# Bayesian inference in an extended SEIR model with nonparametric disease transmission rate: an application to the Ebola epidemic in Sierra Leone

GIANLUCA FRASSO, PHILIPPE LAMBERT*

*Faculté des sciences sociales, Méthodes quantitatives en sciences sociales, Université de Liège, Liège, Belgium.*
*Institut de statistique, biostatistique et sciences actuarielles, Université catholique de Louvain, Louvain-la-Neuve, Belgium*
Corresponding author: p.lambert@ulg.ac.be

Summary

The 2014 Ebola outbreak in Sierra Leone is analysed using a susceptible-exposed-infectious-removed (SEIR) epidemic compartmental model. The discrete time stochastic model for the epidemic evolution is coupled to a set of ordinary differential equations describing the dynamics of the expected proportions of subjects in each epidemic state. The unknown parameters are estimated in a Bayesian framework by combining data on the number of new (laboratory confirmed) Ebola cases reported by the Ministry of Health and prior distributions for the transition rates elicited using information collected by the WHO during the follow-up of specific Ebola cases. The time-varying disease transmission rate is modeled in a flexible way using penalized B-splines. Our framework represents a valuable stochastic tool for the study of an epidemic dynamic even when only irregularly observed and possibly aggregated data are available. Simulations and the analysis of the 2014 Sierra Leone Ebola data highlight the merits of the proposed methodology. In particular, the flexible modeling of the disease transmission rate makes the estimation of the effective reproduction number robust to the misspecification of the initial epidemic states and to under-reporting of the infectious cases.

*Key words*: Bayesian inference; differential equations; effective reproduction number ; Ebola 2014 epidemic; infectious diseases ; penalized B-splines; SEIR model.

## 1. Introduction

It is believed that the Ebola epidemic in Sierra Leone started at the end of May 2014 (WHO, 2015; WHO Ebola Response Team, 2014) with a healer claiming that she could cure people of an illness that turned to be Ebola. Her patients came from Guinea where the epidemic claimed its first victim in December 2013. At least twelve women from Sierra Leone were reported to be infected at the occasion of her burial and started to spread the disease within the country (Vogel, 2014). The absence of efficient treatment for Ebola and the large fatality rate in past and current outbreaks

make the spread of that infectious disease particularly worrying if it takes place in a country where central authorities do not react promptly by adopting adequate treatment policies (Dowell *and others*, 1999). This was the case in Sierra Leone where the government imposed a criticized 3-day population containment only five months after the first Ebola cases (Ozer *and others*, 2014). Mathematical models are useful to understand and to forecast the evolution of an epidemic, to measure or to simulate the effects of public health interventions and to forecast the future course of the disease in a population. Many approaches have been proposed to describe Ebola epidemics and the 2014 outbreak has stimulated new scientific contributions. Most of them aim to estimate an epidemic compartmental model (see e.g. Anderson and May, 1992; Hens *and others*, 2012) to the observed disease records. The virus transmission rate is usually assumed to decline over time along a known (usually exponential) functional form. The idea is to divide the population under study in disjoint classes (or *compartments*) of subjects according to their epidemic status and to describe the epidemic by modeling the transition mechanism between the epidemic compartments. This is usually achieved by defining suitable (deterministic or stochastic) dynamic mathematical models (e.g. differential or state space equations) involving (unknown) parameters regulating the transitions of subjects among the compartments. However, inference using such epidemic models turns out to be very challenging: one or more of the epidemic states are often not observed and when they are, it often happens in an irregular or unreliable way.

In order to describe the Ebola virus dynamics, a deterministic (based on ordinary differential equations, ODE) susceptible-exposed-infectious-removed (SEIR) framework has been adopted in many contributions. For example, Chowell *and others* (2004) use a SEIR model to study the Congo 1995 and Uganda 2000 Ebola outbreaks and estimate the unknown parameters involved in the ODE system by minimizing the sum of squared differences between the approximated state function and observed numbers of new cases. Rivers *and others* (2014) suggest an extended SEIR model to study the impact of central government interventions on the 2014 Ebola epidemic in Sierra Leone and Liberia. It includes specific compartments for hospitalized patients and for subjects taking part to the funeral of an Ebola victim. Althaus (2014) analyzes the West African 2014 Ebola epidemic by estimating the unknown parameters of a SEIR system under the assumption that the cumulative numbers of cases and deaths are Poisson distributed and considering an exponentially declining (with time) disease transmission rate.

Models with stochastic innovations are sometimes more appropriate than deterministic (ODE-based) ones, in particular when dealing with small populations or at the very beginning of an epidemic with an analysis based on a very small number of infected subjects (see Britton, 2010, for more details). For example, Lekone and Finkenstädt (2006) analyze the Congo 1995 Ebola outbreak by fitting a stochastic SEIR model in a Bayesian framework. In particular, they suggest to model the daily number of new cases and removed with binomial distributions. The disease transmission rate is assumed constant in the first days of the epidemic and exponentially declining after public health intervention. This hypothesis about the form of the transmission rate time appears restrictive but is quite popular in the literature on compartmental epidemic models.

Here, we focus on the 2014 Ebola epidemic in Sierra Leone using an extension of a compartmental SEIR model in Bayesian framework. That paradigm enables an efficient and coherent combination of relevant prior information and data evidence. Our approach is a flexible setting for the estimation of parameters in dynamic models from incomplete or aggregated reports on the number of disease state transitions. Furthermore, it enables to quantify the uncertainty in the estimation of key epidemic quantities such as the effective reproduction number $\mathcal{R}_e(t)$ (see e.g. Heffernan *and others*, 2005) and the disease transmission rate modeled in a flexible way using P-splines (Eilers and Marx, 1996).

Simulations show that this framework makes the estimation of $\mathcal{R}_e(t)$ robust to initial conditions and to underreporting.

The paper is organized as follows. In Section 2, we introduce an extended SEIR (susceptible-exposed-infectious-recovered) model involving nonlinear differential equations to describe the dynamics of the Ebola epidemic and briefly review the concept of effective reproduction number. The likelihood used in parameter estimation is described in Section 3. In Section 4, we introduce our Bayesian framework and the results obtained by analyzing the epidemic data observed in Sierra Leone during the 2014 Ebola outbreak. We conclude the paper with a discussion in Section 5.

## 2. SEIR-D MODEL FOR THE EBOLA EPIDEMIC

When an infection spreads across a given population of size $N$, one can define at a given time $t$ a set of disjoint groups (or *compartments*) of subjects according to their disease status. Here, following a SEIR framework, we distinguish five compartments: *Susceptible* (i.e. healthy subjects at risk to get the disease), *Exposed* (i.e. infected subjects but not yet infectious), *Infectious* (i.e. subjects able to transmit the disease), *Recovered* and *Dead*. Susceptible persons (in proportion $s(t) = S(t)/N$) have contacts with a given number of subjects. Some of these contacts occur with infectious subjects (in proportion $i(t) = I(t)/N$) and lead to new infections and, hence, transitions to the exposed compartment (in proportion $e(t) = E(t)/N$). After some time, exposed subjects develop Ebola symptoms and move to the infectious compartment with the ability to transmit the disease. With time, infectious subjects will finally join the Recovered (in proportion $r(t) = R(t)/N$) or the Dead (in proportion $d(t) = D(t)/N$) groups. In addition, we assume that all cases are symptomatic and that the total population size (including the persons who recovered or died from Ebola), $N$, is fixed over the time range of the study (thereby neglecting deaths from other causes and emigration, as well births or immigration).

### 2.1   *Transitions between states*

In order to model the virus transmission mechanism, we define a discrete time Markov chain framework and associate it to a set of ordinary differential equations which solution(s) describe the expected proportions of subjects in each epidemic compartment at any time $t$ (see e.g. Anderson and May, 1992; Hens *and others*, 2012). We assume that the five epidemic states are homogeneously mixed in the population. This enables us to define a simple model to describe the course of the epidemic in the overall population without taking into account regional specificities. The impact of a violation of this hypothesis is evaluated through simulations (see Section 2.2 in the online Supplementary Materials).

The possible transitions between the five epidemic states are summarized in Fig. 1 and can be described as follows in a time interval $(t, t + \mathrm{d}t)$, where $\mathrm{d}t$ is small enough to ensure that a single person can only experience at most one epidemic state transition:

$\underline{S \longrightarrow E}$ : a susceptible person has, on average, contacts with $\beta(t)\mathrm{d}t$ persons during $(t, t + \mathrm{d}t)$. If $\pi_{\mathcal{E}}(t)$ denotes the time-varying probability that a contact between a susceptible and an infectious subject effectively leads to a new infection, then, using the mass-action principle and defining the *force of infection* as $\psi(t) = \beta(t)\pi_{\mathcal{E}}(t)$, one can show (see Section 1.2 in the online Supplementary Materials) that the expected number of transitions $(\mathrm{d}E^+(t) = S(t) - S(t + \mathrm{d}t))$ from the susceptible

to the exposed state in $(t, t + \mathrm{d}t)$ is

$$\mathbb{E}\left(\mathrm{d}E^+(t)|S(t), I(t), \psi(t)\right) = S(t)\frac{I(t)}{N}\psi(t)\,\mathrm{d}t + o(\mathrm{d}t). \tag{2.1}$$

$\underline{E \longrightarrow I}$ : denote by $\sigma$ (in days$^{-1}$) the transition rate of a person in the exposed state to the infectious one. Then, the expected number of transitions $E \to I$ during $(t, t + \mathrm{d}t)$ is

$$\mathbb{E}\left(\mathrm{d}I^+(t)|E(t), \sigma\right) = \sum_{i=1}^{E(t)} \mathbb{E}\left(\mathrm{d}I_i^+(t)|\sigma\right) = \sigma E(t)\,\mathrm{d}t + o(\mathrm{d}t), \tag{2.2}$$

where $\mathrm{d}I_i^+(t)$ is 1 if the $i$th exposed person at $t$ becomes infectious during $(t, t+\mathrm{d}t)$, and 0 otherwise.

$\underline{I \longrightarrow R, D}$ : when turning infectious, a person has probability $\pi_d$ to die from Ebola. At each time $t$, we can distinguish two groups of infectious subjects: $I_d(t)$ of them will die at rate $\gamma_d$ (i.e. on average $1/\gamma_d$ units of time after the appearance of the first symptoms) and the other $I_r(t)$ $(= I(t) - I_d(t))$ will recover (and become immune) at a slower rate $\gamma_r$. One can show (see Section 1.2 in the online Supplementary Materials) that the expected numbers of new recoveries and deaths in $(t, t + \mathrm{d}t)$ are

$$\begin{array}{rcl}
\mathbb{E}\left(\mathrm{d}R^+(t)|I(t), \pi_d, \gamma_r\right) & = & \gamma_r(1 - \pi_d)I(t)\,\mathrm{d}t + o(\mathrm{d}t), \\
\mathbb{E}\left(\mathrm{d}D^+(t)|I(t), \pi_d, \gamma_d\right) & = & \gamma_d\pi_d I(t)\,\mathrm{d}t + o(\mathrm{d}t).
\end{array} \tag{2.3}$$

Finally, note that the distinction between the two categories of infectious ($I_d$ and $I_r$) enables us to use the available prior information about recovery and death rates (see Section 4).

## 2.2   Stochastic and ODE models for the Ebola epidemic

Consider a time interval $(t, t + \mathrm{d}t)$ with $\mathrm{d}t$ small. Using basic theory on Poisson processes (see Section 1.1 in the online Supplementary Materials), the conditional expectations in Equations (2.1) to (2.3) and the transitions between states described above, one concludes that

$$\begin{array}{rcl}
S(t + \mathrm{d}t) & = & S(t) - \mathrm{d}E^+(t) \quad \text{with} \quad (\mathrm{d}E^+(t)|S(t), I(t), \psi(t)) \sim \mathrm{Pois}\left(S(t)\frac{I(t)}{N}\psi(t)\,\mathrm{d}t\right) \\
E(t + \mathrm{d}t) & = & E(t) + \mathrm{d}E^+(t) - \mathrm{d}I^+(t) \quad \text{with} \quad (\mathrm{d}I^+(t)|E(t), \sigma) \sim \mathrm{Pois}\left(\sigma E(t)\,\mathrm{d}t\right) \\
I(t + \mathrm{d}t) & = & I(t) + \mathrm{d}I^+(t) - \left(\mathrm{d}D^+(t) + \mathrm{d}R^+(t)\right) \\
R(t + \mathrm{d}t) & = & R(t) + \mathrm{d}R^+(t) \quad \text{with} \quad (\mathrm{d}R^+(t)|I(t), \pi_d, \gamma_r) \sim \mathrm{Pois}\left(\gamma_r(1 - \pi_d)I(t)\,\mathrm{d}t\right) \\
D(t + \mathrm{d}t) & = & D(t) + \mathrm{d}D^+(t) \quad \text{with} \quad (\mathrm{d}D^+(t)|I(t), \pi_d, \gamma_d) \sim \mathrm{Pois}\left(\gamma_d\pi_d I(t)\,\mathrm{d}t\right).
\end{array} \tag{2.4}$$

This is a discrete time Markov chain (DTMC) model when $\mathrm{d}t$ is a fixed unit of time (Allen, 2010).

