

Data-Based Analysis, Modelling and Forecasting of the COVID-19 outbreak

Cleo Anastassopoulou^{1*}, Lucia Russo², Athanasios Tsakris¹, Constantinos Siettos^{3*}

1 Department of Microbiology, Medical School, University of Athens, Athens, Greece

2 Consiglio Nazionale delle Ricerche, Science and Technology for Energy and Sustainable Mobility, Napoli, Italy

3 Dipartimento di Matematica e Applicazioni “Renato Caccioppoli”, Università degli Studi di Napoli Federico II, Napoli, Italy

* constantinos.siettos@unina.it * cleoa@med.uoa.gr

Abstract

Since the first suspected case of coronavirus disease-2019 (COVID-19) on December 1st, 2019, in Wuhan, Hubei Province, China, a total of 40,235 confirmed cases and 909 deaths have been reported in China up to February 10, 2020, evoking fear locally and internationally. Here, based on the publicly available epidemiological data [1, 2] (WHO, CDC, ECDC, NHC and DXY), for Hubei from January 11 to February 10, 2020, we provide estimates of the main epidemiological parameters, i.e. the basic reproduction number (R_0) and the infection, recovery and mortality rates, along with their 90% confidence intervals. As the number of infected individuals, especially of those with asymptomatic or mild courses, is suspected to be much higher than the official numbers, which can be considered only as a sample of the actual numbers of infected and recovered cases in the population, we have repeated the calculations under a second scenario that considers twenty times the number of confirmed infected cases and forty times the number of recovered, leaving the number of deaths unchanged. Our computations and analysis were based on a mean field Susceptible-Infected-Recovered-Dead (SIRD) model. Based on the reported data, the expected value of R_0 as computed considering the period from the 11th of January until the 18th of January, using the official counts of confirmed cases was found to be 4.6, while the one computed under the second scenario was found to be 3.2. Thus, based on the SIRD simulations, the estimated average value of R_0 under both scenarios was found to be 2.4. Furthermore, using the estimated parameters from both scenarios, we provide tentative three-week forecasts of the evolution of the outbreak at the epicenter. Our forecasting flashes a note of caution for the presently unfolding outbreak in China. Based on the official counts for confirmed cases, the simulations suggest that the cumulative number of infected will surpass 62,000 (as a lower bound) and could reach 140,000 (with an upper bound of 300,000) by February 29. Regarding the number of deaths, simulations forecast that on the basis of the up to the 10th of February data and estimations of the actual numbers of infected and recovered in the population, the death toll might exceed 5,500 by February 29. However, our analysis further reveals a significant decline of the mortality rate to which various factors may have contributed, such as the severe control measures taken in Hubei, China (e.g. quarantine and hospitalization of infected individuals), but mainly because of the fact that the actual cumulative numbers of infected and recovered cases in the population are estimated to be of the order of twenty times higher for the infected and forty times higher for the

recovered than reported, thus resulting in a much lower mortality rate, which according to our computations is of the order of 0.15%.

Introduction

An outbreak of “pneumonia of unknown etiology” in Wuhan, Hubei Province, China in early December 2019 has spiraled into an epidemic that is ravaging China and threatening to reach a pandemic state [3]. The causative agent soon proved to be a new betacoronavirus related to the Middle East Respiratory Syndrome virus (MERS-CoV) and the Severe Acute Respiratory Syndrome virus (SARS-CoV). The novel coronavirus SARS-CoV-2 disease has been named “COVID-19” by the World Health Organization (WHO) and on January 30, the COVID-19 outbreak was declared to constitute a Public Health Emergency of International Concern by the WHO Director-General [4]. Despite the lockdown of Wuhan and the suspension of all public transport, flights and trains on January 23, a total of 40,235 confirmed cases, including 6,484 (16.1%) with severe illness, and 909 deaths (2.2%) had been reported in China by the National Health Commission up to February 10, 2020; meanwhile, 319 cases and one death were reported outside of China, in 24 countries [5]. The origin of COVID-19 has not yet been determined although preliminary investigations are suggestive of a zoonotic, possibly of bat, origin [6,7]. Similarly to SARS-CoV and MERS-CoV, the novel virus is transmitted from person to person principally by respiratory droplets, causing such symptoms as fever, cough, and shortness of breath after a period believed to range from 2 to 14 days following infection, according to the Centers for Disease Control and Prevention (CDC) [3,8,9]. Preliminary data suggest that older males with comorbidities may be at higher risk for severe illness from COVID-19 [8,10,11]. However, the precise virologic and epidemiologic characteristics, including transmissibility and mortality, of this third zoonotic human coronavirus are still unknown. Using the serial intervals (SI) of the two other well-known coronavirus diseases, MERS and SARS, as approximations for the true unknown SI, Zhao et al. estimated the mean basic reproduction number (R_0) of SARS-CoV-2 to range between 2.24 (95% CI: 1.96-2.55) and 3.58 (95% CI: 2.89-4.39) in the early phase of the outbreak [12]. Very similar estimates, 2.2 (95% CI: 1.4-3.9), were obtained for R_0 at the early stages of the epidemic by Imai et al. 2.6 (95% CI: 1.5-3.5) [13], as well as by Li et al., who also reported a doubling in size every 7.4 days [3]. Wu et al. estimated the R_0 at 2.68 (95% CI: 2.47–2.86) with a doubling time every 6.4 days (95% CI: 5.8–7.1) and the epidemic growing exponentially in multiple major Chinese cities with a lag time behind the Wuhan outbreak of about 1–2 weeks [14]. Amidst such an important ongoing public health crisis that also has severe economic repercussions, we reverted to mathematical modelling that can shed light to essential epidemiologic parameters that determine the fate of the epidemic [15]. Here, we present the results of the analysis of time series of epidemiological data available in the public domain (WHO, CDC, ECDC, NHC and DXY) from January 11 to February 10, 2020, and attempt a three-week forecast of the spreading dynamics of the emerged coronavirus epidemic in the epicenter in mainland China.

Methodology

Our analysis was based on the publicly available data of the new confirmed daily cases reported for the Hubei province from the 11th of January until the 10th of February [1,2]. Based on the released data, we attempted to estimate the mean values of the main epidemiological parameters, i.e. the basic reproduction number R_0 , the infection (α), recovery (β) and mortality rate (γ), along with their 90% confidence

intervals. However, as suggested [16], the number of infected, and consequently the number of recovered, people is likely to be much higher. Thus, in a second scenario, we have also derived results by taking twenty times the number of reported cases for the infected and forty times the number for the recovered cases, while keeping constant the number of deaths that is more likely to be closer to the real number. Based on these estimates we also provide tentative forecasts until the 29th of February. Our methodology follows a two-stage approach as described below.

The basic reproduction number is one of the key values that can predict whether the infectious disease will spread into a population or die out. R_0 represents the average number of secondary cases that result from the introduction of a single infectious case in a totally susceptible population during the infectiousness period. Based on the reported data of confirmed cases, we provide estimations of the R_0 from the 16th up to the 20th of January in order to satisfy as much as possible the hypothesis of $S \approx N$ that is a necessary condition for the computation of R_0 .

We also provide estimations of β , γ over the entire period using a rolling window of one day from the 11th of January to the 16th of January to provide the very first estimations.

Furthermore, we exploited the SIRD model to provide an estimation of the infection rate α . This is accomplished by starting with one infected person on the 16th of November, which has been suggested as a starting date of the epidemic [6], run the SIR model until the 10th of February and optimize α to fit the reported confirmed cases from the 11th of January to the 10th of February. Below, we describe analytically our approach.

Let us start by denoting with $S(t)$, $I(t)$, $R(t)$, $D(t)$, the number of susceptible, infected, recovered and dead persons respectively at time t in the population of size N . For our analysis we assume that the total number of the population remains constant. Based on the demographic data for the province of Hubei $N = 59m$. Thus, the discrete SIRD model reads:

$$S(t) = S(t - 1) - \frac{\alpha}{N} S(t - 1) I(t - 1) \quad (1)$$

$$I(t) = I(t - 1) + \frac{\alpha}{N} S(t - 1) I(t - 1) - \beta I(t - 1) - \gamma I(t - 1) \quad (2)$$

$$R(t) = R(t - 1) + \beta I(t - 1) \quad (3)$$

$$D(t) = D(t - 1) + \gamma I(t - 1) \quad (4)$$

The above system is defined in discrete time points $t = 1, 2, \dots$, with the corresponding initial condition at the very start of the epidemic: $S(0) = N - 1$, $I(0) = 1$, $R(0) = D(0) = 0$.

Initially, when the spread of the epidemic starts, all the population is considered to be susceptible, i.e. $S \approx N$. Based on this assumption, by Eq.(2),(3),(4), the basic reproduction number can be estimated by the parameters of the SIRD model as:

$$R_0 = \frac{\alpha}{\beta + \gamma} \quad (5)$$

A problem with the approximation of the epidemiological parameters α , β and γ and thus R_0 , from real-world data based on the above expressions is that in general, for large scale epidemics, the actual number of infected $I(t)$ persons in the whole population is unknown. Thus, the problem is characterized by high uncertainty. However, one can attempt to provide some coarse estimations of the epidemiological parameters based on the reported confirmed cases using the approach described next.

