the difference equations are linear. It can only help the user to tell if there is a peak ahead or behind, depending on  $r_n \approx R_0$  being larger or smaller than 1. If  $r_n$  is kept below one, the Confirmed Infectious will not increase, causing no new threats to the health system. Then the  $S/N$  factor will not decrease substantially, and a full SIR model is not necessary.

As long as countries keep  $r_n$  clearly below one, e.g. below  $1/2$ , this would mean that  $R_0 \approx r_n \frac{N}{S_{SIR}}$  stays below one if  $S_{SIR} \geq N/2$ , i.e. as long as the majority of the population has not been in contact with the SARS-CoV-2 virus. This is good news. But observing a small  $r_n$  can conceal a situation with a large  $R_0$  if  $S_{SIR}/N$  is small. This is one reason why countries need to get a grip on the Susceptibles nationwide.

So far, the above argument cannot replace a SIR model. It only interprets the available data. However, monitoring the Johns Hopkins data in the above way will be very useful when it comes to evaluate the effectiveness of certain measures taken by politicians. It will be highly interesting to see how the data of Figure 5 continue, in particular when countries relax their contact restrictions.

### 3.4 Extension Towards a SIR Model

For cases where one still has to expect  $R_0 > 1$ , e.g. US and Brazil on day 120 (see Figure 5), the challenge remains to predict a possible peak. Using the estimates from

the previous section is impossible, because they concern the sub-population of Confirmed and are systematically underestimating  $R_0$ . The “real” SIR model will have different parameters, a possibly large amount of undetected Infectious, and it needs the Susceptibles to model a peak and to make the  $r_n$  estimates realistic.

For an unrealistic scenario, consider *Total Registration*, i.e. all Infected are automatically confirmed. Then the Susceptibles in the Johns Hopkins model would be  $S_n = N - C_n = N - I_n - R_n - D_n$ . Now the estimate for  $R_0$  must be corrected to

$$r_n \frac{N}{S_n} = r_n \frac{N}{N - C_n} = r_n \left( 1 + \frac{C_n}{N - C_n} \right)$$

but this change will not be serious during an early outbreak.

If the time series  $\beta_n = b_n \frac{N}{S_n} = b_n \frac{N}{N - C_n}$  for  $\beta$  and  $\gamma_n$  for  $\gamma$  are boldly used as predictors for  $\beta$  and  $\gamma$  in a SIR model, and if the model is started using  $S_n = N - C_n = N - I_n - D_n - R_n$  in the discretised form

$$S_{n+1} - S_n = -\beta \frac{S_n}{N} I_n,$$

$$I_{n+1} - I_n = +\beta \frac{S_n}{N} I_n - \gamma I_n,$$

$$(R + D)_{n+1} - (R + D)_n = -\gamma I_n,$$

one gets a crude prediction of the peak in case  $R_0 = \beta/\gamma > 1$ .

Figure 6 shows results for two cases. The **left plot shows** the United States, using data from day 109 (May 10th) and estimating  $\beta$  and  $\gamma$  from the data one week before. The peak is predicted at day 473 (May 9th, 2021) with a total rate of 33% Infectious, i.e. about 124 million people. With an infection fatality rate of 0.5%, this means about 600,000 casualties in the two weeks around the peak. To see how crude the technique is, the second plot shows Germany using data up to day 75 (April 6th, 2020), i.e. before the peak, and the peak is predicted at day 230 (Sept. 8th, 2020) with about 16% Infected. This would imply about 65,000 casualties around the peak. At day 75,  $R_0$  was estimated at 2.01, but a few days later the estimate went below 1 (Figure 5) by political intervention changing  $b_n$  considerably. See Figure 10 for a much better prediction using data only up to day 67.

#### 4 Extended SIR Model

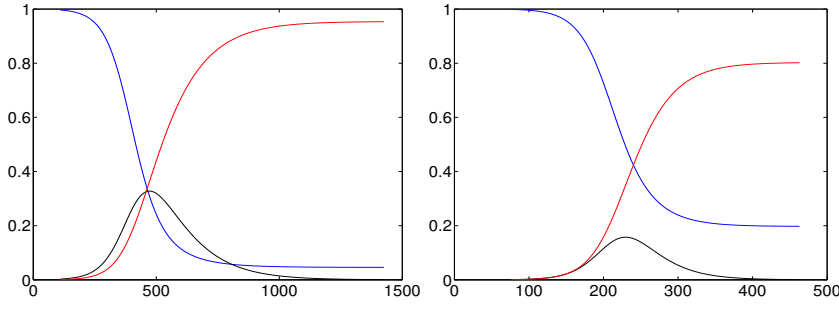
To get closer to reality, one should combine the data-oriented Johns Hopkins Data Model with a SIR model that accounts for what happens outside of the Confirmed. We introduce the time series

$S$  for the Susceptibles like in the SIR model,

$M$  for the Infectious, not yet confirmed, ( $M$  standing for *mysterious*),

$H$  for the unconfirmed Recovered ( $H$  standing for *healed*).

This implies that all deaths occur within the Confirmed, though this is a highly debatable issue. It assumes that persons with serious symptoms get confirmed, and nobody dies of COVID-19 without prior confirmation.



**Fig. 6** Brute force SIR modelling for US and Germany using last week's data, at days 109 and 75, with  $R_0 = 3.22$  and  $R_0 = 2.01$ , respectively.

#### 4.1 The Hidden Model

The Removed from the viewpoint of a global SIR model including  $H$  and  $M$  are  $H + C$ , and thus the SIR model is

$$\begin{aligned} S_{n+1} - S_n &= -\beta \frac{S_n}{N} M_n, \\ M_{n+1} - M_n &= \beta \frac{S_n}{N} M_n - \gamma M_n, \\ (H + C)_{n+1} - (H + C)_n &= \gamma M_n. \end{aligned} \quad (21)$$

To run this *hidden* model with constant  $N = S + M + H + C$ , one needs initial values and good estimates for  $\beta$  and  $\gamma$ , which are not the ones of the Johns Hopkins Data Model of section 3.3. We need other ways to get them.

#### 4.2 The Observable Model

The Johns Hopkins variables  $D$  and  $R$  are linked to the hidden model via  $C = I - R - D$ . They follow an *observable* model

$$\begin{aligned} I_n &= C_n - R_n - D_n, \\ D_{n+1} - D_n &= \gamma_{CD} I_n, \\ R_{n+1} - R_n &= \gamma_{CR} I_n \end{aligned} \quad (22)$$

with *instantaneous case death and recovery rates*  $\gamma_{CD}$  and  $\gamma_{CR}$  for the Confirmed Infectious. These rates can be estimated separately from the available Johns Hopkins data, and we shall do this below. We call these rates *instantaneous*, because they artificially attribute the new deaths or recoveries at day  $n + 1$  to the Infectious of the previous day, not of earlier days. They are *case* rates, because they concern the Confirmed. The difference between standard and instantaneous case rates will be treated in sections 4.3.1 and 4.3.2.

The observable model is coupled to the hidden model only by  $C_n$ . Any data-driven  $C_n$  from the observable model can be used to enter the  $H + C$  variable of the

hidden model, but in an unknown ratio. Conversely, any version of the hidden model produces  $H + C$  values that do not determine the  $C$  part. Summarising, there is no way to fit the hidden model to the data without additional assumptions.

Various possibilities were tried to connect the Hidden to the Observable. Two will be presented now.

### 4.3 Fatality Rates

#### 4.3.1 Infection Fatality Rate

Recall that the parameter  $\gamma_{CD}$  in the observable model (22) relates case fatalities to the confirmed Infectious of the previous day. In contrast to this, the *infection fatality rate* in the standard literature, denoted by  $\gamma_F$  here, is relating to the infection directly, independent of the confirmation, and gives the probability to die of COVID-19 after infection with the SARS-CoV-2 virus, whatever the delay between infection and death is. It was estimated as  $\gamma_F = 0.56\%$  by an der Heiden/Buchholz [10] and  $0.66\%$  by Verity et al. [25], but specialised for China. Recent data of Streeck et. al. [24] gives a value of  $0.36\%$  for the Heinsberg population in Germany. For the UK, Ferguson et al. [6] arrive at  $0.9\%$ . We shall later use  $0.5\%$  for our predictions. But it is very desirable to get more information on infection fatality rates, in particular for different countries. So far, we use a single value globally.

The idea to use the infection fatality rate for information about the hidden system comes from Bommer/Vollmer [1]. The infection fatality rate will be used below in (26) and (28) together with case fatality rates that we consider next.

#### 4.3.2 Estimation of Case Fatality Rates

We now focus on probabilities to die either after an infection or after confirmation of an infection. The first is the infection fatality rate given in the literature, but what is latter, the *case fatality rate*  $\gamma_{CF}$  when using the Johns Hopkins data? It is clearly not the  $\gamma_{CD}$  in (22), giving the ratio of new deaths at day  $n + 1$  as a fraction of the confirmed Infectious at day  $n$ . The deaths at day  $n + 1$  must be assigned to various earlier days instead.

Case fatality rates in the literature vary strongly, and they are country-dependent. Countries have different ways to detect cases, and because the mortality is age-dependent, different age structures will have a serious influence. The Robert-Koch-Institute [21] mentions  $10.5\%$  for Europe and  $4.6\%$  for Germany, while De Brouwer et al. [3] has  $10.0\%$  for Italy,  $4.0\%$  for China,  $6.0\%$  for Spain, and  $4.3\%$  worldwide. According to Streeck et al.[24], the current estimate of the case fatality rate in Germany by the World Health Organization (WHO) is between  $2.2\%$  and  $3.4\%$ .

We cannot clean up these inconsistencies. Instead, we now describe a way to estimate case fatality rates per country from the Johns Hopkins data. The basic idealistic assumption is that COVID-19 diseases end after  $k$  days from confirmation with either death or recovery. Let us call this the *k-day rule*. Suggested values for  $k$  start from 14 days for mild cases (an der Heiden/Buchholz [10] WHO [26]) and go up to 30

days, composed of an incubation time of about 5 days and various values between 11 and 25 days for hospitalisation, depending on the amount of intensive care (an der Heiden/Buchholz [10], Robert Koch-Institut [21], Verity et al. [25], Mohring et al. [19]).

Following Schaback [22], one can estimate the probability to survive on day  $k+1$  after confirmation, and this works in a stable way per country, based only on  $C$  and  $D$ , not on the unstable  $R$  data. In [22] this approach was used to produce  $R$  values that comply with the  $k$ -day rule, but here we use it for estimating the case fatality.

