

RESEARCH

# Data-Driven Prediction of COVID-19 Cases in Germany for Decision Making

Lukas Refisch<sup>1,2†</sup>, Fabian Lorenz<sup>1,3†</sup>, Torsten Riedlinger<sup>4</sup>, Hannes Taubenböck<sup>4,5</sup>, Martina Fischer<sup>6</sup>, Linus Grabenhenrich<sup>6,7</sup>, Martin Wolkewitz<sup>1</sup>, Harald Binder<sup>1,8</sup> and Clemens Kreutz<sup>1,3,8\*</sup>

\*Correspondence:

ckreutz@imbi.uni-freiburg.de

<sup>1</sup>Institute of Medical Biometry and Statistics, Faculty of Medicine and Medical Center, University of Freiburg, Stefan Meier Str. 26, 79104 Freiburg, Germany

Full list of author information is available at the end of the article

<sup>†</sup>Equal contributor

## Abstract

**Background:** The COVID-19 pandemic has led to a high interest in mathematical models describing and predicting the diverse aspects and implications of the virus outbreak. Model results represent an important part of the information base for the decision process on different administrative levels. The Robert-Koch-Institute (RKI) initiated a project whose main goal is to predict COVID-19-specific occupation of beds in intensive care units:

*Steuerungs-Prognose von Intensivmedizinischen COVID-19 Kapazitäten (SPoCK)*. The incidence of COVID-19 cases is a crucial predictor for this occupation.

**Methods:** We developed a model based on ordinary differential equations for the COVID-19 spread with a time-dependent infection rate described by a spline. Furthermore, the model explicitly accounts for weekday-specific reporting and adjusts for reporting delay. The model is calibrated in a purely data-driven manner by a maximum likelihood approach. Uncertainties are evaluated using the profile likelihood method. The uncertainty about the appropriate modeling assumptions can be accounted for by including and merging results of different modelling approaches.

**Results:** The model is calibrated based on incident cases on a daily basis and provides daily predictions of incident COVID-19 cases for the upcoming three weeks including uncertainty estimates for Germany and its subregions. Derived quantities such as cumulative counts and 7-day incidences with corresponding uncertainties can be computed. The estimation of the time-dependent infection rate leads to an estimated reproduction factor that is oscillating around one. Data-driven estimation of the dark figure purely from incident cases is not feasible.

**Conclusions:** We successfully implemented a procedure to forecast near future COVID-19 incidences for diverse subregions in Germany which are made available to various decision makers via an interactive web application. Results of the incidence modeling are also used as a predictor for forecasting the need of intensive care units.

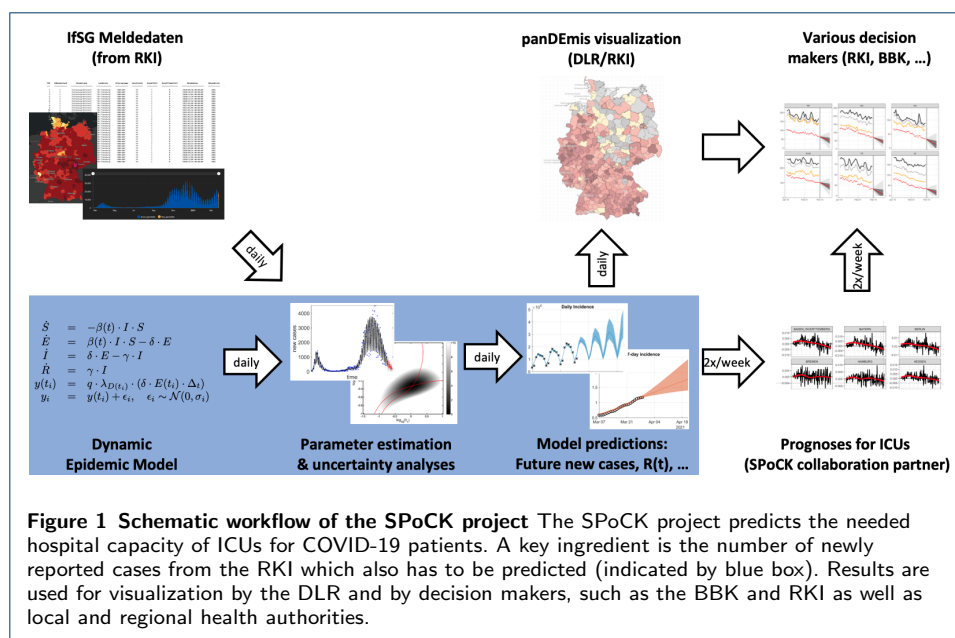
**Keywords:** COVID-19; infectious disease models; input estimation; ordinary differential equations; parameter estimation; nonlinear systems; SEIR models

## 1 Background

Mathematical models of infectious disease epidemiology have experienced a boost of attention since the beginning of the COVID-19 pandemic. One can divide these models into three categories according to their purpose: scenario simulation, now-casting, and forecasting.

Scenario simulation focuses on different assumptions about some aspects of the model in order to compare and illustrate differences between several scenarios of in principle conceivable progressions of the transmission and other dynamics, which do not allow for proper uncertainty assessment. These approaches are used to examine the impact of changing certain parameters in the system, e.g. social behaviour, vaccination rate, etc. Nowcasting focuses on the precise description of the present situation based on incomplete, noisy and/or systematically biased data about the current state ([1], [2]). Forecasting tries to make predictions about the near future providing policy makers with reliable estimates of advancing developments. Similar to nowcasting, forecasting is strongly oriented towards realistic settings. The work presented in this publication focuses on a near-future prediction and can therefore be classified as forecasting.

### 1.1 The SPoCK Project



In Germany, local health authorities collect data about the infection dynamics on population level as mandated by the *Infektionsschutzgesetz* (IfSG) and report it to the national public health institute, the *Robert Koch-Institut* (RKI). In this paper, We describe the fitting and short term forecasting of this quantity, i.e. the newly reported cases of COVID-19 in Germany.

In a second step, which is not covered in this publication, the data about COVID-19-specific occupation of beds in intensive care units which is collected and reported daily by the *DIVI Intensivregister* run by RKI with support of the *Deutschen Interdisziplinären Vereinigung für Intensiv- und Notfallmedizin* (DIVI), is fitted and forecasted by our cooperation partners. The results of the first step are utilized as a predictor to obtain short-term future predictions on the level of intensive care unit (ICU) capacities. This two-step procedure is referred to as the *Steuerungs-Prognose von Intensivmedizinischen COVID-19 Kapazitäten* (SPoCK) project. Several decision makers including the *Federal Ministry of Health* (BMG), the *Robert Koch Institute* (RKI), the *Federal Office of Civil Protection and Disaster Assistance* (BBK),

the local planners of ICU capacities as well as the *Bundesamt für Bevölkerungsschutz und Katastrophenhilfe* (BBK) incorporate these predictions into their risk assessment of the current COVID-19 situation. In addition, the predicted incidences are visualized on an interactive web application provided by the *Deutsches Luft- und Raumfahrtzentrum* (DLR) called *Pandemic Mapping and Information System for Germany* (panDEmis).

