



HAL
open science

Hierarchical Bayesian Estimation of COVID-19 Reproduction Number

Patrice Abry, Juliette Chevallier, Gersende Fort, Barbara Pascal

► **To cite this version:**

Patrice Abry, Juliette Chevallier, Gersende Fort, Barbara Pascal. Hierarchical Bayesian Estimation of COVID-19 Reproduction Number. 2024. hal-04695138

HAL Id: hal-04695138

<https://hal.science/hal-04695138v1>

Preprint submitted on 12 Sep 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Hierarchical Bayesian Estimation of COVID-19 Reproduction Number

Patrice Abry
*Laboratoire de Physique
CNRS, ENS Lyon, France
patrice.abry@ens-lyon.fr*

Juliette Chevallier
*IMT, UMR 5219, INSA Toulouse
Université de Toulouse, France
firstname.lastname@math.univ-toulouse.fr*

Gersende Fort
*IMT, UMR 5219, CNRS
Université de Toulouse, France*

Barbara Pascal
*Nantes Université, ECN
CNRS, LS2N, France
barbara.pascal@cnrs.fr*

Abstract—Assessing the intensity of an epidemic, such as the COVID-19 pandemic, during the epidemic outbreak, constitutes a significant technical challenge with high societal stakes. Elaborating on classical epidemiological models, this work aims to define a hierarchical Bayesian model that permits the robust estimation of the temporal evolution of the pandemic intensity despite highly corrupted daily new infection counts. It also outputs uncertainty assessment, in the form of credibility intervals robust to the priors choice, accounting for uncertainties on model parameters. The estimation is performed by carefully designed Monte Carlo samplers. The relevance of the proposed estimation procedure is illustrated on real COVID-19 pandemic data for several countries and periods, made available from the Johns Hopkins University repository.

Index Terms—Bayesian statistics, Epidemiology, COVID-19, Monte Carlo Sampling

I. INTRODUCTION

Context. The COVID-19 pandemic stroke heavily societies worldwide and over the long term. Assessing pandemic intensity from highly corrupted daily new infection counts turned critical, notably for designing sanitary policies and countermeasures [1], [2], and prompted several statistical signal processing approaches, e.g., [3]–[5]. This work aims to address this technical challenge by investigating the relevance of hierarchical Bayesian procedures for the estimation of the pandemic intensity temporal evolution, robust to both the limited quality of available pandemic data and uncertainties in model parameters, while providing credibility intervals.

Related works. Compartmental models, such as the classical Susceptible-Infectious-Recovered reference, are uneasy to use during epidemic outbreaks, as they heavily rely on accurate information matching precisely social realities (social groups, contact, etc.) [6]–[11]. Instead, during epidemic outbreaks or within active epidemic phases, the time-varying reproduction number, R_t , is considered a relevant proxy for quantifying pandemic intensity [12]–[15]. Classical epidemiological models rely on time-varying Poisson distributions or negative binomial distributions for daily new infection counts [14], [16], [17]. Though, some of these models yield estimates, such as the straightforward maximum likelihood estimators derived in [16], that may be statistically non-consistent and might not provide confidence assessment. In addition, they may prove insufficiently robust given the limited quality (missing counts, misreports, pseudo-seasonalities, etc.) of available epidemic data [14]. The Bayesian-inspired estimator proposed in [14, Web appendix 1] significantly gained robustness by enforcing a smooth temporal behavior, making it an essential reference for viral epidemic monitoring, but at the price of a lack of statistical soundness. Recently, both consistency and robustness against the low quality of COVID-19 data have been achieved by solving inverse problems stemming from regularizing the likelihood model, yet requiring arbitrarily parameter selection and lacking

confidence assessment [18], [19]. Another approach that also provides confidence assessment in estimation has been intended via Bayesian models, at the price of arbitrary a priori assumptions on the pandemic temporal dynamics and model parameters choice [17], [20], [21]. To explore the intractable a posteriori distribution, Bayesian samplers have been derived in [21]–[23], still based on arbitrarily chosen model parameters. This allows to derive an automated data-driven model parameter selection through the Expectation-Maximization (EM) algorithm [24]. Achieving jointly parameters selection with uncertainty quantification is the main challenge addressed here.

Goals, contributions and outline. The overarching goal of this work is to study the relevance of the hierarchical Bayesian framework, proposed to perform pandemic intensity estimation, with jointly confidence assessment and robustness against, at the same time, available low-quality pandemic data and unknown model parameters. To that end, Section II briefly recalls the classical epidemiological models and estimators, serving as a base for this work. Then, Section III devises, as a first contribution, the proposed hierarchical model, detailing the constraints ensuring robustness to low-quality data as well as the choice of priors enabling robustness to unknown model parameters. In Section IV, as a second contribution, the proposed strategy is illustrated at work on real COVID-19 data, extracted from the Johns Hopkins University repository, for several countries. The convergence of the resulting Monte Carlo algorithm exploiting the proposed model is illustrated and its outcome and performance are compared against other estimation procedures proposed in the literature, that did not ensure robustness against arbitrary choices pertaining to unknown model parameters. Matlab codes implementing the proposed robust estimation procedure will be shared publicly at the time of publication.