0.1 Identification of epidemiological parameters from the reported data of confirmed cases

Lets us denote with $\Delta I(t) = I(t) - I(t - 1)$, $\Delta R(t) = R(t) - R(t - 1)$, $\Delta D(t) = D(t) - D(t - 1)$, the reported new cases of infected, recovered and deaths at time t , with $C\Delta I(t)$, $C\Delta R(t)$, $C\Delta D(t)$ the cumulative numbers of confirmed cases at time t . Thus:

$$C\Delta X(t) = \sum_{i=1}^t \Delta X(t), \quad (6)$$

where, $X = I, R, D$.

Let us also denote by $C\Delta \mathbf{X}(t) = [C\Delta X(1), C\Delta X(2), \dots, C\Delta X(t)]^T$, the $t \times 1$ column vector containing the corresponding cumulative numbers up to time t . On the basis of Eqs.(2), (3), (4), one can provide a coarse estimation of the parameters R_0 , β and γ as follows.

Starting with the estimation of R_0 , we note that as the province of Hubei has a population of 59m, one can reasonably assume that for any practical means, at least at the beginning of the outbreak, $S \approx N$. By making this assumption, one can then provide an approximation of the expected value of R_0 using Eq.(5) and Eq.(2), Eq.(3), Eq.(4). In particular, substituting in Eq.(2), the terms $\beta I(t - 1)$ and $\gamma I(t - 1)$ with $\Delta R(t) = R(t) - R(t - 1)$ from Eq.(3), and $\Delta D(t) = D(t) - D(t - 1)$ from Eq.(4) and bringing them into the left-hand side of Eq.(2), we get:

$$I(t) - I(t - 1) + R(t) - R(t - 1) + D(t) - D(t - 1) = \frac{\alpha}{N} S(t - 1) I(t - 1) \quad (7)$$

Adding Eq.(3) and Eq.(4), we get:

$$R(t) - R(t - 1) + D(t) - D(t - 1) = \beta I(t - 1) + \gamma I(t - 1) \quad (8)$$

Finally, assuming that for any practical means at the beginning of the spread that $S(t - 1) \approx N$ and dividing Eq.(7) by Eq.(8) we get:

$$\frac{I(t) - I(t - 1) + R(t) - R(t - 1) + D(t) - D(t - 1)}{R(t) - R(t - 1) + D(t) - D(t - 1)} = \frac{\alpha}{\beta + \gamma} = R_0 \quad (9)$$

Note that one can use directly Eq.(9) to compute R_0 with regression, without the need to compute first the other parameters, i.e. β , γ and α .

At this point, the regression can be done either by using the differences per se, or by using the corresponding cumulative functions (instead of the differences for the calculation of R_0 using Eq.(9)). Indeed, it is easy to prove that by summing up both sides of Eq.(7) and Eq.(8) over time and then dividing them we get the following equivalent expression for the calculation of R_0 .

$$\frac{C\Delta I(t) + C\Delta R(t) + C\Delta D(t)}{C\Delta R(t) + C\Delta D(t)} = \frac{\alpha}{\beta + \gamma} = R_0 \quad (10)$$

Here, we have least squares using Eq. (10) to estimate R_0 in order to reduce the noise included in the differences. Note that the above expression is a valid approximation only at the beginning of the spread of the disease.

Thus, based on the above, a coarse estimation of R_0 and its corresponding confidence intervals can be provided by solving a linear regression problem using least-squares problem as:

$$\hat{R}_0 = ([\mathbf{C}\Delta\mathbf{R}(t) + \mathbf{C}\Delta\mathbf{D}(t)]' [\mathbf{C}\Delta\mathbf{R}(t) + \mathbf{C}\Delta\mathbf{D}(t)])^{-1} [\mathbf{C}\Delta\mathbf{R}(t) + \mathbf{C}\Delta\mathbf{D}(t)]' [\mathbf{C}\Delta\mathbf{I}(t) + \mathbf{C}\Delta\mathbf{R}(t) + \mathbf{C}\Delta\mathbf{D}(t)], \quad (11)$$

The prime (') is for the transpose operation. 120

The estimation of the mortality and recovery rates was based on the simple formula used by the National Health Commission (NHC) of the People's Republic of China [17] (see also [18]) that is the ratio of the cumulative number of recovered/deaths and that of infected at time t . Thus, a coarse estimation of the mortality can be calculated solving a linear regression problem for the corresponding cumulative functions by least squares as: 121
122
123
124
125

$$\hat{\gamma} = [(\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))' (\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))]^{-1} \quad (12)$$

$$(\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))' \mathbf{C}\Delta\mathbf{D}(t), \quad (13)$$

In a similar manner, the recovery rate can be computed as: 126

$$\hat{\beta} = [(\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))' (\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))]^{-1} \quad (14)$$

$$(\mathbf{C}\Delta\mathbf{I}(t) - \mathbf{C}\Delta\mathbf{D}(t) - \mathbf{C}\Delta\mathbf{R}(t))' \mathbf{C}\Delta\mathbf{R}(t), \quad (15)$$

As the reported data are just a subset of the actual number of infected and recovered cases including the asymptomatic and/or mild ones, we have repeated the above calculations considering twenty times the reported number of infected and forty times the reported number of recovered in the population, while leaving the reported number of deaths the same given that their cataloguing is close to the actual number of deaths due to COVID-19. 127
128
129
130
131
132

0.2 Estimation of the infection rate from the SIRD model 133

Having estimated the expected values of the parameters β and γ , an approximation of the infected rate α , that is not biased by the assumption of $S = N$ can be obtained by using the SIRD simulator. In particular, in the SIRD model we set $\hat{\beta}$ and $\hat{\gamma}$, and set as initial conditions one infected person on the 16th of November and run the simulator until the last date for which there are available data (here up to the 10th of February). Then, the value of the infection rate α can be found by “wrapping” around the SIRD simulator an optimization algorithm (such as a nonlinear least-squares solver) to solve the problem: 134
135
136
137
138
139
140
141

$$\operatorname{argmin}_{\alpha} \left\{ \sum_{t=1}^M (w_1 f_t(\alpha; \hat{\beta}, \hat{\gamma})^2 + w_2 g_t(\alpha; \hat{\beta}, \hat{\gamma})^2 + w_3 h_t(\alpha; \hat{\beta}, \hat{\gamma})^2) \right\}, \quad (16)$$

where

$$f_t(\alpha; \hat{\beta}, \hat{\gamma}) = C\Delta I^{SIRD}(t) - C\Delta I(t),$$

$$g_t(\alpha; \hat{\beta}, \hat{\gamma}) = C\Delta R^{SIRD}(t) - C\Delta R(t),$$

$$h_t(\alpha; \hat{\beta}, \hat{\gamma}) = C\Delta D^{SIRD}(t) - C\Delta D(t)$$

where, $C\Delta X^{SIRD}(t)$, ($X = I, R, D$) are the cumulative cases resulting from the SIRD simulator at time t ; w_1, w_2, w_3 correspond to scalars serving in the general case as weights to the relevant functions. For the solution of the above optimization problem we used the function “lsqnonlin” of matlab [19] using the Levenberg-Marquard algorithm. 142
143
144
145

1 Results

As discussed, we have derived results using two different scenarios (see in Methodology). For each scenario, we first present the results obtained by solving the least squares problem as described in section 0.1 using a rolling window of an one-day step. The first time window was that from the 11th up to the 16th of January i.e. we used the first six days to provide the very first estimations of the epidemiological parameters. We then proceeded with the calculations by adding one day in the rolling window as described in the methodology until the 10th of February. We also report the corresponding 90% confidence intervals instead of the more standard 95% because of the small size of the data. For each window, we also report, the corresponding coefficients of determination (R^2) representing the proportion of the variance in the dependent variable that is predictable from the independent variables, and the root mean square of error (RMSE). The estimation of R_0 was based on the data until January 20, in order to satisfy as much as possible the hypothesis underlying its calculation by Eq.(9).

Then, we used the SIRD model to provide an estimation of infection rate by “wrapping” around the SIRD simulator the optimization algorithm as described in section 0.2. Finally, we provide tentative forecasts for the evolution of the outbreak based on both scenarios until the end of February.

1.1 Scenario I: Results obtained using the exact numbers of the reported confirmed cases

Figure1 depicts an estimation of R_0 for the period January 16-January 20. Using the first six days from the 11th of January, \hat{R}_0 results in ~ 4.80 (90% CI: 3.36-6.67); using the data until January 17, \hat{R}_0 results in ~ 4.60 (90% CI: 3.56-5.65); using the data until January 18, \hat{R}_0 results in ~ 5.14 (90%CI: 4.25-6.03); using the data until January 19, \hat{R}_0 results in ~ 6.09 (90% CI: 5.02-7.16); and using the data until January 20, \hat{R}_0 results in ~ 7.09 (90% CI: 5.84-8.35)

Figure2 depicts the estimated values of the recovery (β) and mortality (γ) rates for the period January 16 to February 10. The confidence intervals are also depicted with dashed lines. Note that the large variation in the estimated values of β and γ should be accounted to the small size of the data and data uncertainty. This is also reflected in the corresponding confidence intervals. As more data are taken into account, this variation is significantly reduced. Thus, using all the available data from the 11th of January until the 10th of February, the estimated value of the mortality rate γ is $\sim 3.2\%$ (90% CI: 3.1%-3.3%) and that of the recovery rate is ~ 0.054 (90% CI: 0.049-0.060) corresponding to ~ 18 days (90% CI: 16-20). It is interesting to note that as the available data become more, the estimated recovery rate increases significantly from the 31th of January (see Fig.2).