The workflow within the SPoCK project is depicted in Figure 1. In this paper, we describe the daily analysis and prediction of incident cases of COVID-19 in different regions in Germany which are, in addition to the entire country, the 16 federal states (*Bundesländer*) and their 413 counties (*Land- und Stadtkreise*) summing to a total of 430 regions.

## 2 Methods

A standard approach when describing infectious disease transmission are compartmental models or SIR-like models [3]. In general, both approaches divide the population into subpopulations with disjoint properties. Transition rates allow for flows between the subpopulations and define, in combination with the initial values of the subpopulations, the time evolution of the system. The ordinary differential equation (ODE) representation of the compartmental scheme we use is the well-known SEIR model [4]:

$$\begin{aligned} \dot{S} &= -\beta(t) \cdot I \cdot S/N \\ \dot{E} &= \beta(t) \cdot I \cdot S/N - \delta \cdot E \\ \dot{I} &= \delta \cdot E - \gamma \cdot I \\ \dot{R} &= \gamma \cdot I \end{aligned} \quad (1)$$

with  $N = S + E + I + R$  and where the dot notation is used to indicate time derivatives. A special characteristic of the current pandemic is the massive political and social reaction. In contrast to, e.g. the annual influenza season during which the social and professional life used to proceed pretty much as usual, the COVID-19 pandemic has led to vast political interventions and personal restrictions aiming mainly at the reduction of infections [5]. Within the SEIR scheme these changes over time can be described by a time-dependent infection rate  $\beta(t)$ . There are several studies dealing with this problem in different manners. For example, at the beginning of the COVID-19 pandemic the impact of different non-pharmacological interventions (NPIs) was examined via step functions that implement  $\beta(t)$  via different variants of (smoothed) step functions, e.g. to examine the impact of different NPIs [6, 7, 8, 9]. Often, these approaches are restricted to time ranges in which the infection rate is assumed to be constant or monotonously decreasing or increasing, respectively.

In contrast, we aim for a more general approach which enables the infection rate to vary flexibly, i.e. to decrease and/or increase repeatedly within the considered time range. This is necessary for an accurate description of the COVID-19 transmission dynamics since it is influenced by many factors that may vary over the course of the ongoing COVID-19 pandemics:

- 1 Various NPIs are implemented, repealed and reintroduced iteratively.

- 2 The population's compliance to regulative measures changes over time.
- 3 Seasonal effects, e.g. weather conditions, lead to changes in infection risk.
- 4 Mutations alter the physiological mechanisms underlying the disease transmission and other aspects.
- 5 Vaccinations reduce the population's susceptible fraction.

In order to fit a strictly positive and time-dependent infection rate simultaneously with the SEIR model's parameters, we introduce the following parametrization for the infection rate:

$$\beta(t) = b \cdot \frac{1}{1 + e^{-f(t)}} \quad , \quad (2)$$

where the argument of the exponential function is given by an interpolating cubic spline

$$f(t) = \text{cubic\_spline} \left( t, \{ \tau_i, u_i \}_{i \in \{1, \dots, n\}} \right) \quad . \quad (3)$$

We utilize joint estimation of input spline and ODE parameters as introduced for biological systems in [10]. The composition of the interpolating spline (3) with the logistic function (2) allows for a nearly arbitrary time dependence, while still ensuring that the infection rate  $\beta(t)$  is strictly positive, smooth and restricted to a maximal value  $b$ . The cubic spline curve is determined by estimated parameters  $u_i = \text{cubic\_spline}(\tau_i)$  that represent its values at fixed and evenly spaced dates  $\tau_i$  for  $i \in \{1, \dots, n-2\}$  which cover the time range of observed data. In our model, the last two spline knots are placed after the date  $t_{\text{Last}}$  of the last data point:  $\tau_{n-1} = t_{\text{Last}} + 50\text{d}$  and  $\tau_n = t_{\text{Last}} + 300\text{d}$ . The value  $u_{n-1}$  is fitted to allow for some flexibility in the most recent regime, whereas  $u_n = 0$  is fixed for numerical stability and reflecting the end of the pandemic in at least 300 days.

The predictions for the infection dynamics are primarily determined by the time-dependent infection rate  $\beta(t)$ . In general, assumptions for the future development of  $\beta(t)$  are difficult to justify as many different factors contribute to it. For illustrative purposes, several different assumptions could be made and visualised as done e.g. in various online simulator tools [11]. For example, one such scenario study nicely illustrates the effectiveness of a Test-Trace-Isolate strategy [12].

For a data-driven approach focused on short-term forecasts, we need to be more practical: For extrapolation purposes, we fix

$$\beta(t > t_{\text{Last}}) = \beta(t_{\text{Last}}) \quad (4)$$

i.e. we assume the infection rate to be constant starting from the day where the last data point is reported.

## 2.1 Data-Driven Approach

Typically, there exist a multitude of model classes and structures which can be used to describe the same phenomenon. However, it is generally not possible to transfer results about estimated parameters between different models in a straightforward

manner due to their differing mechanistic structures. To circumvent this problem, we here rely on a purely data-driven approach meaning that no prior knowledge about parameter values is incorporated into the optimization procedure. The only three *a priori* fixed parameters are the initial number of individuals in the susceptible, the exposed and the recovered state:  $S_{\text{init}}$ ,  $E_{\text{init}}$  and  $R_{\text{init}}$ . Time point zero  $t_0$  is set to the first day that has at least a total of 100 reported cases to ensure the well-mixing assumption of ODE modeling.  $S_{\text{init}}$  was set to the total population of the respective region as given by the Federal Statistical Office of Germany [13].  $E_{\text{init}}$  was set to  $\gamma \cdot I_{\text{init}}/\delta$ , which is motivated by the assumption that  $\dot{I} \approx 0$  at the beginning of an epidemic reflecting a slow onset.  $R_{\text{init}}$  is set to zero. The only remaining initial occupation number  $I_{\text{init}}$  is estimated from the data.