II. COVID-19 EPIDEMIOLOGICAL MODEL

To model the propagation of viral epidemics, such as the COVID-19 pandemic, it has been proposed in [14] to consider the daily new infection count at day t , Z_t , as a Poisson random variable conditionally to past counts $\mathbf{Z}_{1:t-1} := (Z_1, \dots, Z_{t-1})$, and to initial values $(Z_{-\tau_\phi+1}, \dots, Z_0)$. The conditional distribution writes:

$$Z_t \mid \mathbf{Z}_{-\tau_\phi+1:t-1}; R_t \sim \mathcal{P}(R_t \Phi_t^Z). \quad (1)$$

The time-varying Poisson parameter, $R_t \Phi_t^Z$, combines the reproduction number at day t , R_t , defined as the expected number of secondary infections stemming from one typical contagious individual, with $\Phi_t^Z := \sum_{s=1}^{\tau_\phi} \phi(s) Z_{t-s}$ that accounts for the *global infectiousness* in the population, computed as a weighted sum of past counts involving the *serial interval distribution* ϕ , modeling the random delay between primary and secondary infections. For the COVID-19 pandemic, we model ϕ as a Gamma distribution of mean and standard deviation

of respectively 6.6 and 3.5 days [25], [26], with temporal support truncated to $\tau_\phi = 25$. Fig. 2 (first and third rows, solid and dashed black curves) displays counts $\mathbf{Z}_{1:t-1}$ and global infectiousness Φ_t^Z for different countries and pandemic stages.

Maximum Likelihood Estimator (MLE) of Model (1). A straightforward estimator of the reproduction number R_t is obtained by maximizing the likelihood of $\mathbf{Z} := \mathbf{Z}_{1:T}$ associated to the model (1):

$$\mathbf{R}_t^{\text{MLE}^{(1)}} := Z_t / \Phi_t^Z \text{ if } \Phi_t^Z > 0, \quad \mathbf{R}_t^{\text{MLE}^{(1)}} := 0 \text{ otherwise;} \quad (2)$$

see [16]. Fig 2 (second and fourth rows) displays $\mathbf{R}^{\text{MLE}^{(1)}}$ (gray curve). First, $\mathbf{R}^{\text{MLE}^{(1)}} := \mathbf{R}_{1:T}^{\text{MLE}^{(1)}}$ is not consistent as there are as many unknown parameters as observations. Second, the low quality of COVID-19 data induces far too large and rapid fluctuations over time, making it unusable by epidemiologists [14].

Reference epidemiological estimator. To tackle robustness against low quality data, the widely used EpiEstim software¹ implements the Bayesian-inspired framework devised in [14, Web appendix 1]. To estimate each R_t , a *local* Poisson likelihood model is introduced, involving $R_{t,\tau}$, assumed constant over τ days, together with a conjugate Gamma prior $\Gamma(a, 1/b)$; $\Gamma(\alpha, \beta)$ denotes a Gamma distribution with shape parameter α and rate parameter β . Hence, the a posteriori distribution of R_t is a Gamma distribution with explicit parameters. This yields a point estimate defined as the a posteriori mean:

$$\mathbf{R}_{t,\tau}^{\text{EpiEstim}} := \frac{\sum_{s=t-\tau+1}^t Z_s + a}{\sum_{s=t-\tau+1}^t \Phi_s^Z + 1/b}, \quad (3)$$

together with credibility intervals (CIs). Examples of $t \mapsto \mathbf{R}_{t,\tau}^{\text{EpiEstim}}$, with the classical choice of $\tau = 7$ days, and CIs are displayed in Fig. 2 (second and fourth rows, solid green curves, and light green areas). They are smoother than $\mathbf{R}^{\text{MLE}^{(1)}}$, and hence more acceptable from an epidemiological point of view. The main bottleneck in the practical use of EpiEstim for emerging epidemics is the choice of the parameters a, b and τ . In the literature, they are chosen by experts with $a \times b$ *large enough* to prevent underestimation of the reproduction number at the onset of an epidemic [14]. Furthermore, assuming a *locally* constant model excludes rigorous statistical ground. In the following, this estimator is referred to as $\mathbf{R}^{\text{EpiEstim}} := (\mathbf{R}_{1,\tau}^{\text{EpiEstim}}, \dots, \mathbf{R}_{T,\tau}^{\text{EpiEstim}})$, omitting the dependence in τ .

III. ROBUST HIERARCHICAL BAYESIAN MODEL

This section aims to derive a novel Bayesian statistical model and to explain how the a posteriori distribution of $\mathbf{R} := \mathbf{R}_{1:T}$ is numerically investigated. This new, hierarchical, model will provide uncertainty quantification while taking into account the self-selection of some parameters of the model. Point estimates and credibility intervals will be derived from empirical expectations and empirical quantiles of the derived a posteriori distributions.

Throughout this section $\pi(\mathbf{X}|\mathbf{Y})$ stands for the density distribution of the random vector \mathbf{X} given the random vector \mathbf{Y} , with respect to the Lebesgue measure.

Penalized Poisson model. To increase robustness against the low-quality data, it was proposed in [18], [19], [21] to model daily new infection counts Z_t as a Poisson distribution of intensity $R_t \Phi_t^Z + O_t$, where O_t denotes errors, unknown by construction. For consistency, we set $\mathbf{O} := \mathbf{O}_{1:T}$ hereafter.

Following [24], we consider the extended model where, on the set $(R_t, O_t) \in \mathcal{D}_t := \{(R_t, O_t) : R_t \geq 0 \text{ and } R_t \Phi_t^Z + O_t \geq 0\}$

and conditionally to R_t, O_t , past counts $\mathbf{Z}_{1:t-1}$ and initial values $\mathcal{I} := (R_{-1}, R_0, Z_{-\tau_\phi+1}, \dots, Z_0)$, Z_t follows a Poisson distribution:

$$Z_t | R_t, O_t, \mathbf{Z}_{1:t-1}, \mathcal{I} \sim \mathcal{P}\left(R_t \Phi_t^Z + O_t\right); \quad (4)$$

by convention $\mathcal{P}(0)$ is a Dirac mass at 0. Conditionally to the past, R_t and O_t are independent. The R_t are modeled through a second order autoregressive process with Laplace noise:

$$\pi(R_t | R_{t-1}, R_{t-2}, \lambda_R) \propto \lambda_R e^{-\frac{\lambda_R}{4} |R_t - 2R_{t-1} + R_{t-2}|}. \quad (5)$$