The basic argument lets the new Confirmed  $C_n - C_{n-1}$  at day  $n$  enter into the new deaths  $D_{n+1} - D_n$  at day  $n+1$  with probability  $p_1 =: q_1$ , into  $D_{n+2} - D_{n+1}$  with probability  $p_2(1-p_1) =: q_2$  and so on. The rest enters into the new Recovered at day  $n+k$  with probability  $q_{k+1}$  if we set  $p_{k+1} = 1$  and define

$$q_i = p_i \prod_{j=1}^{i-1} (1 - p_j), \quad 1 \leq i \leq k+1. \quad (23)$$

Then the estimated case fatality rate is  $1 - q_{k+1}$ , while the case recovery rate is  $q_{k+1}$ . Therefore the technique of [22] performs a fit

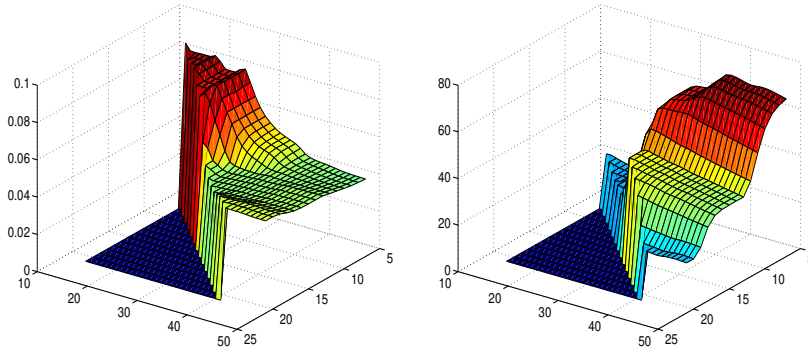
$$D_n - D_{n-1} \approx \sum_{i=1}^k q_i (C_{n-i} - C_{n-i-1}), \quad (24)$$

over all possible probabilities  $p_i$  with sum **bounded by one** connected to the  $q_i$  by (23). It assigns all new deaths at day  $n$  to previous new infections on previous days in a hopefully consistent way, minimising the error in the above formula under variation of the probabilities  $p_i$  to die on day  $i$  after confirmation, and it delivers case fatality and case recovery rates per country. It formally assigns all recoveries to day  $k+1$  after confirmation. Before that day, a living Confirmed cannot be declared to be recovered.

At this point, there is a hidden assumption. The change  $C_{n+1} - C_n$  to the Confirmed is understood as the number of new registered infections, i.e. it is treated like  $I_{n+1} - I_n$ , disregarding short-time death or recovery. But replacing  $C_{n-i} - C_{n-i-1}$  by  $I_{n-i} - I_{n-i-1}$  in (24) would connect a cumulative function to a non-cumulative function. Furthermore, this requires the unsafe data of the Recovered.

In fact, the estimation via the fit (24) is unexpectedly reliable, provided one looks at  $1 - q_{k+1}$  or  $q_{k+1}$ , not at single probabilities  $p_j$ , and if sufficiently many  $n$  are used. This follows from a series of experiments that we do not document fully here, except for Figure 7. In [22], data for  $2k$  days backwards were used for the estimation, and results did not change much when more or less data were used or when  $k$  was modified. Here, the range  $7 \leq k \leq 21$  was tested, and backlogs of up to 50 days from day 109. See Figure 7 below for an example. It is typical here and for many other cases that a value of  $k = 14$  performs well, with a backlog of  $2k = 28$  days for the fit in (24). Using larger  $k$  needs a larger backlog, but then the estimation is not time-local enough to produce up-to-date estimates, because outdated values are used. Figure 7 shows the variation of the case fatality rate estimation when  $k$  and the backlog are varied. The rates usually do not vary much and have plateaus for  $k \geq 14$ , but of course the errors decrease when  $k$  is taken larger, because there are more days to assign deaths to.





**Fig. 7** Left: case fatality rate for Germany based on data at day 109, as functions of  $k$  (right axis) and the data backlog  $B \geq 2k$  (left axis). Right: Root mean-square error for (24).

Country	Death rate	Detection rate
Germany	0.047	0.106
Brazil	0.094	0.053
Italy	0.138	0.036
Spain	0.085	0.059
Sweden	0.157	0.032
Austria	0.052	0.096
France	0.122	0.041
UK	0.145	0.035
US	0.067	0.075

**Table 1** Case fatality and detection rates, estimated on day 109 using the 14-day rule and a backlog of 28 days.

See the first column of Table 1 for estimates of case fatality rates for different countries, calculated on day 109 (May 10th) for  $k = 14$  and a backlog of 28 days. They comply with the values from the literature cited above. Their interpretation depends strongly on the strategy for confirmation. In particular, they are high when only serious cases are confirmed, e.g. cases that need hospital care. If many more people are tested, confirmations will contain plenty of much less serious cases, and then the case fatality rates are low.

The instantaneous case death rate  $\gamma_{iCD}$  of (22) for the Johns Hopkins data comes out around 0.004 for Germany on day 109 by direct inspection of the data via

$$\gamma_{iCD} \approx \frac{D_{n+1} - D_n}{I_n}, \quad (25)$$

while the Case Fatality Rate  $\gamma_{CF}$  in Table 1 is about 0.047. The deaths have to be attributed to different days using the  $k$ -day rule, they cannot easily be assigned to the previous day without making the rate smaller.

### 4.3.3 The Detection Rate

A simple way to understand the quotient  $\frac{\gamma_F}{\gamma_{CF}}$  of the infection fatality rate  $\gamma_F$  and the case fatality rate  $\gamma_{CF}$  as a *detection rate* is to ask for the probability  $p(C)$  for Confirmation. If the probability to die after Confirmation is  $\gamma_{CF}$ , and if there are no deaths outside confirmation, then

$$p(D) = p(C) \cdot p(D|C),$$

by conditional probabilities, and

$$p(C) = \frac{p(D)}{p(D|C)} = \frac{\gamma_F}{\gamma_{CF}}.$$

See the second column of Table 1, prepared for  $\gamma_F = 0.005$ . The rate depends on good estimates of the infection fatality rate, and the new value 0.0036 by Streeck et al. [24] will decrease the detection rate for Germany from 10.6% to 7.7% for the Heinsberg subpopulation.

All of this is comparable to the findings of Bommer/Vollmer [1] and uses the basic idea from there, but with a somewhat different technique and different results. There, the values were 7% for March 23rd and 9% for March 30th, while Mohring et al. [19] assume 20% on April 29th.

### 4.3.4 Using Fatality Rates for the Hidden Model

If the case fatality rates  $\gamma_{CF}$  of Table 1 are used with a known infection fatality rate  $\gamma_F$ , one should obtain an estimate of the total Infectious. If the formula (24) is written as

$$\sum_{i=1}^k q_i (C_{n-i} - C_{n-i-1}) \approx D_n - D_{n-1} \approx \sum_{i=1}^k \tilde{q}_i (S_{n-i-1} - S_{n-i})$$

in terms of the previous new infections  $S_{n-i-1} - S_{n-i}$  in terms of Susceptibles with daily infection fatality probabilities  $\tilde{q}_i$ , one should maintain

$$\gamma_{CF} = \sum_{i=1}^k q_i \text{ and } \gamma_F = \sum_{i=1}^k \tilde{q}_i,$$

and this works by setting

$$C_n - C_{n-1} = \frac{\gamma_F}{\gamma_{CF}} (S_{n-1} - S_n) \quad (26)$$

in general, without using the unstable  $p_i$ . This is the first connection of the Observable to the Hidden, namely  $C$  to  $S$ . Like in the discussion following (24) one can argue to use  $M$  instead of  $S$  in (26), but this would again connect a cumulative variable to a non-cumulative one.

### 4.3.5 Local Estimation of Fatality Rates

Because politicians change testing strategies and the parameters  $\beta$  and  $\gamma$ , the estimation of the Case Fatality Rate should be made locally, not globally. Using the experience of Schaback [22] and section 4.3.2, we shall use a fixed  $k = 14$  for the  $k$ -day rule and data for a fixed backlog of  $2k$  days. Then the formula (26) has  $\gamma_{CF}$  varying with  $n$  as far as Johns Hopkins data are available.

## 4.4 Recovery Rates

We need another parameter to connect the hidden to the observable model. There are many choices, and after some failures we selected the constant  $\gamma_{iIR}$  in a model equation

$$H_{n+1} - H_n = \gamma_{iIR} M_n.$$

Following what was mentioned about *instantaneous* rates in section 4.2,  $\gamma_{iIR}$  is an *instantaneous Infection Recovery Rate*, relating the new unregistered Recovered to the unregistered Infections the day before.

### 4.4.1 Estimation of Recovery Rates

A good value of  $\gamma_{iIR}$  can come out of a field experiment that produces time series for  $M$  and  $H$ , i.e. for unregistered Infectious and unregistered Recovered. Then the instantaneous Infection Recovery rate  $\gamma_{iIR}$  can be obtained directly by

$$\frac{H_{n+1} - H_n}{M_n} \approx \gamma_{iIR}.$$

The Infection Recovery rate  $\gamma_{IR} = 1 - \gamma_{IF}$  does not help, because we need an instantaneous rate that has no interpretation as a probability.

With the risk of using unstable data of the Recovered, we can look at the instantaneous Case Recovery rate

$$\frac{R_{n+1} - R_n}{I_n} \approx \gamma_{iCR} \quad (27)$$

that is available from the Johns Hopkins data, and comes out experimentally to be rather stable, provided that countries have useful data for the Recovered. Otherwise, we have to use the technique of Schaback [22] for estimating them. The rate  $\gamma_{iIR}$  must be larger than  $\gamma_{iCR}$  because we now are not in the subpopulation of the Confirmed, and nobody can die without going first into the population of the Confirmed. As long as no better data are available, we shall use the formula

$$\gamma_{iIR} = \frac{1 - \gamma_{IF}}{1 - \gamma_{CF}} \gamma_{iCR} = \frac{\gamma_{IR}}{\gamma_{CR}} \gamma_{iCR} = \frac{\gamma_{iCR}}{\gamma_{CR}} \gamma_{IR} \quad (28)$$

that implements two meaningful arguments:

1. the value  $\gamma_{iCR}$  is increased by the ratio  $\frac{\gamma_{IR}}{\gamma_{CR}}$  of Recovered probabilities for the Infected and the Confirmed,

2. the value  $\gamma_R$  is multiplied by a factor  $\frac{\gamma_{CR}}{\gamma_{CR}}$  for transition to immediate rates, and this factor is the transition factor for the Confirmed Recovered.

The above strategy is debatable and may be the weakest point of this approach. However, others turned out to be worse, mainly due to instability of results. On the positive side, the final prediction will not need it, see (33) below. It enters only the intermediate step when  $S$ ,  $M$ , and  $H$  are calculated in the time range of the available Johns Hopkins data, see (29) in section 4.5. And, finally, there is hope that there will be field experiments that yield reliable values directly.

#### 4.4.2 Practical Approximation of Recovery Rates

In (28) the rate  $\gamma_R$  is fixed, and the rate  $\gamma_{CR}$  is determined locally via section 4.3.5. The rate  $\gamma_{CR}$  follows from the time series

$$\frac{R_{n+1} - R_n}{I_n} \approx \gamma_{CR}$$

as in (22). This works for countries that provide useful data for the Recovered. In that case, and in others to follow below, we can take the time series itself as long as we have data. For prediction, we estimate the constant from the time series using a fixed backlog of  $m$  days from the current day, i.e. we take the mean of the last  $m + 1$  values. Since many data have a weekly oscillation, due to data being not properly delivered during weekends, the backlog should not be less than 7.

But for certain countries, like the United Kingdom, the data for the Recovered are useless. In such cases, we employ the technique of Schaback [22] to estimate the Recovered using the  $k$ -day rule and a backlog of  $2k$  days, like in section 4.3.5 for the case fatality rates.