## 2.2 Link between Model and Observed Data

In order to calibrate the ODE model, it needs to be linked to the observed data. The data we use for calibration is the daily incidence  $y_i$  published by the reporting date (*Meldedatum*)  $t_i$  at the local health authority. Therefore, we introduce the observation function

$$y(t_i) = q \cdot \lambda_{D(t_i)} \cdot (\delta \cdot E(t_i) \cdot \Delta) \quad , \quad (5)$$

where the parameters can be interpreted as follows:

- $q \in [0, 1]$  is the fraction of all infectious individuals that are detected and reported.
- $D(t_i) \in \{1, \dots, 7\}$  is an index for the weekday at date  $t_i$  where  $\{1, \dots, 7\}$  are naturally identified with the weekdays  $W = \{\text{Monday}, \dots, \text{Sunday}\}$ .
- $\lambda_D$  is a factor for the weekday  $D$  that adjusts for the weekly modulation occurring in the IfSG data (see 2.2.1).
- $(\delta \cdot E(t) \cdot \Delta)$  approximates the influx into the state  $I(t)$  of equation (1). As the considered data represents daily incidences, we set  $\Delta$  to 1 day. This approximation of the true incidence quantity  $\int_{t-1}^t \delta \cdot E(t') dt'$  is exact if the state  $E(t)$  remains constant within that day. Comparison with this exact but computationally much more expensive approach showed minor deviations for real data applications.

The observable function (5) connects the model's predictions to the reported data. The observations are assumed to scatter around this mean according to a normal distribution:

$$y_i = y(t_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_i^2) \quad . \quad (6)$$

As we are dealing with a count process we use the standard deviation inspired by a Poisson model

$$\sigma_i = C \cdot \sqrt{1 + y(t_i)} \quad . \quad (7)$$

The error parameter  $C$  is fitted jointly with all others.

### 2.2.1 Weekly Modulation Factors

The IfSG data shows an oscillatory pattern with a period of one week. The main reason for this is the reporting procedure, displaying a major delay during weekends, instead of actual infection dynamics. Therefore, we account for this effect within the observation function via seven weekday-specific factors  $\lambda_D$  with the integer  $D \in \{1, \dots, 7\}$ . In order to

- 1 guarantee that the factors  $\lambda_D$  essentially do not change the 7-day-incidence and
- 2 separate the weekly modulation from a global scaling of the observation function, which is realized via the factor  $q$ ,

we, furthermore, set the constraint that

$$\sum_{D \in \{1, \dots, 7\}} \lambda_D = 7. \quad (8)$$

As a consequence, we are left with six degrees of freedom to describe the weekly effects. For a convenient implementation in the used software, we introduce a Fourier series with six parameters  $\Theta_{\text{weekly}} = \{A_1, A_2, A_3, \phi_1, \phi_2, \phi_3\}$ :

$$\psi(t) = A_0 + \sum_{k=1}^3 A_k \cdot \cos(k\omega t + \phi_k) \quad (9)$$

where offset and frequency are fixed to

$$A_0 = 1, \quad \omega = \frac{2\pi}{7 \text{ days}}. \quad (10)$$

Instead of fitting the factors  $\lambda_D$  directly, we rewrite them in terms of equation (9) as

$$\lambda_D = \frac{\psi(D)}{\sum_{j=1}^7 \psi(j)} \quad (11)$$

and calibrate the parameters  $\Theta_{\text{weekly}}$ . Doing so allows to set the amplitudes  $A_1$ ,  $A_2$  and  $A_3$  to zero in order to get an adjusted curve that does not feature the weekly oscillations and therefore reflects the ideal case of no reporting artifacts in the data.

### 2.2.2 Correction of Last Data Points

The IfSG data published on date  $t_n$  contains information about the reported cases at all past dates  $t_n, t_{n-1}, \dots, t_1$  since the beginning of reporting. However, due to reporting delays between the test facilities, the local health authorities and the RKI, the data update from date  $t_{n-1}$  to  $t_n$  contains not only cases that were reported to the local health authorities at date  $t_{n-1}$ , but also before that at dates  $t_{n-2}, t_{n-3}, \dots$  and so on. This means that the number of reported cases on day  $t_n$  will be underestimated especially for the most recent dates. For some regions this correction factor be as big as three for the most recent day.

Meaningful handling of this data artifact can be done in at least two ways: For instance, one could choose to ignore some of the latest data points, since they are

most prominently affected by this data artifact. An alternative is to estimate the systematic deviation from historically published data sets. In order to avoid the bias towards smaller incidences in the prediction, the data can be adjusted accordingly. Therefore, one assumes, that the future data sets of  $t_n$  will not change reported counts older than four weeks  $t_{n-28}$ . Let  $N_{t_1}^{t_2}$  denote the number of reported cases, that were published at time point  $t_1$  to be reportedly infected at date  $t_2$  where  $N_{t_1}^{t_2 > t_1} = 0$  as future cases cannot be reported. Then, one can learn from this history of published data sets the correction factor  $CF_k$

$$CF_k = \frac{\sum_{\hat{t}} N_{\hat{t}}^{\hat{t}-k}}{\sum_{\hat{t}} N_{t_{\text{Last}}}^{\hat{t}-k}} \quad (12)$$

the initial publication of  $k$  day old counts had to be corrected to obtain the number in the latest data set  $t_n$ . The factors  $CF_k$  can then be applied to the newest data set.

This was done for Germany and all the federal states separately. We showcase the resulting differences of these two data preprocessing strategies in section 2.4.2.

For the county level, this adjustment is not as crucial for two reasons: 1) the count numbers are much lower, so the stochasticity can lead to wrong correction factors and 2) the shape of the estimated dynamics is inherited from the federal states in our model.

### 2.3 Parameter Estimation

In general, we follow the maximum likelihood estimation (MLE) approach. As there are a total of 429 regions for which the data has to be fitted and predictions are calculated, we rely on a two-step procedure to reduce computation time which is described in the following paragraphs.

#### 2.3.1 Federal States and Germany

The parameter estimation problem given by the above defined ODE model and the IfSG daily incidence data is solved separately for Germany and each federal state by an MLE approach. The latter has been well established for ODE models [14]. The deviation between data and the model's observation function as specified in equation (5) is minimized, taking into account the error model of equations (6) and (7). The simultaneous parameter estimation of the spline parameters  $u_i$  follows the lines of [10]. In particular, no explicit regularization term is implemented that penalizes non-vanishing spline curvatures.

#### 2.3.2 County Level

Analysis at the rural and urban county level (*Land- and Stadtkreise*) is important to obtain a spatially resolved picture of the infection dynamics in Germany. The previously described approach is computationally not feasible because the analysis of 429 regions cannot be performed within 24 hours without access to a sufficiently large computing cluster which can be used 24/7 without queuing. Moreover, the number of infected individuals can generally be so small at the county level that

inference and prediction based on a purely deterministic model is not appropriate. Therefore, we used the results on the higher-level administrative structure, i.e. the fitted model of the federal state, as prior information about the dynamics, and scaled it down to the county level for predictions.