In Figures3,4,5, we show the coefficients of determination (R^2) and the root of mean squared errors ($RMSE$), for \hat{R}_0 , $\hat{\beta}$ and $\hat{\gamma}$, respectively.

As described in the methodology, we have also used the SIRD simulator to provide an estimation of the infection rate by optimization with $w_1=1$, $w_2=2$, $w_3=2$. Thus, we performed the simulations by setting $\beta=0.054$ and $\gamma=0.032$, and as initial conditions one infected, zero recovered and zero deaths on November 16th 2019, and run until the 10th of February. The optimal, with respect to the reported confirmed cases from the 11th of January to the 10th of February, value of the infected rate (α) was ~ 0.206 (90% CI: 0.204-0.208). This corresponds to a mean value of the basic reproduction number $\hat{R}_0 \approx 2.4$. Note that this value is different compared to the value that was estimated using solely the reported data.

Finally, using the derived values of the parameters α , β , γ , we have run the SIRD

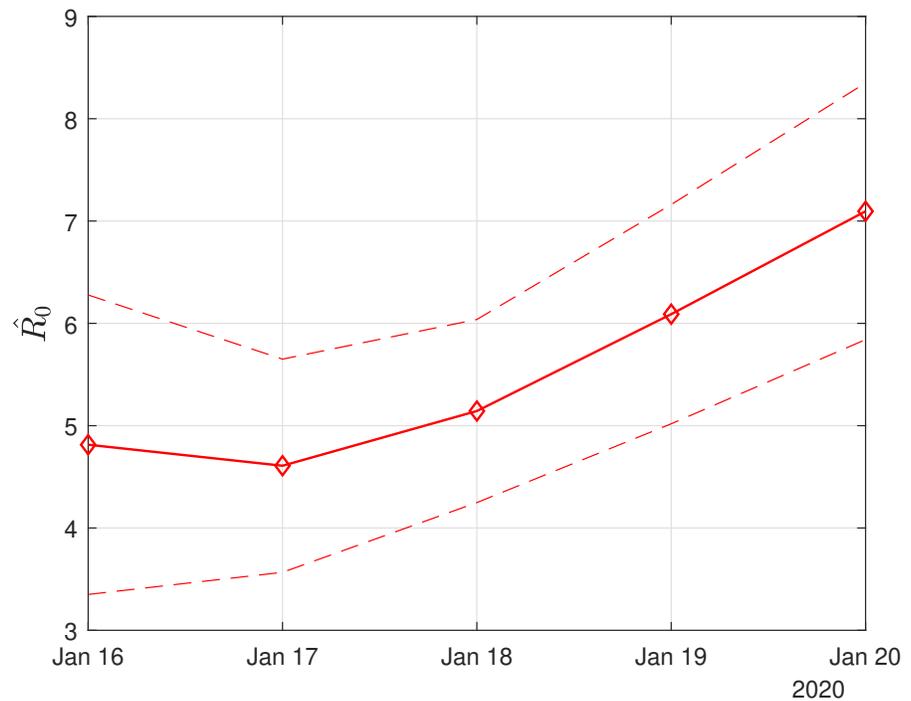


Fig 1. Scenario I. Estimated values of the basic reproduction number (R_0) as computed by least squares using a rolling window with initial date the 11th of January. The solid line corresponds to the mean value and dashed lines to lower and upper 90% confidence intervals.

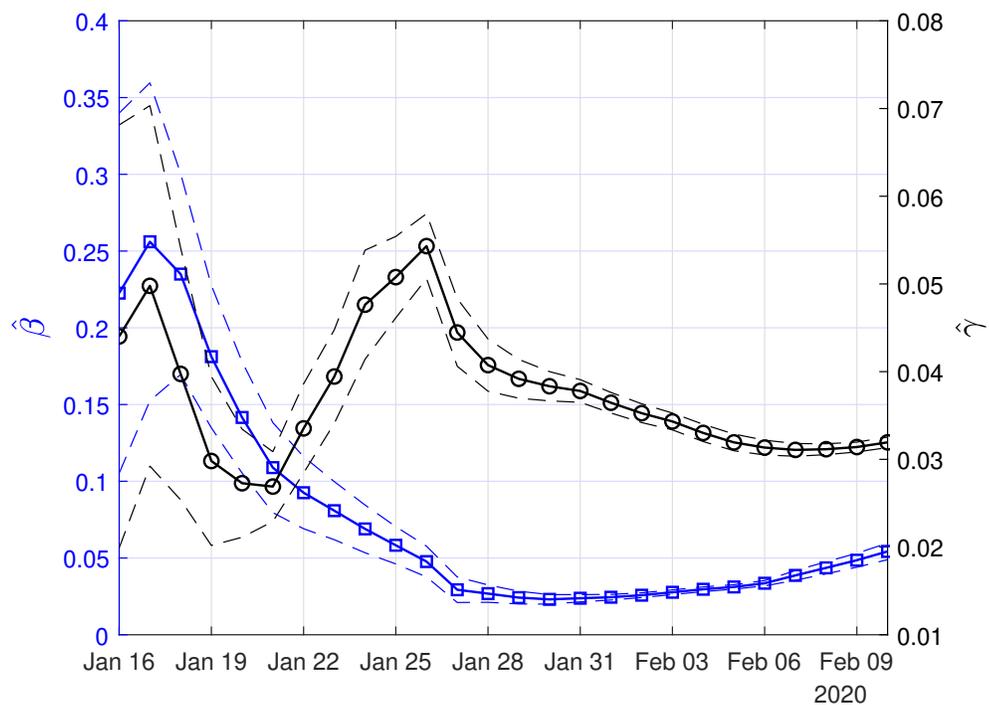


Fig 2. Scenario I. Estimated values of the recovery ($\hat{\beta}$) and mortality ($\hat{\gamma}$) rate as computed by least squares using a rolling window (see section 0.1). Solid lines correspond to the mean values and dashed lines to lower and upper 90% confidence intervals

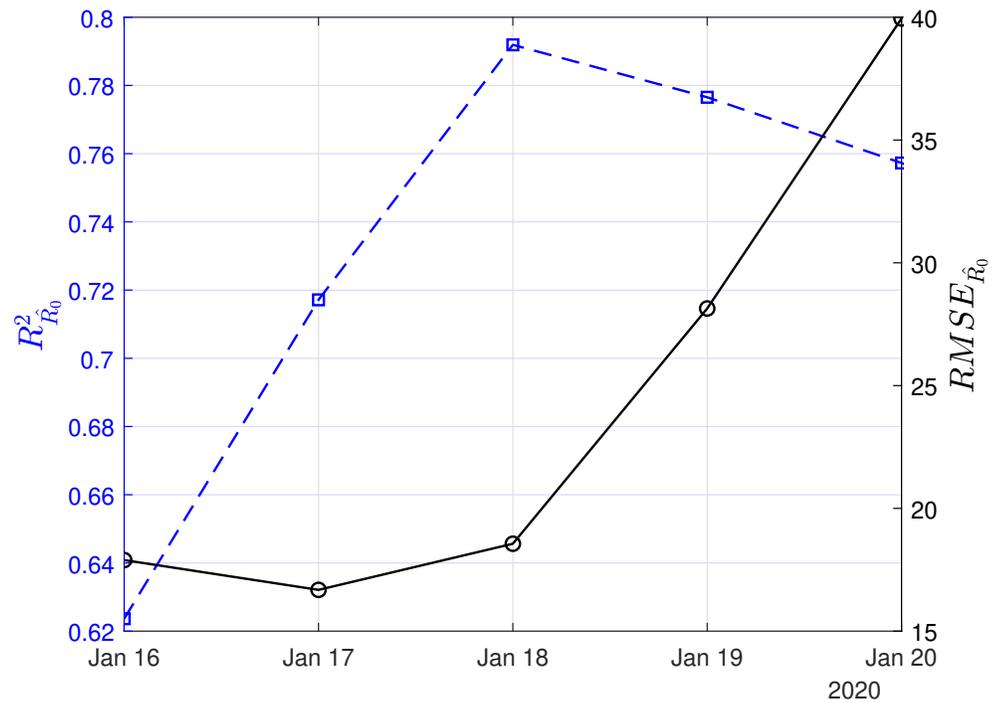


Fig 3. Scenario I. Coefficient of determination (R^2) and root mean square error ($RMSE$) resulting from the solution of the linear regression problem with least-squares for the basic reproduction number (R_0)