### 4.5 Model Calibration

We now have everything to run the hidden model, but we do it first for days with Johns Hopkins data, delaying predictions to section 5. This is a *calibration step* that leads to estimations of  $S$ ,  $M$ , and  $H$  from the observed data of the Johns Hopkins source, without any need for sophisticated fitting algorithms. With the parameters from above, we use the new relations

$$\begin{aligned} C_{n+1} - C_n &= \frac{\gamma_{IF}}{\gamma_{CF}}(S_n - S_{n+1}), \\ H_{n+1} - H_n &= \gamma_{IR}M_n \end{aligned} \quad (29)$$

in a specific way. We set up the second model equation in (21) for  $M$  as

$$\begin{aligned} M_{n+1} - M_n &= S_n - S_{n+1} - \gamma_n M_n \\ &= \frac{\gamma_{CF}}{\gamma_{IF}}(C_{n+1} - C_n) - \gamma_n M_n \\ &= \frac{\gamma_{CF}}{\gamma_{IF}}(C_{n+1} - C_n) - (C_{n+1} - C_n + H_{n+1} - H_n) \\ &= \left( \frac{\gamma_{CF}}{\gamma_{IF}} - 1 \right) (C_{n+1} - C_n) - \gamma_{IR}M_n \end{aligned} \quad (30)$$

that can be solved if an initial value  $\tilde{M}_0$  is prescribed. Then (29) is run to produce the  $S_n$  and  $H_n$ , with starting values that we describe in section 4.5.1. If  $\beta_n$  and  $\gamma_n$  are calculated by

$$\begin{aligned}\beta_n \frac{S_n}{N} M_n &= S_n - S_{n+1}, \\ \gamma_n M_n &= C_{n+1} - C_n + H_{n+1} - H_n,\end{aligned}\quad (31)$$

respectively, the balance equation  $N = S + M + H + C$  follows from (30) and (31).

#### 4.5.1 Starting Values

Since the populations are large, the starting values for  $S$  are not important. Beginning at the full population  $N$  from a very early day, the  $S$  values are calculated from (29) first, just to get values  $S_j$  for actually starting at later days.

Then the first day  $j$  is taken where  $C_j$  is at least 10, and  $k$  days later the start value for  $H$  is set as

$$H_{j+k} = C_j \frac{\gamma_{CF}}{\gamma_F} \quad (32)$$

using the  $k$ -day rule with  $k = 14$ . This divides the  $C_j > I_j$  value by the detection rate, i.e. roughly all estimated undetected Infectious at time  $j$  are assumed to be healed  $k$  days later, i.e. at day  $j+k$ . Then the starting value for  $M_{j+k}$  is calculated via the balance equation  $N = S + M + H + C$  from the  $S_{j+k}$  value calculated by the previous paragraph. Finally, the calibration starts at day  $j+k$  by the above formulae. Unfortunately, this is a serious limit preventing application to very short time series.

The starting value for  $H$  is irrelevant for  $H$  itself, because only differences enter into the model, but it determines the starting value for  $M$  due to the balance equation. Anyway, it turns out experimentally that the starting values do not matter much, if the model is started early. The hidden model (21) depends much more strongly on  $C$  than on the starting values.

Figure 10 contains a wide variation of the starting value (32) for  $H$  at the starting point, by multipliers between  $1/32$  and  $32$ . This has hardly any effect on the results, the lines getting somewhat thicker. The variation in starting values get more visible in other cases, see the right-hand plot in Figure 10 for the United States. But the influence on predictions is negligible.

#### 4.5.2 Examples

The figures to follow in section 5.2 show the original Johns Hopkins data together with the hidden variables  $S$ ,  $M$ , and  $H$  that are calculated by the above technique. The calibration runs up to the vertical line where predictions start. Note that the only ingredients beside the Johns Hopkins data are the number  $k = 14$  for the  $k$ -day rule, the Infection Fatality Rate  $\gamma_F$  from the literature, equations (29), and the backlog of  $m = 7$  days for estimation of constants from time series.

## 5 Predictions using the Full Model

To let the combined model predict the future, or to check what it would have predicted if used at an earlier day, we take the calibrated model of the previous sections up to a day  $n$  and use the values  $S_n, M_n, H_n, C_n, I_n, R_n$  and  $D_n$  for starting the prediction. With the variable  $HC := H + C$ , we use the recursion

$$\begin{aligned}
 S_{i+1} &= S_i - \beta \frac{S_i}{N} M_i, \\
 M_{i+1} &= M_i + \beta \frac{S_i}{N} M_i - \gamma M_i, \\
 HC_{i+1} &= HC_i + \gamma M_i, \\
 C_{i+1} &= C_i + \gamma_F (S_i - S_{i+1}) / \gamma_{CF}, \\
 R_{i+1} &= R_i + \gamma_{CR} I_i, \\
 D_{i+1} &= D_i + \gamma_{CD} I_i, \\
 I_{i+1} &= C_{i+1} - R_{i+1} - D_{i+1}, \\
 H_{i+1} &= HC_{i+1} - C_{i+1}.
 \end{aligned} \tag{33}$$

This needs fixed values of  $\beta$  and  $\gamma$  that we estimate from the time series for  $\beta_n$  and  $\gamma_n$  by using a backlog of 7 days, following Section 4.5. The instantaneous rates  $\gamma_{CR}$  and  $\gamma_{CD}$  can be calculated via their time series, as in (27) and (25), using the same backlog. We do this at the starting point of the prediction, and then the model runs in a *no political change* mode. Examples will follow in section 5.2.

### 5.1 Properties of the Full Model

The first part of the full model (33) is a standard SIR model for the variables  $S, M$  and  $H + C$ , and inherits the properties of these as described in section 2. It does not use the  $\gamma_{IR}$  parameter of the second equation in (29), and it uses the first the other way round, now determining  $C$  from  $S$ , not  $S$  from  $C$ .

The balances  $N = S + M + H + C$  and  $C = I + D + R$  are maintained automatically, and the time series for  $S, C, R, H + C$ , and  $D$  stay monotonic as long as  $M$  and  $I$  are non-negative. To check the monotonicity of  $H$ , consider

$$\begin{aligned}
 H_{i+1} - H_i &= HC_{i+1} - HC_i - C_{i+1} + C_i \\
 &= \gamma M_i - \frac{\gamma_F}{\gamma_{CF}} (S_i - S_{i+1}) \\
 &= \left( \gamma - \beta \frac{\gamma_F}{\gamma_{CF}} \frac{S_i}{N} \right) M_i.
 \end{aligned}$$

The bracket is positive if

$$R_0 = \frac{\beta}{\gamma} < \frac{\gamma_{CF}}{\gamma_F} \frac{N}{S_i} \geq \frac{\gamma_{CF}}{\gamma_F},$$

which is enough for practical purposes as long as detection rates  $\frac{\gamma_F}{\gamma_{CF}}$  are low and  $R_0$  is not excessively large. Anyway,  $H$  should be monitored.

The slopes of  $S$  and  $C$  are always connected by (26), and those of  $R$  and  $D$  are connected by

$$R_{i+1} - R_i = \frac{\gamma_{iCR}}{\gamma_{iCD}} (D_{i+1} - D_i) \quad (34)$$

in the prediction part. But the figures below will show logarithms, and therefore the slope parallelism will not be visible.

By section 2.11, the hidden Infectious  $M$  will always go to zero, and the variables  $S$  and  $H + C$  will level out in the long run. Since  $C$  is increasing, it must level out as well, and  $I$  must level out because  $R$  and  $D$  do. But due to the equations for  $R$  and  $D$ , the final level of  $I$  must be zero.

The asymptotic levels of  $S$  and  $H + C$  follow from 2.11, but not the interesting level of  $D$ , the total death toll. If the prediction is started at day  $n$ , then

$$R_\infty - R_n = \frac{\gamma_{iCR}}{\gamma_{iCD}} (D_\infty - D_n),$$

obtained by summation of (34), connects the asymptotic deaths and confirmed recoveries. From the connection of  $S$  and  $C$  we likewise get

$$C_\infty - C_n = \frac{\gamma_{iIF}}{\gamma_{iCF}} (S_n - S_\infty).$$

With  $C_\infty = R_\infty + D_\infty$  we now have three independent equations for the unknowns  $C_\infty, D_\infty, R_\infty$ . Because the theory of Section 2.11 yields  $S_\infty$  and  $H_\infty + C_\infty$  in terms of  $\beta$  and  $\gamma$ , we know  $S_\infty$  and can get  $H_\infty$  from  $C_\infty$ . But if the simulation is run long enough, one can easily read the asymptotic values off the plots.

## 5.2 Examples of Predictions

Figure 8 shows predictions on day 122, May 23rd, for Germany, Brazil, France, and USA, from the top. The plots for countries behind their peak are rather similar to those for Germany and France. The other two countries are selected because they still have to face their peak, if no action is taken to change the parameters.

The plots show that Germany can expect to get away with no more than 10000 casualties in the long run, while Brazil goes for a peak of about 20 million hidden Infectious in fall 2020 ( $M$ , symbol  $\square$ ) and a final death toll of about 1 million ( $D$ , symbol  $+$ ). The United States would have to face a peak of hidden Infectious of about 25 million in mid-January 2021, and more than 1 million COVID-19 deaths in October 2021, and still rising. But of course, these predictions assume that reality follows the model and that there are no parameter changes by political action.

The estimated  $R_0$  values are 0.65, 2.19, 0.42, and 1.75, respectively. Note that these are not directly comparable to Figure 5, because they are the fitted constants to the backlog of a week, and using (31) instead of (20), avoiding the systematic underestimation of the latter. The hidden  $M$  and  $H$  (symbols  $\square$  and  $\diamond$ ) follow roughly the observable  $I$  and  $C$  (symbols  $O$  and  $x$ ), but with a factor due to the detection rate that is different between countries, see Table 1. To enhance visibility, not all data points in the plots are marked with symbols. The  $C, R, I$  and  $D$  data left of the vertical

line are the original Johns Hopkins data. The  $S$ ,  $M$ ,  $H$  data there are calculated by section 4, while to the right the data are predictions for all variables by the full model (33).

All test runs were made for the infection fatality rate  $\gamma_F = 0.005$ , the delay  $k = 14$  for estimating case fatalities, and a backlog of 7 days when estimating constants out of recent values of time series. The choice  $\gamma_F = 0.005$  is somewhat between 0.56% from an der Heiden/Buchholz [10], 0.66% from Verity et al. [25], and 0.36% from Streeck et al. [24]. New information on infection fatality rates should be included as soon as they are available, and if possible per country, not global.

When used within estimation routines, the Johns Hopkins data were smoothed by a double application of the 1/4, 1/2, 1/4 filter on the logarithms, like for Figure 5. But the plots show the original Johns Hopkins and prediction data.

### 5.3 Evaluation of Predictions

To evaluate the prediction quality, one should go back and start the predictions on earlier days, to compare with what happened later. Figure 9 shows over-plots of predictions for days 94, 108, and 122, each a fortnight apart, though there may be parameter changes in the meantime. The starting points of the predictions are marked by vertical lines again. For better visibility, only the death count  $D$  (symbol  $+$ ) and the two non-cumulative variables  $M$  and  $I$  for the hidden and confirmed Infectious (symbols  $\square$  and  $O$ ) are shown. In particular, the case fatality rates and detection rates of Table 1 change with the starting point of the prediction, and they determine  $S$ ,  $M$ , and  $H$  in the calibration step of section 4.5. This is why the  $S$ ,  $M$ , and  $H$  values differ left of the starting points.