More specifically, the county-level data was used to merely estimate two parameters in a county-specific manner: the scaling parameter  $q$  from equation (5), which in this context can be related to the proportion of current infections occurring in the county  $c$ , and the error parameter  $C$  from equation (7) which quantifies the stochasticity of county-level observations analogous to its meaning on the level of federal states. All other parameter values for a county  $c$  are taken from the estimated set of parameters  $\hat{\Theta}_{FS(c)}$  for the corresponding federal state  $FS(c)$ .

The county-level dynamics might change rapidly as new clusters of infection emerge. For predictions, it is important that such rapid changes are detected by the model calibration procedure, i.e. fitting of  $q$  and  $C$  has to account for such rapid changes. We implemented this requirement by exponentially weighting down the county level data observed in the past by increasing the standard deviations via

$$\sigma_i^2 \leftarrow \frac{\sigma_i^2}{w_i}, \quad w_i = A \cdot \sqrt{(\exp(t_i - t_{\text{Last}})/\tau)^2 + (w_{\text{min}}/A)^2}. \quad (13)$$

Here,  $w_{\text{min}} = 0.01 \cdot A$  denotes the minimal weight factor used for data observed in the past.  $A = 7.56$  denotes the normalization factor that ensures that the sum of all weights  $w_i$  is equal to one. Moreover, we chose  $\tau = 7$  as time-constant of this weighting step. To be clear, on the county-level,  $\sigma_i$  from equation (7) should be thought of as first being transformed according to the mapping (13) before entering equation (6) as the standard deviation of Gaussian observation errors.

Just as the analysis for the federal states, the described scaling procedure for the counties is updated on a daily basis, i.e. the county-specific parameters  $q$  and  $C$  are updated every day. This accounts for time-dependent deviations of the local infection history on the federal state level, i.e. each county has an individual kinetics.

## 2.4 Calculation of Uncertainties

To quantify the uncertainty in the predictions of the model, our forecasting tool provides confidence intervals along with proposed predictions. Here, we describe two main sources of uncertainties: parameter uncertainty and approach uncertainty. The first is captured by simulating all parameter combinations that agree with the observed data as will be explained in section 2.4.1, the second is incorporated by running the analysis with several models as detailed in section 2.4.2.

### 2.4.1 Profile Likelihood Analysis

For non-linear models, uncertainties for estimated parameters can be determined using the *profile likelihood* (PL) method which estimates parameter values that still ensure agreement of model and data to a certain confidence level in a pointwise and iterative manner [15]. This approach has been showcased for infectious disease models [16]. Parameter uncertainties naturally translate to prediction uncertainties which can be analyzed systematically [17]. Following the given references, we



simulate the data-compatible parameter combinations from the parameter profiles and then take the envelope of the resulting family of curves to obtain confidence intervals.

One could also analyze the uncertainty of a model prediction directly via the *prediction profile likelihood* method [18]. Prediction profiles need to be computed via a costly iterative fitting procedure for each predicted quantity and time point separately. However, by using the parameter combinations from the profile likelihood method, we can calculate uncertainties for any desired model quantities and time points only by simulation, thus rendering this method more efficient for our purposes.

#### 2.4.2 Averaging of Approaches

When utilizing ODE models to describe certain aspects of reality, a multitude of assumptions are implicitly made, which include (but are not limited to) the selected model structure, the noise model of the data, the appropriate data preprocessing. All these decisions result in a certain *approach*. These necessary decisions along the modeling process impact the space of possibly described and therefore also predicted dynamics. To account for this origin of uncertainty, we perform the procedure described so far simultaneously for several approaches and merge their results into one comprehensive result. The latter is done by taking the mean / minimum / maximum of the different approaches' MLE / lower bound / upper bound curves. Accounting for different modeling decisions prevents overconfidence in the results.

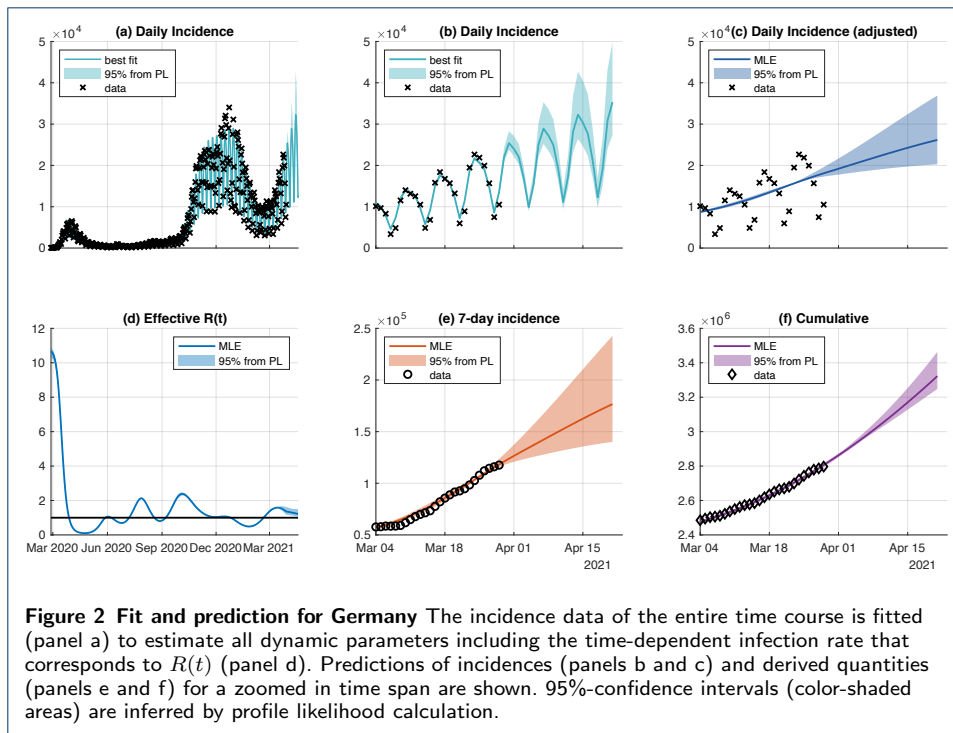
### 3 Results

Since April 2020, the described methodology has delivered daily predictions and the *ansatz* has evolved and several changes and refinements have been implemented. Currently, the resulting predictions for ICU bed capacity, which use the here presented results of estimated incidences as a main predictor, are reported bi-weekly to public health decision makers. The presented methodology and results were generated on April 1st, 2021. The data fitted had therefore registered infections up to March 31st, 2021.