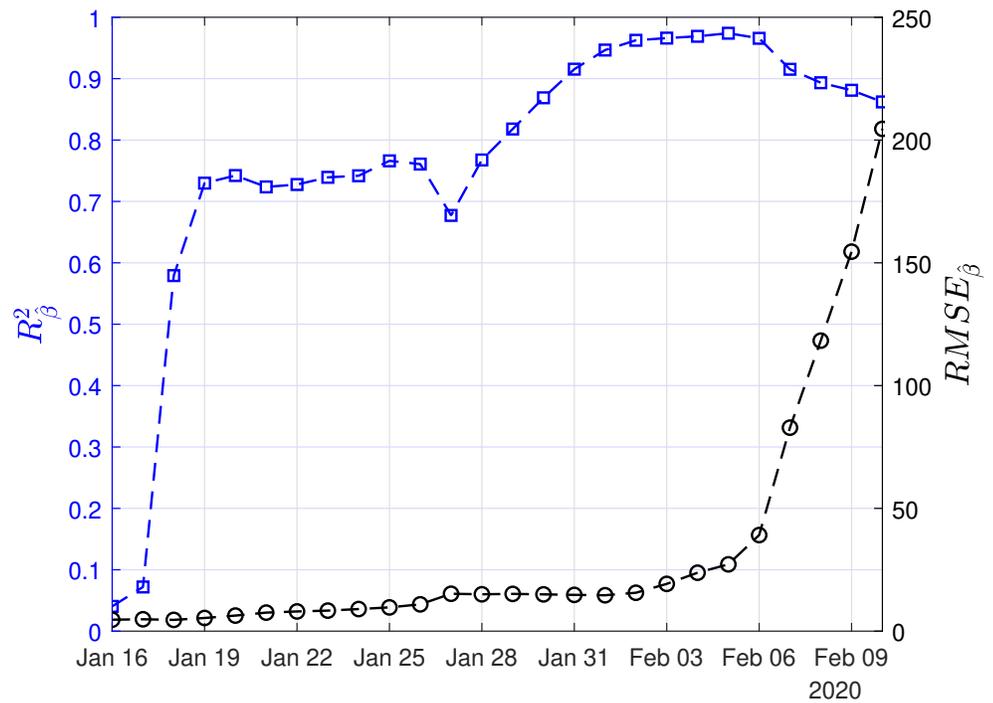


Fig 4. Scenario I. Coefficient of determination (R^2) and root mean square error ($RMSE$) resulting from the solution of the linear regression problem with least-squares for the recovery rate (β)

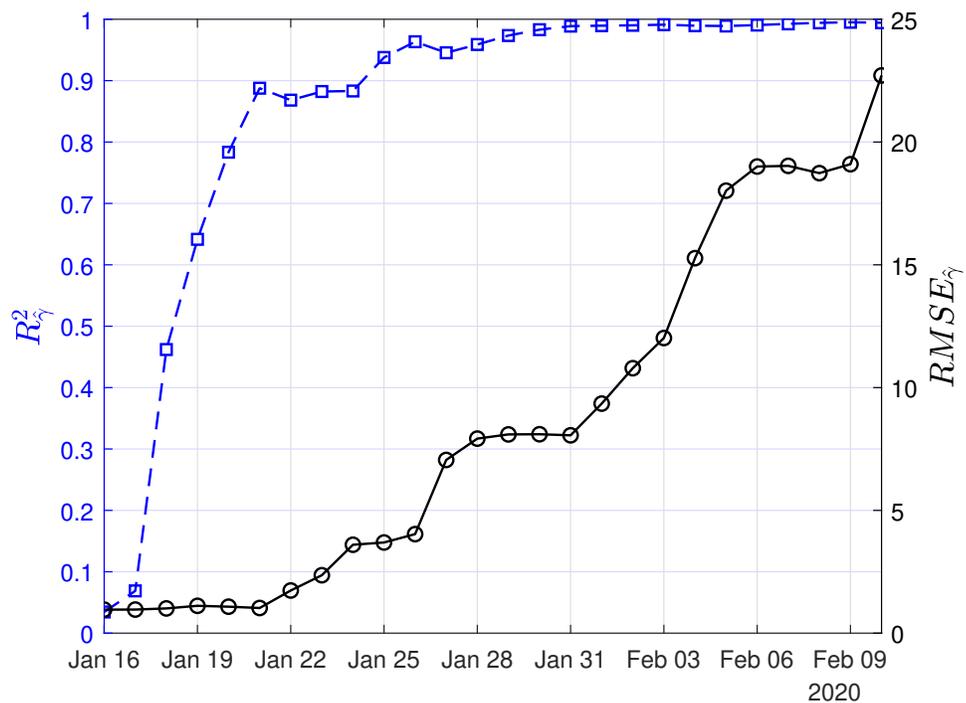


Fig 5. Scenario I. Coefficient of determination (R^2) and root mean square error (RMSE) resulting from the solution of the linear regression problem with least-squares for the mortality rate (γ)

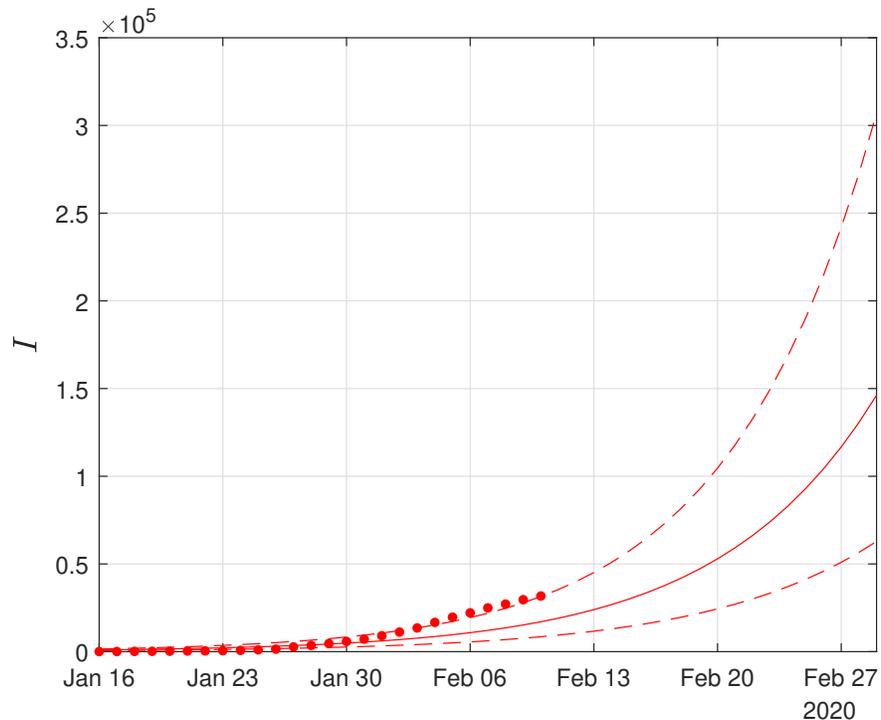


Fig 6. Scenario I. Simulations until the 29th of February of the cumulative number of infected as obtained using the SIRD model. Dots correspond to the number of confirmed cases from the 16th of January to the 10th of February. The initial date of the simulations was the 16th of November with one infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.206$, $\beta = 0.054$, $\gamma = 0.032$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

simulator until the end of February. The results of the simulations are given in Figures 6, 7, 8. Solid lines depict the evolution, when using the expected (mean) estimations and dashed lines illustrate the corresponding lower and upper bounds as computed at the limits of the confidence intervals of the estimated parameters. 195

As Figures 6, 7 suggest, the forecast of the outbreak at the end of February, through the SIRD model is characterized by high uncertainty. In particular, simulations result in an expected number of $\sim 140,000$ infected cases but with a high variation: the lower bound is at $\sim 62,000$ infected cases while the upper bound is at $\sim 300,000$ cases. Similarly for the recovered population, simulations result in an expected number of $\sim 65,000$, while the lower and upper bounds are at $\sim 40,000$ and $\sim 118,000$, respectively. Finally, regarding the deaths, simulations result in an average number of $\sim 39,000$, with lower and upper bounds, $\sim 20,000$ and $\sim 67,000$, respectively. 196
197
198
199
200
201
202
203
204
205
206

However, as more data are released it appears that the mortality rate is much lower than the predicted with the current data and thus the death toll is expected, to be significantly less compared with the predictions. 207
208
209

Furthermore, simulations reveal that the confirmed cumulative number of deaths is significantly smaller than the lower bound of the simulations. This suggests that the mortality rate is considerably lower than the estimated one based on the officially reported data. Thus, it is expected that the actual numbers of the infected cases, and 210
211
212
213

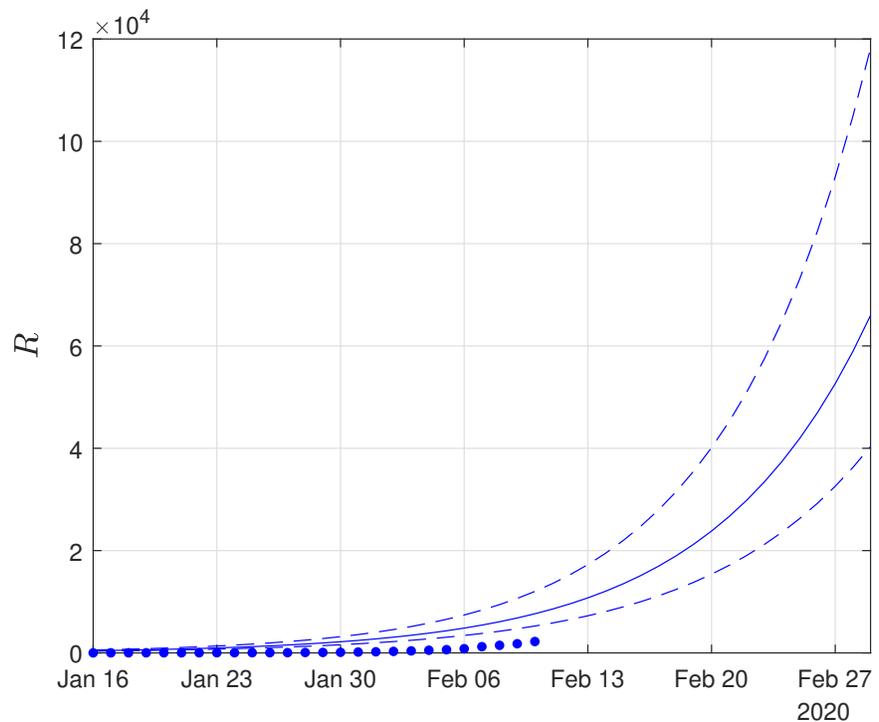


Fig 7. Scenario I. Simulations until the 29th of February of the cumulative number of recovered as obtained using the SIRD model. Dots correspond to the number of confirmed cases from the 16th of January to the 10th of February. The initial date of the simulations was the 16th of November with one infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.206$, $\beta = 0.054$, $\gamma = 0.032$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