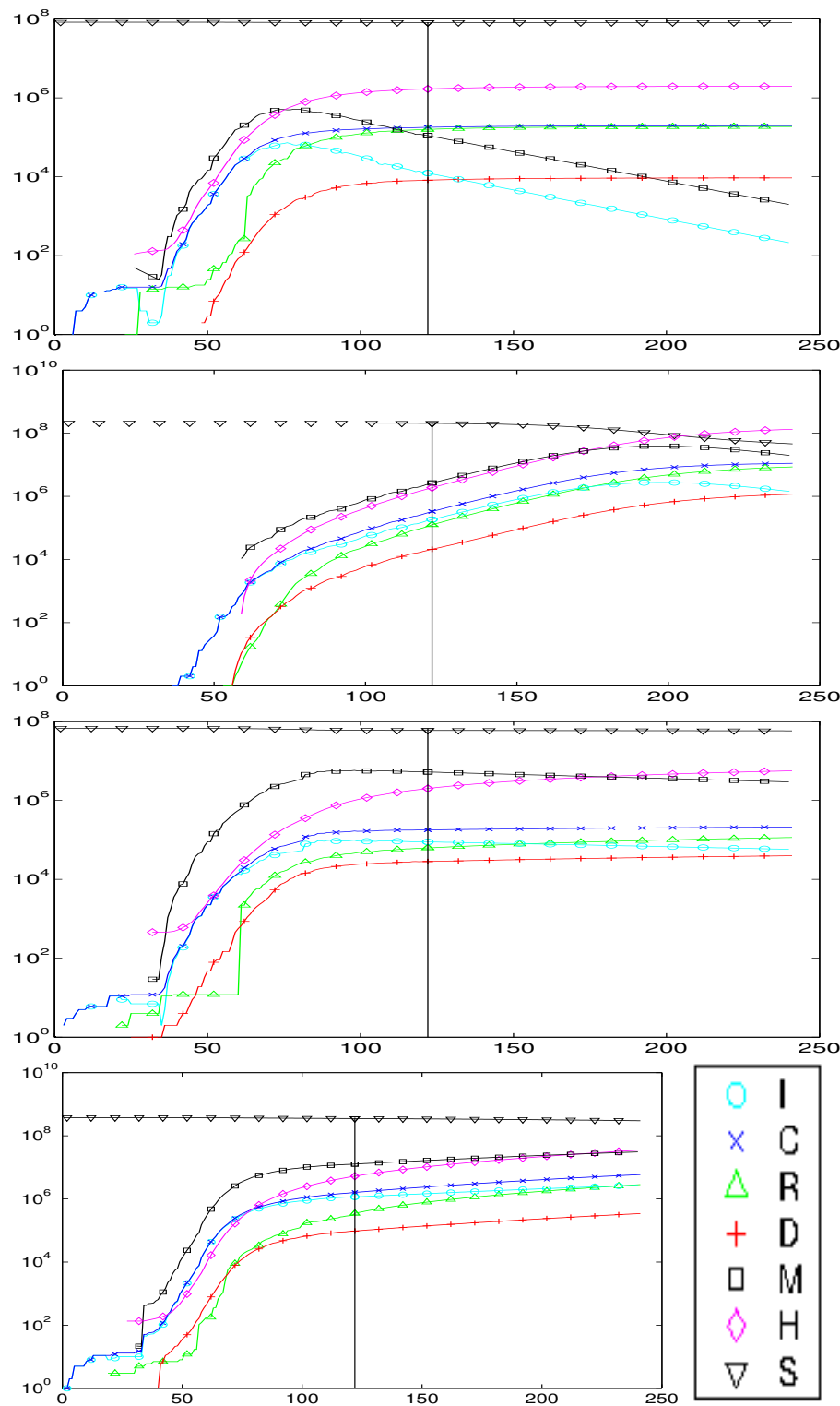
The leftmost prediction on day 94 roughly matches the data available up to day 122 in all cases. It has to be taken into account that errors in such models must proliferate exponentially, and then linearly in logarithmic plots. One can see that the Brazil parameters do not change much, while the three predictions for the United States get better. This might be used to assess effectivity of administrative efforts to handle the pandemics.

For an early case in Germany, Figure 10 shows the prediction based on data of day 67, March 27th. The peak of about 35 million hidden and 3.2 million confirmed Infected is predicted on day 121, May 22nd, with about 82,000 casualties at the peak and about 250,000 finally. A good reason to act politically. Note that the real death count is about 8300 on May 23rd, and the prediction of the day, in Figure 8, targets a final count of below 10,000.

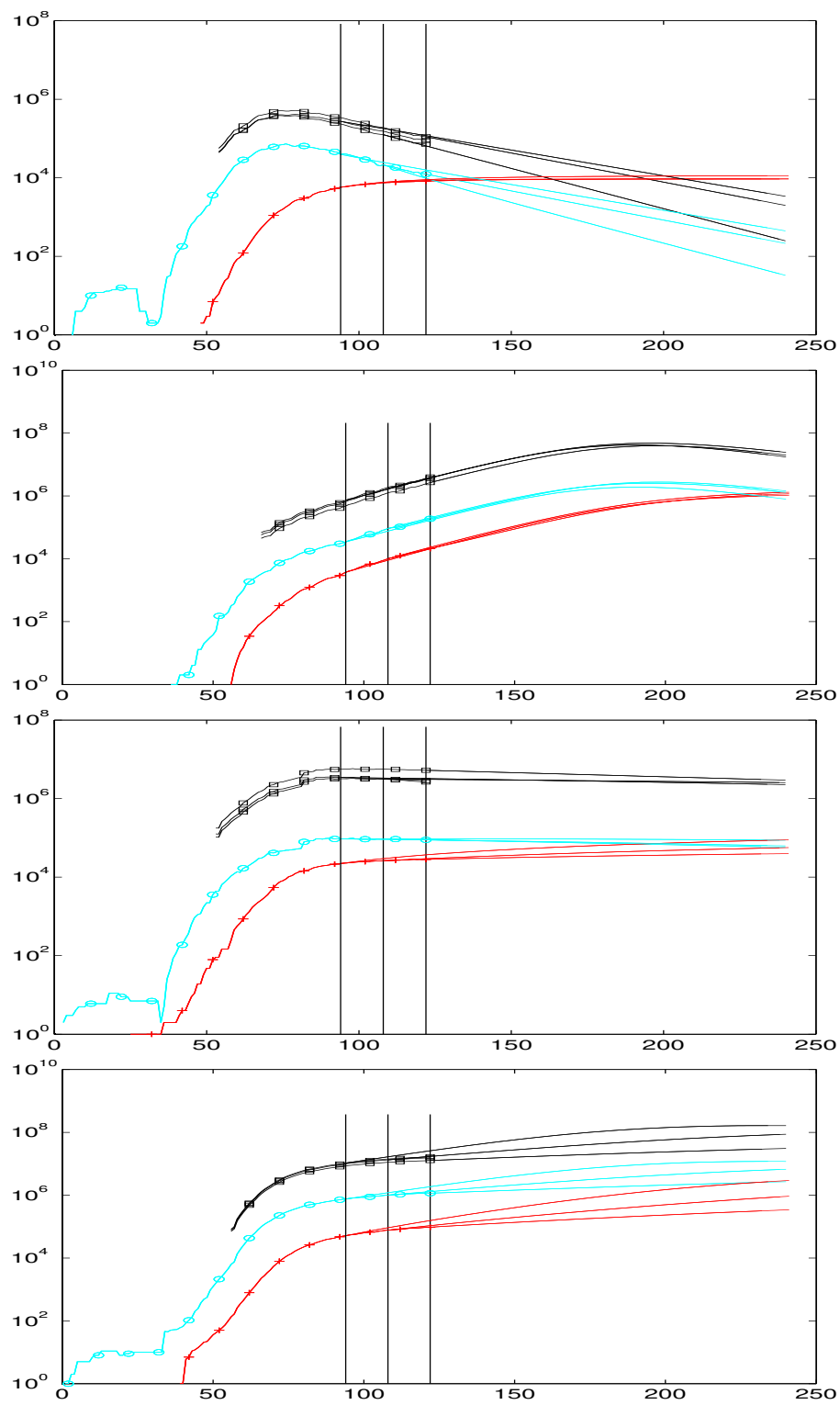
Quantitative commitments to predictions are rare in the literature, except for rough estimations of dramatic outbreak scenarios. On April 3rd, after the last public restrictions in Germany of March 22nd, 2020, Germany had 1107 deaths and Khailaie et al. [16] predicted “an order of 10,000 deaths” for the next four weeks. This model predicts 15,500 for May 3rd when run on data of April 3rd, while the true deaths were 6812 on May 3rd, after the interventions worked.

On March 16th, day 54, Ferguson et al. [6] predicted deaths in the order of 250,000 in Great Britain, and 1.1 to 1.2 million in the USA “in the most effective





**Fig. 8** Predictions for countries Germany, Brazil, France, and US on day 122 marked by the vertical line. The  $S, M, H$  values to the left are obtained by calibration, the  $C, R, D, I$  values there are the original Johns Hopkins data. Not all data points have marks.



**Fig. 9** Predictions for countries Germany, Brazil, France, and USA on days 122, 108, and 94, marked by vertical lines. Legend as in Figure 8, but only  $M$ ,  $I$ , and  $D$  shown ( $M=\square, I=O, D=+$ ).

mitigation strategy examined”, but not based on the data of the day. In an “unmitigated epidemic” 520,000 deaths in the UK and 2.2 million in the US were predicted, under assumption of  $R_0 = 2.4$  and a range of  $R_0$  tested between 2.2 and 2.6. Unfortunately, the model (33) cannot be safely run on day 54 for these countries because there are not enough reliable backlog data. The model can be run if the amount of data used is cut down by choosing  $k = 7$  for the  $k$ -day rule. Then the predictions on day 54 are more than 30 million deaths for the US and 801,000 for the UK, with a data-based estimation of  $R_0 = 6.06$  for the USA and 4.55 for the UK. There is no reasonable data-driven estimate for  $R_0$  that comes close to  $R_0 = 2.4$  used by Ferguson et. al. [6] for both countries. They had a much more serious outbreak than assumed by Ferguson et al. on March 16th. See Figure 5 for much later data-based estimates for the US that still are very large.

The use of the Infection Fatality Rate is somewhat different from Streeck et al. [24] and Bommer/Vollmer [1], but results are similar. If the rate 0.0036 of [24] is used in a test run based on data of May 2nd, the estimated number  $M_n + C_n$  of total Infected comes out as 1.7 million, while [24] gets 1.8 million by the formula  $D_n/0.0036$  for the same day.

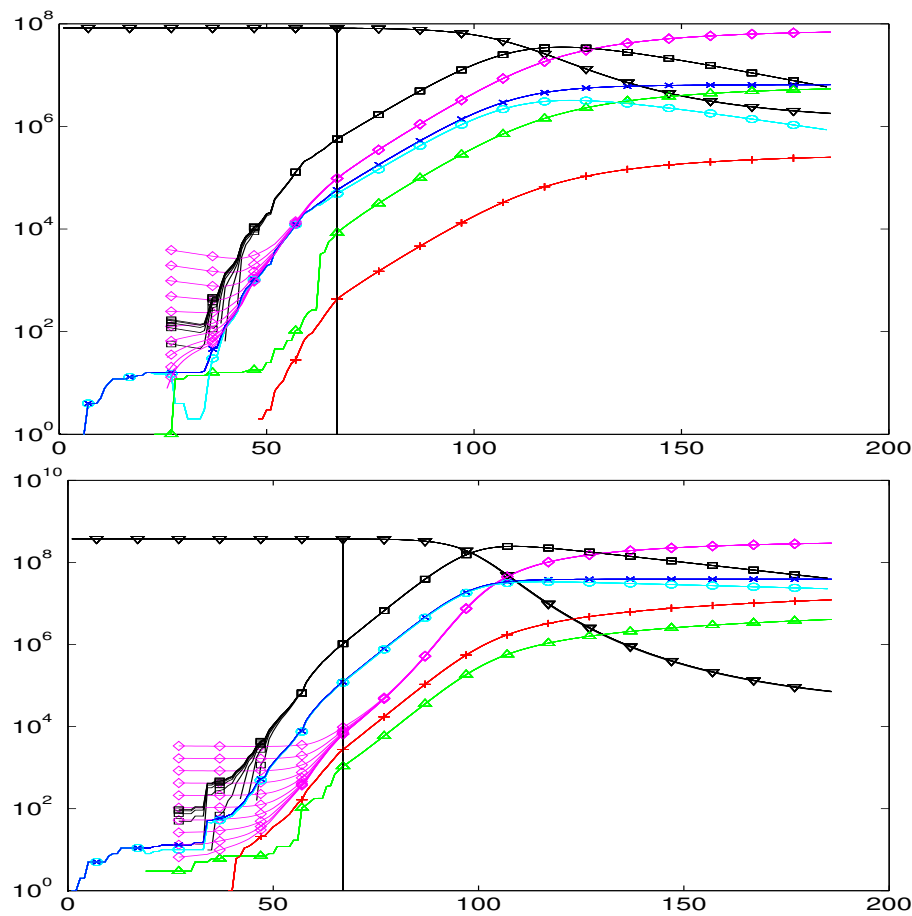
The parameter changes by political measures turned out to be rather effective, like in many countries that applied similar strategies. But since parts of the population want to go back to their previous lifestyle, all of this is endangered, and the figures should be monitored carefully.