Such an a priori distribution favors limited changes in the (second order differences) of $t \mapsto R_t$ [18], [19], [21]. The O_t are assumed to be independent and distributed according to a Laplace distribution:

$$\pi(O_t | \lambda_O) \propto \lambda_O e^{-\lambda_O |O_t|}. \quad (6)$$

Denote $\underline{\lambda} := (\lambda_R, \lambda_O)$ and $\mathbf{D} \in \mathbb{R}^{T \times T}$ the discrete-time second order derivative matrix acting on $\mathbf{R} \in \mathbb{R}^T$ as

$$\forall t \in \{2, \dots, T\}, \quad (\mathbf{DR})_t := \frac{1}{4} (R_t - 2R_{t-1} + R_{t-2}), \quad (7)$$

with initial conditions $(\mathbf{DR})_1 := 1/4 R_1$, $(\mathbf{DR})_2 := 1/4 R_2 - 1/2 R_1$. This description implies that the logarithm of the joint distribution of (\mathbf{R}, \mathbf{O}) given $(\mathbf{Z}, \mathcal{I}, \underline{\lambda})$ is, up to an additive constant,

$$\begin{aligned} \ln \pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda}) &= -\lambda_R \|\mathbf{DR} + \delta\| - \lambda_O \|\mathbf{O}\| + T \ln \lambda_R \\ &+ T \ln \lambda_O - \sum_{t=1}^T ((R_t \Phi_t^Z + O_t) - Z_t \ln(R_t \Phi_t^Z + O_t)), \end{aligned} \quad (8)$$

where $4\delta := (R_{-1} - 2R_0, R_0, 0, \dots, 0)^\top \in \mathbb{R}^T$ and where $\|\cdot\|$ denotes the 1-norm. Estimators $\hat{\mathbf{R}}_t^{\text{Mean}^{(8)}}$ and $\hat{\mathbf{O}}_t^{\text{Mean}^{(8)}}$ are defined as the *empirical* means of the distribution $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda})$ defined in (8), which is intractable. A Metropolis-within-Gibbs sampler (See [27, Section 10.3]) using a proximal Langevin proposal mechanism targeting $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda})$ has been designed in [21] (See also [22], [23]). Hereafter, a one-step iteration of this Gibbs sampler starting from (r, o) is denoted $\text{PGdual}(r, o; \underline{\lambda})$.

Understanding the role of parameters $\underline{\lambda}$. Eq. (8) suggests that $\underline{\lambda}$ plays a balancing role between a data fidelity term and penalty terms on \mathbf{R} and \mathbf{O} . Therefore, particular attention must be paid to the choice of $\underline{\lambda}$. The function $(\mathbf{R}, \mathbf{O}) \mapsto \ln \pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda})$ is concave and admits a maximizer $(\mathbf{R}^{\text{MAP}^{(8)}}, \mathbf{O}^{\text{MAP}^{(8)}})$ under conditions [19], [21], which is the maximum of the a posteriori distribution (8) (MAP). It satisfies from the Fermat's rule [28, Theorem 16.3]:

$$\forall t, \quad \left| \frac{Z_t}{\mathbf{R}_t^{\text{MAP}^{(8)}} \Phi_t^Z + \mathbf{O}_t^{\text{MAP}^{(8)}}} - 1 \right| \leq \min\left(\frac{\lambda_R}{\Phi_t^Z}, \lambda_O\right). \quad (9)$$

Upon noting that the MLE of model (4) satisfies for all t

$$Z_t = \mathbf{R}_t^{\text{MLE}^{(4)}} \Phi_t^Z + \mathbf{O}_t^{\text{MLE}^{(4)}},$$

Eq. (9) shows that λ_O and λ_R / Φ_t^Z quantify how much $(\mathbf{R}^{\text{MAP}^{(8)}}, \mathbf{O}^{\text{MAP}^{(8)}})$ is allowed to differ from $(\mathbf{R}^{\text{MLE}^{(4)}}, \mathbf{O}^{\text{MLE}^{(4)}})$, thus avoiding overfitting and providing consistency of the estimators. Furthermore, small values of $\underline{\lambda}$ imply that the a priori distributions for O_t and $(\mathbf{DR})_t$ are Laplace distributions with large variance. Such a choice allows a broader range of possible values for O_t and the second order derivatives of $t \mapsto R_t$. Finally, Eq. (9) also suggests that λ_O and λ_R / Φ_t^Z scale one with the other. Examples of $\mathbf{R}^{\text{MAP}^{(8)}}$ and corrected counts $\mathbf{Z} - \mathbf{O}^{\text{MAP}^{(8)}}$, with selected parameters fixed to (See [19], [21])

$$\lambda_R = 3.5 \times \text{std}(\mathbf{Z}), \quad \lambda_O = 50 \times 10^{-3}, \quad (10)$$

¹<https://github.com/mrc-ide/EpiEstim>

Algorithm 1 MCMC sampler targeting $\pi(\mathbf{R}, \mathbf{O}, \underline{\lambda} | \mathbf{Z}, \mathcal{I}; \theta)$ (See (12))