#### 3.1 COVID-19 Spread in Germany

For the aggregated data over all of Germany, we obtained a fit and predictions with uncertainties as shown in Figure 2. The data can be described by the model. Adjusting for weekday effects turned out to be beneficial and the prediction is a reasonable continuation of the last data points. The most interesting model quantity is the time-dependent infection rate  $\beta(t)$  which translates to an effective time-dependent reproduction number  $R(t) = \frac{\beta(t) \cdot S}{\gamma \cdot N}$ . The latter quantifies how many other people are infected on average by a single infectious individual and determines at which rate the number of currently infectious individuals is growing ( $R(t) > 1$ ) or decaying ( $R(t) < 1$ ). It should be noted that, despite the fact that  $\beta(t)$  is extrapolated as remaining constant (see equation (4)),  $R(t)$  is not necessarily constant. This is because  $R(t)$  includes the monotonously decreasing susceptible density  $\frac{S(t)}{N}$ .

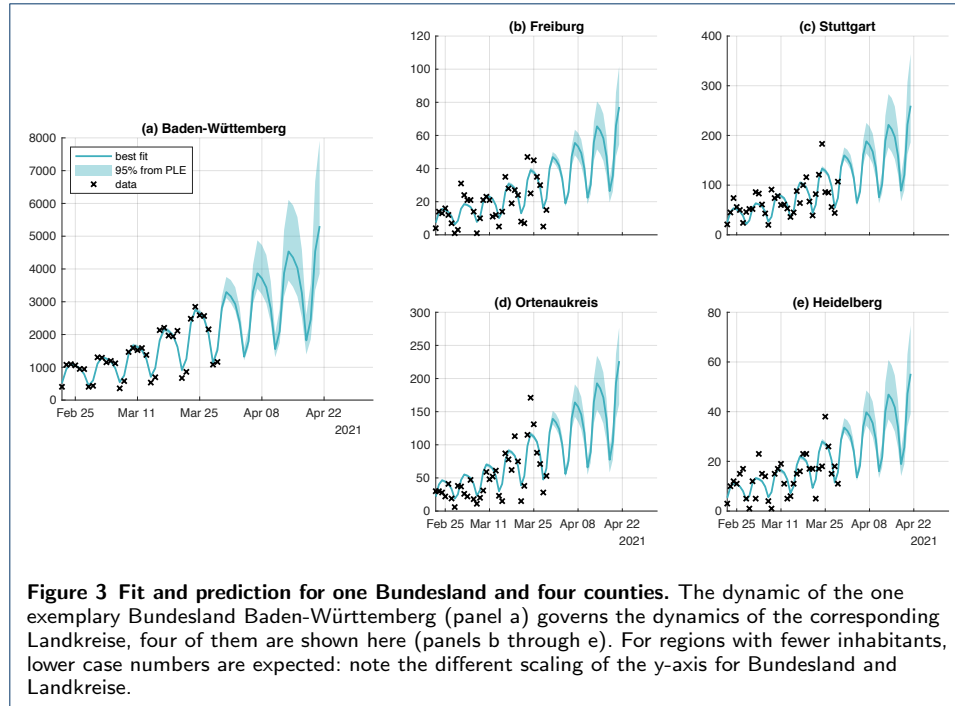
The estimated reproduction number  $R(t)$  oscillates around a value of 1 and illustrates the effect of describing the politics' countermeasures and the population's



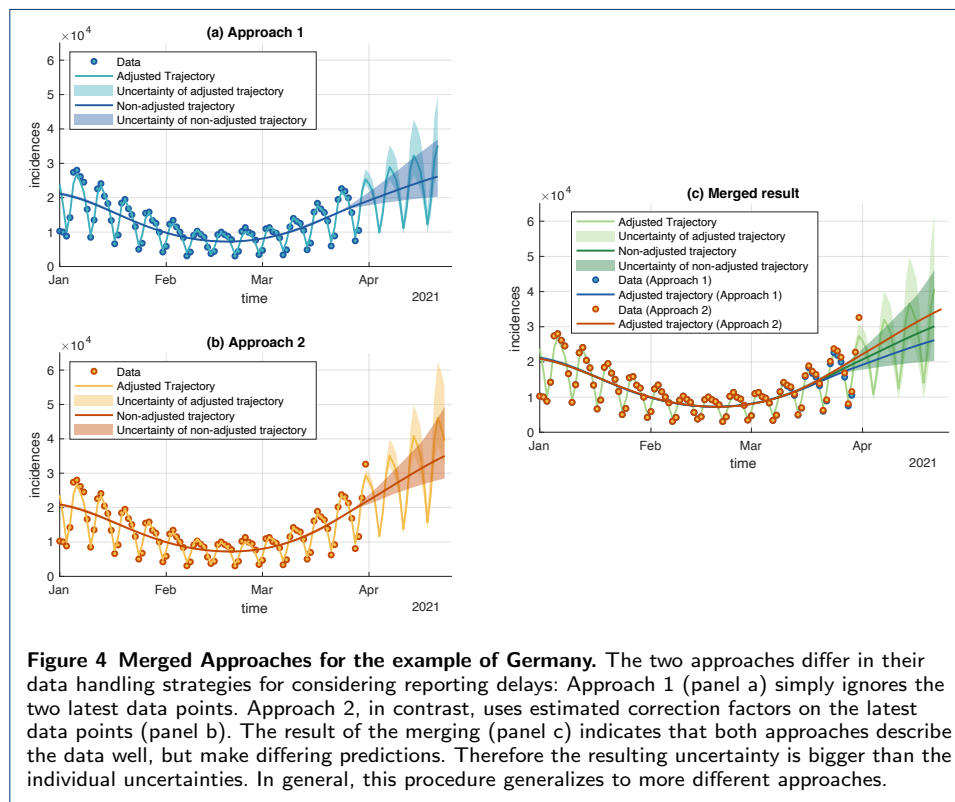
compliance to them (Figure 2, Panel (d)). This is in line with several publications [9], [19], [20] reporting similar behavior of the reproduction number. In general, oscillations in dynamical systems often are attributed to a feedback with delay, which is also the case here for the reproduction number  $R(t)$ . Several additional quantities of interest, such as the 7-day incidence or the cumulative number of cases can be computed from the model's predictions. In addition, the associated confidence intervals of these quantities can be determined using the parameter sets below the 95% threshold of likelihood profiles. We stress here again, that only the incidence data was used for model calibration (Figure 2, Panel a).

### 3.2 COVID-19 Spread in Subregions of Germany

For the county-level (*Landkreise*) we obtain results by the scaling approach described in section 2.3.2. The shape of dynamics is preserved and describes the latest data. Due the exponential scaling on later data points, it is unlikely that the entire time course is described well by the scaled dynamics. As we are primarily interested in the forecast, we display only the latest time interval. The data is more noisy due lower numbers of cases and inhabitants (Figure 3). Here, we show already merged results for clarity (see section 3.3).