Of course, all of this makes sense only under the assumption that reality follows the model, in spite of all attempts to design a model that follows reality.

## 6 Conclusion and Open Problems

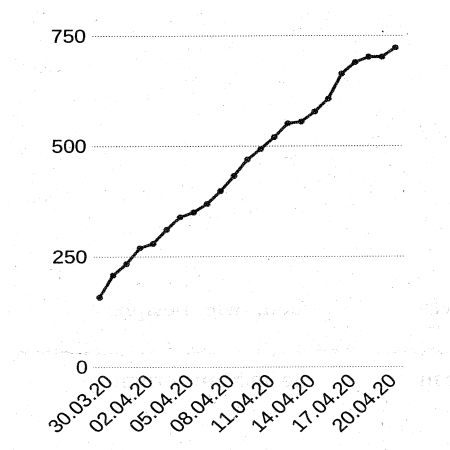
So far, the model presented here seems to be useful, combining theory and practically available data. It is data-driven to a very large extent, using only the infection fatality rate from outside for prediction, and the approximation (28) for calibration. On the downside, there is quite a number of shortcomings:

- Like the data themselves, the model needs regular updating. As far as the Johns Hopkins data are concerned, the model updates itself by using the latest data for its internal parameter estimation, but it needs changes as soon as new information on the hidden infections come in.
- There may be better ways of estimating the hidden part of the epidemics. However, it will be easy to adapt the model to other parameter choices. If time series for the unknown variables get available, the model can easily be adapted to being data-driven by the new data.
- The treatment of delays is unsatisfactory. In particular, infected persons get infectious immediately, and the  $k$ -day rule is not followed at all places in the model. But the rule is violated as well in the data (Schaback [22]).
- There is no stochastics involved, except for simple things like estimating constants by means, or for certain probabilistic arguments on the side, e.g. in section 4.3.2. But it is not at all clear whether there are enough data to do a proper probabilistic analysis.



**Fig. 10** Predictions for Germany and USA on day 67, March 27th, with varying starting values. Legend as in Figure 8.

- As long as there is no probabilistic analysis, there should be more simulations under perturbations of the data and the parameters. A few were included, e.g. for section 4.3.2 and Figures 7 and 10, but a large number was performed in the background when preparing the paper, making sure that results are stable. However, there are never too many test simulations.
- Totally different models were not considered, e.g. the classical ones with delays (Hoppenstaedt/Waltman [13,14]), and agent-based approaches (Ferguson et al. [6]) that model infections via contacts and can care for spatial distributions.
- The model needs quite an amount of backward data, making it useless at the very beginning of an outbreak.
- Under certain circumstances, epidemics do not show an exponential outbreak, in particular if they hit only locally and a prepared population. See Figure 11 for the COVID-19 cases in Göttingen and vicinity.



**Fig. 11** Infectious in Göttingen city and county, as of April 22nd, 2020 in the local newspaper “Göttinger Tageblatt”. No exponential outbreak.

**Acknowledgements** MATLAB programs and more recent predictions will be on the research website <http://num.math.uni-goettingen.de/schaback/research/group.html> of the author. Special thanks go to Tara Fickle, Reiner Kree, Viola Priesemann, Jalda Schaback, and Wolfgang Warth for various forms of input. All links in the references were verified on June 2nd, 2020.

### Conflict of interest

The author declares no conflict of interest.

### References

1. Bommer, C., Vollmer, S.: Average detection rate of SARS-CoV-2 infections is estimated around six percent (April 2nd, 2020). <https://reason.com/wp-content/uploads/2020/04/Bommer-Vollmer-2020-COVID-19-detection-April-2nd.pdf>
2. Dandekar, R., Barbastathis, G.: Quantifying the effect of quarantine control in Covid-19 infectious spread using machine learning. <https://www.medrxiv.org/content/10.1101/2020.04.03.20052084v1> (April 6th, 2020). DOI:10.1101/2020.04.03.20052084
3. De Brouwer, E., Raimondi, D., Moreau, Y.: Modeling the COVID-19 outbreaks and the effectiveness of the containment measures adopted across countries (April 19th, 2020). DOI: 10.1101/2020.04.02.20046375
4. Dehning, J., Spitzner, P., Linden, M., Mohr, S., Zierenberg, J., Wibral, M., Wilczek, M., Priesemann, V.: Model-based and model-free characterization of epidemic outbreaks - Technical notes on Dehning et al., Science, 2020 (June 22, 2020)
5. Dehning, J., Zierenberg, J., Spitzner, P., Wibral, M., Pinheiro Neto, J., Wilczek, M., Priesemann, V.: Inferring covid-19 spreading rates and potential change points for case number forecasts (May 4th, 2020). ArXiv:2004.01105v2
6. Ferguson, N., Laydon, D., Nedjati-Gilani, G., Imai, N., Ainslie, K., Baguelin, M., Bhatia, S., Boonyasiri, A., Cucunubá, Z., Cuomo-Dannenburg, G., Dighe, A., Dorigatti, I., Fu, H., Gaythorpe, K., Green, W., Hamlet, A., Hinsley, W., Okell, L., van Elsland, S., Thompson, H., Verity, R., Volz, E.,

- Wang, H., Wang, Y., Walker, P., Walters, C., Winskill, P., Whittaker, C., Donnelly, C., Riley, S., Ghani, A.: Impact of non-pharmaceutical interventions (NPIs) to reduce COVID-19 mortality and healthcare demand. DOI: 10.25561/77482. Imperial College London (16-03-2020)
7. Friston, K.J., Parr, T., Zeidman, P., Razi, A., Flandin, G., Daunizeau, J., Hulme, O.J., Billig, A.J., Litvak, V., Moran, R.J., Price, C.J., Lambert, C.: Dynamic causal modelling of covid-19 (April 9th, 2020)
  8. Giordano, G., Blanchini, F., Bruno, R., Colaneri, P., Di Filippo, A., Di Matteo, A., Colaneri, M.: Modelling the COVID-19 epidemic and implementation of population-wide interventions in Italy. <https://doi.org/10.1038/s41591-020-0883-7> (April 22nd, 2020). [www.nature.com/naturemedicine](http://www.nature.com/naturemedicine)
  9. GitHub: Covid-19 repository at github. [https://github.com/CSSEGISandData/COVID-19/tree/master/csse\\_covid\\_19\\_data/csse\\_covid\\_19\\_time\\_series](https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data/csse_covid_19_time_series) (2020)
  10. An der Heiden, M., Buchholz, U.: Modellierung von Beispielszenarien der SARS-CoV-2-Epidemie 2020 in Deutschland. Bekanntmachungen des Robert Koch-Instituts (March 3rd, 2020). DOI: 10.25646/6571.2, <https://edoc.rki.de/handle/176904/6547.2>
  11. An der Heiden, M., Hamouda, O.: Schätzung der aktuellen Entwicklung der SARS-CoV-2-Epidemie in Deutschland - Nowcasting. *Epidemiologisches Bulletin* **17**, 10–15 (2020)
  12. Höhle, M., an der Heiden, M.: Bayesian nowcasting during the STEC 0104:H4 outbreak in Germany, 2011. *Biometrics* pp. 993–1002 (2014)
  13. Hoppensteadt, F., Waltman, P.: A problem in the theory of epidemics I. *Math. Biosci.* **9**, 71–91 (1970)
  14. Hoppensteadt, F., Waltman, P.: A problem in the theory of epidemics II. *Math. Biosci.* **12**, 133–145 (1971)
  15. Kermack, W., McKendrick, A.: A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society A* **115**, 700–721 (1927)
  16. Khailaie, S., Mitra, T., Bandyopadhyay, A., Schips, M., Mascheroni, P., Vanella, P., Lange, B., Binder, S., Meyer-Hermann, M.: Estimate of the development of the epidemic reproduction number  $r_t$  from Coronavirus SARS-CoV-2 case data and implications for political measures based on prognostics (April 7th, 2020). DOI: 10.1101/2020.04.04.20053637, <https://www.medrxiv.org/content/10.1101/2020.04.04.20053637v1>
  17. Kucharski, A., Russell, T., Diamond, C., Liu, Y., Edmunds, J., Funk, S., Eggo, R.: Early dynamics of transmission and control of COVID-19: a mathematical modelling study. *Lancet Infect Dis.* **20**, 553–558 (May 1st, 2020). DOI: 10.1016/S1473-3099(20)30144-4
  18. Maier, B.F., Brockmann, D.: Effective containment explains subexponential growth in recent confirmed COVID-19 cases in China. *Science* **368**(6492), 742–746 (April 2020). DOI 10.1126/science.abb4557. URL <http://dx.doi.org/10.1126/science.abb4557>
  19. Mohring, J., Wegener, R., Gramsch, S., Schöbel, A.: Prognosemodelle für die Corona-Pandemie (April 29th, 2020). Fraunhofer-Institut für Techno- und Wirtschaftsmathematik ITWM Kaiserslautern, [https://www.itwm.fraunhofer.de/content/dam/itwm/de/documents/PressemitteilungenPDF/2020/20200429\\_Bericht\\_Prognosemodelle-für-die-Coronapandemie.pdf](https://www.itwm.fraunhofer.de/content/dam/itwm/de/documents/PressemitteilungenPDF/2020/20200429_Bericht_Prognosemodelle-für-die-Coronapandemie.pdf)
  20. Robert-Koch-Institut: Erläuterung der Schätzung der zeitlich variierenden Reproduktionszahl  $r_t$ . Tech. rep., Robert-Koch-Institut (2020). 15.05.2020, [https://www.rki.de/DE/Content/InfAZ/N/Neuartiges\\_Coronavirus/Projekte\\_RKI/R-Wert-Erlaeuterung.html](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/R-Wert-Erlaeuterung.html)
  21. Robert-Koch-Institut: SARS-CoV-2 Steckbrief zur Coronavirus-Krankheit-2019 (COVID-19). Tech. rep., Robert-Koch-Institut (2020). 22.05.2020, [https://www.rki.de/DE/Content/InfAZ/N/Neuartiges\\_Coronavirus/Steckbrief.html](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Steckbrief.html)
  22. Schaback, R.: Modelling recovered cases and death probabilities for the COVID-19 outbreak (March 26th, 2020). URL <https://arxiv.org/abs/2003.12068>
  23. Sentker, A.: Bloß raus hier! *DIE ZEIT* (2020). April 16th, 2020, page 20
  24. Streeck, H., Schulte, B., Kuemmerer, B., Richter, E., Hoeller, T., Fuhrmann, C., Bartok, E., Dolscheid, R., Berger, M., Wessendorf, L., Eschbach-Bludau, M., Kellings, A., Schwaiger, A., Coenen, M., Hoffmann, P., Noethen, M., Eis-Huebinger, A.M., Exner, M., Schmithausen, R., Schmid, M., Hartmann, G.: Infection fatality rate of SARS-CoV-2 infection in a German community with a super-spreading event (May 8th, 2020). DOI 10.1101/2020.05.04.20090076
  25. Verity, R., Okell, L.C., Dorigatti, I., Winskill, P., Whittaker, C., Imai, N., Cuomo-Dannenburg, G., Thompson, H., Walker, P.G.T., Fu, H., Dighe, A., Griffin, J.T., Baguelin, M., Bhatia, S., Boonyasiri, A., Cori, A., Cucunubá, Z., FitzJohn, R., Gaythorpe, K., Green, W., Hamlet, A., Hinsley, W., Laydon, D., Nedjati-Gilani, G., Riley, S., van Elsland, S., Volz, E., Wang, H., Wang, Y., Xi, X., Donnelly, C.A.,