When the total population $N$ is large, the proportions of persons in each of the (sub-)states at time $t$, $\boldsymbol{p}(t) = (s(t), e(t), i_r(t), i_d(t), r(t), d(t))$, can be approximated by time-dependent continuous functions. We describe the change in these proportions by a set of ordinary differential equations (with given initial value conditions) derived from (2.4). Indeed, by rewriting e.g. the 1st equation, dividing by $N$, taking conditional expectations when $\mathrm{d}t \to 0^+$ and using (2.1), one gets

$$\begin{array}{rcl}
\displaystyle\lim_{\mathrm{d}t \to 0^+} \mathbb{E}\left(\frac{s(t + \mathrm{d}t) - s(t)}{\mathrm{d}t}\Big| S(t), I(t), \psi(t)\right) & = & -\displaystyle\lim_{\mathrm{d}t \to 0^+} \frac{1}{N\,\mathrm{d}t}\mathbb{E}\left(\mathrm{d}E^+(t)\Big| S(t), I(t), \psi(t)\right) \\
& = & -\psi(t)\left(i_r(t) + i_d(t)\right)s(t).
\end{array}$$

Repeating similar operations on the other four equations in (2.4) suggests the following set of ODE:

$$
\begin{array}{rcl}
s'(t) & = & -\psi(t)\,i(t)\,s(t) \;\; ; \;\; e'(t) = \psi(t)\,i(t)\,s(t) - \sigma\,e(t) \\
i_r'(t) & = & \sigma\,(1 - \pi_d)\,e(t) - \gamma_r\,i_r(t) \;\; ; \;\; i_d'(t) = \sigma\,\pi_d\,e(t) - \gamma_d\,i_d(t) \\
r'(t) & = & \gamma_r\,i_r(t) \;\; ; \;\; d'(t) = \gamma_d\,i_d(t)
\end{array}
\tag{2.5}
$$

Given initial conditions $\boldsymbol{p}(0) = \boldsymbol{p}_0$ and the ODE parameters $\boldsymbol{\theta} = (\psi(t), \sigma, \pi_d, \gamma_r, \gamma_d)$, the solution $\boldsymbol{p}(t)$ is a deterministic vectorial function of $t$ providing the expected proportions of subjects in each of the states at time $t$. Discrete (DTMC) or continuous time Markov chain (CTMC) models such as in (2.4) are possible stochastic alternatives (Allen, 2010), but the proposed deterministic ODE model for $\boldsymbol{p}(t)$ combined with a stochastic component for the number of newly observed Ebola cases and a flexible specification of $\psi(t)$ appropriately describes the epidemic propagation (see Section 4).

### 2.3    *Possible model extensions and nonparametric time-varying disease transmission rate*

Our ODE model can be extended or simplified in many ways (see e.g. Section 1.3 in the online Supplementary Materials). Given that the monitoring data on the number of new transitions between states are only reasonably reliable for the entries in the $I$ state, we decided to focus our efforts on a flexible specification of the disease transmission rate. According to (2.5), the transition rate between states $S$ and $E$ is modeled to vary with time not only because of its dependence on the proportion $i(t)$ of infectious persons in the population, but also through $\psi(t)$. We do not specify a restrictive parametric form for $\psi(t)$, but, instead, assume that it changes in a smooth way over time with its logarithm modeled using a linear combination of basis functions $\{b_k(t) : k = 1, \ldots, K\}$,

$$
\log \psi(t) = \sum_{k=1}^{K} b_k(t)\alpha_k.
\tag{2.6}
$$

Here, we shall take a large number of cubic B-splines (see e.g. Dierckx, 1995) associated to equidistant knots over the range of observation times, and penalize for change in the successive $\alpha_k$'s during the parameter estimation procedure (Eilers and Marx, 1996), see Fig. 2 for an illustration. That P-spline framework is particularly convenient for our estimation purposes since the number of B-splines in the basis (and their degree) hardly affects the shape of the final estimates which is mainly regulated by the strength of the (roughness) penalty (Eilers and Marx, 2010).

### 2.4    *The effective reproduction number*

One key quantity measuring the potential spread of an epidemic is the *effective reproduction number* $\mathcal{R}_e(t)$, see e.g. Heffernan *and others* (2005) for a review paper on that concept and its computation. It indicates the expected number of secondary cases caused by an infected individual during the course of the disease. A $\mathcal{R}_e(t)$ less than one suggests that the disease will die out, whereas it is in an epidemic state otherwise. The *next generation method* was proposed by Diekmann *and others* (1990) to compute $\mathcal{R}_e(t)$ for mathematical models based on differential equations. By using it in combination with (2.5), one can show (see Section 1.4 in the Online Supplementary Materials) that

$$
\mathcal{R}_e(t) = \psi(t)s(t) \times \left( \frac{1 - \pi_d}{\gamma_r} + \frac{\pi_d}{\gamma_d} \right).
\tag{2.7}
$$

Intuitively, the first factor corresponds to the expected number of new infections caused by an infectious subject during one unit of time (e.g. one day) at time $t$, while the second one is the expected time spent by that person in the infectious state.

While the sequence of $\mathcal{R}_e(t)$ values provides a first summary of the epidemic evolution, parameter estimates and the posterior distributions can also be used together with model (2.4) for prediction purposes. However, caution is needed when making forecasts, in particular with the flexible form assumed for $\psi(t)$. Indeed, its extrapolated behavior crucially depends on the order of the roughness penalty (Eilers and Marx, 2010) in the P-spline model. For this reason, we recommend to restrict the focus on short term predictions or simply to examine $\mathcal{R}_e(t)$ to anticipate future qualitative evolution of the number of cases.

## 3. LIKELIHOOD, OVERDISPERSION AND UNDER-REPORTING

Assume that the epidemic in the area of interest started during an identified day with a single infected and non infectious person ($E(t_0) = 1$) at $t_0$. Then, one has $S(t_0) = N - 1$ and $I_r(t_0) = I_d(t_0) = R(t_0) = D(t_0) = 0$ (with $N$ constant by hypothesis). The corresponding epidemic state proportions at $t_0$ are directly obtained by dividing the preceding quantities by $N$. These could be taken as starting values $\boldsymbol{p}_0$ for the state probabilities $\boldsymbol{p}(t)$ involved in the dynamic system described by (2.5). Further assume that reliable reports on the number of new (infectious) cases are available on a subset of non necessarily consecutive days, $\mathcal{I}^+ = \{\mathrm{d}I^+(t_\ell) : \ell = 1, \ldots, L\}$ with $\mathrm{d}t > 0$ and $t_{\ell_2} > t_{\ell_1}$ whenever $\ell_2 > \ell_1$, and where $\mathrm{d}I^+(t)$ refers to the number of transitions between the $E$ and the $I$ compartments during $(t, t + \mathrm{d}t)$. The modeled state proportions $\boldsymbol{p}(t)$ can be calculated at any time $t$ by solving the set of nonlinear differential equations (2.5) for given values of $\boldsymbol{\theta}$ and the preceding initial conditions $\boldsymbol{p}_0$. This can be done numerically using, for example, a Runge-Kutta scheme. Then, the proportion $i(t)$ of infectious subjects at time $t$ is the sum of $i_r(t)$ and $i_d(t)$.

As discussed above, at time $t$, conditionally on the total number $E(t)$ of exposed subjects, the expected number of new (infectious) cases occurring during $(t, t + \mathrm{d}t)$ is given by (2.2). The DTMC model implicitly assumes that an individual leaving state $E$ during $(t, t + \mathrm{d}t)$ only effectively becomes infectious at $t + \mathrm{d}t$. The potential impact of such an approximation is limited when taking $\mathrm{d}t = 1$ *day* with the Ebola epidemic. Since $E(t)$ is not observed, we further approximate it by its expected value $Ne(t)$ (given $\boldsymbol{\theta}$ and $\boldsymbol{p}_0$) as obtained from the (deterministic) solution of the ODE system, yielding $\mathbb{E}\big(\mathrm{d}I^+(t)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) = \mathbb{E}\big(\mathbb{E}(\mathrm{d}I^+(t)|E(t), \sigma)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) = \sigma Ne(t)\,\mathrm{d}t + o(\mathrm{d}t)$. The likelihood contribution for the number of new cases over a single day, $\mathrm{d}I^+(t_\ell)$, is obtained using

$$\big(\mathrm{d}I^+(t_\ell)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) \sim \mathrm{Pois}\big(Ne(t_\ell)\sigma\,\mathrm{d}t\big), \tag{3.1}$$

as can be shown by applying the general result on Poisson processes in the online Supplementary Materials (Section 1.1) with $M(t, t + \mathrm{d}t) = \mathrm{d}I^+(t)$, $\rho = Ne(t)\sigma$ and $\mathrm{d}t$ sufficiently small to handle $\rho$ as if it were constant in $(t, t + \mathrm{d}t)$. We suggest to use the following log-likelihood in inferential procedures,

$$\ell\big(\boldsymbol{\theta}; \mathcal{I}^+, \boldsymbol{p}_0)\big) = \sum_{\ell=1}^{L} \left\{ -N\sigma e\big(t_\ell|\boldsymbol{\theta}, \boldsymbol{p}_0)\big) + \mathrm{d}I^+(t_\ell) \log\big(\sigma e(t_\ell|\boldsymbol{\theta}, \boldsymbol{p}_0)\big) \right\},$$

where we assume for simplicity that the elements in $\mathcal{I}^+$ are conditionally independent although the number of new infectious during $(t, t + \mathrm{d}t)$ indirectly affects future values of the disease transmission rate and, hence, future numbers of new cases. Simulations (in the online Supplementary Materials)

show that this hypothesis has a negligible impact on the reliability of the estimation procedure. When the number of new cases is reported in an aggregated way and refers to a $d$-day time span, $(t, t + d\,\mathrm{d}t)$, the likelihood contribution is obtained using similar arguments with

$$
\begin{aligned}
\mathbb{E}\big(\mathrm{d}I^{+d}(t)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) &= \mathbb{E}\big(\mathrm{d}I^+(t) + \mathrm{d}I^+(t + \mathrm{d}t) + \ldots + \mathrm{d}I^+(t + (d-1)\,\mathrm{d}t)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) \\
&= \big(\sigma Ne(t) + \ldots + \sigma Ne(t + d - 1)\big)\mathrm{d}t + o(\mathrm{d}t).
\end{aligned}
$$

One might also have to deal with overdispersion, a common feature in count data (Breslow, 1984). It occurs when the conditional variance of the count response tends to be larger than its conditional mean, usually as a consequence of missing unidentified or non measured covariates. Then, the homogeneous mixing of the population assumed in the preceding models does not hold: more flexible count data distributions are required to provide a reasonable description of the dispersion of the number of persons entering daily in the infectious state. For that reason, we suggest to replace the Poisson distribution in (3.1) by a negative binomial with identical conditional mean but larger variance, $\big(\mathrm{d}I^+(t)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) \sim \mathrm{NB}\big(Ne(t)\sigma\,\mathrm{d}t, \phi_I\big)$, where $\mathrm{NB}(\mu, \phi)$ denotes a negative binomial distribution with mean $\mu$ and variance $\mu(1 + \mu/\phi)$. This is equivalent to a Poisson sampling process with a gamma distributed rate with mean $\mu$ and shape parameter $\phi$. That parameter quantifies overdispersion with the negative binomial reducing to the Poisson when $\phi \to +\infty$. Finally, a negative binomial distribution can also model the number of new recoveries and deaths in situations where reliable reports on these quantities are also available.

Under-reporting is also expected as some Ebola cases may remain unnoticed or not have been confirmed by a laboratory analysis. One can account for that by replacing the mean in the preceding negative binomial by a fraction of them, $\big(\mathrm{d}I^+(t)|\boldsymbol{\theta}, \boldsymbol{p}_0\big) \sim \mathrm{NB}\big(\rho_I(t)Ne(t)\sigma\,\mathrm{d}t, \phi_I\big)$, where the $\rho_I(\cdot)$ function takes values in $(0, 1)$. Of course, some model parameters might not be identifiable. Therefore, some of them should be fixed arbitrarily (e.g. $\rho_I(t)$) or estimated by combining informative priors (such as on $\sigma$, $\gamma_r$ and $\gamma_d$) from historical studies. The impact of ignoring under-reporting on parameter estimation was assessed using simulations in Section 2.5 of the online Supplementary materials. A synthesis of our conclusions can be found at the end of Section 4.

## 4. Application on the Ebola data for Sierra Leone

Observations are partially available since the start of the epidemic. Reports by the Ministry of Health of Sierra Leone provide, since 18 July 2014, daily summaries on the number of cases, recovered and dead people due to Ebola (http://health.gov.sl/). Sparse and non detailed information is available before that date and up to the start of the epidemic (end of May 2014). We focus our interest on the reported number of new entries (confirmed by a laboratory analysis) for the infectious, recovered or dead persons from Ebola (see Fig 3). That restriction follows from the similarity of Ebola symptoms with other endemic diseases like malaria, typhoid fever and meningitis. Therefore, the selected data just report on a (likely growing) portion of the epidemic. A non negligible and unknown proportion of persons suffering or having suffered from Ebola are not reported, and even when they are, may not have been confirmed by a laboratory analysis.

*Model and likelihood* : The SEIR-D model described in Section 2 was fitted to the number of new Ebola cases available at different calendar times of the epidemic, $\mathcal{I}^+ = \{dI^+(t_\ell) : \ell = 1, \ldots, L\}$. Data were aggregated on a weekly basis (from Sunday to Saturday) starting on August 10 to account for the irregularities in the administrative report of the laboratory analyzes (see upper panel of Fig. 3). This mitigates the impact of outlying reported numbers of cases probably recorded

with some delay on specific days due to administrative reasons and the time required to obtain and to collect lab results. The (approximated) likelihood contribution for $\mathrm{d}I^+(t_\ell)$ was obtained along the lines of Section 3. When a report follows several days of silence or results from a weekly aggregation, $\mathrm{d}I^+(t_\ell)$ is interpreted as the cumulative number of new cases over $(t_\ell - t_{\ell-1})$ days. Then the likelihood contribution is obtained from $\left(\mathrm{d}I^+(t_\ell)|\boldsymbol{\theta}, \boldsymbol{p}_0\right) \sim \mathrm{Pois}\left(\sigma N\left(e(t_{\ell-1}) + \ldots + e(t_\ell - 1)\right)\right)$ and its extension accounting for possible overdispersion.

*Priors* : Some of the model parameters referring to the disease course cannot be estimated from the number of new cases only. Therefore, priors (in Table 1) were elicited for these quantities using information gathered during the follow-up of a subset of Ebola patients in Sierra Leone from June 08 till September 14, 2014. The information reported in WHO Ebola Response Team, 2014 cover all the states defined in the ODE system (2.5). A uniform prior on $(0, 100)$ was taken for the overdispersion parameter $\phi_I$ in the negative binomial distribution for the observed counts (see e.g. Jackman, 2009). Alternatively, one can use a uniform prior on $(0, 1)$ for $1/\phi_I$ (see e.g. Gelman and Hill, 2007). A prior penalizing changes in successive spline parameters was assumed for $\boldsymbol{\alpha}$ in (2.6), following the Bayesian translation (Jullion and Lambert, 2007) of the frequentist P-spline approach in Eilers and Marx (1996):

$$p(\boldsymbol{\alpha}|\lambda) \propto \exp\left(-\frac{\lambda}{2}\sum_k (D\boldsymbol{\alpha})_k^2\right) = \exp\left(-\frac{\lambda}{2}\boldsymbol{\alpha}' D^\top D \boldsymbol{\alpha}\right) \;\; ; \;\; (\lambda|\xi) \;\; \sim \mathcal{G}\left(\frac{\nu}{2}, \frac{\nu\xi}{2}\right) \;\; ; \;\; \xi \sim \mathcal{G}\left(\epsilon, \epsilon\right),$$

where $D$ denotes the 3rd order difference matrix and $\lambda$ the roughness penalty parameter with a continuous mixture of gammas (with $\epsilon = 10^{-4}$, say) as prior. Notice that the (B-spline based) P-spline framework has valuable advantages over alternative flexible (penalized or not) regression approaches. The combination of a basis of cubic B-splines associated to equidistant knots and difference penalties on the spline coefficients ensures a flexible and smooth behavior of $\psi(t)$ and of the reproduction number between the knot locations. On the other hand, as discussed in Eilers and Marx (2010), the choice of the number $K$ of B-splines is not crucial in regulating the shape of the final estimates provided that it is large enough. Indeed, the smoothness of the final estimate is, then, regulated by the penalty parameter $\lambda$ and does not depend on $K$.