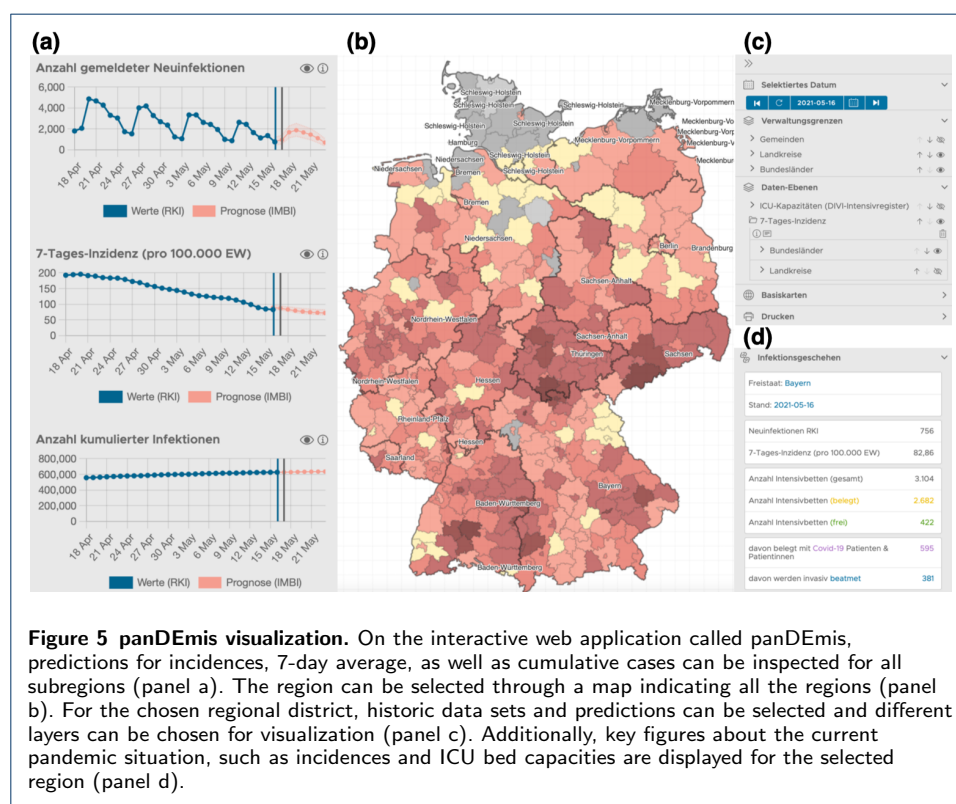
### 3.3 Approach Averaging



The analyses can be carried out for different approaches representing a variety of *a priori* equally feasible modeling strategies. To account for the uncertainty

that arises from (possibly over-)simplifying modelling assumptions, those different approaches are analyzed independent from each other. After results for all regional entities, i.e. federal states (as in 2 and counties 3) have been obtained for each approach, the results are merged into one comprehensive prediction, which features by construction (see 2.4.2) a higher uncertainty, now including both the uncertainty in the data and the uncertainty which modeling strategy is used. We illustrate this for two different approaches which differ only in the handling of the most recent data points (Figure 4). In general, this methodology generalizes to an arbitrary number of different approaches with the available computing resources as the only limiting factor.

### 3.4 Availability of Results



Sound political or social decisions are based on an empirical or prognostic foundation. To make the daily generated predictions available to various stakeholders, the forecasts are integrated into a web-application called panDEmis: In this interactive application, the recent infection situation is analyzed and displayed. For all registered users of the *DIVI Intensivregister* the tool is available at <https://pandemis.dlr.de/de/#/overview>. Current capacities of hospital beds and intensive care units, exposed population in the catchment areas of hospitals are merged with the forecast data. The combined display of all available data sets allows a situation picture for each day including also for past and future time steps. Figure 5 shows different features of the web-application from May 17th, 2021 for the occurrence of infection in the map entire Germany (panel b), as well as for the

selected administrative district of Bayern (panel a). Here, the blue graphs represent 1) the daily reported new infections by RKI, 2) the incidence of COVID-19 cases in the past 7 days per 100,000 people and 3) the cumulative infections. The prognosis is displayed as red curve, including a 95% confidence interval. All data can be interactively analyzed and visualized for different administrative units, i.e. federal states and county level.

The results of this incidence modeling approach are also a main predictor for a prediction analysis of ICU beds. The results of this second analysis step which is not detailed within this paper, is available for all registered users of the *DIVI Intensivregister* at <https://www.intensivregister.de/#/aktuelle-lage/prognosen>.

## 4 Discussion

Different model classes as ODE models or stochastic differential equation (SDE) models with or without mixed effects could be used for a data-driven parameter estimation approach. An SDE approach might be beneficial for small regions with low infection numbers or during times with very low total infection numbers. In these cases local outbreaks dominate the infection dynamics and the population is not *well-mixed* which renders an ODE approach ineffective. For the presented regional entities, the underlying assumptions for ODE modeling are reasonable and the ODE model was successfully adapted. We here focused on a pragmatic procedure that allows daily analysis and reliably calculates predictions.

When fitting data about the number of reported cases of an infectious disease outbreak, it is beneficial to fit incidences (or fluxes) instead of the total (or cumulative) number of cases [21]. The residuals of a fit on cumulative data will be correlated by construction. Most noise models assume independent measurement errors. Thus, the uncertainty will be underestimated in these cases and obtained results will be overly confident. By fitting the model to incidence data, the measurement errors remain uncorrelated.

The presented modeling approach heavily relies on the time-dependent infection rate  $\beta(t)$ . We assume dynamic processes to be continuously differentiable which leads to a smoothing of possible steps in the real infection rate which might occur due to rapid policy changes. Also, the temporal change  $\beta(t)$  incorporates many different mechanisms, which include but are not limited to: vaccinations, NPIs, changes in compliance to NPIs, viral mutations, seasonality and testing frequency. For an assumed constant vaccination rate, we saw that our approach delivers the same results when omitting the explicit vaccination state since  $\beta(t)$  is flexible enough to compensate the vaccination effect.

In general, it is *a priori* unclear how much flexibility this function should have. In the presented procedure, this corresponds to the number of knots employed in the spline. The spline's freedom should allow for a good fit of the dynamics, but also prevent overfitting.

Furthermore, the dynamics of the prediction are primarily determined by the value of  $R(t)$  at the latest data point. Hence, this value should not be estimated by too few data points meaning that the last spline knot should not be too close to the end of the time series.