Input: $\mathbf{Z}, \Phi^{\mathbf{Z}}, \mathcal{I}$; **Parameters:** $(\alpha_{\mathbf{R}}, \beta_{\mathbf{R}}, \alpha_{\mathbf{O}}, \beta_{\mathbf{O}}), k_{\max}$
Initialize: $(\mathbf{R}^{(0)}, \mathbf{O}^{(0)}, \underline{\lambda}^{(0)})$
for $k = 0, 1, \dots, k_{\max} - 1$ **do**
 # Sample \mathbf{R}, \mathbf{O} at fixed $\underline{\lambda}^{(k)}$
 $\mathbf{R}^{(k+1)}, \mathbf{O}^{(k+1)} \sim \text{PGdual}(\mathbf{R}^{(k)}, \mathbf{O}^{(k)}; \underline{\lambda}^{(k)})$
 # Sample $\underline{\lambda}$ at fixed $\mathbf{R}^{(k+1)}, \mathbf{O}^{(k+1)}$
 $\lambda_{\mathbf{R}}^{(k+1)} \sim \Gamma(T + \alpha_{\mathbf{R}}, \|\mathbf{DR}^{(k+1)} + \delta\| + \beta_{\mathbf{R}})$
 $\lambda_{\mathbf{O}}^{(k+1)} \sim \Gamma(T + \alpha_{\mathbf{O}}, \|\mathbf{O}^{(k+1)}\| + \beta_{\mathbf{O}})$
end for
Output: $\{\mathbf{R}^{(k)}, \mathbf{O}^{(k)}, \underline{\lambda}^{(k)}\}_{k=1, \dots, k_{\max}}$

$\text{std}(\mathbf{Z})$ denoting the standard deviation of the observations \mathbf{Z} , are provided for different countries and pandemic stages in Fig. 2 (second and fourth rows, yellow curves, most of the time covered by red curves). The computation of this MAP relies on a primal-dual algorithm, suited to nonsmooth convex objective functions [19].

A novel hierarchical Bayesian model. In [24], a Stochastic Approximation of the EM algorithm was derived to perform an automated data-driven selection of $\underline{\lambda}$. Yet, estimating \mathbf{R} and \mathbf{O} given a fixed value of $\underline{\lambda}$ does not yield robustness quantification with respect to this choice. Instead, the novel model proposed below addresses this issue. Since the Gamma distribution is the conjugate prior for the Laplace distribution [29], we assume Gamma prior on both $\lambda_{\mathbf{R}}$ and $\lambda_{\mathbf{O}}$: $\Gamma(\alpha_{\mathbf{R}}, \beta_{\mathbf{R}})$ and $\Gamma(\alpha_{\mathbf{O}}, \beta_{\mathbf{O}})$. Set $\theta := (\alpha_{\mathbf{R}}, \beta_{\mathbf{R}}, \alpha_{\mathbf{O}}, \beta_{\mathbf{O}})$. This yields a novel a posteriori distribution for \mathbf{R} and \mathbf{O} , denoted $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}; \theta)$ and which satisfies:

$$\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}; \theta) = \int_{\mathbb{R}^+ \times \mathbb{R}^+} \pi(\mathbf{R}, \mathbf{O}, \underline{\lambda} | \mathbf{Z}, \mathcal{I}; \theta) d\underline{\lambda}, \quad (11)$$

where, on the set $\bigcap_{t=1}^T \mathcal{D}_t \times \mathbb{R}^+ \times \mathbb{R}^+$,

$$\begin{aligned} \ln \pi(\mathbf{R}, \mathbf{O}, \underline{\lambda} | \mathbf{Z}, \mathcal{I}; \theta) &:= \ln \pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda}) \\ &- \beta_{\mathbf{R}} \lambda_{\mathbf{R}} + (\alpha_{\mathbf{R}} - 1) \ln \lambda_{\mathbf{R}} - \beta_{\mathbf{O}} \lambda_{\mathbf{O}} + (\alpha_{\mathbf{O}} - 1) \ln \lambda_{\mathbf{O}}. \end{aligned} \quad (12)$$

Novel point estimates and CIs are defined as expectations and quantiles of the distribution $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}; \theta)$ (See (11)). Unfortunately, this distribution is intractable: we derive hereafter a novel Markov Chain Monte Carlo (MCMC) sampler (See Alg. 1) and replace exact expectations and quantiles by their empirical counterpart. We denote the new point estimates by $\hat{\mathbf{R}}^{\text{Mean}^{(H)}}$ and $\hat{\mathbf{O}}^{\text{Mean}^{(H)}}$.

Sampling from $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}; \theta)$ with MwG-Hierarchical. This distribution (11) is not explicit, and must be approximated by points $(\mathbf{R}^{(k)}, \mathbf{O}^{(k)})_k$, produced by a Monte-Carlo sampler. From Eq. (11), such points are obtained from the first two components of $(\mathbf{R}^{(k)}, \mathbf{O}^{(k)}, \underline{\lambda}^{(k)})_k$ themselves produced by a Metropolis-within-Gibbs sampler targeting $\pi(\mathbf{R}, \mathbf{O}, \underline{\lambda} | \mathbf{Z}, \mathcal{I}; \theta)$. A draw approximating the first conditional distribution $\pi(\mathbf{R}, \mathbf{O} | \mathbf{Z}, \mathcal{I}, \underline{\lambda})$, is obtained from a call to PGdual. From (8) and (12), the second conditional distribution $\pi(\underline{\lambda} | \mathbf{Z}, \mathbf{R}, \mathbf{O}, \mathcal{I}; \theta)$, are independent Gamma distributions:

$$\Gamma(T + \alpha_{\mathbf{R}}, \|\mathbf{DR} + \delta\| + \beta_{\mathbf{R}}) \quad \text{and} \quad \Gamma(T + \alpha_{\mathbf{O}}, \|\mathbf{O}\| + \beta_{\mathbf{O}}). \quad (13)$$

The sampling algorithm is sketched in Alg. 1 and hereafter, it is named MwG-Hierarchical.