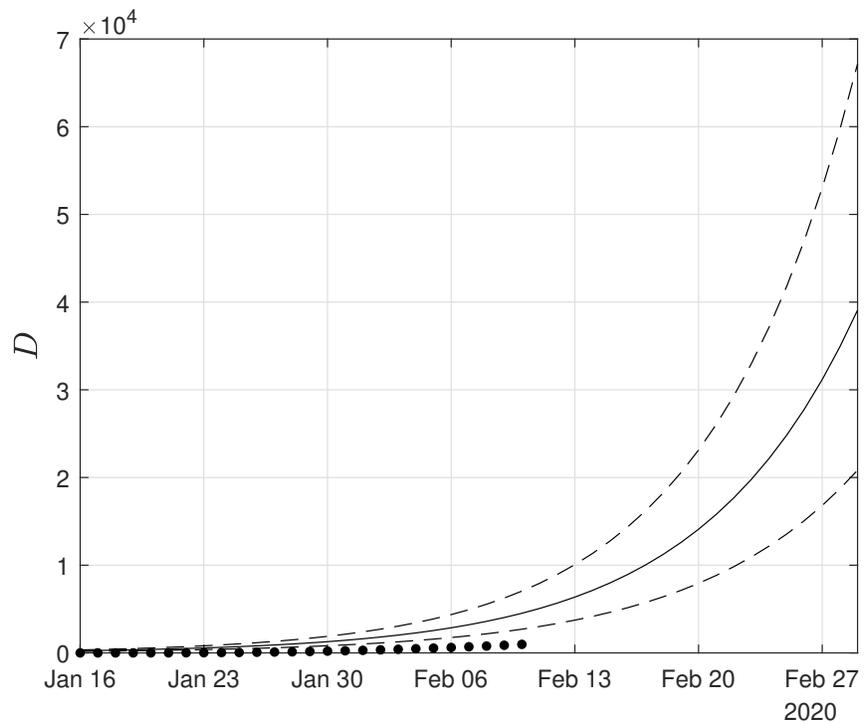


Fig 8. Scenario I. Simulations until the 29th of February of the cumulative number of deaths as obtained using the SIRD model. Dots correspond to the number of confirmed cases from 16th of January to the 10th of February. The initial date of the simulations was the 16th of November with one infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.206$, $\beta = 0.054$, $\gamma = 0.032$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

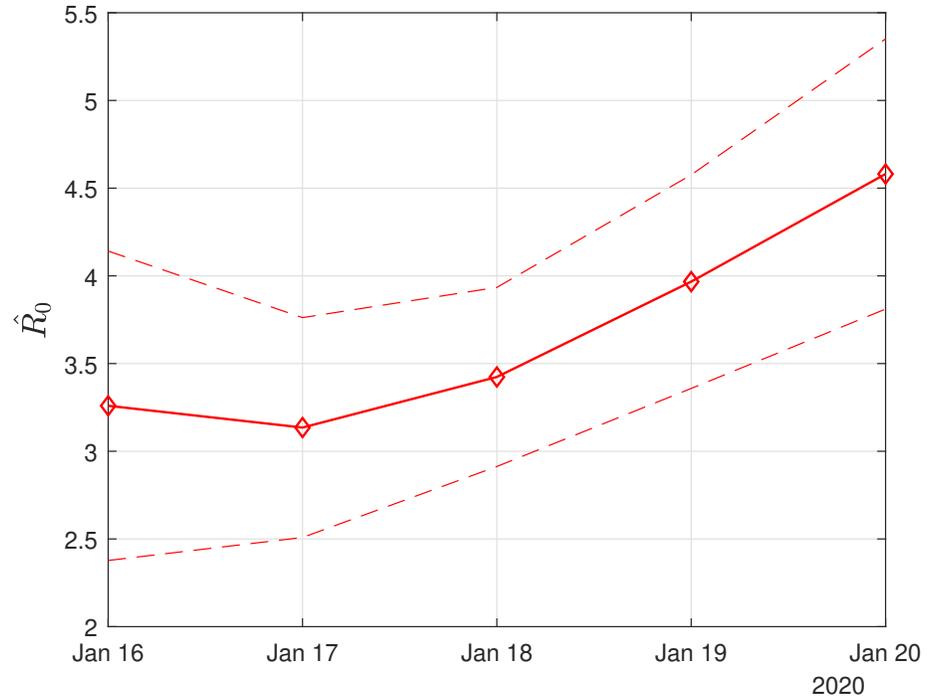


Fig 9. Scenario II. Estimated values of the basic reproduction number (R_0) as computed by least squares using a rolling window with initial date the 11th of January. The solid line corresponds to the mean value and dashed lines to lower and upper 90% confidence intervals.

consequently of the recovered ones too, are considerably larger than reported. Hence, we assessed the dynamics of the outbreak by considering a different scenario that we present in the following subsection.

1.2 Scenario II. Results obtained based by taking twenty times the number of infected and forty times the number of recovered people with respect to the confirmed cases

For our illustrations, we assumed that the number of infected is twenty times the number of the confirmed infected and forty times the number of the confirmed recovered people. Figure9 depicts an estimation of R_0 for the period January 16-January 20. Using the first six days from the 11th of January to 16th of January, \hat{R}_0 results in 3.25 (90% CI: 2.37-4.14); using the data until January 17, \hat{R}_0 results in 3.13 (90% CI: 2.50-3.76); using the data until January 18, \hat{R}_0 results in 3.42 (90% CI: 2.91-3.93); using the data until January 19, \hat{R}_0 results in 3.96 (90% CI: 3.36-4.57) and using the data until January 20, \hat{R}_0 results in 4.58 (90% CI: 3.81-5.35).

It is interesting to note that the above estimation of R_0 is close enough to the one reported in other studies (see in the Introduction for a review).

Figure10 depicts the estimated values of the recovery (β) and mortality (γ) rates for the period January 16 to February 10. The confidence intervals are also depicted with dashed lines. Note that the large variation in the estimated values of β and γ should be accounted to the small size of the data and data uncertainty. This is also reflected in the corresponding confidence intervals. As more data are taken into account, this variation

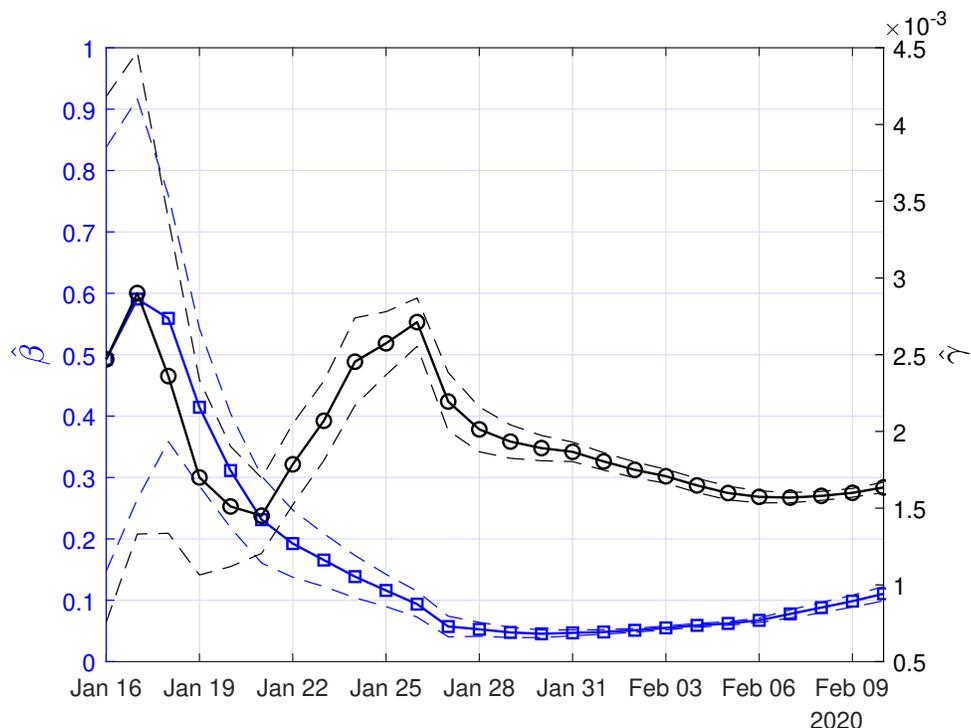


Fig 10. Scenario II. Estimated values of the recovery ($\hat{\beta}$) rate and mortality ($\hat{\gamma}$) rate, as computed by least squares using a rolling window (see section 0.1). Solid lines correspond to the mean values and dashed lines to lower and upper 90% confidence intervals

is significantly reduced. Thus, using all the (scaled) data from the 11th of January until the 10th of February, the estimated value of the mortality rate γ now drops to $\sim 0.163\%$ (90% CI: 0.160%-0.167%) while that of the recovery rate is ~ 0.11 (90% CI: 0.099-0.122) corresponding to ~ 9 days (90% CI: 8-10 days). It is interesting also to note, that as the available data become more, the estimated recovery rate increases slightly (see Fig.10), while the mortality rate seems to be stabilized at a rate of $\sim 0.16\%$.