Any prediction model used for forecasting should not exceed a certain time period as the future infection rate is hard to determine. But even at a short prediction

time span, it is unclear how recent political measures and the population's resulting behavior will alter the future infection rate. Therefore, we assume  $\beta(t)$  to be constant starting at the last data points. By additional precise knowledge about the effect of planned or recently made political decisions or other effects like weather conditions, this assumption could be further refined.

In contrast to other modeling approaches, we do not feed the actual NPIs into the model, but can instead correlate the estimated time development in a second step of the infection rate to NPIs. Quantifying the NPIs' effect and time lag on  $R(t)$  is difficult as most NPIs are not imposed or lifted independently of each other and estimates will therefore be highly correlated [22].

Whenever discussing the required amount of flexibility to obtain a good model fit, one should be aware of bias-variance-tradeoff: The introduction of more parameters included to explain a certain time dependence (reducing the bias), the bigger the resulting prediction uncertainty will be (increasing the variance). Similar arguments can be made when discussing the amount of utilized spline parameters or accounting for age structure. More available and consistent data can help.

There are no explicit states in our model to distinguish between recovered and dead people, mainly for the reason that there is no reliable data over the entire time course for those quantities. Recovered individuals are not tested to be non-sick anymore, and people who died were not consistently assessed in real-time in Germany.

Furthermore, the unobserved infected and infectious individuals are not in an explicit state. This fact is compensated by two aspects: Firstly, the used data does not contain information about the duration from beginning of infectivity to reporting to the local health authority. Thus, since the additional state would not help to better describe the used data, it is omitted. Secondly, the factor  $q$  introduced in the observation function in section 2.2 accounts for individuals that are overseen at all times. The estimated dark figure from equation (5) when fitting only incidence data is in the presented modeling approach in most regions compatible with a broad set of values ranging from 0.1 to 1 within the confidence level. This means that anywhere between 10% to 100% of all cases are detected by local authorities and both edge cases still agree sufficiently with the data. Therefore, the dark figure can not be estimated solely based on reported incidence cases. For reliable determination of the dark figure, additional testing in pre-specified cohorts is necessary.

## 5 Conclusions

We presented a data-driven ODE approach to fit and predict incidences of COVID-19 cases for different subregions of Germany. The key ingredients in doing so are 1) likelihood-based estimation and uncertainty quantification and 2) a time-dependent infection rate which is estimated by utilizing a cubic spline. All parameters are estimated from data and uncertainty in parameter estimates are translated to prediction uncertainty. As many different modeling assumptions will affect the outcomes, we average over similarly plausible approaches to account for this source of uncertainty. A major constraint for a feasible analysis strategy is a maximum runtime of 24 hours as the analysis should be repeated on a daily basis in an automated manner including the respectively newest data set.

In the future, more work for validation of competing modeling approaches and comparison of the various efforts undertaken in the currently highly dynamic field of mathematical modeling of infectious diseases is needed and will certainly be seen.

#### Competing interests

All authors have completed the ICMJE uniform disclosure form at [www.icmje.org/coi\\_disclosure.pdf](http://www.icmje.org/coi_disclosure.pdf) and declare: no support from any organization for the submitted work; no financial relationships with any organizations that might have an interest in the submitted work in the previous three years; no other relationships or activities that could appear to have influenced the submitted work.

#### Author's contributions

Lead conception of the overall study design: LG, HB; Conceived and designed the analysis: LR, FL, CK; Collected the data: MF, RKI; Contributed data or analysis tools: LR, FL, CK; Performed the analysis: LR, FL, CK; Integration and description of information system components (panDEmis): HT, TR; Wrote the paper: LR, FL, CK.

#### Acknowledgements

The project was funded by the Bundesministerium für Gesundheit (BMG).

We thank Matthäus Lottes, Janina Esins and the team from DIVI Intensivregister at the RKI.

Also, we thank the RKI's statisticians responsible for processing of the raw and routine data.

We thank Mario Menk, Steffen Weber-Carstens, Christian Karagiannidis, Uwe Janssens for fruitful discussions during project planning and implementation.

Thanks to Rafael Aruntjunjan for critically revising the manuscript.

#### Author details

<sup>1</sup>Institute of Medical Biometry and Statistics, Faculty of Medicine and Medical Center, University of Freiburg, Stefan Meier Str. 26, 79104 Freiburg, Germany. <sup>2</sup>Institute of Physics, University of Freiburg, Hermann-Herder-Str. 3, 79104 Freiburg, Germany. <sup>3</sup>Centre for Integrative Biological Signalling Studies (CIBSS), Schänzlestr. 18, 79104 Freiburg, Germany. <sup>4</sup>German Aerospace Center, Earth Observation Center, Münchener Str. 20, 82234 Weßling, Germany. <sup>5</sup>Institute for Geography and Geology, Julius-Maximilians-Universität Würzburg, Am Hubland, 97074 Würzburg, Germany. <sup>6</sup>Robert-Koch-Institute, Department for Methodology and Research Infrastructure, Nordufer 20, 13353 Berlin, Germany. <sup>7</sup>Charité - Universitätsmedizin Berlin, Department of Dermatology, Venerology and Allergology, Luisenstraße 2, 10117 Berlin, Germany. <sup>8</sup>Freiburg Center for Data Analysis and Modelling (FDM), University of Freiburg, Ernst-Zermelo-Str. 1, 79104 Freiburg, Germany.