*Posterior and results* : The joint posterior is obtained using Bayes' theorem by multiplying the approximated likelihood (see Section 3) and the preceding priors. A random-walk Metropolis-within-Gibbs algorithm is used to sample the posterior. Metropolis steps based on multivariate normal proposals with adaptive steps (Haario *and others*, 2001) are made for each of the three blocks of parameters $\boldsymbol{\alpha}$, $(\sigma^{-1}, \pi_d, \gamma_r^{-1}, \gamma_d^{-1})$ and $\phi_I$, while Gibbs steps are used for the penalty parameters by sampling from the following conditional posteriors:

$$(\lambda|\xi, \boldsymbol{\alpha}, \mathcal{I}^+) \sim \mathcal{G}\left(\frac{\nu + \rho(P)}{2}, \frac{\nu\xi + \boldsymbol{\alpha}' P \boldsymbol{\alpha}}{2}\right) \;\; \text{where} \;\; P = D^\top D \;\; ; \;\; (\xi|\lambda, \mathcal{I}^+) \sim \mathcal{G}\left(\epsilon + \frac{\nu}{2}, \; \epsilon + \frac{\nu\lambda}{2}\right).$$

The sampling algorithm is written in R (R Core Team, 2013) by exploiting parallel computing (we run several independent chains, one per virtual core of the computer chipset). The `deSolve` package (Soetaert *and others*, 2010) is used to solve the set of differential equations governing $\left(s(t), e(t), i_d(t), i_r(t), d(t), r(t)\right)$ for a given value of $\boldsymbol{\theta}$ with, as initial values for the state proportions, $s(0) = 1 - 13/N$, $e(0) = 12/N$, $i_r(0) = 0$, $i_d(0) = 1/N$, $r(0) = 0$, $d(0) = 0$. The first Ebola case was reported on 26 May 2014 at $t = t_0 = 0$ (see Section 1) in a population of approximately $N = 6.2$ million inhabitants. Thanks to the flexible specification of $\psi(t)$, our results were found to

be robust to the choice of the initial number of infected people, $E(t = 0)$ (see Section 2.3 of the online Supplementary Materials).

A chain of length $500k$ with a $10k$ burn-in was considered to explore the joint posterior. Summary statistics on the marginal posterior for $\phi_I$, $\sigma$, $\pi_d$, $\gamma_d$ and $\gamma_r$ can be found in Table 2, while histograms of the sampled values are displayed in the online Supplementary Materials (Section 3). Not surprisingly, the information on $\pi_d$, $\gamma_d$ and $\gamma_r$ mainly comes from their respective priors (since it is not possible to infer on these parameters by using the number of new infectious only). The marginal posterior for $\sigma^{-1}$ is slightly left-skewed but has a mean and a variance similar to its prior values. The estimation of the overdispersion $\phi_I$ (with a uniform prior on $(0,100)$) in the negative binomial distribution used for the response suggests that there is no clear indication of overdispersion and that the Poisson assumption is a legitimate simplification. The effective reproduction number, see (2.7), can also be computed at each MCMC iteration. Then, the posterior sample can be used to estimate the posterior mean of $\mathcal{R}_e(t)$ and to compute pointwise 95% credible intervals for it, see the upper right panel of Fig. 4.

Our results suggest that, at the end of December 2014, Sierra-Leone was most likely still in an epidemic state, although there were strong indications that the disease propagation steadily stepped back since the end of September 2014 after an increase of about two months (starting mid-July 2014) corresponding to the propagation of the epidemic Westwards to the densely populated urban districts. Note that the reproduction number at the first day of the epidemic (also named the *basic reproduction number*) is estimated to be much larger than reported in the literature (see e.g. Althaus, 2014). This is most likely due to an under-evaluation of the number of persons truly exposed to Ebola (here assumed to be $E(0) = 12$) at the start of the epidemic (WHO, 2015).

The fitted number of new cases, see the upper left panel of Fig. 4, shows that the SEIR-D model does an excellent job in describing the dynamics of the epidemic. The involved parameters have a meaningful interpretation and the model enables qualitative and quantitative understanding of the epidemic propagation. The fitted number of deaths and recoveries can also be visualized and compared to the less reliable official reports of these quantities, see Fig. 4 (lower panels). It appears that the dynamics in the officially reported numbers of deaths and recoveries are totally incompatible with the observed numbers of (confirmed) Ebola cases: deaths and recoveries appear strongly under-reported given the large fatality rate among symptomatic Ebola cases. Furthermore, the number of Ebola cases confirmed by a laboratory analysis likely underestimates the true number of cases, and, thus, the under-reporting of deaths and recoveries is probably even more severe than what the solid lines in the lower panels of Fig. 4 suggest.

The impact of an under-reported number of cases on the estimation of the ODE parameters and of the dynamics in the reproduction number has been evaluated through simulations (see Section 2.5 of the online Supplementary Materials). It is shown that, in the presence of under-reporting of infectious cases, our framework ensures an accurate description of the dynamic in the disease transmission rate and in the reproduction number $\mathcal{R}_e(t)$ even when the proportion of under-reported infectious subjects, $\rho_I(t)$, is wrongly assumed equal to 1. This is made possible by an automatic over-estimation of the expected time spent by an Ebola subject in the infectious state.

## 5. Discussion

In this paper, we propose a stochastic framework for the analysis of the 2014 Ebola epidemic in Sierra Leone. We model the virus transmission dynamic by subdividing the subjects belonging to the population under study into five disjoint classes: susceptible, exposed, infectious, dead and recov-

ered. The transition of individuals between states is described by a system of differential equations obtained by extending a SEIR compartmental model (see e.g. Anderson and May, 1992, Chap. 6). The parameters governing the transition between compartments are estimated via Bayesian inferential techniques. This makes possible an efficient and coherent combination of prior information (WHO Ebola Response Team, 2014) and data evidence on the epidemic course.

In Section 4, we focused on the analysis of the number of new confirmed Ebola cases reported by the Ministry of Health and Sanitation of Sierra Leone between 26 May 2014 and 1 January 2015. Our approach efficiently deals with coarse and irregular data and can take overdispersion into account. Furthermore, the proposed P-spline (Eilers and Marx, 1996; Jullion and Lambert, 2007) definition of the disease transmission rate delivers a smooth and convincing fit to the number of new Ebola cases. There is strong evidence that the official numbers of deaths and recoveries in Ebola subjects were largely under-reported and, given the large fatality rate among symptomatic cases, incoherent with the dynamic in the number of confirmed Ebola cases.

The proposed Bayesian framework enables to estimate the transition parameters and the effective reproduction number $\mathcal{R}_e(t)$ with their uncertainties. The study of $\mathcal{R}_e(t)$ is of crucial interest to evaluate the evolution of the epidemic. On the basis of the analysis in Section 4, we conclude that the rate of the disease transmission started to decrease end of September 2014, with $\mathcal{R}_e(t)$ approaching 1.0 by the end of December 2014.

A sensitivity analysis showed that, thanks to the flexible specification of $\psi(t)$ in (2.6), the transition parameter estimates are robust to the specification of the initial value conditions $E(t = 0)$ (see Section 2.3 of the online Supplementary Materials). On the other hand, the simulation studies in Sections 2.2 and 2.5 of the online Supplementary Materials highlight the robustness of the estimation of the effective reproduction number to the violation of the homogeneous mixing assumption (cf. Section 2) and to under-reporting of the number of new cases.

We also tested the appropriateness of our framework through a large simulation study (with results summarized in Section 2.1 of the online Supplementary Materials). Considering four different scenarios of data incompleteness and aggregation, our approach was found to be very accurate in the estimation of the model parameters and of the evolution of the effective reproduction number.

Some extensions of the presented framework are possible. If more reliable prior information or data on some extra disease states were available, the SEIR-D model could be extended by defining a more sophisticated virus transmission dynamic as proposed, for example, by Rivers *and others* (2014). On the other hand, a spatial description of the virus spread and its connection to neighbor countries would also be insightful. Finally, it could be interesting to include a time-varying under-reporting mechanism into our proposal. A starting point is given by the method of Gamado *and others* (2014) where the under-reporting rate is assumed piecewise constant with known jumps.

## 6. SUPPLEMENTARY MATERIAL

Supplementary material is available online at http://biostatistics.oxfordjournals.org.

## ACKNOWLEDGMENTS

REFERENCES

ALLEN, L.J.S. (2010). *An Introduction to Stochastic Processes with Applications to Biology (2nd edition)*. Chapman and Hall / CRC.

ALTHAUS, C.L. (2014). Estimating the reproduction number of Ebola virus (EBOV) during the 2014 outbreak in West Africa. *PLOS Currents Outbreaks* **1**.

ANDERSON, R.M. AND MAY, R.M. (1992). *Infectious Diseases of Humans: Dynamics and Control*, Oxford science publications. OUP Oxford.

BRESLOW, N E. (1984). Extra-Poisson variation in log-linear models. *Applied Statistics* **33**(1), 38–44.

BRITTON, T. (2010). Stochastic epidemic models: a survey. *Mathematical Biosciences* **225**(1), 24–35.

CHOWELL, G., HENGARTNER, N.W., CASTILLO-CHAVEZ, C., FENIMORE, P.W. AND HYMAN, J.M. (2004). The basic reproductive number of Ebola and the effects of public health measures: the cases of Congo and Uganda. *Journal of Theoretical Biology* **229**(1), 119–126.

DIEKMANN, O., HEESTERBEEK, J.A.P. AND METZ, J.A.J. (1990). On the definition and the computation of the basic reproduction ratio $R_0$ in models for infectious diseases in heterogeneous populations. *Journal of Mathematical Biology* **28**, 365–382.

DIERCKX, P. (1995). *Curve and Surface Fitting with Splines*. Oxford University Press.

DOWELL, S.F., MUKUNU, R., KSIAZEK, T.G., KHAN, A.S., ROLLIN, P. E. AND PETERS, C.J. (1999). Transmission of Ebola hemorrhagic fever: a study of risk factors in family members, Kikwit, Democratic Republic of the Congo, 1995. *Journal of Infectious Diseases* **179**(Supplement 1), S87–S91.

EILERS, P.H.C. AND MARX, B.D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science* **11**(2), 115–121.

EILERS, P.H.C. AND MARX, B.D. (2010). Splines, knots, and penalties. *Wiley Interdisciplinary Reviews: Computational Statistics* **2**(6), 637–653.

GAMADO, K.M., STREFTARIS, G. AND ZACHARY, S. (2014). Modelling under-reporting in epidemics. *Journal of Mathematical Biology* **69**(3), 737–765.

GELMAN, A. AND HILL, J. (2007). *Data Analysis Using Regression and Multilevel/Hierarchical Models*, Analytical Methods for Social Research. Cambridge University Press.

HAARIO, H., SAKSMAN, E. AND TAMMINEN, J. (2001, 04). An adaptive Metropolis algorithm. *Bernoulli* **7**(2), 223–242.

HEFFERNAN, J.M., SMITH, R.J. AND WAHL, L.M. (2005). Perspectives on the basic reproductive ratio. *Journal of the Royal Society Interface* **2**(4), 281–293.

HENS, N., SHKEDY, Z., AERTS, M., FAES, C., VAN DAMME, P. AND BEUTELS, P. (2012). *Modeling Infectious Disease Parameters Based on Serological and Social Contact Data*, Statistics for Biology and Health. Springer.

JACKMAN, S. (2009). *Bayesian Analysis for the Social Sciences*, Wiley Series in Probability and Statistics. Wiley.

JULLION, A. AND LAMBERT, P. (2007). Robust specification of the roughness penalty prior distribution in spatially adaptive Bayesian P-spline models. *Computational Statistics & Data Analysis* **51**(5), 2542 – 2558.

LEKONE, P.E. AND FINKENSTÄDT, B.F. (2006). Statistical inference in a stochastic epidemic SEIR model with control intervention: Ebola as a case study. *Biometrics* **62**(4), 1170–1177.

OZER, P., THIRY, A., FALLON, C., BLOCHER, J. AND DE LONGUEVILLE, F. (2014). Containment in sierra leone: the inability of a state to confront Ebola ? *The Lancet* **384**(9950), e47.

R CORE TEAM. (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

RIVERS, C.M., LOFGREN, E.T., MARATHE, M., EUBANK, S. AND LEWIS, B.L. (2014). Modeling the impact of interventions on an epidemic of Ebola in Sierra Leone and Liberia. *PLOS Currents Outbreaks* **1**.

SOETAERT, K., PETZOLDT, T. AND SETZER, R.W. (2010). Solving differential equations in R: Package desolve. *Journal of Statistical Software* **33**(9), 1–25.

VOGEL, G. (2014). Genomes reveal start of Ebola outbreak. *Science* **345**(6200), 989–990.

WHO. (2015). Sierra Leone: a traditional healer and a funeral. http://www.who.int/csr/disease/ebola/ebola-6-months/sierra-leone/en/. Accessed: 2015-05-12.