In Figures 11,12,13, we show the coefficients of determination (R^2) and the root of mean squared errors ($RMSE$), for \hat{R}_0 , $\hat{\beta}$ and $\hat{\gamma}$, respectively.

Again, we used the SIRD simulator to provide estimation of the infection rate by optimization setting $w_1 = 1$, $w_2 = 400$, $w_3 = 1$ to balance the residuals of deaths with the scaled numbers of the infected and recovered cases. Thus, to find the optimal infection run, we used the SIRD simulations with $\beta = 0.11$, and $\gamma = 0.00163$ and as initial conditions one infected, zero recovered, zero deaths on November 16th 2019, and run until the 10th of February. The optimal, with respect to the reported confirmed cases from the 11th of January to the 10th of February value of the infected rate (α) was ~ 0.258 (90% CI: 0.256-0.26). This corresponds to a mean value of the basic reproduction number $\hat{R}_0 \approx 2.31$.

Finally, using the derived values of the parameters α , β , γ , we have run the SIRD simulator until the end of February. The simulation results are given in Figures 14,15,16. Solid lines depict the evolution, when using the expected (mean) estimations and dashed lines illustrate the corresponding lower and upper bounds as computed at the limits of the confidence intervals of the estimated parameters.

Again as Figures 15,16 suggest, the forecast of the outbreak at the end of February,

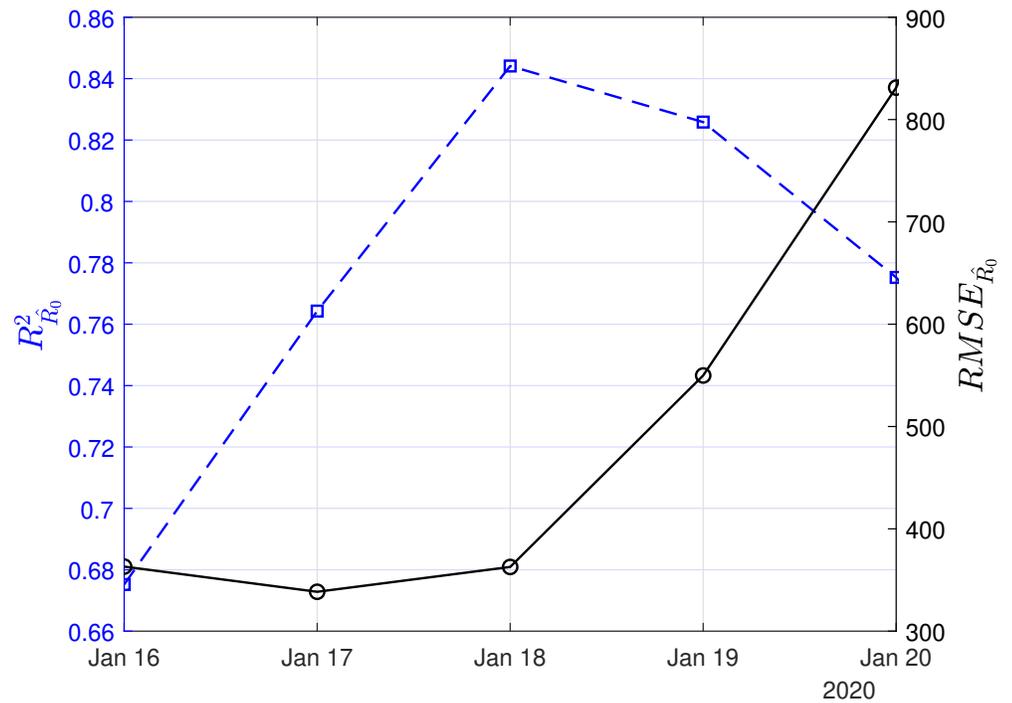


Fig 11. Scenario II. Coefficient of determination (R^2) and root mean square error ($RMSE$) resulting from the solution of the linear regression problem with least-squares for the basic reproduction number (R_0).

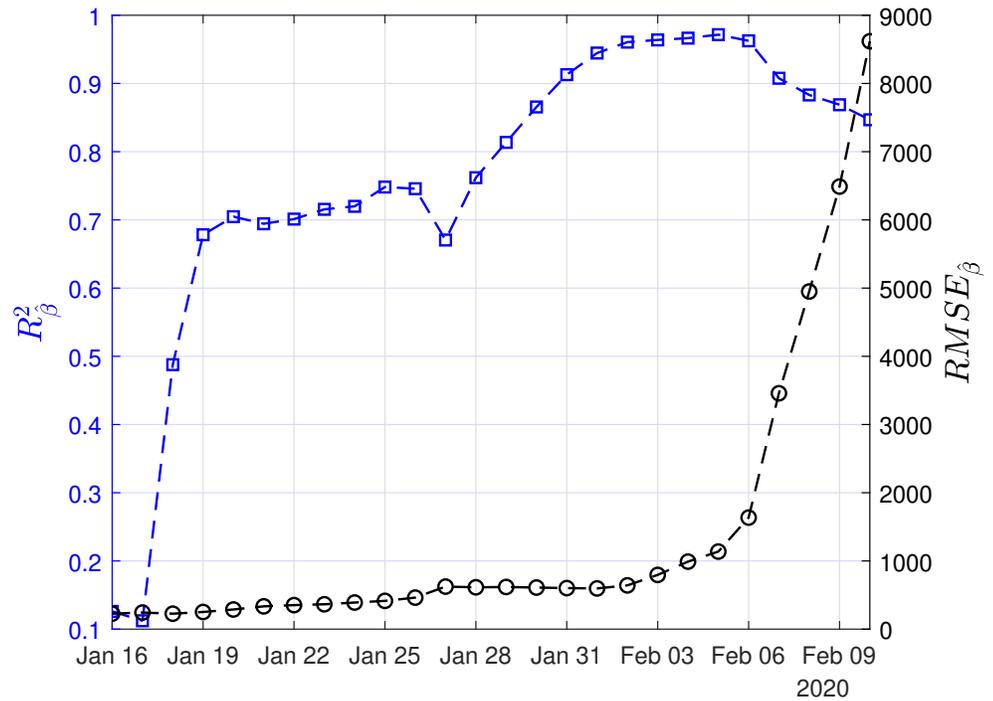


Fig 12. Scenario II. Coefficient of determination (R^2) and root mean square error ($RMSE$) resulting from the solution of the linear regression problem with least-squares for the recovery rate (β).

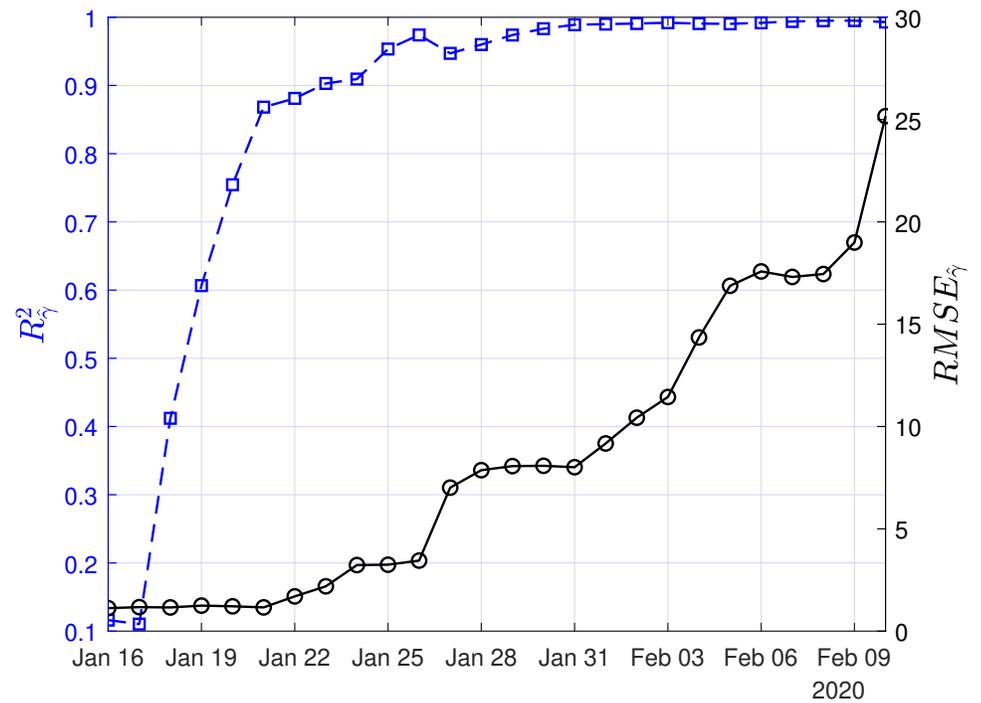


Fig 13. Scenario II. Coefficient of determination (R^2) and root mean square error ($RMSE$) resulting from the solution of the linear regression problem with least-squares for the mortality rate (γ).

