

# Towards inferring network properties from epidemic data

Istvan Z. Kiss<sup>1,2</sup>, Luc Berthouze<sup>3</sup>, & Wasiur R. KhudaBukhsh<sup>4</sup>

<sup>1</sup> Department of Mathematics, University of Sussex, Falmer, Brighton BN1 9QH, UK

<sup>2</sup> Network Science Institute, Northeastern University London, London E1W 1LP, UK

<sup>3</sup> Department of Informatics, University of Sussex, Falmer, Brighton BN1 9QH, UK

<sup>4</sup> School of Mathematical Sciences, University of Nottingham, University Park, Nottingham NG7 2RD, UK

## Abstract

Epidemic propagation on networks represents an important departure from traditional mass-action models. However, the high-dimensionality of the exact models poses a challenge to both mathematical analysis and parameter inference. By using mean-field models, such as the pairwise model (PWM), the high-dimensionality becomes tractable. While such models have been used extensively for model analysis, there is limited work in the context of statistical inference. In this paper, we explore the extent to which the PWM with the susceptible-infected-recovered (SIR) epidemic can be used to infer disease- and network-related parameters. Data from an epidemics can be loosely categorised as being population level, e.g., daily new cases, or individual level, e.g., recovery times. To understand if and how network inference is influenced by the type of data, we employed the widely-used MLE approach for population-level data and dynamical survival analysis (DSA) for individual-level data. For scenarios in which there is no model mismatch, such as when data are generated via simulations, both methods perform well despite strong dependence between parameters. In contrast, for real-world data, such as foot-and-mouth, H1N1 and COVID19, whereas the DSA method appears fairly robust to potential model mismatch and produces parameter estimates that are epidemiologically plausible, our results with the MLE method revealed several issues pertaining to parameter unidentifiability and a lack of robustness to exact knowledge about key quantities such as population size and/or proportion of under reporting. Taken together, however, our findings suggest that network-based mean-field models can be used to formulate approximate likelihoods which, coupled with an efficient inference scheme, make it possible to not only learn about the parameters of the disease dynamics but also that of the underlying network.

Keywords: Epidemics, Networks, Inference.

## 1 Introduction

Exact mathematical models for describing the spread of epidemics on networks are often insoluble or intractable for large networks [19, 16]. ‘Mean-field’ models provide a solution by introducing approximations and focusing on quantities at the population level, such as the expectation of the number of infected or susceptible individuals, or the number of direct connections between two such groups [18]. Many mean-field models exist to describe the dynamics of epidemic processes on networks. They usually take the form of a system of ODEs describing these processes [9]. Such models typically involve applying a ‘closure’ to exact models. Closures rely on assumptions about the underlying contact network and/or even the dynamics (usually simplifying ones), and these assumptions bring the complexity of a given system to manageable levels [23, 4].

Modelling epidemics on networks using mean-field approximations is a well studied and active area of research [20, 1]. In both theoretical and applied settings, it is used for parameter estimation, prediction and informing intervention or policy making [3], as recently demonstrated during the COVID-19 global pandemic [21]. However, there is a lack of understanding as to how such models operate in combination with the explicit inclusion of contact structures via networks, especially when placed in the context of statistical parameter inference. As such an investigation is warranted into whether current methods could be improved upon, or otherwise better informed, by incorporating models of epidemics on networks and by including structured population-level information and/or assumptions.

As previously mentioned, existing mean-field models are characterised by varying levels of complexity based on the assumptions used to close the exact system. This often requires making a statement about the links in the network, e.g., the number of edges that form [SI] (susceptible-infected) pairs, or [ISI] (infected-susceptible-infected) triples. For example, contact homogeneity – that is, a fixed number of links between each node in the network – is a common assumption [10, 16]. In this work, we use the ‘pairwise’ mean-field model, closed at the level of triples. Pairwise models are based on a bottom-up approach starting at node-level and building towards links and thereafter triples. This makes them very intuitive and the ‘go-to choice’ in many different areas. Moreover, pairwise models extend naturally to networks with heterogeneous degrees, weighted networks or even more complex epidemic dynamics.

The aim of this paper is to investigate to what extent this model can be used for inference purposes, and more specifically, for gaining insights about both the value of the parameters of the disease dynamics and that of the contact network, thus expanding the current body of work in the field (a review of which can be found in [17]).

In Section 2, we outline the principle of epidemics on networks as stochastic processes before detailing the pairwise system of ODEs constituting the so-called mean-field SIR model. Section 3 describes simulated data – namely, the output from the forward model with noise and Gillespie simulations, which we used to benchmark the performance of our inference schemes – as well as three real-world datasets: (i) the 2001 UK foot-and-mouth disease outbreak, (ii) The A(H1N1) outbreak in Washington State University (WSU) campus at Pullman, and (iii) the third wave of COVID-19 in India. Section 4 details the two inference schemes we considered, namely, maximum likelihood estimation and dynamical survival analysis. Section 5 presents a comparative analysis of these two schemes, both when ground-truth data is available (simulated data) and when it is not (real-world datasets). An interpretation of these results is provided in Section 6, along with potential new research directions.