WHO EBOLA RESPONSE TEAM. (2014). Ebola virus disease in West Africa. the first 9 months of the epidemic and forward projections. *New England Journal of Medicine* **371**(16), 1481–1495.
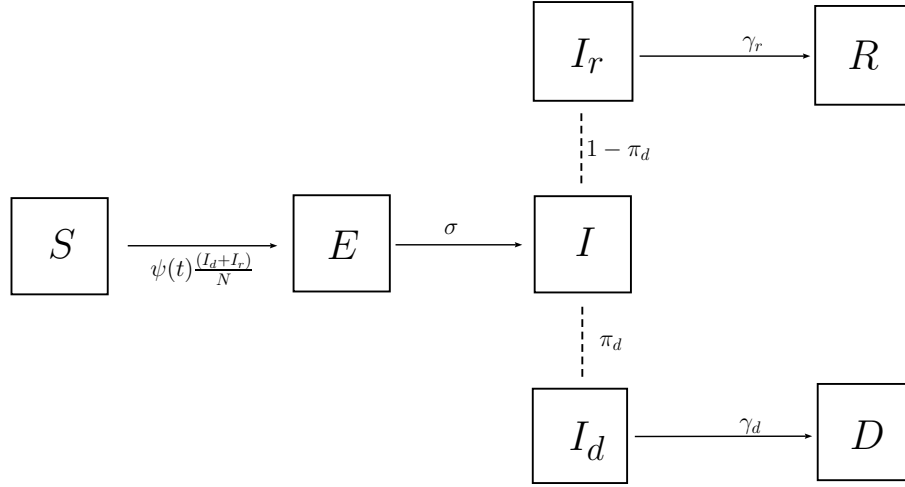
Fig. 1. SEIR-D model for the Ebola epidemic.

Table 1. Prior elicitation on SEIR-D model parameters using data (WHO Ebola Response Team, 2014) gathered during the individual follow-up of a small portion of subjects suffering from Ebola in Sierra Leone.

|  | Sample mean $\pm$ s.d. | # Ebola subjects followed-up | Specified prior |
|---|---|---|---|
| Number of deaths | 307 deaths | 445 | $\pi_d \sim \text{Beta}(308, 139)$ |
| Duration of incubation | $9.0 \pm 8.1$ days | 201 | $\sigma^{-1} \sim \mathcal{N}\left(9.0,\ 8.1^2/201\right)$ |
| Time from 1st symptoms |  |  |  |
| – to death | $8.6 \pm 6.9$ days | 128 | $\gamma_d^{-1} \sim \mathcal{N}\left(8.6,\ 6.9^2/128\right)$ |
| – to recovery | $17.2 \pm 6.2$ days | 70 | $\gamma_r^{-1} \sim \mathcal{N}\left(17.2,\ 6.2^2/70\right)$ |

Table 2. Descriptive summaries of the samples produced by $500k$ iterations of the random-walk Metropolis-within-Gibbs algorithm (after a burn-in of $10k$) for parameters $\phi_I$, $\sigma^{-1}$, $\pi_d$, $\gamma_d^{-1}$ and $\gamma_r^{-1}$ in the SEIR-D model.

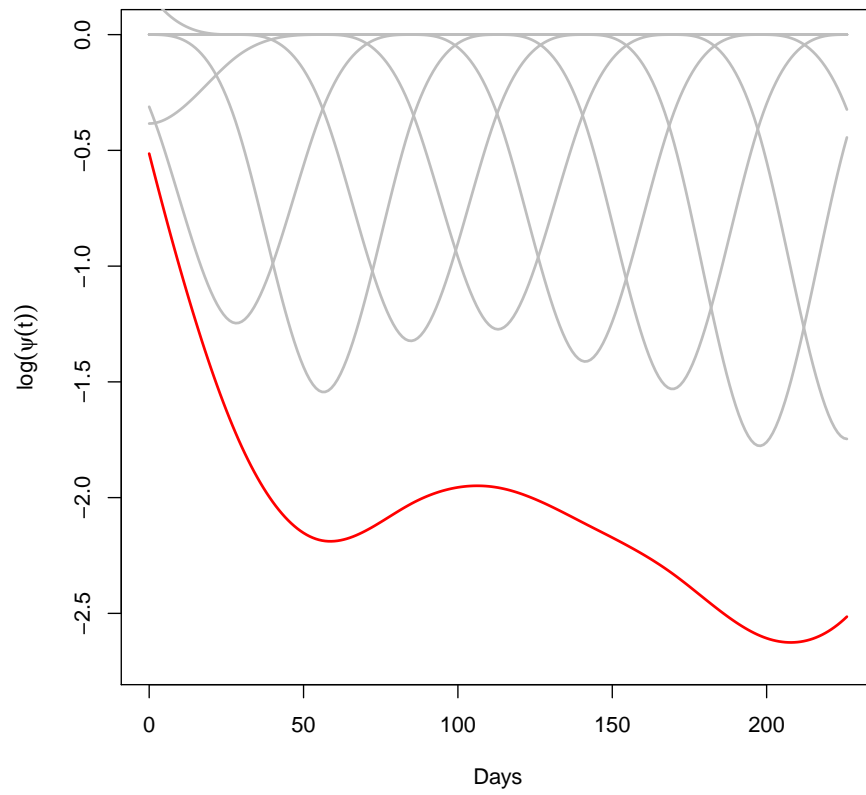|  | Mean | SD | Quantiles 2.5% | 50% | 97.5% |
|---|---|---|---|---|---|
| $\phi_I$ | 65.817 | 19.080 | 29.653 | 66.082 | 97.765 |
| $1/\sigma$ | 9.310 | 0.558 | 8.219 | 9.306 | 10.417 |
| $\pi_d$ | 0.689 | 0.021 | 0.645 | 0.689 | 0.731 |
| $1/\gamma_d$ | 8.552 | 0.617 | 7.339 | 8.550 | 9.777 |
| $1/\gamma_r$ | 17.203 | 0.742 | 15.747 | 17.204 | 18.676 |

Fig. 2. B-spline model for $\log \psi(t)$. The successive terms in (2.6) are plotted in the same order as grey lines adding up to $\log \psi(t)$ in red.

**Daily new Ebola cases in Sierra Leone**

**Daily new Ebola deaths in Sierra Leone**
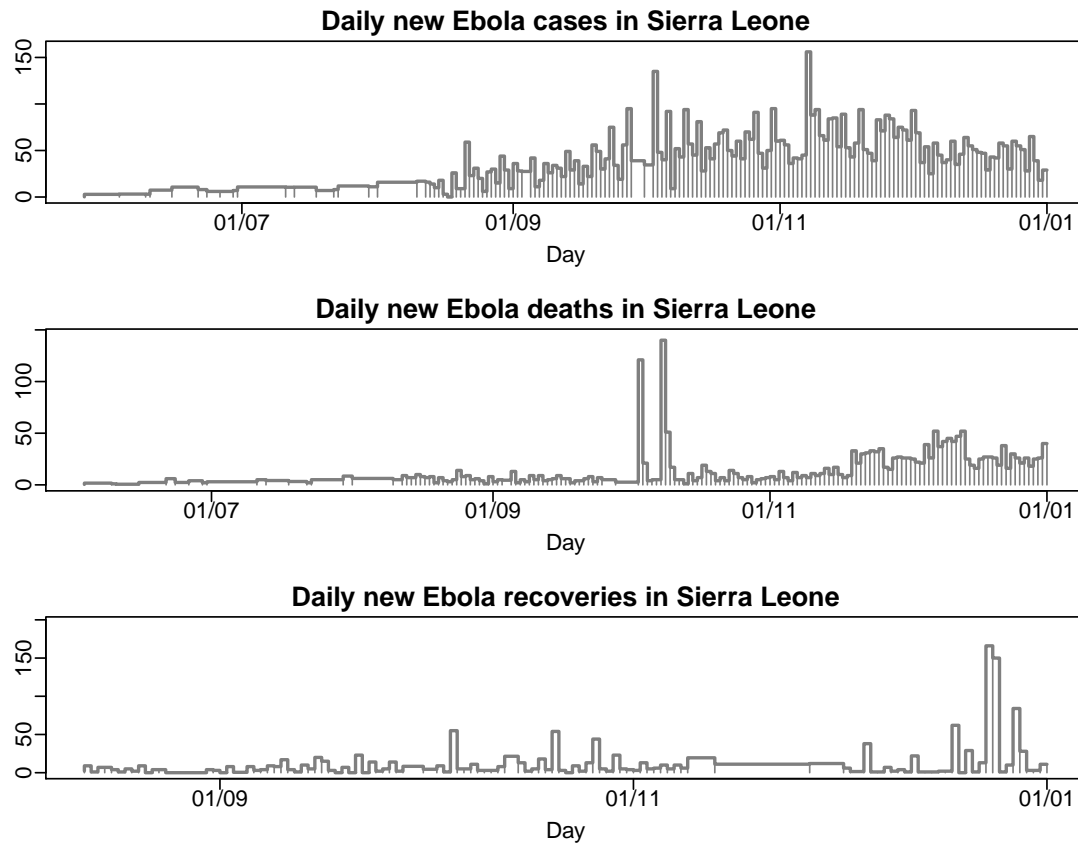
**Daily new Ebola recoveries in Sierra Leone**

Fig. 3. Histograms of the reported numbers of new Ebola cases, deaths and recoveries confirmed by a laboratory analysis in Sierra Leone in 2014.
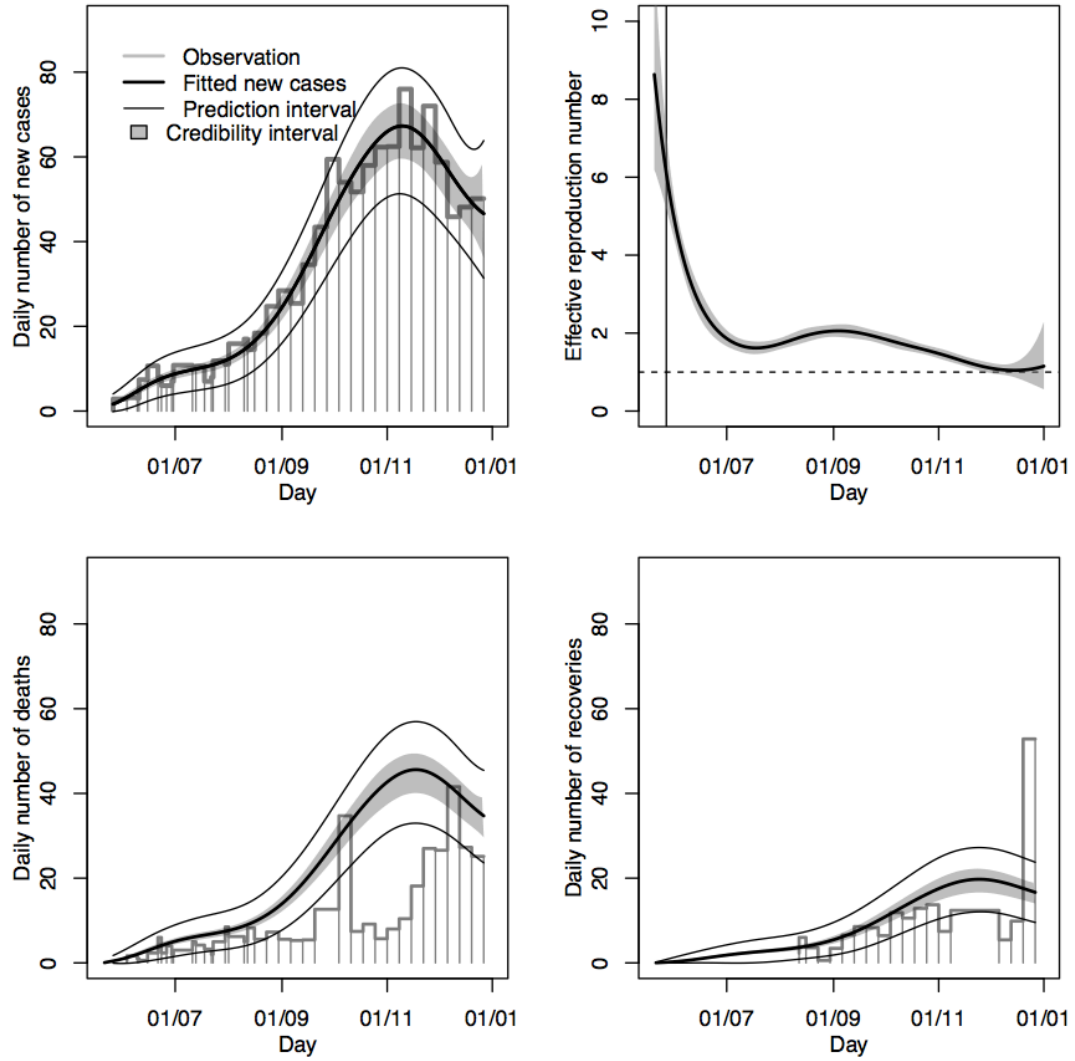
Fig. 4. Upper left panel: reported (histogram), fitted (thick solid line) and predicted (thin solid line) number of new confirmed Ebola cases in Sierra-Leone using the SEIR-D model. Upper right panel: posterior mean and pointwise 95% credible intervals for the effective reproduction number (where the vertical black line indicates the date of the first available observation: 26 May 2014). Lower left panel: reported (histogram), fitted (thick solid line), and predicted (thin solid lines) numbers of deaths in Sierra-Leone using the SEIR-D model. Lower right panel: reported (histogram), fitted (thick solid line), and predicted (thin solid lines) numbers of recovered in Sierra-Leone using the SEIR-D model

# Bayesian inference in an extended SEIR model with nonparametric disease transmission rate: an application to the Ebola epidemic in Sierra Leone – Supplementary Materials –

GIANLUCA FRASSO, PHILIPPE LAMBERT[*]

*Faculté des sciences sociales, Méthodes quantitatives en sciences sociales, Université de Liège, Liège, Belgium.*

*Institut de statistique, biostatistique et sciences actuarielles (ISBA), Université catholique de Louvain, Louvain-la-Neuve, Belgium.*

p.lambert@ulg.ac.be

## 1. SEIR-D MODEL FOR THE EBOLA EPIDEMIC

### 1.1  *Poisson process*

Basic theory on Markov processes with discrete states in continuous time (see e.g. Cox and Miller, 1965) provide the following results. Consider point events occurring singly in time. If $M(t, t + \mathrm{d}t)$ is the number of events occurring in $(t, t + \mathrm{d}t)$ such that, when $\mathrm{d}t \to 0^+$,

$$\Pr\big(M(t, t + \mathrm{d}t) = 0\big) = 1 - \rho\,\mathrm{d}t + o(\mathrm{d}t),$$
$$\Pr\big(M(t, t + \mathrm{d}t) = 1\big) = \rho\,\mathrm{d}t + o(\mathrm{d}t),$$
$$\Pr\big(M(t, t + \mathrm{d}t) > 1\big) = o(\mathrm{d}t),$$

then
1. $M(t, t + h) \sim \mathrm{Pois}(\rho h)$ for any $h > 0$. In particular, $\rho$ is the expected number of events per unit time.
2. The intervals between successive events are independently distributed with an exponential distribution of mean $1/\rho$.