IV. REPRODUCTION NUMBER ESTIMATION

COVID-19 data. From the early stages of the COVID-19 pandemic, the Johns Hopkins University set up an impressive reposi-

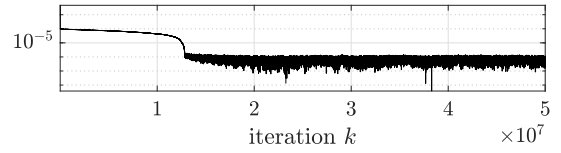


Fig. 1: **Convergence of Alg. 1** on French dataset monitoring $(\max \log \pi - \pi^{(k)}) / \pi^{(0)}$ where $\pi^{(k)} := \pi(\mathbf{R}^{(k)}, \mathbf{O}^{(k)}, \underline{\lambda}^{(k)} | \mathbf{Z}, \mathcal{I}; \theta)$ (See Eq. (12)). The maximum is computed across three Markov chains generated by Alg. 1 with three different initializations.

tory² by collecting on a daily basis new infection counts reported by National Health Authorities worldwide and organizing them to be shared with the research community. In the present work, we make use of these data for several different countries and periods of time, spanning $T = 70$ days, and chosen as representative of various pandemic phases: India (2020/08/09-2020/10/17), Germany (2021/02/21-2021/05/01), France (2022/02/20-2022/04/30), South Korea (2022/06/22-2022/08/30).

MCMC algorithms setup. For each dataset, considering the standard choice $\tau = 7$ days, the reported infection counts during the two days before the monitored time period are temporarily included to compute $\mathbf{R}_{-1, \tau}^{\text{EpiEstim}}, \mathbf{R}_{0, \tau}^{\text{EpiEstim}}$, so that PGdual and MwG-Hierarchical use as past values (stored in \mathcal{I} in Alg. 1):

$$\mathbf{R}_{-1, \tau} = \mathbf{R}_{-1, \tau}^{\text{EpiEstim}}, \quad \mathbf{R}_{0, \tau} = \mathbf{R}_{0, \tau}^{\text{EpiEstim}}.$$

Setting $k_{\max} = 5 \times 10^7$, $\hat{\mathbf{R}}^{\text{Mean}^{(S)}}, \hat{\mathbf{O}}^{\text{Mean}^{(S)}}$ (resp. $\hat{\mathbf{R}}^{\text{Mean}^{(H)}}, \hat{\mathbf{O}}^{\text{Mean}^{(H)}}$) are estimated from k_{\max} iterations of PGdual (resp. by running the MwG-Hierarchical Alg. 1 for k_{\max} iterations). Then, in both cases, the empirical means and quantiles are computed after discarding a burnin phase of $3 \cdot 10^7$ iterations, ensuring the convergence of the Markov chains (See Fig. 1 for MwG-Hierarchical Alg. 1). The \mathbf{R}, \mathbf{O} Markov chains are initialized at:

$$\mathbf{R}_t^{(0)} := (\mathbf{R}_{t, \tau}^{\text{EpiEstim}} + 1) / 2, \quad \mathbf{O}_t^{(0)} := (\mathbf{Z}_t - \mathbf{R}_{t, \tau}^{\text{EpiEstim}} \Phi_t^{\mathbf{Z}}) / 2,$$

which constitutes an average between the EpiEstim estimate and the raw initialization $\mathbf{R}_t = 1, \mathbf{O}_t = 0$ used in [21]. The design parameters of PGdual are chosen as in [21] and the estimates $\hat{\mathbf{R}}^{\text{Mean}^{(S)}}, \hat{\mathbf{O}}^{\text{Mean}^{(S)}}$ are computed with a fixed $\underline{\lambda}$, set as in Eq. (10). For MwG-Hierarchical, the hyperparameters θ are set so that

$$\alpha_{\mathbf{R}} / \beta_{\mathbf{R}} = 3.5 \times \text{std}(\mathbf{Z}), \quad \alpha_{\mathbf{O}} / \beta_{\mathbf{O}} = 50 \times 10^{-3},$$

thus imposing that the expectations of the Gamma priors on $\underline{\lambda}$ coincide with the values of (10), but letting the MwG-Hierarchical sampling scheme explore around these values by fixing the standard deviations of the priors to a preset percentage of their expectations:

$$\sqrt{\alpha_{\mathbf{R}}} / \beta_{\mathbf{R}} = 0.02 \times \alpha_{\mathbf{R}} / \beta_{\mathbf{R}}, \quad \sqrt{\alpha_{\mathbf{O}}} / \beta_{\mathbf{O}} = 0.015 \times \alpha_{\mathbf{O}} / \beta_{\mathbf{O}}.$$

$\underline{\lambda}$ is initialized at the historical parameters (See Eq. (10)).

Robust epidemiological indicator estimations. Fig. 2 compares, for all countries and periods, the reproduction number point estimates $\mathbf{R}^{\text{MLE}^{(1)}}, \mathbf{R}^{\text{MAP}^{(S)}},$ and $\mathbf{R}^{\text{EpiEstim}}, \hat{\mathbf{R}}^{\text{Mean}^{(S)}}, \hat{\mathbf{R}}^{\text{Mean}^{(H)}}$ accompanied with 95% CIs, and the corrected counts point estimates $\mathbf{Z} - \hat{\mathbf{O}}^{\text{MAP}^{(S)}},$ and $\mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(S)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(H)}}$ with 95% CIs. The $(\mathbf{R}^{\text{MAP}^{(S)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{MAP}^{(S)}})$ estimates (yellow) are barely visible as they superimpose almost perfectly with the $(\hat{\mathbf{R}}^{\text{Mean}^{(S)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(S)}})$ estimates obtained with PGdual (blue). This shows that the PGdual sampling scheme complements the estimates $\mathbf{R}^{\text{MAP}^{(S)}},$ and $\mathbf{O}^{\text{MAP}^{(S)}}$ with CIs at fixed