#### References

1. an der Heiden, M., Hamouda, O.: Erfassung der SARS-CoV-2-Testzahlen in Deutschland - Nowcasting. *Epidemiologisches Bulletin* **17**, 10–17 (2020)
2. Günther, F., Bender, A., Katz, K., Küchenhoff, H., Höhle, M.: Nowcasting the COVID-19 pandemic in bavaria. *Biometrical Journal* **63**(3), 490–502 (2020). doi:10.1002/bimj.202000112
3. Kermack, W.O., McKendrick, A.G.: A contribution to the mathematical theory of epidemics. *Proc R Soc A* **115**(772), 700–721 (1927)
4. Keeling, M.J., Rohani, P.: *Modeling Infectious Diseases in Humans and Animals*, pp. 41–44. Princeton University Press, Princeton, NJ, USA (2008). doi:10.1515/9781400841035. <https://doi.org/10.1515/9781400841035>
5. Maier, B.F., Brockmann, D.: Effective containment explains subexponential growth in recent confirmed COVID-19 cases in China. *Science* **368**(6492), 742–746 (2020). doi:10.1126/science.abb4557. Publisher: American Association for the Advancement of Science Section: Research Article
6. Dehning, J., Zierenberg, J., Spitzner, F.P., Wibral, M., Neto, J.P., Wilczek, M., Priesemann, V.: Inferring change points in the spread of COVID-19 reveals the effectiveness of interventions. *Science* **369**(6500) (2020). doi:10.1126/science.abb9789. Publisher: American Association for the Advancement of Science Section: Research Article
7. Linka, K., Peirlinck, M., Kuhl, E.: The reproduction number of COVID-19 and its correlation with public health interventions. *Computational Mechanics* **66**(4), 1035–1050 (2020). doi:10.1007/s00466-020-01880-8
8. Flaxman, S., Mishra, S., Gandy, A., Unwin, H.J.T., Mellan, T.A., Coupland, H., Whittaker, C., Zhu, H., Berah, T., Eaton, J.W., Monod, M., Ghani, A.C., Donnelly, C.A., Riley, S., Vollmer, M.A.C., Ferguson, N.M., Okell, L.C., Bhatt, S.: Estimating the effects of non-pharmaceutical interventions on COVID-19 in Europe. *Nature* **584**(7820), 257–261 (2020). doi:10.1038/s41586-020-2405-7. Number: 7820 Publisher: Nature Publishing Group
9. Dings, C., Götz, K., Och, K., Sihnevich, I., Selzer, D., Werthner, Q., Kovar, L., Marok, F., Schräpel, C., Fuhr, L., Türk, D., Britz, H., Smola, S., Volk, T., Kreuer, S., Rissland, J., Lehr, T.: *Mathematische Modellierung und Vorhersage von COVID-19 Fällen, Hospitalisierung (inkl. Intensivstation und Beatmung) und Todesfällen in dendeutschen Bundesländern* (2021). [https://covid-simulator.com/wp-content/uploads/2021/04/Report\\_2021\\_03\\_31.pdf](https://covid-simulator.com/wp-content/uploads/2021/04/Report_2021_03_31.pdf) Accessed April 1st 2021
10. Schelker, M., Raue, A., Timmer, J., Kreutz, C.: Comprehensive estimation of input signals and dynamics in biochemical reaction networks. *Bioinformatics* **28**(18), 529–534 (2012). doi:10.1093/bioinformatics/bts393
11. Noll, N.B., Aksamentov, I., Druelle, V., Badenhorst, A., Ronzani, B., Jefferies, G., Albert, J., Neher, R.A.: COVID-19 Scenarios: an interactive tool to explore the spread and associated morbidity and mortality of SARS-CoV-2. medRxiv, 2020–050520091363 (2020). doi:10.1101/2020.05.05.20091363. Publisher: Cold Spring Harbor Laboratory Press. Accessed 2021-05-06

12. Contreras, S., Dehning, J., Loidolt, M., Zierenberg, J., Spitzner, F.P., Urrea-Quintero, J.H., Mohr, S.B., Wilczek, M., Wibral, M., Priesemann, V.: The challenges of containing SARS-CoV-2 via test-trace-and-isolate. *Nature Communications* **12**(1), 378 (2021). doi:10.1038/s41467-020-20699-8. Number: 1 Publisher: Nature Publishing Group. Accessed 2021-05-06
13. Kreisfreie Städte und Landkreise nach Fläche, Bevölkerung und Bevölkerungsdichte am 31.12.2019 - Statistisches Bundesamt. <https://www.destatis.de/DE/Themen/Laender-Regionen/Regionales/Gemeindeverzeichnis/Administrativ/04-kreise.html>
14. Raue, A., Schilling, M., Bachmann, J., Matteson, A., Schelker, M., Kaschek, D., Hug, S., Kreutz, C., Harms, B.D., Theis, F.J., Klingmüller, U., Timmer, J.: Lessons Learned from Quantitative Dynamical Modeling in Systems Biology. *PLoS ONE* **8**(12) (2013). doi:10.1371/annotation/ea0193d8-1f7f-492a-b0b7-d877629fdcee
15. Kreutz, C., Raue, A., Kaschek, D., Timmer, J.: Profile likelihood in systems biology. *FEBS Journal* **280**(11), 2564–2571 (2013). doi:10.1111/febs.12276
16. Tönsing, C., Timmer, J., Kreutz, C.: Profile likelihood-based analyses of infectious disease models. *Statistical Methods in Medical Research*, 962280217746444 (2017). doi:10.1177/0962280217746444
17. Steiert, B., Raue, A., Timmer, J., Kreutz, C.: Experimental Design for Parameter Estimation of Gene Regulatory Networks. *PLOS ONE* **7**(7), 40052 (2012). doi:10.1371/journal.pone.0040052. Publisher: Public Library of Science
18. Kreutz, C., Raue, A., Timmer, J.: Likelihood based observability analysis and confidence intervals for predictions of dynamic models. *BMC Systems Biology* **6**(1), 120 (2012). doi:10.1186/1752-0509-6-120
19. Khailaie, S., Mitra, T., Bandyopadhyay, A., Schips, M., Mascheroni, P., Vanella, P., Lange, B., Binder, S.C., Meyer-Hermann, M.: Development of the reproduction number from coronavirus SARS-CoV-2 case data in Germany and implications for political measures. *BMC Medicine* **19**(1), 32 (2021). doi:10.1186/s12916-020-01884-4
20. Abbott, S., Hellewell, J., Thompson, R.N., Sherratt, K., Gibbs, H.P., Bosse, N.I., Munday, J.D., Meakin, S., Doughty, E.L., Chun, J.Y., Chan, Y.-W.D., Finger, F., Campbell, P., Endo, A., Pearson, C.A.B., Gimma, A., Russell, T., CMMID COVID modelling group, Flasche, S., Kucharski, A.J., Eggo, R.M., Funk, S.: Estimating the time-varying reproduction number of SARS-CoV-2 using national and subnational case counts. *Wellcome Open Research* **5**, 112 (2020). doi:10.12688/wellcomeopenres.16006.1. Accessed 2021-05-06
21. King, A.A., Domenech de Cellès, M., Magpantay, F.M.G., Rohani, P.: Avoidable errors in the modelling of outbreaks of emerging pathogens, with special reference to Ebola. *Proceedings of the Royal Society B: Biological Sciences* **282**(1806) (2015). doi:10.1098/rspb.2015.0347
22. Haug, N., Geyrhofer, L., Londei, A., Dervic, E., Desvars-Larrive, A., Loreto, V., Pinior, B., Thurner, S., Klimek, P.: Ranking the effectiveness of worldwide COVID-19 government interventions. *Nature Human Behaviour* **4**(12), 1303–1312 (2020). doi:10.1038/s41562-020-01009-0. Number: 12 Publisher: Nature Publishing Group. Accessed 2021-05-06