### 1.2  *Expected number of state transitions*

We provide details for the derivation of some of the expected numbers of transitions between several of the disease states in Section 2.1 of the main paper:

$\underline{S \longrightarrow E}$ : the $i$th susceptible person at time $t$ has contacts with $\mathrm{d}C_i(t)$ subjects in a time interval $(t, t + \mathrm{d}t)$. Denote by $\beta(t)\mathrm{d}t$ the mean of that random variable and assume that $\mathrm{d}t$ is small enough to ensure that a single person can only experience at most one epidemic state transition in $(t, t + \mathrm{d}t)$.

We want to compute the expected number of new exposed, $dE^+(t) = S(t) - S(t + dt)$, i.e. the expected number of transitions from the susceptible to the exposed state during this time interval. According to the state definitions, a susceptible subject can only get infected through contact(s) with (the body fluids of) one of the $I(t)$ infectious persons at time $t$. We denote by $\pi_{\mathcal{E}}(t)$ the time-varying probability that a contact between a susceptible and an infectious subject effectively leads to a new infection. This probability is expected to vary over time as a consequence of prevention actions (such as quarantine or population containment policies for example) or due to change in social behavior as a reaction to the epidemic threat. From our definitions, it follows that, over the time span $(t, t + dt)$, the expected number of transitions from state $E$ to $I$ is

$$\mathbb{E}\left(dE^+(t)|S(t), I(t), \beta(t), \pi_{\mathcal{E}}(t)\right) = \sum_{i=1}^{S(t)} \mathbb{E}\left(dE_i^+(t)|S(t), I(t), \beta(t), \pi_{\mathcal{E}}(t)\right), \qquad (1.1)$$

where $dE_i^+(t)$ is 1 if the $i$th susceptible subject at time $t$ changes state by $t + dt$ and becomes exposed after contacts with $dC_i(t)$ persons, or 0 if s/he remains susceptible. One has

$$\Pr\left(dE_i^+(t) = 1|S(t), I(t), \beta(t), \pi_{\mathcal{E}}(t)\right)$$
$$= \sum_{j \in \mathbb{N}} \Pr\left(dE_i^+(t) = 1|dC_i(t) = j, S(t), I(t), \pi_{\mathcal{E}}(t)\right) \Pr\left(dC_i(t) = j|\beta(t)\right), \qquad (1.2)$$

On the other hand, as already mentioned, a contact with an infectious subject does not necessarily lead to an infection. Let $dK_i(t)$ be the number of contacts (for subject $i$) in $(t, t + dt)$ leading to a virus transmission. Given $dC_i(t) = j$, the preceding definition of $\pi_{\mathcal{E}}(t)$, the homogeneous mixing assumption and the proportion of infectious subjects in the population, we conclude that

$$\left(dK_i(t)|dC_i(t) = j, I(t), \pi_{\mathcal{E}}(t)\right) \sim \text{Bin}\left(j, \pi_{\mathcal{E}}(t)\frac{I(t)}{N} dt\right).$$

Therefore, the conditional probability that the $i$th susceptible subject becomes infected during $(t, t + dt)$ through contacts with $j$ persons is

$$\Pr\left(dE_i^+(t) = 1|dC_i(t) = j, S(t), I(t), \pi_{\mathcal{E}}(t)\right) = \Pr\left(dK_i(t) \geqslant 1|dC_i(t) = j, I(t), \pi_{\mathcal{E}}(t)\right)$$
$$= 1 - \left(1 - \pi_{\mathcal{E}}(t)\frac{I(t)}{N} dt\right)^j$$
$$= j\pi_{\mathcal{E}}(t)\frac{I(t)}{N} dt + o(dt). \qquad (1.3)$$

Note that the first term of the last expression as an approximation to the probability of interest also holds for non necessarily small time interval provided that the number of infectious is small compared to the population size (as in the Ebola epidemic) or if $\pi_{\mathcal{E}}(t)$ is small. Combining (1.1), (1.2) and (1.3), one gets

$$\mathbb{E}\left(dE^+(t)|S(t), I(t), \beta(t), \pi_{\mathcal{E}}(t)\right) = \sum_{i=1}^{S(t)} \sum_{j \in \mathbb{N}} \left(j\pi_{\mathcal{E}}(t)\frac{I(t)}{N} dt\right) \Pr\left(dC_i(t) = j|\beta(t)\right) + o(dt)$$
$$= \beta(t)\pi_{\mathcal{E}}(t)\frac{I(t)}{N} S(t) dt + o(dt). \qquad (1.4)$$

The mean $\beta(t)\mathrm{d}t$ and the disease transmission probability $\pi_{\mathcal{E}}(t)$ during a contact with an infectious subject are unknown and need to be estimated. In order to avoid identifiability issues, we estimate their product $\psi(t) = \beta(t)\pi_{\mathcal{E}}(t)$. Then, the expected number of new exposed in a time interval $(t, t + \mathrm{d}t)$ becomes

$$\mathbb{E}\left(\mathrm{d}E^+(t)|S(t), I(t), \psi(t)\right) = S(t)\frac{I(t)}{N}\psi(t)\,\mathrm{d}t + o(\mathrm{d}t).$$

$I \longrightarrow R, D$ : when turning infectious, a person has probability $\pi_d$ to die from Ebola. At each time $t$, we can distinguish two groups of infectious subjects: $I_d(t)$ of them will die at rate $\gamma_d$ (i.e. on average $1/\gamma_d$ days after the appearance of the first symptoms) and the other $I_r(t)$ ($= I(t) - I_d(t)$) will recover (and become immune) at a slower rate $\gamma_r$ (i.e. on average $1/\gamma_r$ days after the appearance of the first symptoms). We want to compute the expected number of new recoveries and deaths in a time interval $(t, t + \mathrm{d}t)$. Denote by $\delta_i^R(t)$ the indicator variable taking value 1 (with probability $1 - \pi_d$) if the $i$th infectious subject will move to the recovery compartment $I_r$ and zero otherwise. Then, given that $\delta_i^R(t) = 0$ necessarily implies that $\mathrm{d}R_i^+(t)$ is also zero, one has

$$\mathbb{E}\left(\mathrm{d}R^+(t)|I(t), \pi_d, \gamma_r\right) = \sum_{i=1}^{I(t)}\left\{\mathbb{E}\left(\mathrm{d}R_i^+(t)|\delta_i^R(t) = 0, \gamma_r\right)\pi_d + \mathbb{E}\left(\mathrm{d}R_i^+(t)|\delta_i^R(t) = 1, \gamma_r\right)(1 - \pi_d)\right\},$$

$$= \sum_{i=1}^{I(t)}\gamma_r(1 - \pi_d)\,\mathrm{d}t + o(\mathrm{d}t) = \gamma_r(1 - \pi_d)I(t)\,\mathrm{d}t + o(\mathrm{d}t).$$

Following a similar reasoning, the expected number of new deaths in the same time interval is

$$\mathbb{E}\left(\mathrm{d}D^+(t)|I(t), \pi_d, \gamma_d\right) = \gamma_d\pi_d I(t)\,\mathrm{d}t + o(\mathrm{d}t).$$

### 1.3 *Other possible model extensions*

Our differential equation model can be extended or simplified in many ways. Among the possible extensions, one could think of

- Forcing the risk $\pi_{\mathcal{E}}(t)$ of being infected during a contact with an infectious person to decrease over time;

- Assuming that the probabilities of death and recovery of infectious persons (and the speed of their transitions to these absorbing states) are changing over time thanks to e.g. an improving handling of symptomatic patients;

- Acknowledging that an increasing number of symptomatic persons are sent to an hospital, thereby reducing the risk of transmission of the disease. By ignoring this feature (as we lack reliable information on the time varying proportion of hospitalized Ebola patients), we implicitly entrust that task to $\psi(t)$;

- Adding a 'Funeral' state between the $I_d$ and the $D$ compartments where persons deceased due to Ebola have a temporary increased chance to transmit the disease during the funeral rites. That feature was particularly relevant at the start of the epidemic. However, like for hospitalized persons, the lack of reliable statistics on this aspect led us to prefer the proposed SEIR-D model;

- Using the ODE model on each district of Sierra Leone and of the neighbor countries with the inclusion of interactions between them. Indeed, internal population movements and emigration have an impact on the virus transmission: the epidemic started in Sierra Leone in the East through contacts with Guinean residents (WHO, 2015) and progressively spread Westwards to the capital city and the more densely populated Western districts.

Given that the monitoring data on the number of new transitions between states are only reasonably reliable for the entries in the $I$ state, we decided to focus our efforts in the main paper (see Section 2.3) on a flexible specification of the disease transmission rate.

### 1.4    *The effective reproduction number*

In order to derive an expression for the reproduction number we need to compute the so called *new generation matrix* $\mathcal{G}$ (Diekmann *and others*, 1990). This is a square matrix with element $(i, j)$ giving the expected number of secondary infections of type $i$ caused by a single infected individual of type $j$ in the dynamical system under consideration. In other words, for the computation of $\mathcal{G}$, we need to consider all the subjects that can spread the disease including the exposed ones since they will become infectious after the incubation period.

At each time $t$, according to the stochastic model (2.4) (in the main paper) with associated ODE equations (2.5) (in the main paper), we can distinguish three classes of infected subjects: persons in states $E$, $I_r$ and $I_d$. Hence, in order to compute $\mathcal{G}$, we need to focus on the following subset of equations,

$$e'(t) = \psi(t) \left(i_r(t) + i_d(t)\right) s(t) - \sigma\, e(t) \; ; \; i'_r(t) = \sigma\left(1 - \pi_d\right) e(t) - \gamma_r\, i_r(t) \; ; \; i'_d(t) = \sigma\, \pi_d\, e(t) - \gamma_d\, i_d(t).$$

Letting $\boldsymbol{x}(t) = (e(t), i_r(t), i_d(t))^T$ be the vector of proportions of subjects in each of the infected states (since we assume that only exposed and infectious people can spread the disease), these equations can be rewritten as

$$\boldsymbol{x}'(t) = (\mathcal{F} - \mathcal{V})\boldsymbol{x}(t),$$

with

$$\mathcal{F} = \left[\frac{\partial F_k}{\partial x_\ell}\right] = \begin{bmatrix} 0 & \psi(t)s(t) & \psi(t)s(t) \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \; ; \; \mathcal{V} = \left[\frac{\partial V_k}{\partial x_\ell}\right] = \begin{bmatrix} \sigma & 0 & 0 \\ -\sigma(1 - \pi_d) & \gamma_r & 0 \\ -\sigma\pi_d & 0 & \gamma_d \end{bmatrix},$$

where $F_k$ denotes the rate of change in $x_k(t)$ due to new infected subjects and $V_k$ the transfer rate of already infected subjects to the $k$th compartment (of infected subjects).

Following Diekmann *and others* (1990), the next generation matrix is equal to $\mathcal{G} = \mathcal{F}\mathcal{V}^{-1}$ and the reproduction number is equal to the spectral radius (i.e. the largest eigenvalue) of $\mathcal{G}$. In our case it is easy to verify that

$$\mathcal{G} = \psi(t)s(t) \begin{bmatrix} (1 - \pi_d)\gamma_r^{-1} + \pi_d\gamma_d^{-1} & \gamma_r^{-1} & \gamma_d^{-1} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

and, hence, that

$$\mathcal{R}_e(t) = \psi(t)s(t) \times \left(\frac{1 - \pi_d}{\gamma_r} + \frac{\pi_d}{\gamma_d}\right).$$

## 2. SIMULATION STUDIES

### 2.1 *Properties of the parameter estimates obtained using the ODE method*

In this section, we evaluate the performances of our estimation strategy based on ODE to estimate transition parameters from data generated using the associated discrete time stochastic SEIR-D model. Five hundred datasets were generated using a discrete time stochastic SEIR-D model

$$S(t + \mathrm{d}t) = S(t) - \mathrm{d}E^+(t)$$
$$E(t + \mathrm{d}t) = E(t) + \mathrm{d}E^+(t) - (\mathrm{d}I_r^+(t) + \mathrm{d}I_d^+(t))$$
$$I_r(t + \mathrm{d}t) = I_r(t) + \mathrm{d}I_r^+(t) - \mathrm{d}R^+(t) \ ; \ \ I_d(t + \mathrm{d}t) = I_d(t) + \mathrm{d}I_d^+(t) - \mathrm{d}D^+(t)$$
$$R(t + \mathrm{d}t) = R(t) + \mathrm{d}R^+(t) \ ; \ \ D(t + \mathrm{d}t) = D(t) + \mathrm{d}D^+(t),$$

with hourly ($\mathrm{d}t = 1/24$ day) Poisson innovations over 360 days,

$$\mathrm{d}E^+(t) \sim \mathrm{Pois}\Big(S(t)\frac{I_r(t) + I_d(t)}{N}\psi(t)\,\mathrm{d}t\Big)$$
$$\mathrm{d}I_r^+(t) \sim \mathrm{Pois}\Big(\sigma\pi_d E(t)\,\mathrm{d}t\Big) \ ; \ \ \mathrm{d}I_d^+(t) \sim \mathrm{Pois}\Big(\sigma(1 - \pi_d)E(t)\,\mathrm{d}t\Big)$$
$$\mathrm{d}R^+(t) \sim \mathrm{Pois}\Big(\gamma_r I_r(t)\,\mathrm{d}t\Big) \ ; \ \ \mathrm{d}D^+(t) \sim \mathrm{Pois}\Big(\gamma_d I_d(t)\,\mathrm{d}t\Big),$$

and considering the following initial conditions

$$S(0) = N - 23 \ ; \ E(0) = 13 \ ; \ I_r(0) = 0 \ ; \ I_d(0) = 10 \ ; \ R(0) = 0 \ ; \ D(0) = 0 \ ; \ N = 10^6.$$

Values consistent with the rate values measured during the follow-up of Ebola patients in Sierra Leone (cf. Table 1 in the main paper ; see also WHO Ebola Response Team (2014)) were chosen for the parameters

$$\sigma^{-1} = 9.0 \ ; \ \gamma_d^{-1} = 8.6 \ ; \ \gamma_r^{-1} = 17.2 \ ; \ \pi_d = 0.69.$$

The time-varying parameter $\psi(t)$ was assumed to decline smoothly over time:

$$\psi(t) = \exp\big(-1.25 + t/150 + \cos(t/50)\mu(t)\big) \ \text{ with } \ \mu(t) = 0.25 + 0.25t/360.$$

The corresponding evolution of the reproduction number is presented in Fig. Sup.1 (solid black line).