- 
- Ghani, A.C., Ferguson, N.M.: Estimates of the severity of coronavirus disease 2019: a model-based analysis (June 1st, 2020). URL [https://doi.org/10.1016/S1473-3099\(20\)30243-7](https://doi.org/10.1016/S1473-3099(20)30243-7)
26. WHO: Report of the WHO-China Joint Mission on Coronavirus Disease 2019 (COVID-19). [https://www.who.int/publications-detail/report-of-the-who-china-joint-mission-on-coronavirus-disease-2019-\(covid-19\)](https://www.who.int/publications-detail/report-of-the-who-china-joint-mission-on-coronavirus-disease-2019-(covid-19)) (Feb. 28th, 2020)
  27. Wikipedia: Compartmental models in epidemiology. [https://en.wikipedia.org/wiki/Compartmental\\_models\\_in\\_epidemiology](https://en.wikipedia.org/wiki/Compartmental_models_in_epidemiology) (June 2nd, 2020)

## 7 Additions

This is a very informal appendix containing additions, remarks, and changes to the original version of the paper<sup>2</sup> for the DMV. It will change regularly, and this version is of **August 1, 2020**. Changes in these additions will **not** be marked in red. Earlier versions of these additions are on my webpage.<sup>3</sup>

As the only addition to the version of July 19th, a **new** section 9 below deals with the situation of roughly constant new infections, as is the case for various countries in mid-July 2020. This is a modification of the Johns Hopkins Data Model that allows to assess the roughly constant death rate implied by a roughly constant rate of new infections.

### 7.1 Errors and Typos

... if not already marked in **red** in sections 1 to 6 of the original text ...

1. line after (23): replace “with sum 1” by “with sum bounded by one”
2. Comment to Fig. 6 in 3.4: replace “The top shows” by “the left plot shows”.

### 7.2 Updated Predictions

In the meantime, new superspreading events took place in Germany, and the Infections were increasing again, but decreasing after fast intervention by the administration. Various countries seem to have a second outbreak. This calls for updates of the figures in the paper. But also, the data quality gets worse over time, and predictions get more difficult. For example, the UK figures of Recovered always were useless, but on July 19th, even the registration of deaths will be suspended, making any predictions impossible.<sup>4</sup>Data of other states are suspected to be not only questionable, but even manipulated.<sup>5,6</sup>

The updated predictions use certain changes to the programs, see section 7.5.

#### 7.2.1 Predictions

The full model predictions replacing Figure 8 are in Figure 12, now based on data up to day 179 (July 19th). Germany has overcome the small intermediate increase of the registered Infectious  $I$  (black  $\square$ ) and the hidden Infectious  $M$  (cyan  $\square$ ) around day 150-160 due to outbreaks e.g. in the Tönnies factory on day 146. But lowering the

<sup>2</sup> <https://link.springer.com/article/10.1365/s13291-020-00219-9>

<sup>3</sup> <http://num.math.uni-goettingen.de/schaback/research/group.html>

<sup>4</sup> July 19th: <https://www.independent.co.uk/news/uk/politics/coronavirus-uk-death-toll-nhs-phe-covid-19-government-england-scotland-a9626336.html>

<sup>5</sup> June 9th: <https://www.telegraph.co.uk/news/2020/06/09/coronavirus-world-round-up-jair-bolsonaro-accused-manipulating/>

<sup>6</sup> May 19th: <https://kutv.com/news/coronavirus/states-accused-of-manipulating-covid-19-statistics-to-make-situation-look-better>



Country	Death rate	Detection rate	Death rate	Detection rate
Germany	0.039	0.128	0.017	0.300
Brazil	0.040	0.125	0.039	0.170
Italy	0.117	0.043	0.074	0.068
Spain	0.147	0.034	0.017	0.300
Sweden	0.033	0.151	0.023	0.215
Austria	0.041	0.123	0.011	0.449
France	0.083	0.060	0.054	0.092
UK	0.114	0.044	0.133	0.038
US	0.036	0.141	0.014	0.346

**Table 2** Case fatality and detection rates, estimated on day 159 (June 29th) and day 179 (July 19th) using the 14-day rule and a backlog of 28 days.

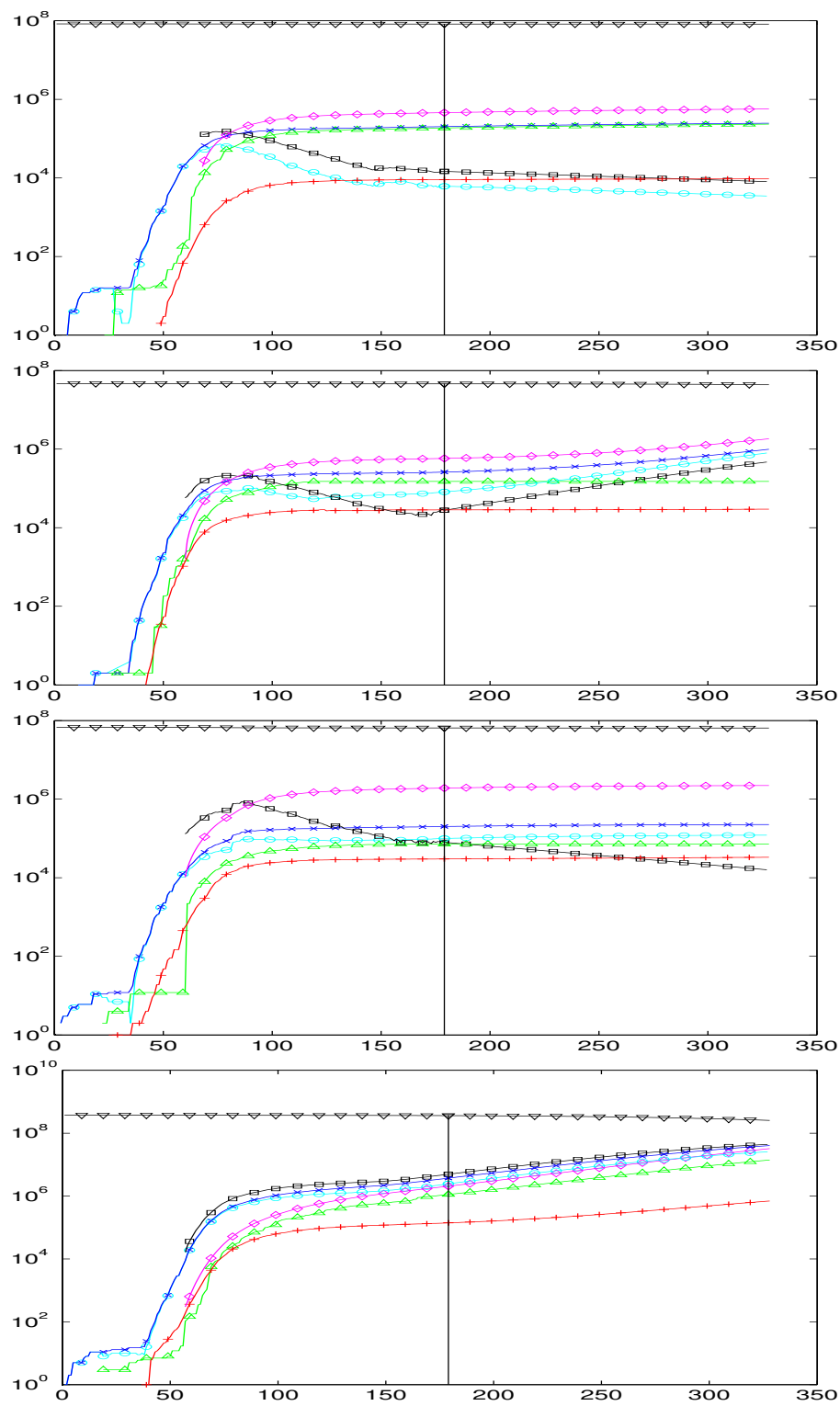
restrictions has stopped the rapid decline around day 100 and replaced it by a very gradual decline. The same is visible for France, and in the plot for US one can see that the increase of Infectious got larger. Anyway, Germany still has good chances to stay below a total of 10,000 deaths, while the long-term prediction of the full model for the US has an  $R_0$  of 1.93 and targets a final death toll of 2 million with a peak of the Infectious around day 380 (Feb. 5th, 2021) if no actions are taken. Brazil also showed a strong increase, but was taken out due to unreliable data. Spain was included to show how fast the model reacts to new small outbreaks. But the results for deaths are very questionable shortly after a new increase of Infectious, because the instantaneous case fatality rate is near to zero until the deaths following the new infections show up. This applies to Spain and Germany. For reasons to be explained below, the results for France are somewhat questionable.

The update of Table 1 is Table 2, based on data up to day 159 (June 29th) and day 179 (July 19th), using the backlog of 28 days. Most death rates are smaller now, and detection rates are mostly higher. The Recovered of Sweden and the UK had to be estimated using the 14-day rule by [22].

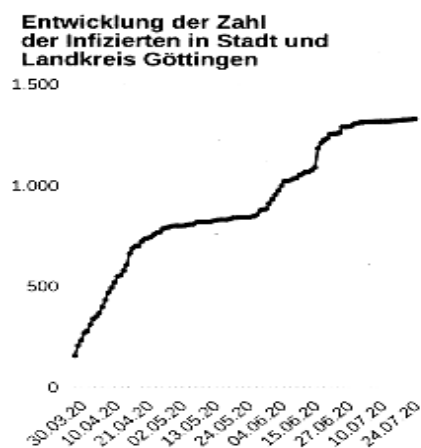
For readers interested in what happened to Figure 11 describing the COVID-19 situation in Göttingen, here is Figure 13. The two superspreading events in rundown apartment houses are clearly visible, and how the local authorities regained control.

### 7.3 Media Coverage in German TV

As of early July, the standard broadcast does not mention terms like Reproduction Number or doubling time anymore. They give the 7-day mean of *new* infections and the increase or decrease of it. This is information about the second derivative of the Infectious, and therefore useful. The strong variations in the data provided by the Robert Koch Institute are smoothed away by the 7-day mean. Daily values are not provided anymore, avoiding any false positive interpretation of lower values on Mondays due to the transmission delays incurred by the weekend. See section 7.6.1 below on how the RKI changes the data to be able to deal with transmission delays.