²<https://coronavirus.jhu.edu/>

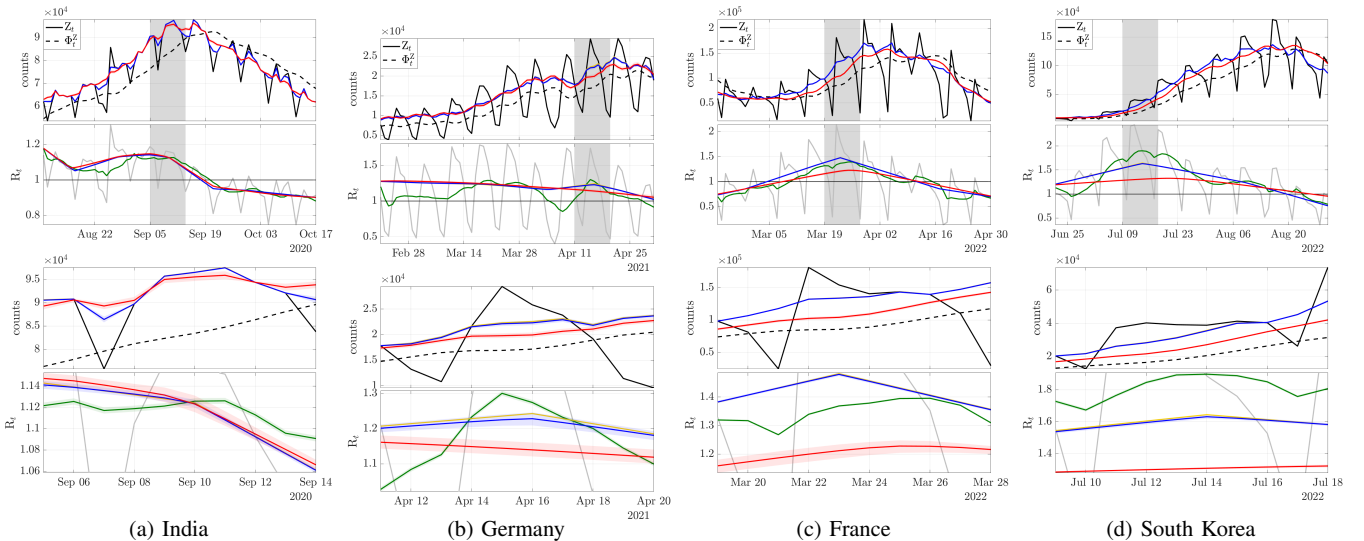


Fig. 2: **Compared reproduction number estimations:** $\mathbf{R}^{\text{MLE}^{(1)}}$ (Eq. (2), gray); $\mathbf{R}^{\text{EpiEstim}}$ and 95% CIs (Eq. (3), green); $\mathbf{R}^{\text{MAP}^{(8)}}$ and $\mathbf{Z} - \mathbf{O}^{\text{MAP}^{(8)}}$ (Eq. (8), yellow); $\hat{\mathbf{R}}^{\text{Mean}^{(8)}}$, $\mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(8)}}$ and 95% CIs (Eq. (8), blue); $\hat{\mathbf{R}}^{\text{Mean}^{(H)}}$, $\mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(H)}}$ and 95% CIs (Eq. (12), red). First and third rows: COVID-19 daily new infection counts and infectiousness (solid and dashed black curves respectively). Second and fourth rows: estimated reproduction number. Top rows: $T = 70$ -day whole period. Bottom rows: zoom on the ten-day period shaded in gray on top rows.

Mean; CI	$\lambda_{\mathbf{R}}/\text{std}(\mathbf{Z})$	$\lambda_{\mathbf{O}} \times 10^3$
India	2.6; [2.5, 2.7]	19.6; [19.1, 20.2]
Germany	3.4; [3.3, 3.6]	14.7; [14.4, 15.1]
France	1.9; [1.9, 2.0]	3.2; [3.1, 3.3]
South Korea	2.8; [2.7, 2.9]	4.9; [4.8, 5.0]
Historical, Eq. (10)	3.5	50

TABLE I: Empirical a posteriori mean estimates of λ and 95% CIs, computed from the MwG-Hierarchical sampler, Alg. 1.

λ , without modifying much the pointwise estimates. For India and Germany, in which the reporting errors correspond to a limited fraction of the true counts, Figs 2a and 2b, $(\hat{\mathbf{R}}^{\text{Mean}^{(H)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(H)}})$ (red) is also in very good agreement with $(\mathbf{R}^{\text{MAP}^{(8)}}, \mathbf{Z} - \mathbf{O}^{\text{MAP}^{(8)}})$. To the contrary, when the reporting errors are of the order of magnitude of true counts, as in the France and South Korea datasets, Figs 2c and 2d, $(\hat{\mathbf{R}}^{\text{Mean}^{(H)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(H)}})$ (red) show significantly smoother dynamics than $(\hat{\mathbf{R}}^{\text{Mean}^{(8)}}, \mathbf{Z} - \hat{\mathbf{O}}^{\text{Mean}^{(8)}})$ (blue). Tab. I further reports a Monte-Carlo approximation deduced from Alg. 1 of the expectation and CIs of λ under the distribution $\pi(\mathbf{R}, \mathbf{O}, \lambda | \mathbf{Z}, \mathcal{I}; \theta)$ (See (12)), which are observed to significantly depart from the historical initial values (See (10)), thus illustrating the benefit of considering λ as a random variable with a nontrivial prior distribution, described jointly with the quantities of interest \mathbf{R}, \mathbf{O} .