Synthetic observations were simulated on a hourly basis over approximately one year (360 days) under four possible scenarios. In scenario 1, exact reports were generated on a daily basis for the number of new cases, deaths and recoveries. In scenario 2, we assume that data are available for all the epidemic states (cases, deaths and recoveries), but only through weekly reports of aggregated counts. Finally, in the last two scenarios, the impact of incomplete information on the epidemic has been assessed: parameter estimation was performed using only the reported numbers of new cases (either on a daily or on a weekly basis) by treating the other two states as latent (unobserved).

The results of our simulation study are presented in Table Sup.1 and Fig. Sup.1. Parameter estimates were obtained by maximizing the joint posterior distribution described in Section 4 with the priors of Table 1 in the main paper. As expected, the quality of the estimates tends to be higher with daily observations, even if good results are already obtained with weekly aggregated data. By comparing the results for scenarios 1-2 and 3-4, it appears that observing only the new cases (and

not deaths and recoveries) has a limited impact on bias, and tends to slightly increase the standard deviations and the relative RMSEs of the parameter estimators.

Finally, for all simulation settings, Fig. Sup.1 highlights the quality of the estimation of the true reproduction number with the proposed ODE model combined with a P-spline specification of the disease transmission rate. These results were obtained by using cubic B-splines built on 8 equally spaced internal knots and a third order difference penalty.

## 2.2   *Impact of a violation of the homogeneous population hypothesis*

We want to evaluate the impact of a violation of the homogeneous mixing hypothesis on the quality of the estimates obtained by using the ODE-based estimation strategy. To reach this goal, we have simulated 500 datasets by taking into account a patch SEIR-D model including two coupled subpopulations ($k = 1, 2$) with a model specification generalizing Eq. (2.5) in the main manuscript:

$$\frac{\mathrm{d}s^k(t)}{\mathrm{d}t} = -s^k(t) \sum_{j=1}^{2} \psi_{kj}(t) \left( i_r^j(t) + i_d^j(t) \right)$$

$$\frac{\mathrm{d}s^k(t)}{\mathrm{d}t} = s^k(t) \sum_{j=1}^{2} \psi_{kj}(t) \left( i_r^j(t) + i_d^j(t) \right) - \sigma e^k(t)$$

$$\frac{\mathrm{d}i_r^k(t)}{\mathrm{d}t} = (1 - \pi_d)\sigma e^k(t) - \gamma_r i_r^k(t)$$

$$\frac{\mathrm{d}i_d^k(t)}{\mathrm{d}t} = \pi_d \sigma e^k(t) - \gamma_d i_d^k(t)$$

$$\frac{\mathrm{d}r^k(t)}{\mathrm{d}t} = \gamma_r i_r^k(t)$$

$$\frac{\mathrm{d}d^k(t)}{\mathrm{d}t} = \gamma_d i_d^k(t).$$

In this system, the force of infection in patch $k$ is given by

$$\sum_{j=1}^{2} \psi_{kj}(t) \left( i_r^j(t) + i_d^j(t) \right).$$

Hence, infectious individuals belonging to one subpopulation can infect susceptible individuals from the other one, but the contact mechanism between the two patches is left unspecified. An analogous model has been proposed by Lloyd and May (1996) to incorporate population heterogeneity in a SEIR framework. Analogously to the simulation settings presented in the previous subsection, the counts for the simulated compartments were generated by discrete time stochastic simulation of the patch ODE model. In particular we used the following system

$$S^k(t + \mathrm{d}t) = S^k(t) - \mathrm{d}E^{k+}(t)$$
$$E^k(t + \mathrm{d}t) = E^k(t) + \mathrm{d}E^{k+}(t) - (\mathrm{d}I_r^{k+}(t) + \mathrm{d}I_d^{k+}(t))$$
$$I_r^k(t + \mathrm{d}t) = I_r^k(t) + \mathrm{d}I_r^{k+}(t) - \mathrm{d}R^{k+}(t) \ ; \ I_d^k(t + \mathrm{d}t) = I_d^k(t) + \mathrm{d}I_d^{k+}(t) - \mathrm{d}D^{k+}(t)$$
$$R^k(t + \mathrm{d}t) = R^k(t) + \mathrm{d}R^{k+}(t) \ ; \ D^k(t + \mathrm{d}t) = D^k(t) + \mathrm{d}D^{k+}(t),$$

with innovations $dE^{k+}(t), dI_r^{k+}(t), dI_d^{k+}(t), dR^{k+}(t), dD^{k+}(t)$ drown from Poisson distributions with expected values as described in (2.4) (see main manuscript). The two patches were coupled using

$$\mathbb{E}(dE^{k+}(t)|\boldsymbol{\theta}, \boldsymbol{p}_0) = S^k(t) \sum_{j=1}^{2} \psi_{kj}(t) \frac{\left(I_r^j(t) + I_d^j(t)\right)}{N^j} \, dt.$$

The time varying parameters $\psi_{kj}(t)$ were specified as in the previous subsection. We supposed that subpopulation 2 is one fifth of the other one and has a smaller contact rate (equal to 85% of the contact rate for patch 1). Also the probability to die after an Ebola infection has been set differently for the two subpopulations ($\pi_d^1 = 0.65$, $\pi_d^2 = 0.75$) while the other parameters were set as follows, $\sigma^{-1} = 9.0$ ; $\gamma_d^{-1} = 8.6$ ; $\gamma_r^{-1} = 17.2$. Finally, we considered the following initial states (at $t = 0$) to simulate the data:

$N^1 = 4800000, \ S^1(0) = (N^1 - 10), \ E^1(0) = 10, \ I_r^1(0) = 0, \ I_d^1(0) = 0, \ R^1(0) = 0, \ D^1(0) = 0,$

$N^2 = 1200000, \ S^2(0) = (N^2 - 30), \ E^2(0) = 20, \ I_r^2(0) = 0, \ I_d^2(0) = 10, \ R^2(0) = 0, \ D^2(0) = 0.$

The ODE parameters and the effective reproduction number were estimated using the Bayesian approach introduced in the manuscript. In particular, we modeled the daily and weekly national number of new cases (obtained by aggregating the data simulated for the two regions) as negative binomial distributed with unknown overdispersion parameter. The results in terms of parameter estimates are summarized in Table Sup.2, while Fig. Sup.2, shows the simulated and estimated effective reproduction numbers.

From our results it clearly appears that the estimation bias tends to be lower when dealing with daily observations. The same is true for the relative root mean squared errors and the standard deviation of the estimates. As expected, the estimated overdispersion parameter has always been found close to 100 pointing to a Poisson distribution for the observed counts (not shown).

Finally, Fig. Sup.2 suggests that the dynamic of $R_e(t)$ is properly estimated, although when dealing with weekly observations, the reproduction number appears (on average) underestimated during the first few weeks.

## 2.3   *Sensitivity analysis to the choice of $E(t = 0)$*

In this section, we evaluate the robustness of the results presented in Section 4 of the paper to the choice of the initial state conditions $E(t = 0)$. In particular we are interested in evaluating the impact of the number of exposed subjects at time zero on the behavior of estimates of the effective reproduction number and of the ODE parameter estimates.

The results of our sensitivity study, based on the maximization of the posterior distribution defined in the manuscript, are presented in Fig. Sup.3 and Table Sup.3. Fig. Sup.3 shows that the initial number of exposed subjects in the population has a limited impact on the evolution of the estimated effective reproduction number. For all configurations, the reproduction number appears to steadily step back after the end of September and is close to 1.0 in the last days of December. In the period preceding June 2014, when few and coarse data are available, the estimated $\mathcal{R}_e(t)$ are less homogeneous. It is however remarkable that the estimates in Fig. Sup.3 are compatible with the credible region estimated for the effective reproduction number and reported in the paper.

The ODE-parameter estimates under the different scenarios are reported in Table Sup.3. As for the effective reproduction number, they are robust to the choice of $E(t = 0)$ and appear compatible with the 95% credible intervals reported in the paper.

### 2.4   *Evolution of the model based perception of the epidemic*

We also studied how the perception of the Ebola epidemic changed over time by fitting the flexible SEIR-D model on the data available for Sierra Leone at different stages of the outbreak. The estimates obtained by using the proposed Bayesian approach were also compared with those computed using the method advocated by Althaus (2014).

As mentioned in the main paper (Section 1), Althaus (2014) analyzed the 2014 Sierra Leone Ebola outbreak by fitting a SEIR model to the cumulative number of Ebola cases and deaths assumed to be Poisson distributed. The transmission rate (and hence the effective reproduction number) is assumed to be exponentially declining over time. In his analysis, the author estimates the appearance of the first Ebola case in the last days of April 2014 and considers only the data available up to the end of August 2014. The results obtained using this approach on the data available up to three different calendar times are reported in Fig. Sup.4. Considering the data available till 20 August 2014, the estimated effective reproduction numbers suggest that the epidemic is already under control in the first half of August. This message changes if the reproduction number is computed by adding data acquired later in time. The estimated disease transmission dynamic and the mortality rates (see the legends in Fig. Sup.4) seem unrealistic and do not properly assess the severity of the health crisis experienced in Sierra Leone during 2014.

A different perception of the epidemic dynamic is obtained using our proposal, see Fig. Sup.5. Our analysis provides a better fit to the number of new cases and suggests a more realistic behavior of the estimated effective reproduction number $\mathcal{R}_e$. For every observational horizon, the reproduction number shows a fast decline in the first period of the disease propagation. As time passes and data accumulate, $\mathcal{R}_e(t)$ shows a bounce late July till about the end of September 2014 where it starts a slow and steady decline towards 1.0 by the end of 2014.

### 2.5   *Impact of under-reported new cases on $\boldsymbol{\theta}$ and $\mathcal{R}_e(t)$ estimates*

We are interested in evaluating the impact of under-reporting on the estimates of the ODE parameters and of the effective reproduction number. Here, analogously to the data analysis presented in the main manuscript, we suppose to observe only the number of new cases reported on a weekly basis. A set of 100 data series was simulated following the guidelines described in Section 2.1 (Supplementary Materials). The numbers of new infectious cases were simulated using a discrete time stochastic approach with Poisson innovations and, in order to simulate under-reporting, we assume a monotone increasing reporting rate $\rho(t)$ taking vales in $(0.4, 0.8)$: $\rho(t) = (1 + \tanh(t/150))/2.5$ for $t \in (0, 360)$. The results are presented in Figure Sup.6, Figures Sup.7 and Figure Sup.8. Under-reporting of the number of new cases leads to biased parameter estimates (Figure Sup.8) with an over-estimation of the expected time spent by an Ebola subject in the infectious state, while the number of susceptible subjects tends to be over-estimated (upper right panel of Figure Sup.7). These two effects seem to compensate. Indeed, despite under-reporting and thanks to the flexible P-spline representation of $\psi(t)$, the effective reproduction number is estimated without bias (see upper left panel of Figure Sup.7), except at the very start of the epidemic.

## 3. Details on posterior results

In this section we illustrate some details on the posterior results. In particular, we report the posterior distribution of the P-spline coefficients and of the ODE parameters. Furthermore, chain

trace plots and ACF functions of all the unknowns are also shown. The Geweke diagnostic statistics for convergence (not shown) indicate convergence for all the unknown. Finally, the trace plot of the log-posterior is shown in Figure Sup.13.

## References

Althaus, C.L. (2014). Estimating the reproduction number of Ebola virus (EBOV) during the 2014 outbreak in West Africa. *PLOS Currents Outbreaks* **1**.

Cox, D.R. and Miller, H.D. (1965). *The Theory of Stochastic Processes*. Chapman & Hall.

Diekmann, O., Heesterbeek, J.A.P. and Metz, J.A.J. (1990). On the definition and the computation of the basic reproduction ratio $R_0$ in models for infectious diseases in heterogeneous populations. *Journal of Mathematical Biology* **28**, 365–382.

Lloyd, Alun L. and May, Robert M. (1996). Spatial heterogeneity in epidemic models. *Journal of Theoretical Biology* **179**(1), 1 – 11.

WHO. (2015). Sierra Leone: a traditional healer and a funeral. http://www.who.int/csr/disease/ebola/ebola-6-months/sierra-leone/en/. Accessed: 2015-05-12.

WHO Ebola Response Team. (2014). Ebola virus disease in West Africa. the first 9 months of the epidemic and forward projections. *New England Journal of Medicine* **371**(16), 1481–1495.

Table Sup.1. Relative bias (in percentage), standard deviations and relative root mean squared errors (rRMSE) estimated using 500 simulated datasets.

| | | Data for three states | | Data for cases only | |
|---|---|---|---|---|---|
| | | Scenario 1: daily data | Scenario 2: weekly data | Scenario 3: daily data | Scenario 4: weekly data |
| **Bias** | $\phi_c = 100$ | -0.11% | -0.24% | -0.47% | -0.93% |
| | $\phi_r = 100$ | -0.12% | -0.28% | | |
| | $\phi_d = 100$ | -0.11% | -0.17% | | |
| | $1/\sigma = 9$ | -1.22% | 1.31% | -1.52% | -2.17% |
| | $1/\gamma_d = 8.6$ | 1.03% | -1.12% | -1.35% | -1.57% |
| | $1/\gamma_r = 17.2$ | 0.75% | -0.97% | -1.02% | -2.01% |
| | $\pi_d = 0.69$ | 0.31% | 0.72% | -1.26% | -1.83% |
| | | | | | |
| **Std. dev** | $\phi_c = 100$ | 1.59E-01 | 4.07E-01 | 4.41E-01 | 5.74E-01 |
| | $\phi_r = 100$ | 1.72E-01 | 5.12E-01 | | |
| | $\phi_d = 100$ | 1.91E-01 | 5.92E-01 | | |
| | $1/\sigma = 9$ | 4.67E-01 | 5.52E-01 | 6.12E-01 | 7.08E-01 |
| | $1/\gamma_d = 8.6$ | 4.46E-01 | 5.40E-01 | 2.83E-01 | 3.19E-01 |
| | $1/\gamma_r = 17.2$ | 6.69E-01 | 8.09E-01 | 2.49E-01 | 3.13E-01 |
| | $\pi_d = 0.69$ | 2.91E-03 | 6.74E-03 | 4.14E-02 | 7.07E-02 |
| | | | | | |
| **rRMSE** | $\phi_c = 100$ | 0.28E-01 | 0.47E-01 | 2.52E-01 | 3.81E-01 |
| | $\phi_r = 100$ | 0.32E-01 | 1.21E-01 | | |
| | $\phi_d = 100$ | 0.17E-01 | 0.82E-01 | | |
| | $1/\sigma = 9$ | 4.12E-03 | 5.60E-03 | 8.76E-03 | 3.49E-02 |
| | $1/\gamma_d = 8.6$ | 3.19E-02 | 4.68E-02 | 2.51E-02 | 4.14E-02 |
| | $1/\gamma_r = 17.2$ | 1.01E-02 | 6.07E-02 | 5.20E-02 | 7.70E-02 |
| | $\pi_d = 0.69$ | 3.18E-03 | 6.65E-03 | 1.20E-02 | 1.52E-02 |

Table Sup.2. Relative bias (in percentage), standard deviations and relative root mean squared errors (rRMSE) estimated using 500 datasets simulated according to a two regions patch SEIR-D model.