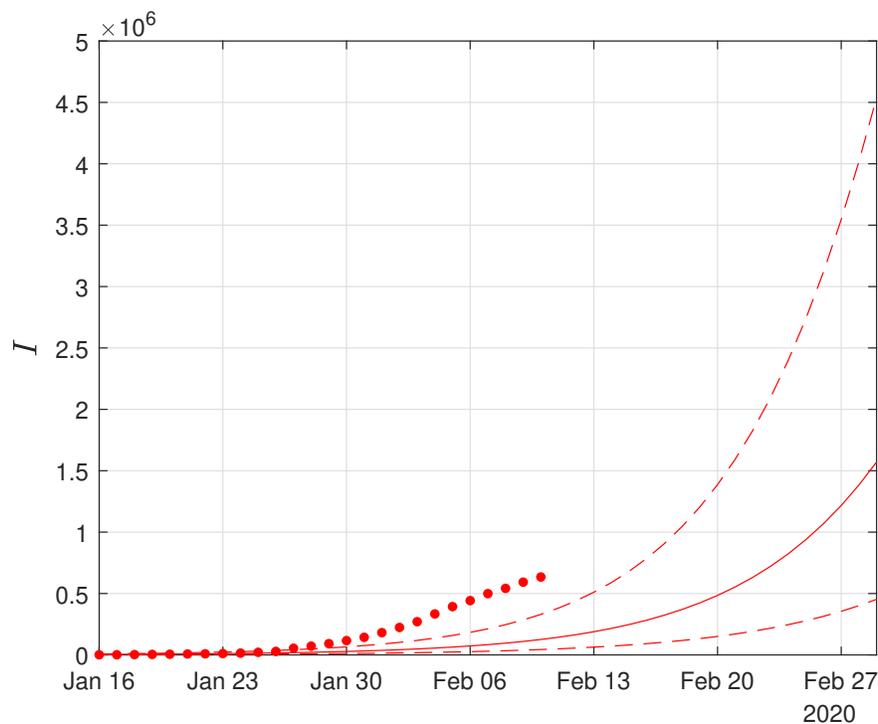


Fig 14. Scenario II. Simulations until the 29th of February of the cumulative number of infected as obtained using the SIRD model. Dots correspond to the number of confirmed cases from 16th of Jan to the 10th of February. The initial date of the simulations was the 16th of November with one infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.258$, $\beta = 0.11$, $\gamma = 0.00163$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

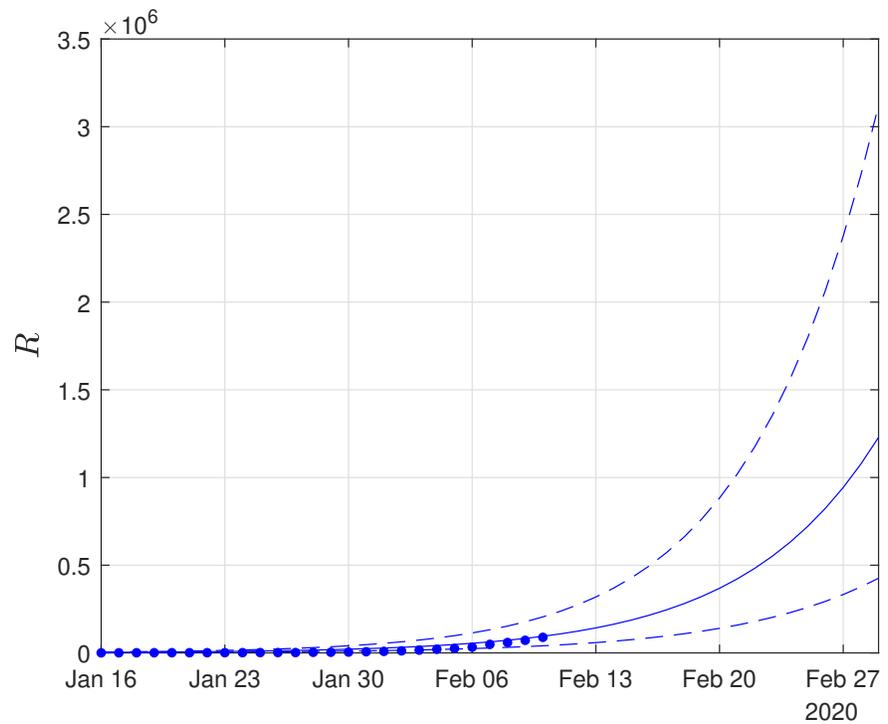


Fig 15. Scenario II. Simulations until the 29th of February of the cumulative number of recovered as obtained using the SIRD model. Dots correspond to the number of confirmed cases from 16th of January to the 10th of February. The initial date of the simulations was the 16th of November, with one infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.258$, $\beta = 0.11$, $\gamma = 0.00163$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

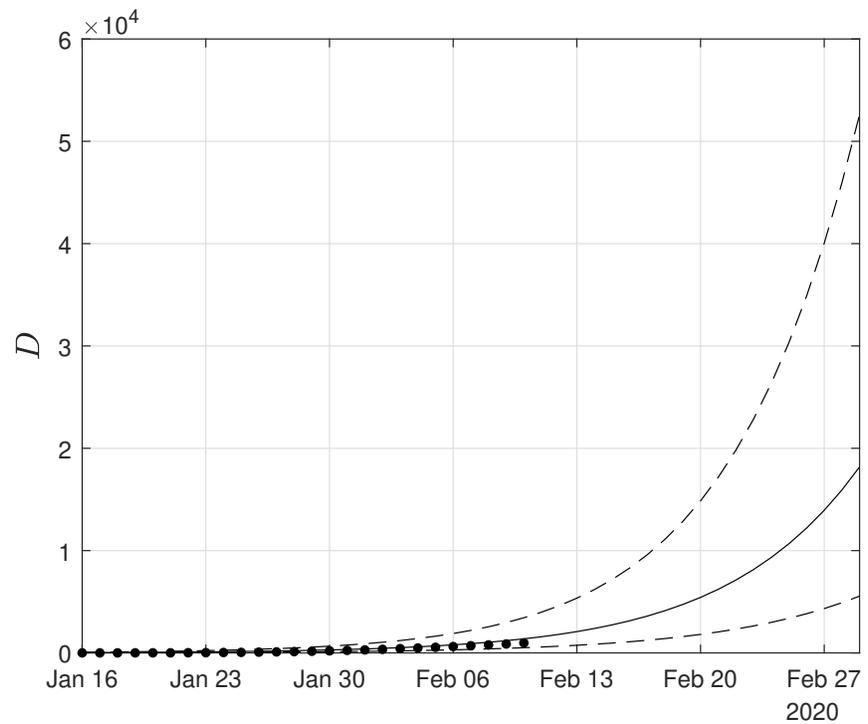


Fig 16. Scenario II. Simulations until the 29th of February of the cumulative number of deaths as obtained using the SIRD model. Dots correspond to the number of confirmed cases from the 16th of November to the 10th of February. The initial date of the simulations was the 16th of November with zero infected, zero recovered and zero deaths. Solid lines correspond to the dynamics obtained using the estimated expected values of the epidemiological parameters $\alpha = 0.258$, $\beta = 0.11$, $\gamma = 0.00163$; dashed lines correspond to the lower and upper bounds derived by performing simulations on the limits of the confidence intervals of the parameters.

through the SIRD model is characterized by high uncertainty. In particular, in Scenario II, by February 29, simulations result in an expected actual number of 1.57m infected cases (corresponding to 78,000 unscaled reported cases) in the total population with a lower bound at 450,000 (corresponding to 23,000 unscaled reported cases) and an upper bound at \sim 4.5m cases (corresponding to 225,000 unscaled reported cases). Similarly, for the recovered population, simulations result in an expected actual number of 1.22m (corresponding to 31,000 unscaled reported cases), while the lower and upper bounds are at 425,000 (corresponding to 11,000 unscaled reported cases) and 3.1m (corresponding to 77,000 unscaled reported cases), respectively. Finally, regarding the deaths, simulations under this scenario result in an average number of 18,000, with lower and upper bounds, 5,500 and 52,000.

Table 1 summarizes the above results for both scenarios.

We note, that the results derived under Scenario II seem to better reflect the actual situation as the reported number of deaths is within the average and lower limits of the SIRD simulations. In particular, as this paper was revised, the reported number of deaths on the 22th February was 2,346, while the lower bound of the forecast is 2,300; This indicates that the mortality rate is 0.16%. Regarding the number of infected and recovered cases by February 20, the cumulative numbers of confirmed reported cases were 64,084 infected and 15,299 recovered. Thus, the corresponding scaled numbers are 1,281,680 infected and 611,960 recovered. Based on Scenario II, for the 22th of February, our simulations give an expected cumulative number of 630,000 infected with 1.8m as an upper bound, and a total cumulative number of 480,000 recovered with a total of 1.2m as an upper bound.