Accurate uncertainty quantification. In all the examples of Fig. 2, the CIs given by the novel hierarchical model $(\mathbf{R}, \mathbf{O}) | \mathbf{Z}, \mathcal{I}; \theta$ (in red) are larger than the CIs in the EpiEstim model (in green) and the CIs in the model $(\mathbf{R}, \mathbf{O}) | \mathbf{Z}, \mathcal{I}, \lambda$ (in blue) notably for India, Germany and France in Figs. 2a, 2b and 2c, demonstrating that accounting for uncertainty on λ yields a wider exploration, leading to more accurate uncertainty assessment. In a pandemic context when drastic sanitary measures with high social and economic impact are to be decided based on estimated epidemiological indicators, avoiding underestimation of the CIs on the time-varying reproduction number R_t , and hence overconfidence, is a critical concern. To further quantify how accounting for uncertainty on λ impacts the a posteriori uncertainty on the estimation of R_t , Tab. II reports the

	EpiEstim Eq. (3)	Eq. (8)	Hierarchical Eq. (12)
India	0.42	0.40 ± 0.01	0.58 ± 0.01
Germany	1.09	1.56 ± 0.02	2.13 ± 0.01
France	0.36	0.46 ± 0.01	1.35 ± 0.07
South Korea	0.82	0.79 ± 0.03	0.89 ± 0.01

TABLE II: Area covered by the CIs for the estimation of the reproduction number R within the whole $T = 70$ days time period.

area covered by the estimated 95% CIs given by the hierarchical model $(\mathbf{R}, \mathbf{O}) | \mathbf{Z}, \mathcal{I}; \theta$ compared to the area of those obtained with EpiEstim and with the model $(\mathbf{R}, \mathbf{O}) | \mathbf{Z}, \mathcal{I}, \lambda$, computed from the trapezoidal numerical integration method [30, Section 5.1] with one-day step size. For the two estimators computed from MCMC samplers, these areas are averaged over 5 runs and accompanied by associated confidence intervals. Tab. II shows that the CIs obtained with the hierarchical model are consistently larger compared to those from both other models. The size of the CIs strongly depends on the dataset, supporting the conclusion that the novel Hierarchical approach leverages the designed hierarchical Bayesian model (See (12)) to capture accurately intrinsic uncertainty contained in real COVID-19 infection counts.

V. CONCLUSIONS AND PERSPECTIVES

In the present work, we have constructed a hierarchical Bayesian epidemiological model, and devised a sampling scheme that permits the estimation of the reproduction number, quantifying the intensity of a pandemic. The achieved estimation benefits jointly from confidence assessment in terms of credibility intervals and from robustness against both the corrupted nature of the available COVID-19 pandemic data and the unknown values of parameters of the model. The relevance of the procedure is assessed on real pandemic data for several countries. Future investigations include exploring online estimation schemes and better-suited models for data corruption and missing data, focusing on aggregated count models [31]–[33]. Matlab codes implementing the proposed robust estimation procedure will be shared publicly at the time of publication.