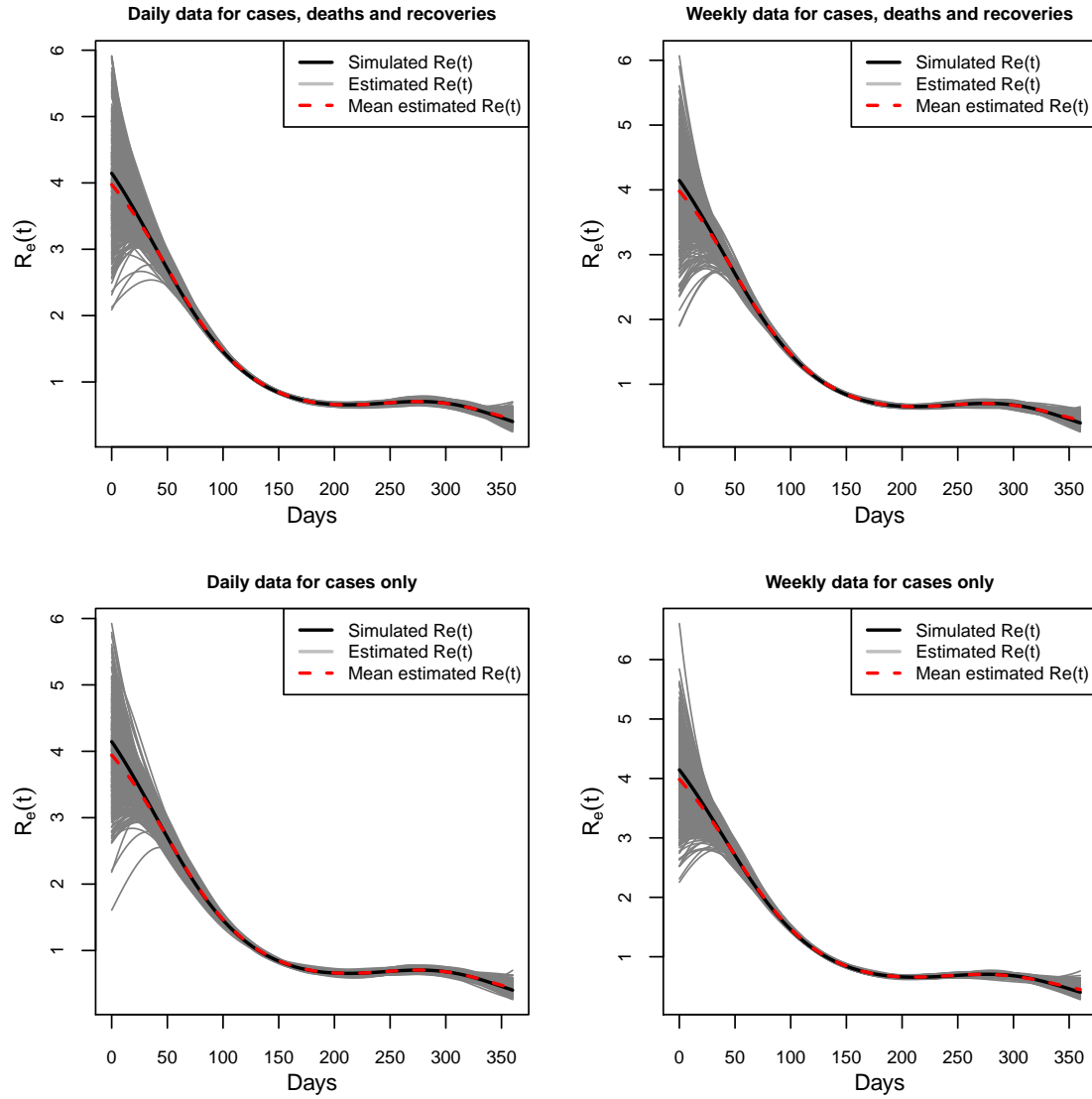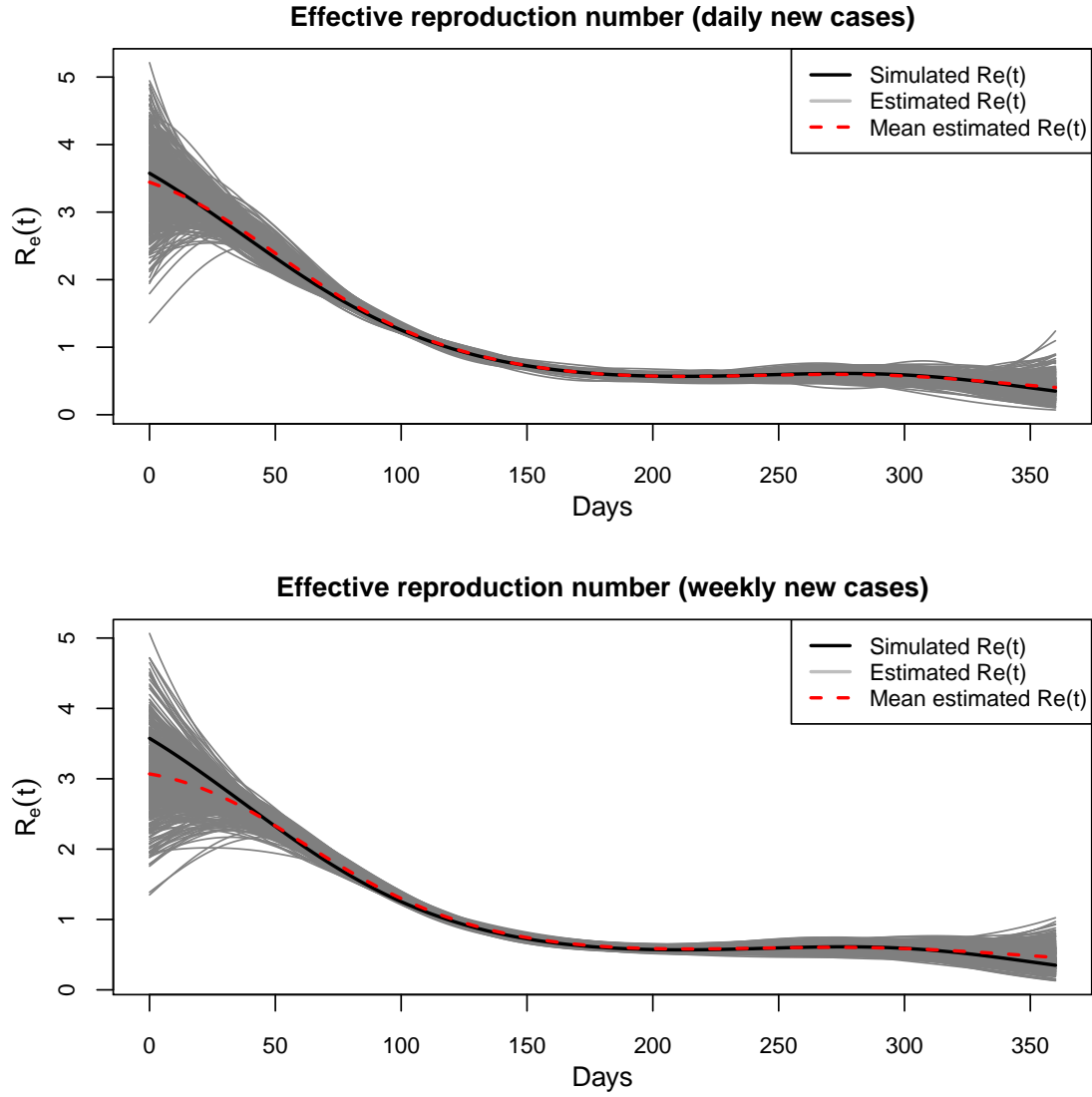| | Parameters | BIAS | Std. Dev. | rRMSE |
|---|---|---|---|---|
| **Daily observations** | | | | |
| | $1/\sigma = 9$ | 2.62% | 2.99E-01 | 1.98E-02 |
| | $1/\gamma_d = 8.6$ | 0.89% | 2.67E-01 | 1.33E-02 |
| | $1/\gamma_r = 17.2$ | 0.35% | 3.01E-01 | 7.76E-03 |
| | $\bar{\pi}_d = 0.7$ | -1.24% | 9.38E-03 | 1.13E-03 |
| | | | | |
| **Weekly observations** | | | | |
| | $1/\sigma = 9$ | 3.79% | 4.99E-01 | 3.52E-02 |
| | $1/\gamma_d = 8.6$ | -1.86% | 4.85E-01 | 2.94E-02 |
| | $1/\gamma_r = 17.2$ | 0.52% | 7.13E-01 | 1.76E-02 |
| | $\bar{\pi}_d = 0.7$ | -2.17% | 2.42E-02 | 9.01E-03 |

Fig. Sup.1. Upper panels: simulated (black lines) and estimated (gray and red lines) effective reproduction numbers over 500 simulation runs with daily (left panel) and weekly (right panel) observations on the number of new cases, deaths and recoveries. Lower panels: simulated (black lines) and estimated (gray and red lines) effective reproduction numbers over 500 simulation runs with daily (left panel) and weekly (right panel) observations of new cases only.

**Effective reproduction number (daily new cases)**



**Effective reproduction number (weekly new cases)**



Fig. Sup.2. Upper panel: simulated (black lines) and estimated (gray and red lines) effective reproduction numbers over 500 simulation runs with daily new cases obtained from a patch stochastic SEIR-D model considering two regions. Lower panels: simulated (black lines) and estimated (gray and red lines) effective reproduction numbers over 500 simulation runs with weekly new cases obtained from a patch stochastic SEIR-D model considering two regions.

Table Sup.3. Estimates (MAP) of the ODE-parameters $\sigma^{-1}, \gamma_d^{-1}, \gamma_r^{-1}$ and $\pi_d$ for different values of $E(t=0)$.

| $E(t=0)$ | $1/\sigma$ | $1/\gamma_d$ | $1/\gamma_r$ | $\pi_d$ |
|---|---|---|---|---|
| 1 | 9.100 | 8.858 | 17.009 | 0.690 |
| 3 | 8.904 | 9.001 | 16.950 | 0.697 |
| 5 | 8.478 | 8.306 | 17.242 | 0.697 |
| 7 | 8.985 | 8.571 | 17.250 | 0.679 |
| 9 | 8.859 | 8.673 | 17.305 | 0.689 |
| 11 | 8.986 | 8.583 | 17.147 | 0.689 |
| 13 | 8.916 | 8.741 | 17.261 | 0.690 |
| 15 | 8.966 | 8.420 | 17.116 | 0.660 |
| 17 | 8.957 | 8.367 | 17.097 | 0.698 |
| 19 | 8.986 | 8.400 | 17.092 | 0.701 |

Table Sup.4. Descriptive summaries of the posterior samples of parameters $\boldsymbol{\alpha}$, $\phi_I$, $\sigma^{-1}$, $\pi_d$, $\gamma_d^{-1}$, $\gamma_r^{-1}$, $\lambda$ and $\xi$ in the SEIR-D model.

| | Mean | SD | Quantiles | | |
|---|---|---|---|---|---|
| | | | 2.5% | 50% | 97.5% |
| $\phi_I$ | 65.817 | 19.080 | 29.653 | 66.082 | 97.765 |
| $1/\sigma$ | 9.310 | 0.558 | 8.219 | 9.306 | 10.417 |
| $\pi_d$ | 0.689 | 0.021 | 0.645 | 0.689 | 0.731 |
| $1/\gamma_d$ | 8.552 | 0.617 | 7.339 | 8.550 | 9.777 |
| $1/\gamma_r$ | 17.203 | 0.742 | 15.747 | 17.204 | 18.676 |
| $\lambda$ | 0.847 | 0.463 | 0.201 | 0.763 | 1.968 |
| $\xi$ | 1.656 | 2.436 | 0.029 | 0.907 | 7.769 |

Fig. Sup.3. Estimated effective reproduction numbers for different values of the initial number $E(t = 0)$ of exposed subjects in the population (lower panel).