**Fig. 12** Predictions for countries Germany, Spain, France, and US on day 179 marked by the vertical line. The  $S, M, H$  values to the left are obtained by calibration, the  $C, R, D, I$  values there are the original Johns Hopkins data. Not all data points have marks. See Figure 8 for explanation of colours and markers.



**Fig. 13** Infectious in Göttingen city and county, as of end of July 2020 in the local newspaper “Göttinger Tageblatt”. Still no exponential outbreak in the large, but two local superspreading events.

#### 7.4 Shortcomings of the Models

The models in the paper rely strongly on the numbers for the deaths, when it comes to estimating the hidden variables or the case fatalities. This has serious consequences for re-infections after peaks. As long as only the Infectious re-increase without any influence on deaths and recoveries, the new situation does not have a strong influence on the hidden system. In particular, re-infections will quickly increase the  $R_0$  estimate by the Johns Hopkins data model, but will only reluctantly increase the  $R_0$  estimate by the full model, because the hidden part does not change as quickly.

This is one of the places where a major shortcoming of the models shows up: the inadequate treatment of delays. Another case is that the prediction of a second peak should work like a restart of the full model under new initial conditions. But like startup data are discarded until at least 10 deaths and 100 Confirmed are present, the restart requires again at least 10 new deaths and 100 new Confirmed. It cannot be expected that the models produce useful results unless this condition is satisfied.

Some formulas rely strongly on monotonicity of the data, e.g. (20), (25), and (27). And, when cumulative time series get near-stationary after a peak, the  $R_0$  estimation of the Johns Hopkins Data Model by (20) comes close to  $0/0$  and thus becomes extremely unstable. In particular, the estimation by (20) will get very large when the Confirmed already re-increase while deaths and Recovered do not yet follow up. This is why the update of Figure 5 is delayed to section 8.1.

#### 7.5 Technical Notes

The original programs for the published paper are frozen, and available via <http://num.math.uni-göttingen.de/schaback/research/papers/OC19M.zip> from the author’s research website. But there are certain changes made in the meantime that will be reported here.

### 7.5.1 Useful Data

Basic Johns Hopkins data plots will still show all data. But for any prediction algorithm, outbreak data are ignored before they reach at least 10 deaths and 100 Confirmed. This rule complies well with Figure 4.

But this has a serious consequence for modelling re-infections. As stated above, one has to wait for 10 new deaths or 100 new Confirmed to get useful data for a restart of the models. See section 8.1 for changes concerning  $R_0$  estimates via the Johns Hopkins data model.

### 7.5.2 Estimation of Case Fatalities

The optimizations in (24) and in (2) of [22] can be simplified by taking the  $q_i$  as variables under a nonnegativity constraint. This allows to use linear least squares routines. The results are the same except for roundoff, for all cases seen, and the solution now is unique. The interpretation of the  $p_j$  as probabilities are lost, and the  $q_i$  describe nonnegative portions of the newly Confirmed. The constraint on the sum of the  $q_i$  turned out to be automatically satisfied in all test runs performed so far.

### 7.5.3 Changes in Model Calibration

This technique is strongly based on (28), but the constants in the formula can be estimated both in the beginning and in the end of the available data, in particular in (27) and in the estimation of  $\gamma_{CF}$  via section 4.3.5. The new software uses “startup” values for using (28) in model calibration, but “final” values in the full model for prediction. The old software always used the “final” values. The calculations for Tables 1 and 2 are not affected, but the plots in Figure 12 use the new strategy.

## 7.6 Additional References

From April 9th on, the Robert Koch Institute (An der Heiden & Hamouda [11]) published its own way of preprocessing its data by *Imputation* and *Nowcasting*. I missed this publication when working through the RKI website in March, and got it in June via Dehning et al. [4].

So far, I tried my best not to change data. The data *publishers* can do that, but not the *users*. Encouraged by the RKI case, an algorithm in section 8.1 will make reasonable changes also to the Johns Hopkins data, eliminating the worst flaws and contradictions, e.g. non-monotonicity.

### 7.6.1 RKI Imputation and Nowcasting

In [20] of May 15th, the Robert Koch Institute published its current way of estimating  $R$ . This is based on the Imputation and Nowcasting data corrections described in [11] and mentioned above. It roughly implements the logic of (11) because it takes delayed ratios of Infectious and applies additional means. This is not estimating the  $R_0$  of the

mathematics of the SIR model. It estimates  $R_t$  in the formulation of Section 2.5 of the paper, namely as the multiplication factor for new Infections, restricted to registered cases. This factor is quite sufficient in practice, since for the German situation the impact on the health system is very strongly connected to the new daily registered infections. See section 9 for a variation of the Johns Hopkins Data Model dealing with this situation and making predictions.

*Imputation* by the RKI concerns the estimation of a fictitious symptom onset from the known registration time, i.e. an estimation of the delay between symptom onset and registration. The RKI found a delay between 5.3 and 9 days, time-varying, and applied this to the 40% of all cases where the symptom onset was not known. For about 60% of the data, the RKI knows the delay. To do this properly for the JH data, one needs information about the correspondent delays in JH data. These are missing, to my knowledge, and will anyway be strongly country-dependent.

*Nowcasting* a time series of counting “cases” with true unknown values  $X_n$  at time step  $n$  uses the observed  $x_j$  at time  $j$  to estimate  $X_n$  by some estimated value  $\hat{X}_n$  that accounts for delays in data acquisition. The observations  $x_j$  at time  $j$  contain true cases  $X_i$  at time  $i$  for  $i \leq j$ , which means that  $X_n$  contains parts of the  $x_k$  for  $k \geq n$ . This requires the portions  $q_j$  that should have been registered  $j$  steps later, such that roughly

$$\hat{X}_n \approx \sum_{j=0}^{D-1} q_j x_{n+j}. \quad (35)$$

The RKI seems to use imputation to assign the “right”  $n$  to uncertain  $x$  cases, and then use the obtained  $x_n$  for nowcasting. To do this properly one needs experimental statistics about observed delays, or Bayesian priors to assume meaningful distributions of such delays. The RKI has such information, but what about Johns Hopkins data?

From a deterministic viewpoint, (35) will have a smoothing effect, and this may be a way to use it when no additional information is given. We shall apply (35) in section 8.1 to mimic imputation and nowcasting without having additional data on delays.

## 8 Repairs in JH Data

The JH data contain cases where cumulativity is violated, and where strong  $\approx 7$  day oscillations occur that may, like in the RKI case, be due to delays in transmission, e.g. by weekends. Furthermore, there are violations of the  $k$ -day rule, and predictions will improve if certain conservative smoothing techniques are applied. “Conservative” means here that counts of cases are kept in the mean. Missing or doubtful data for Recovered can be roughly estimated using [22] and a  $k$ -day rule, if the Confirmed and Deaths are reliable. A general goal is to find algorithms that detect errors and plausibility flaws in the Johns Hopkins data, but this is ongoing work.

## 8.1 Monotonicity and Smoothness

Assume that a time series  $X_n$  is required that is at least weakly monotonic, based on observations  $x_n$  that may be faulty. The goal is to find an algorithm that works like imputation and nowcasting, combined with a smoothing technique. This will be applicable to  $C$ ,  $R$ , and  $D$  of the Johns Hopkins data.

The idea pursued here is to perform a fit like (35) under monotonicity constraints and non-smoothness penalties. For the latter, we may measure smoothness by a vector-valued function  $f$  whose components  $f_j(X)$  give a penalty for non-smoothness of the  $X$  values around step  $j$ . One way is to use a time window around time  $j$ , take the  $X_k$  in the window, fit them to a fixed local model and return the fitting error. Then one might minimize a weighted sum of  $f_j(X)^2$  and squared errors in (35) under all  $X$  with monotonicity constraints and under all nonnegative  $q_j$ . If necessary, one can restrict the sum of the  $q_j$  to be one, enforcing that no cases are lost in the mean.

The implementation works via MATLAB's `lsqnonlin` acting on the logarithms of  $C$ ,  $R$ , or  $D$ , the non-smoothness penalties being the fitting error to a quadratic polynomial on a window of up to five points. The overall optimization problem thus is quadratic with linear constraints, usually leading to a unique solution. If  $M$  data are treated, this minimizes over the  $D + M$  nonnegative quantities  $(q, X)$  under monotonicity constraints on  $X$  and  $2M - D$  penalties consisting of  $M - D$  errors in (35) and  $M$  non-smoothness penalties. Note that the non-smoothness penalties are independent of the observed data. A more sophisticated technique of non-smoothness penalties will not aim at local quadratics, but could fit data on 5 points by two line segments with a breakpoint at one of the points. This is currently not implemented, but would not iron out any breakpoints. Note that equation (5) in Höhle/an der Heiden [12] uses a local quadratic spline for smoothing within a Bayesian nowcasting technique.

For Johns Hopkins data of France up to day 180 (July 20th), the  $C$ ,  $D$ , and  $R$  values have 20, 7, and 12 places of non-monotonicity, respectively. Figure 14 shows the result of the algorithm for  $D = 7$ . Here and in many other cases, the new data (dotted) are to the left of the original data. This is normal, because they contain values that are possibly falsely registered on later days. Table 3 shows the resulting  $q_j$  for the data of France. They indicate that Recovered are possibly registered later and more irregularly than Confirmed and Dead. This calls for further checks, if time permits.

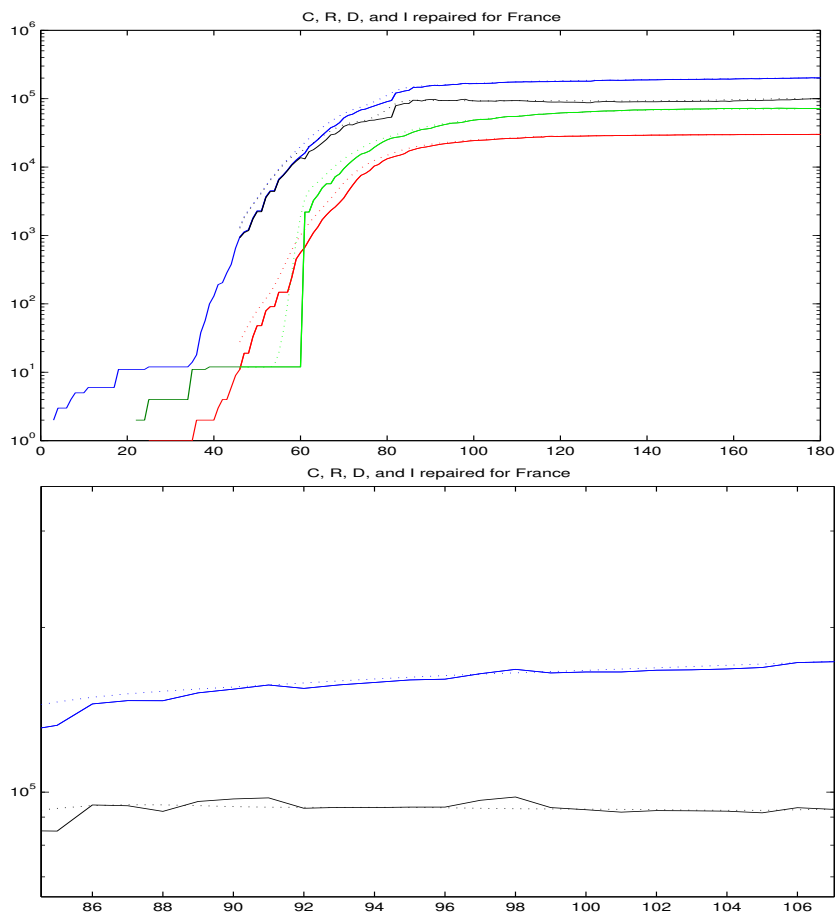
For countries with useless  $R$  data, like the United Kingdom, Denmark, and Sweden, the  $C$  and  $D$  data are processed first, and then the  $R$  data are estimated from the processed  $C$  and  $R$  data under the 14-day rule as proposed in [22].