REFERENCES

- [1] A. Flahault, "Covid-19 cacophony: is there any orchestra conductor?" *The Lancet*, vol. 395, no. 10229, p. 1037, 2020.
- [2] E. Krymova, B. Béjar, D. Thanou, T. Sun, E. Manetti, G. Lee, K. Namigai, C. Choirat, A. Flahault, and G. Obozinski, "Trend estimation and short-term forecasting of covid-19 cases and deaths worldwide," *Proc Natl Acad Sci U S A*, vol. 119, no. 32, p. e2112656119, 2022.
- [3] D. Gaglione, P. Braca, L. M. Millefiori, G. Soldi, N. Forti, S. Marano, P. K. Willett, and K. R. Pattipati, "Adaptive bayesian learning and forecasting of epidemic evolution—data analysis of the covid-19 outbreak," *IEEE Access*, vol. 8, pp. 175 244–175 264, 2020.
- [4] S. Marano and A. H. Sayed, "Decision-making algorithms for learning and adaptation with application to covid-19 data," *Signal Processing*, vol. 194, p. 108426, 2022.
- [5] P. Singh, A. Singhal, B. Fatimah, and A. Gupta, "A novel prfb decomposition for non-stationary time-series and image analysis," *Signal Processing*, vol. 207, p. 108961, 2023.
- [6] Q.-H. Liu, M. Ajelli, A. Aleta, S. Merler, Y. Moreno, and A. Vespignani, "Measurability of the epidemic reproduction number in data-driven contact networks," *Proc Natl Acad Sci U S A*, vol. 115, no. 50, pp. 12 680–12 685, 2018. [Online]. Available: <https://www.pnas.org/content/115/50/12680>
- [7] F. Brauer, C. Castillo-Chavez, and Z. Feng, *Mathematical models in epidemiology*. Springer, New York, 2019.
- [8] I. C.-. F. Team and H. S. I., "COVID-19 scenarios for the United States," *medRxiv*, 2020. [Online]. Available: <https://www.medrxiv.org/content/early/2020/07/14/2020.07.12.20151191>
- [9] N. Bannur, V. Shah, A. Raval, and J. White, "Synthetic Data Generation for Improved covid-19 Epidemic Forecasting," *medRxiv*, pp. 2020–12, 2020.
- [10] P. Singh, A. Singhal, B. Fatimah, and A. Gupta, "An improved data driven dynamic SIRD model for predictive monitoring of COVID-19," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.* IEEE, 2021, pp. 8158–8162.
- [11] S. Ahn and M. Kwon, "Reproduction Factor Based Latent Epidemic Model Inference: A Data-Driven Approach Using COVID-19 Datasets," *IEEE J. Biomed. Health Inform.*, vol. 27, no. 3, pp. 1259–1270, 2022.
- [12] O. Diekmann, J. A. P. Heesterbeek, and J. A. J. Metz, "On the definition and the computation of the basic reproduction ratio R_0 in models for infectious diseases in heterogeneous populations," *J. Math. Biol.*, vol. 28, pp. 365–382, 1990.
- [13] J. Wallinga and P. Teunis, "Different Epidemic Curves for Severe Acute Respiratory Syndrome Reveal Similar Impacts of Control Measures," *Am. J. of Epidemiol.*, vol. 160, no. 6, pp. 509–516, 09 2004. [Online]. Available: <https://doi.org/10.1093/aje/kwh255>
- [14] A. Cori, N. M. Ferguson, C. Fraser, and S. Cauchemez, "A new framework and software to estimate time-varying reproduction numbers during epidemics," *Am. J. Epidemiol.*, vol. 178, no. 9, pp. 1505–1512, 2013.
- [15] R. K. Nash, P. Nouvellet, and A. Cori, "Real-time estimation of the epidemic reproduction number: Scoping review of the applications and challenges," *PLOS Digital Health*, vol. 1, no. 6, p. e0000052, 2022.
- [16] C. Fraser, "Estimating individual and household reproduction numbers in an emerging epidemic," *PLoS one*, vol. 2, no. 8, p. e758, 2007.
- [17] O. Gressani, J. Wallinga, C. L. Althaus, N. Hens, and C. Faes, "EpiIps: A fast and flexible bayesian tool for estimation of the time-varying reproduction number," *PLOS Computational Biology*, vol. 18, no. 10, pp. 1–27, 2022.
- [18] P. Abry, N. Pustelnik, S. Roux, P. Jensen, P. Flandrin, R. Gribonval, C.-G. Lucas, É. Guichard, P. Borgnat, and N. Garnier, "Spatial and temporal regularization to estimate COVID-19 reproduction number $R(t)$: Promoting piecewise smoothness via convex optimization," *PLOS One*, vol. 15, no. 8, p. e0237901, 2020.
- [19] B. Pascal, P. Abry, N. Pustelnik, S. Roux, R. Gribonval, and P. Flandrin, "Nonsmooth convex optimization to estimate the Covid-19 reproduction number space-time evolution with robustness against low quality data," *IEEE Trans. Signal Process.*, vol. 70, pp. 2859–2868, 2022. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-03348154/document>
- [20] P. Abry, G. Fort, B. Pascal, and N. Pustelnik, "Temporal evolution of the covid19 pandemic reproduction number: Estimations from proximal optimization to monte carlo sampling," in *Annu Int Conf IEEE Eng Med Biol Soc.*, 2022, pp. 167–170.
- [21] G. Fort, B. Pascal, P. Abry, and N. Pustelnik, "Covid19 reproduction number: Credibility intervals by blockwise proximal monte carlo samplers," *IEEE Trans. Signal Process.*, vol. 71, pp. 888–900, 2023.
- [22] H. Artigas, B. Pascal, G. Fort, P. Abry, and N. Pustelnik, "Credibility interval design for covid19 reproduction number from nonsmooth langevin-type monte carlo sampling," in *2022 30th European Signal Processing Conference (EUSIPCO)*, 2022, pp. 2196–2200.
- [23] P. Abry, G. Fort, B. Pascal, and N. Pustelnik, "Proximal-langevin samplers for nonsmooth composite posteriors: application to the estimation of covid19 reproduction number," in *2023 31th European Signal Processing Conference (EUSIPCO)*, 2023.
- [24] P. Abry, J. Chevallier, G. Fort, and B. Pascal, "Pandemic Intensity Estimation from Stochastic Approximation-based Algorithms," in *2023 IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, Herradura, Costa Rica, Dec. 2023. [Online]. Available: <https://hal.science/hal-04174245v2/document>
- [25] D. Cereda, M. Tirani, F. Rovida, V. Demicheli, M. Ajelli, P. Poletti, F. Trentini, G. Guzzetta, V. Marziano, A. Barone *et al.*, "The early phase of the COVID-19 outbreak in Lombardy, Italy," *Preprint arXiv:2003.09320*, 2020.
- [26] F. Riccardo, M. Ajelli, X. D. Andriano, A. Bella, M. Del Manso, M. Fabiani, S. Bellino, S. Boros, A. M. Urdiales, V. Marziano *et al.*, "Epidemiological characteristics of COVID-19 cases and estimates of the reproductive numbers 1 month into the epidemic, Italy, 28 January to 31 March 2020." *Euro Surveillance*, 2020.
- [27] C. Robert and G. Casella, *Monte Carlo statistical methods*. Springer Verlag, 2004.
- [28] H. Bauschke and P.-L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. New York: Springer, 2019.
- [29] C. P. Robert, *The Bayesian choice: From Decision-Theoretic Foundations to Computational Implementation*. Springer, 2007.
- [30] K. Atkinson, *An Introduction to Numerical Analysis*. John wiley & Sons, 1991.
- [31] N. M. Ferguson, Z. M. Cucunubá, I. Dorigatti, G. L. Nedjati-Gilani, C. A. Donnelly, M.-G. Basáñez, P. Nouvellet, and J. Lessler, "Countering the Zika epidemic in Latin America," *Science*, vol. 353, no. 6297, pp. 353–354, 2016.
- [32] K. Charniga, Z. M. Cucunubá, M. Mercado, F. Prieto, M. Ospina, P. Nouvellet, and C. A. Donnelly, "Spatial and temporal invasion dynamics of the 2014–2017 Zika and chikungunya epidemics in Colombia," *PLoS Comput. Biol.*, vol. 17, no. 7, p. e1009174, 2021.
- [33] R. K. Nash, S. Bhatt, A. Cori, and P. Nouvellet, "Estimating the epidemic reproduction number from temporally aggregated incidence data: A statistical modelling approach and software tool," *PLoS Comput. Biol.*, vol. 19, no. 8, p. e1011439, 2023.