Fig. Sup.4. Estimates of the cumulative number of cases and deaths and of the effective reproduction number when using the method of Althaus (2014) on data available up to three different calendar times.
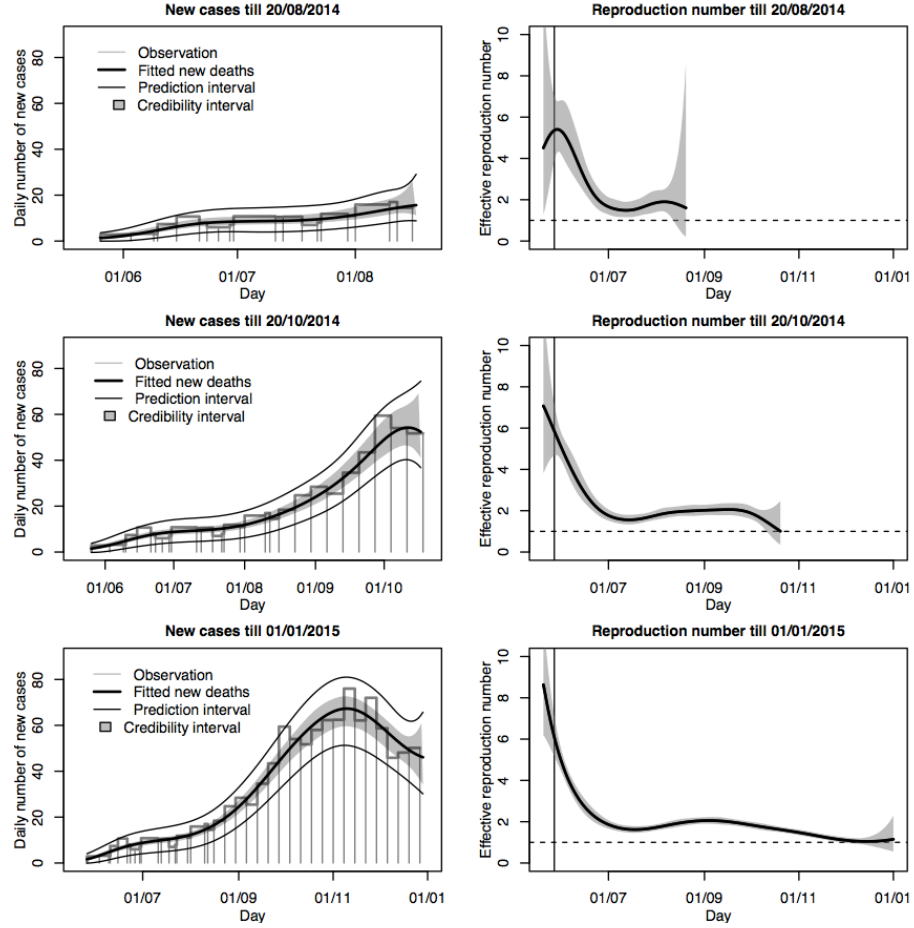
Fig. Sup.5. Estimates of the posterior mean of the expected number of new cases and of the effective reproduction number (with pointwise 95% credible region and prediction interval) from data available up to three calendar times.

Fig. Sup.6. Effective reproduction number estimated for complete number of new cases (100 simulations). Each sub-figure represent a component of the $R_e(t)$ computed as in Section 2.3.

Fig. Sup.7. Effective reproduction number $\mathcal{R}_e(t)$ estimated for time-varying under-reporting of the numbers of new cases (100 simulations). Each sub-figure reports the estimate of one factor entering the formula of $R_e(t)$, see Section 2.6 in the main manuscript.

Fig. Sup.8. Estimated ODE-parameter for complete and under-reported number of new cases. Each wide horizontal lines indicates the true parameter value.
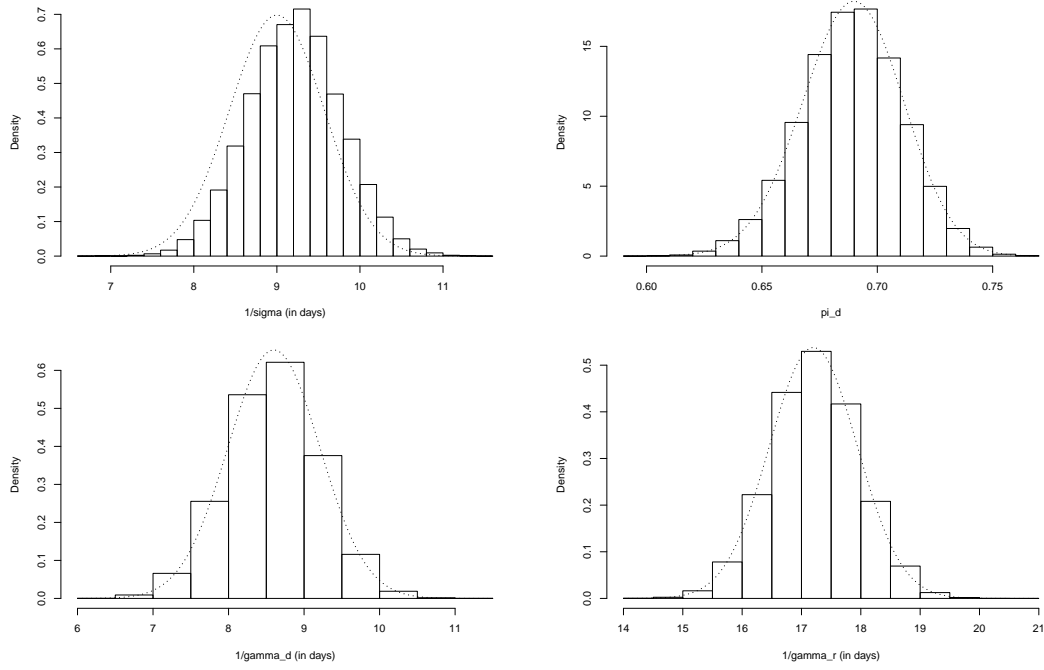
Fig. Sup.9.  Posterior distribution of the spline coefficients.

Fig. Sup.10. Priors (dotted lines) and histograms of the sampled $\sigma^{-1}$, $\pi_d$, $\gamma_d$, $\gamma_r$ parameters.
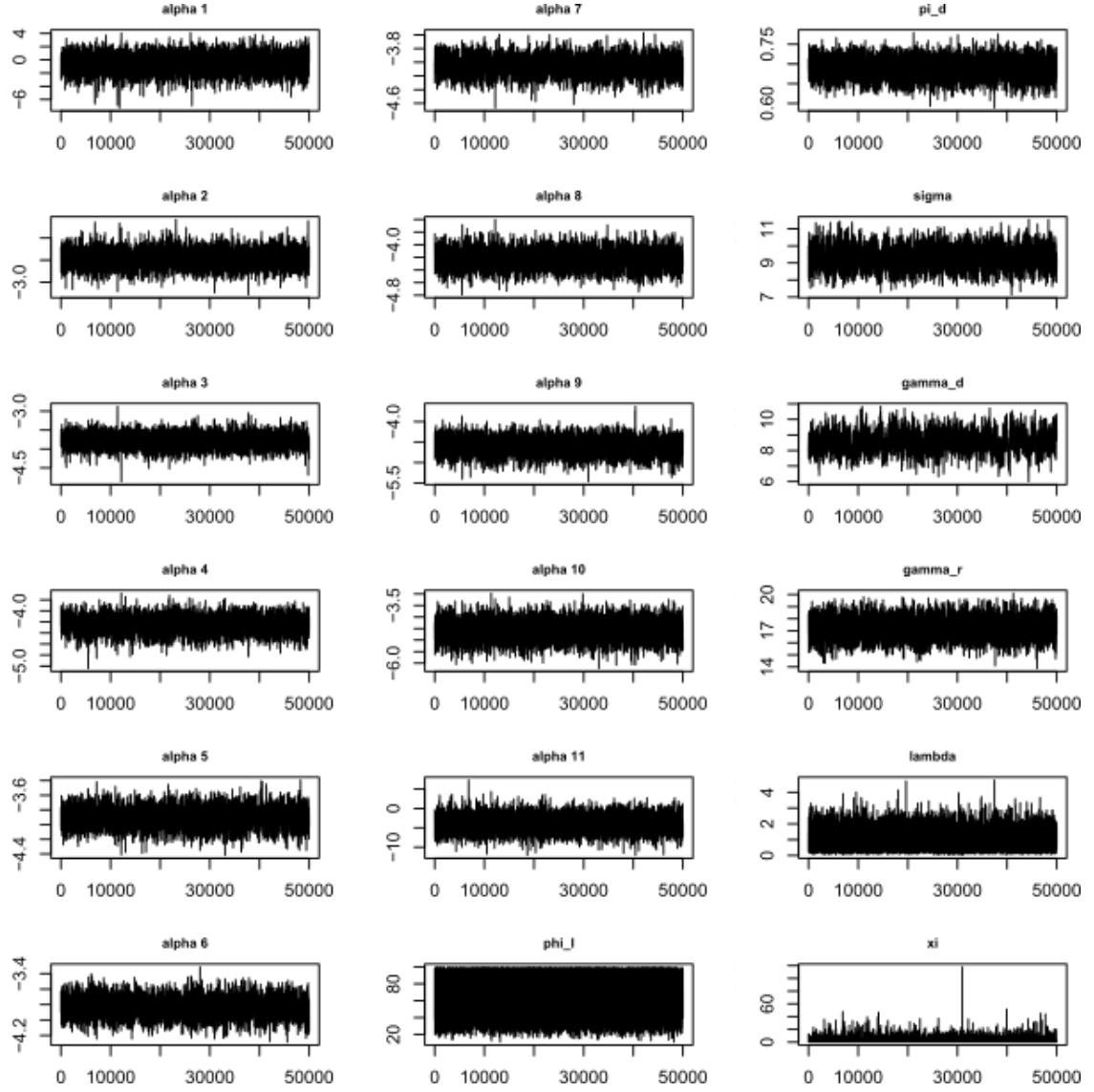
Fig. Sup.11.  Trace plots for P-spline coefficients and ODE parameters.
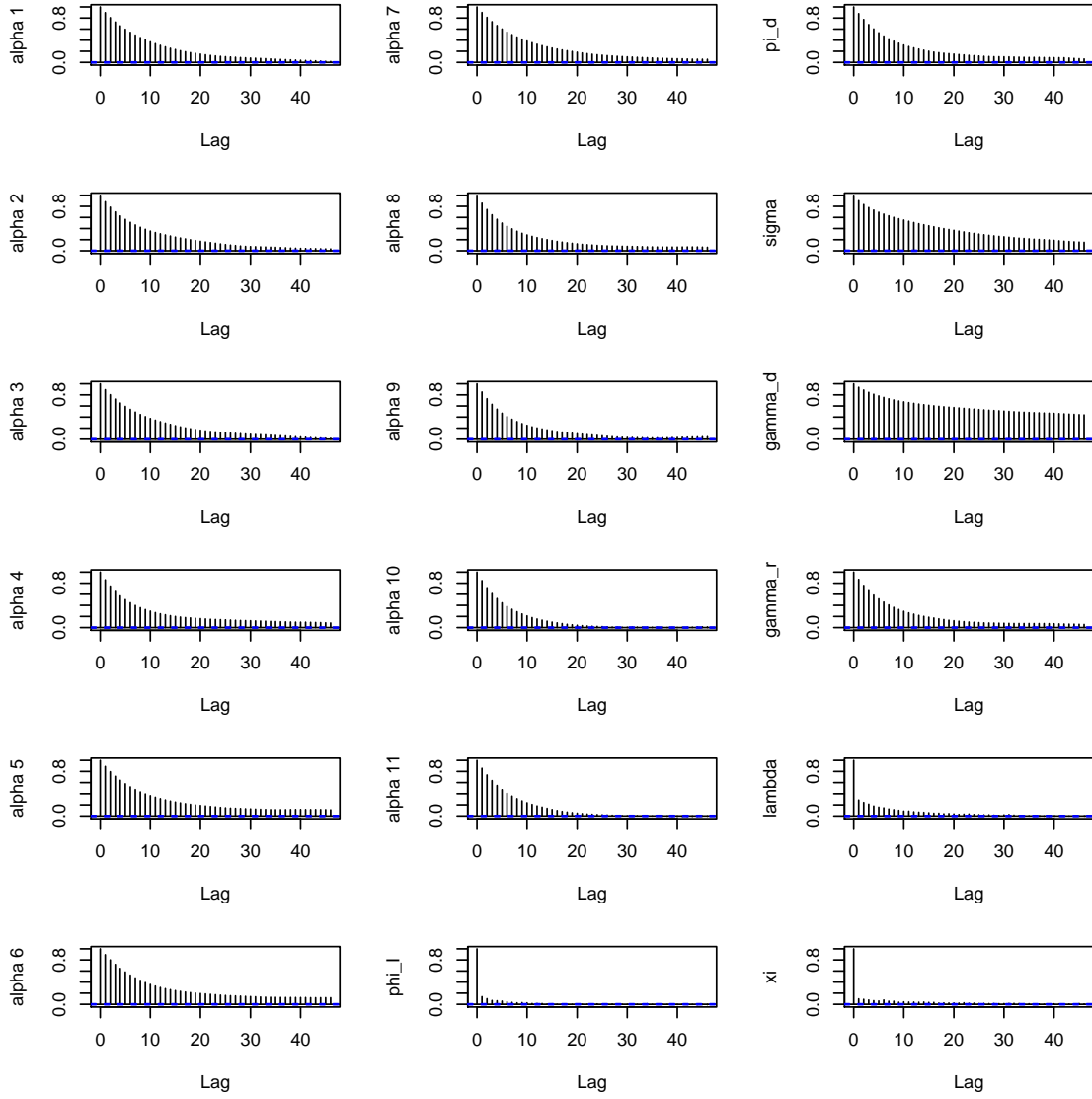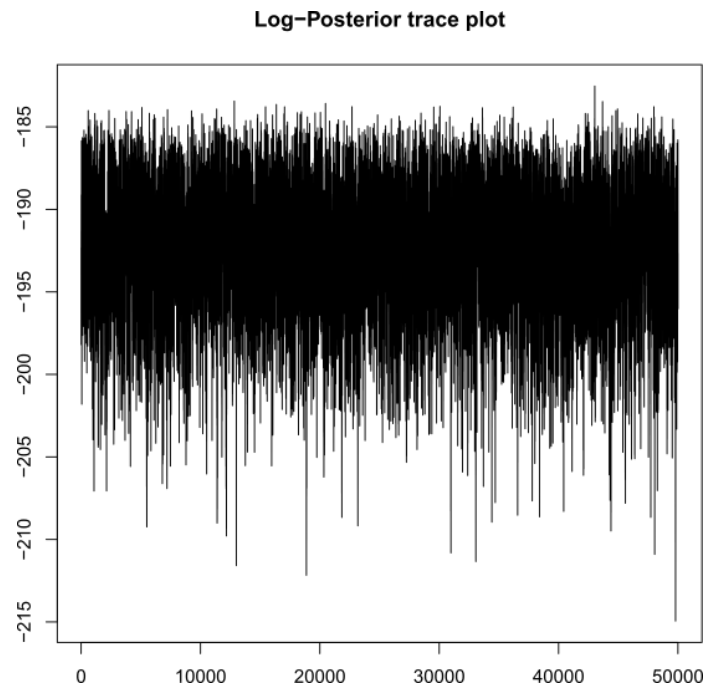
Fig. Sup.12. ACF plots for P-splines and ODE parameters.

Fig. Sup.13. Trace plot of the log-posterior density.