It can be argued that (35) should be used on increments, i.e.

$$\hat{X}_n - \hat{X}_{n-1} \approx \sum_{j=0}^{D-1} q_j (x_{n+j} - x_{n+j-1}),$$

but (35) implies this.

If Figure 5 containing the  $R_0$  estimates from the Johns Hopkins data is produced for data up to day 179 (July 19th), results for France are useless due to non-monotone



**Fig. 14** The JH data for France, with repair, up to day 170 (July 20th). Dotted: Repaired data. Full data and a closeup. Note that the black  $I$  curves are not monotonic.

$j$	$C$	$R$	$D$
0	0.2085	0.0949	0.1261
1	0.1985	0.1369	0.1393
2	0.1965	0.1761	0.2284
3	0.1284	0.1907	0.1577
4	0.1090	0.1738	0.1922
5	0.1167	0.1326	0.0878
6	0.0428	0.0924	0.0685
sum	1.0004	0.9975	0.9998

**Table 3** Quantities  $q_j$  for monotonicizing, nowcasting and smoothing the JH data of France up to day 170. There was no constraint on the sum of the  $q_j$ .

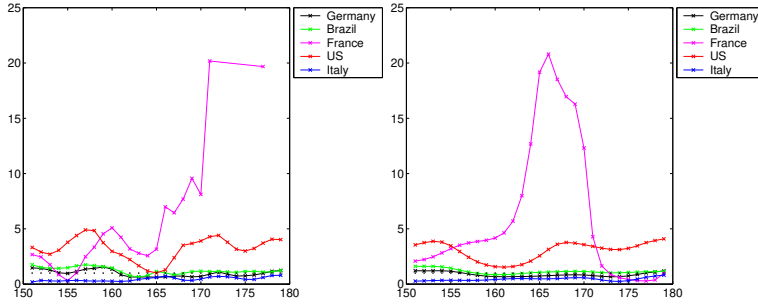


Fig. 15 Estimates of  $R_0$  via the time series  $r_n$  up to day 179, without and with data repair.

values and the 0/0 effect in (20). The left part of Figure 15 shows the results for the standard  $(1/4, 1/2, 1/4)$ -smoothing of raw logarithmic data, while the right-hand side is applied after applying the above technique. The missing  $x$  markers in the left-hand plot are at places where the denominator in (20) was zero or the quotient was negative. The results for France are still useless due to the 0/0 instability in (20), even after monotonicity is enforced by taking the values of Figure 14. Overcoming the 0/0 instability of (20) is still an open problem. The  $r$  estimate there is connected to the derivative of  $C$  reparametrized as a function of  $R + D$ , and this function necessarily behaves badly when there are new confirmations but no new deaths and recoveries.

## 8.2 Testing the $k$ -day Rule

A standard assumption is that the newly Confirmed  $C_n - C_{n-1}$  at day  $n$  will end up dead or alive until day  $n + k$  for  $k$  large enough. This implies

$$C_n - C_{n-1} \leq R_{n+k} - R_{n-1} + D_{n+k} - D_{n-1}. \quad (36)$$

Conversely, the new Removed  $D_n - D_{n-1} + R_n - R_{n-1}$  at day  $n$  must have been confirmed between day  $n$  and day  $n - k$ . This means

$$D_n - D_{n-1} + R_n - R_{n-1} \leq C_n - C_{n-k-1} \quad (37)$$

and both inequalities admit that transition from Confirmed to Removed can be on the same day. This is different from the treatment of the  $k$ -day rule in the paper and in [22], where death and recovery can occur at least one day later than confirmation. Altogether, this differs by just a shift by one day. The extreme case  $k = 0$  now implies that

$$C_n - C_{n-1} = R_n - R_{n-1} + D_n - D_{n-1},$$

i.e. all new Registered at day  $n$  are either dead or recovered. Since the  $k$ -dependent right-hand sides increase with increasing  $k$ , there should be a minimal  $k_{min}$  for which the above relations hold for all reasonable  $n$ .

Table 4 shows the minimal  $k$  for which either (36) or (37) are satisfied for the full range of data with  $D \geq 10$  and  $C \geq 100$  up to day 179 (July 19th). The data were



preprocessed with the technique of section 8.1. Furthermore, for UK, Denmark, and Sweden, the Recovered had to be estimated by the technique of [22] to get reasonable values at all. If the calculation is repeated for the data between days 142 and 170, all values (except for Spain) are considerably smaller. The unusually high value for Spain is due to failing (36) by at most 300 cases out of about 250,000 Confirmed. Altogether, even for the early phase of the outbreak, a 14-day rule is roughly satisfied, supporting the strategy used so far for estimation of Recovered and of the case fatality rate.

Country	Passes (36) for $k$	Passes (37) for $k$
Austria	12	4
Brazil	18	1
Denmark	15	3
France	13	1
Germany	13	1
Italy	6	5
Russia	12	1
Spain	29	4
Sweden	15	1
Switzerland	14	4
UK	13	1
US	11	1

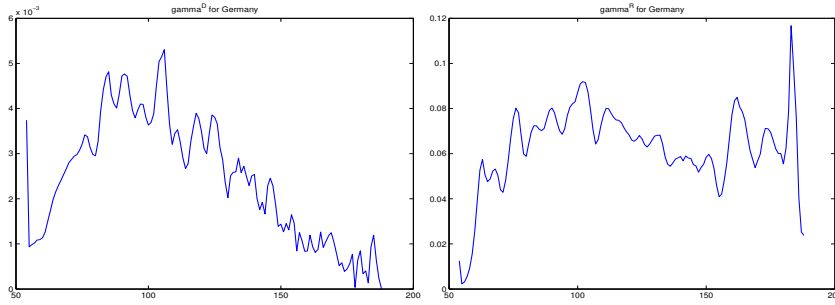
**Table 4** The minimal  $k$  for which the JH data up to day 179 satisfy either (36) or (37).

## 9 Constant New Confirmations

A typical situation in Germany in mid-July 2020 was that there was a roughly constant number  $\Delta C$  of new Confirmed, varying between 200 and 500, slowly increasing. This can be due to more travelling or relaxation of contact constraints, for instance, but here we do not ask for reasons. Instead, we check how the Johns Hopkins data model behaves under such an assumption. Equations (19) turn into

$$\begin{aligned} C_{n+1} - C_n &= \Delta C, \\ I_{n+1} - I_n &= \Delta C - \gamma_n I_n, \\ (R + D)_{n+1} - (R + D)_n &= \gamma_n I_n \end{aligned}$$

maintaining the balance  $C = R + D + I$ . If the  $\gamma_n$  are considered to be constant, the  $I$  values go exponentially to a level  $I = \frac{\Delta C}{\gamma}$ , from below or above, depending on the starting value for  $I$ . The qualitative behaviour of  $I$  now is logistic. There is no peak and no long-term exponential decrease to zero. But a constant number of Infectious means a constant increase of deaths, calling for political changes that stop this scenario. Anyway, it cannot be valid over long time intervals because there is an upper bound on the Confirmed.



**Fig. 16** Time series for estimates of  $\gamma^D$  and  $\gamma^R$  for Germany up to day 188

To get a grip on deaths and recoveries, we should split the last equation into

$$\begin{aligned}\gamma_n^R + \gamma_n^C &= \gamma_n, \\ R_{n+1} - R_n &= \gamma_n^R I_n, \\ D_{n+1} - D_n &= \gamma_n^D I_n.\end{aligned}$$

The  $\gamma_n^R$  and  $\gamma_n^D$  are now time series of the instantaneous case rates  $\gamma_{iCR}$  and  $\gamma_{iCD}$  of (22), and it is easy to get estimates of them by just solving the equations. Figure 16 shows these values as time series, based on the Johns Hopkins data after applying the technique of section 8.1. For Germany on day 188 (July 28th) using a mean over the last 14 days, the values are

$$\gamma^D \approx 0.0005, \quad \gamma^R \approx 0.064.$$

With a number of about  $I = 7600$  Infectious on day 188, the critical value of  $\Delta C$  determining increase or decline of  $I$  is  $\Delta C = I\gamma \approx 500$ . Then a constant  $I$  of about 7600 implies a daily death count of four due to COVID-19. This agrees quite well with what actually happened in mid-July, but this must be expected because we used the real data. If rates  $\gamma^R$  and  $\gamma^D$  stay as they are, doubling  $\Delta C$  implies doubling the asymptotic  $I$  level and the death rate in the long run. The logistic increase to the new  $I$  level is rather fast, see Figure 17 for Germany with an assumed  $\Delta C = 1000$  from day 188 on. The linear increase of  $C$  is still far from reaching saturation, and the total death toll still stays below 10,000 in the period shown.

Italy has about 250 new daily infections around day 188, with values of  $\gamma^D \approx 0.00086$  and  $\gamma^R \approx 0.124$  with Infectious around 12,350. The threshold for  $\Delta C$  is around 125, and a value of 400 new daily infections would imply about 11 daily COVID-19 deaths in the long run.

Summarizing, there are situations where countries with a small and constant daily increase  $\Delta C$  of the Confirmed may get along for quite some time, but at the price of a constant and hopefully low COVID-19 death rate. Societies and their politicians must decide whether this rate is tolerable.

But various other countries like France, Spain, UK, or US cannot be treated that way, because either the current  $\Delta C$  values are not roughly constant or the  $\gamma$  values are completely unreliable because the  $R$  or  $D$  values are. Finally, the  $\gamma^D$  values used in the above examples were much smaller than what was observed in the early phase of

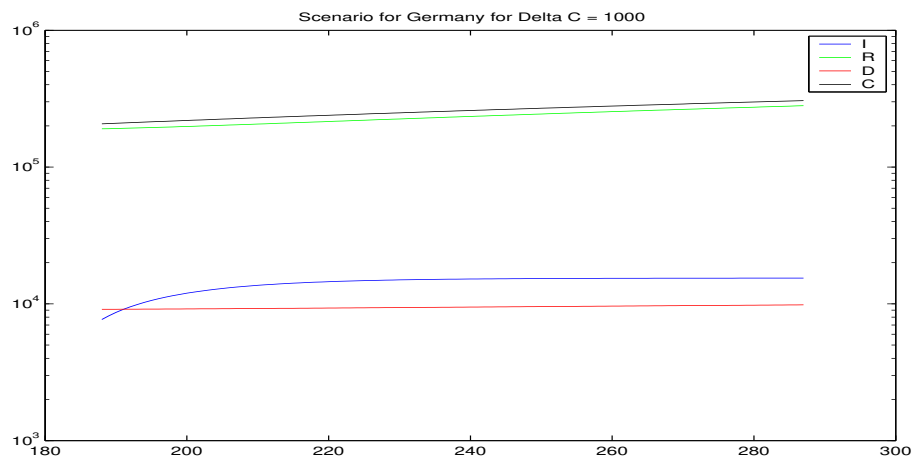


Fig. 17 Prediction for  $\Delta C = 1000$  for Germany based on data up to day 188

the COVID-19 outbreak, see Figure 16. If health systems get under serious stress, and if seniors are not properly protected, the above scenario will become much worse.