Hence, based on this estimation of the actual numbers of infected and recovered in the total population, the evolution of the epidemic is well within the bounds of our forecasting.

Discussion

We have proposed a methodology for the estimation of the key epidemiological parameters as well as the modelling and forecasting of the spread of the COVID-19 epidemic in Hubei, China by considering publicly available data from the 11th of January 2019 to the 10th of February 2020.

At this point we should note that our SIRD modelling approach did not take into account many factors that play an important role in the dynamics of the disease such as the effect of the incubation period in the transmission dynamics, the heterogeneous contact transmission network, the effect of the measures already taken to combat the epidemic, the characteristics of the population (e.g. the effect of the age, people which had already health problems). Of note, COVID-19, which is thought to be principally transmitted from person to person by respiratory droplets and fomites without excluding the possibility of the fecal-oral route [20] had been spreading for at least over a month and a half before the imposed lockdown and quarantine of Wuhan on January 23, having thus infected unknown numbers of people. The number of asymptomatic and mild cases with subclinical manifestations that probably did not present to hospitals for treatment may be substantial; these cases, which possibly represent the bulk of the COVID-19 infections, remain unrecognized, especially during the influenza season [21]. This highly likely gross under-detection and underreporting of mild or asymptomatic cases inevitably throws severe disease courses calculations and death rates out of context, distorting epidemiologic reality. Another important factor that should be taken into consideration pertains to the diagnostic criteria used to determine infection status and confirm cases. A positive PCR test was required to be considered a confirmed case by China's Novel Coronavirus Pneumonia Diagnosis and Treatment program in the

early phase of the outbreak [1]. However, the sensitivity of nucleic acid testing for this novel viral pathogen may only be 30-50%, thereby often resulting in false negatives, particularly early in the course of illness. To complicate matters further, the guidance changed in the recently-released fourth edition of the program on February 6 to allow for diagnosis based on clinical presentation, but only in Hubei province [1]. The swiftly growing epidemic seems to be overwhelming even for the highly efficient Chinese logistics that did manage to build two new hospitals in record time to treat infected patients. Supportive care with extracorporeal membrane oxygenation (ECMO) in intensive care units (ICUs) is critical for severe respiratory disease. Large-scale capacities for such level of medical care in Hubei province, or elsewhere in the world for that matter, amidst this public health emergency may prove particularly challenging. We hope that the results of our analysis contribute to the elucidation of critical aspects of this outbreak so as to contain the novel coronavirus as soon as possible and mitigate its effects regionally, in mainland China, and internationally.

Conclusion

In the digital and globalized world of today, new data and information on the novel coronavirus and the evolution of the outbreak become available at an unprecedented pace. Still, crucial questions remain unanswered and accurate answers for predicting the dynamics of the outbreak simply cannot be obtained at this stage. We emphatically underline the uncertainty of available official data, particularly pertaining to the true baseline number of infected (cases), that may lead to ambiguous results and inaccurate forecasts by orders of magnitude, as also pointed out by other investigators [3, 16, 21].

References

1. for Health Security JHC. Daily updates on the emerging novel coronavirus from the Johns Hopkins Center for Health Security. February 9, 2020; 2020. Available from: <https://hub.jhu.edu/2020/01/23/coronavirus-outbreak-mapping-tool-649-em1-art1-dtd-health/>.
2. for Health Security JHC. Coronavirus COVID-19 Global Cases by Johns Hopkins CSSE; 2020. Available from: <https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6>.
3. Li Q, Guan X, Wu P, et al. Early Transmission Dynamics in Wuhan, China, of Novel Coronavirus-Infected Pneumonia; 2020. Available from: <https://doi.org/10.1088/2F0951-7715/2F16/2F2/2F308>.
4. Organization WH. WHO Statement Regarding Cluster of Pneumonia Cases in Wuhan, China; 2020. Available from: <https://www.who.int/china/news/detail/09-01-2020-who-statement-regarding-cluster-of-pneumonia-cases-in-wuhan-ch>
5. Organization WH. Novel coronavirus(2019-nCoV). Situation report 21. Geneva, Switzerland: World Health Organization; 2020; 2020. Available from: https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200210-sitrep-21-ncov.pdf?sfvrsn=947679ef_2.
6. Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, et al. Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor binding. *The Lancet*. 2020;doi:10.1016/s0140-6736(20)30251-8.

7. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, et al. A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature*. 2020;doi:10.1038/s41586-020-2012-7.
8. Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, et al. Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *The Lancet*. 2020;doi:10.1016/s0140-6736(20)30211-7.
9. Patel A, Jernigan D, nCoV CDC Response Team. Initial Public Health Response and Interim Clinical Guidance for the 2019 Novel Coronavirus Outbreak - United States, December 31, 2019-February 4, 2020. *MMWR Morb Mortal Wkly Rep*. 2020;doi:10.15585/mmwr.mm6905e.
10. Hunag C, Wang Y, Li X, et al. Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet*. 2020;doi:10.1016/S0140-6736(20)30183-5.
11. Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. Clinical Characteristics of 138 Hospitalized Patients With 2019 Novel Coronavirus-Infected Pneumonia in Wuhan, China. *JAMA*. 2020;doi:10.1001/jama.2020.1585.
12. Zhao S, Lin Q, Ran J, Musa SS, Yang G, Wang W, et al. Preliminary estimation of the basic reproduction number of novel coronavirus (2019-nCoV) in China, from 2019 to 2020: A data-driven analysis in the early phase of the outbreak. *Int J Infect Dis*. 2020;doi:10.1101/2020.01.23.916395.
13. Imai N, Cori A, Dorigatti I, et al. Report 3: Transmissibility of 2019-nCoV. *Int J Infect Dis*. 2019;doi:10.1016/j.ijid.2020.01.050.
14. Wu JT, Leung K, Leung GM. Nowcasting and forecasting the potential domestic and international spread of the 2019-nCoV outbreak originating in Wuhan, China: a modelling study. *The Lancet*. 2020;doi:10.1016/s0140-6736(20)30260-9.
15. Siettos CI, Russo L. Mathematical modeling of infectious disease dynamics. *Virulence*. 2013;4(4):295–306. doi:10.4161/viru.24041.
16. Wu P, Hao X, Lau EHY, Wong JY, Leung KSM, Wu JT, et al. Real-time tentative assessment of the epidemiological characteristics of novel coronavirus infections in Wuhan, China, as at 22 January 2020. *Eurosurveillance*. 2020;25(3). doi:10.2807/1560-7917.es.2020.25.3.2000044.
17. NHC. NHS Press Conference, Feb. 4 2020 - National Health Commission (NHC) of the People's Republic of China; 2020.
18. Ghani AC, Donnelly CA, Cox DR, Griffin JT, Fraser C, Lam TH, et al. Methods for Estimating the Case Fatality Ratio for a Novel, Emerging Infectious Disease. *American Journal of Epidemiology*. 2005;162(5):479–486. doi:10.1093/aje/kwi230.
19. MATLAB R2018b; 2018.
20. Gale J. Coronavirus May Transmit Along Fecal-Oral Route, Xinhua Reports; 2020. Available from: <https://www.bloomberg.com/news/articles/2020-02-02/coronavirus-may-transmit-along-fecal-oral-route-xinhua-reports>.
21. Battagay M, Kuehl R, Tschudin-Sutter S, Hirsch HH, Widmer AF, Neher RA. 2019-novel Coronavirus (2019-nCoV): estimating the case fatality rate – a word of caution. *Swiss Medical Weekly*. 2020;doi:10.4414/smw.2020.20203.

	symbol	parameter	computed values	90% CI
Scenario I Exact numbers for confirmed cases	R_0	Basic reproduction number		
	Based on linear regression	11-16 Jan	4.80	3.36-6.67
		11-17 Jan	4.60	3.56-5.65
	Based on the SIRD simulator	11-18 Jan	5.14	4.25-6.03
		Nov 16-Feb 10	2.4	-
	α	infection rate	0.206	0.204-0.208
	β	recovery rate	0.054	0.049-0.060
		recovery time	20 days	18-22 days
γ	mortality rate	3.2%	3.1%-3.3%	
Forecast to Feb 29		infected	140,000	62,000-300,000
		recovered	65,000	40,000 - 118,000
		deaths	39,000	20,000 - 67,000
Scenario II x20 Infected, x40 recovered of confirmed cases	R_0	Basic reproduction number		
	Based on linear regression	11-16 Jan	3.25	2.37-4.14
		11-17 Jan	3.13	2.50-3.76
	Based on the SIRD simulator	11-18 Jan	3.42	2.91-3.93
		Nov 16-Feb 10	2.31	-
	α	infection rate	0.258	0.256-0.26
	β	recovery rate	0.11	0.099-0.122
		recovery time	9 days	8-10 days
γ	mortality rate	0.163%	0.16%-0.167%	
Forecast to Feb 29		infected	1.57m	450,000-4.5m
		recovered	1.22m	425,000 - 3.1m
		deaths	19,000	5,500 - 52,000

Table 1. Model parameters, their computed values and forecasts for the Hubei province under two scenarios: (I) using the exact values of confirmed cases or (II) using estimations for infected and recovered (twenty and forty times the number of confirmed cases, respectively).