## 2 Model

### 2.1 Epidemics on networks as a stochastic process

The starting point is the modelling of population contact structures as a network of nodes connected by links which represent possible routes of disease transmission. The network can be represented by an adjacency matrix  $G = (g_{ij})_{i,j=1,2,\dots,N}$ , where  $N$  is the number of nodes and the entries,  $g_{ij}$ , are either zero, if nodes  $i$  and  $j$  are not connected, or one otherwise. The adjacency matrix is symmetric and all elements on the main diagonal are zero, i.e., no self-loops are allowed. In this paper, we will focus on regular or homogeneous networks where each node has exactly  $n$  links.

When modelled as a continuous-time Markov Chain, a stochastic susceptible-infected-recovered

(SIR) epidemic on a network results in a state space of size  $3^N$  since each of the  $N$  nodes can be independently S, I or R, and each state, that is, a labelled network, needs an equation [9]. This of course makes the model intractable both theoretically and numerically, even at modest values of  $N$ . Of course, Gillespie [8] simulations can help deal with the problem and enable us to produce true stochastic paths of the process, see Figure 1 for example. This is based on the simple principle that in the Markovian framework, infection and recovery are independent Poisson process with rate  $\tau$  and  $\gamma$ .  $\tau$  is the per-link rate of infection and is the rate at which the I (infected) node in an I-S link infects the S (susceptible) node. This process is network-dependent. All infected nodes recover independently of the network and of each other at rate  $\gamma$ .

One way to move beyond simulations while dealing with the challenges of intractable high-dimensional models is to use mean-field models that focus on some expected quantity from the exact system, such as the expected number of infected nodes or the expected number of pairs of various types (e.g., S-S and S-I). One widely used model is the pairwise model [9] which is briefly described below.

## 2.2 Pairwise model as an approximation of epidemics on networks

In essence, the pairwise model focuses on a hierarchical construction where the expected number of nodes in state  $A$  at time  $t$ ,  $[A](t)$ , depends on the expected number of pairs of various types (e.g.,  $[AB]$ ) and these, in turn, depend on triples such as  $[ABC]$ . Here, the counting is done in all possible directions, meaning that  $[SS]$  pairs are counted twice and  $[SI] = [IS]$ . With this in mind, the pairwise model becomes

$$[\dot{S}] = -\tau[SI]; \quad [\dot{I}] = \tau[SI] - \gamma[I]; \quad [\dot{R}] = \gamma[I], \quad (1)$$

$$[\dot{SI}] = -(\tau + \gamma)[SI] + \tau([SSI] - [ISI]); \quad [\dot{SS}] = -2\tau[SSI]. \quad (2)$$

This system is not self-contained as pairs depend on triples and equations for these are needed. This, however, would lead to an explosion in system size as triples will then depend on quadruples connected in ways different from the simple line graphs over four nodes. To tackle this dependency on higher-order moments, the triples in equation (2) are closed using the following relation,

$$[ASB] = \kappa \frac{[AS][SB]}{[S]}, \quad (3)$$

where  $A, B \in \{A, B\}$ . Applying this closure leads to

$$[\dot{S}] = -\tau[SI], \quad (4)$$

$$[\dot{I}] = \tau[SI] - \gamma[I], \quad (5)$$

$$[\dot{R}] = \gamma[I], \quad (6)$$

$$[\dot{SI}] = -(\tau + \gamma)[SI] + \tau\kappa \frac{[SI]([SS] - [SI])}{[S]}, \quad (7)$$

$$[\dot{SS}] = -2\tau\kappa \frac{[SS][SI]}{[S]}, \quad (8)$$

which is now a self-contained system. As it turns out, see Figure (1), this low-dimensional mean-field model is exact in the asymptotic limit of  $N \rightarrow \infty$ , and the numerical solution of the PW model is indistinguishable from the average of stochastic realisations. We note that there are necessary and sufficient conditions which guarantee that the PW model is exact in the limit of large network sizes. In particular, it is true for networks with Binomial (with Regular being a special case of

Binomial), Poisson and Negative Binomial degree distributions [15, 13]. In the present setup, we will use  $\kappa = (n - 1)/n$  which corresponds to a regular network where each node has degree  $n$ . For networks with Poisson degree distribution  $\kappa = 1$ , while for negative binomial with  $NegBin(r, p)$  we have  $\kappa = (r + 1)/r$ . In all cases, the average degree is also needed in defining the initial conditions of the PW system. While we have chosen the closure corresponding to Regular networks, our approach can be extended to the other cases and this is detailed further in the Discussion section.

For a chosen set of parameters  $(n, \tau, \gamma)$  and initial conditions, the system above can be numerically integrated, furnishing us with  $[I](t)$  for example. Using that  $R_0 = \frac{\tau(n-1)}{\tau+\gamma}$ , the closed pairwise equations can be re-parameterised to include  $R_0$  explicitly. Keeping in mind that  $\kappa = (n - 1)/n$ , the re-parameterised system now reads

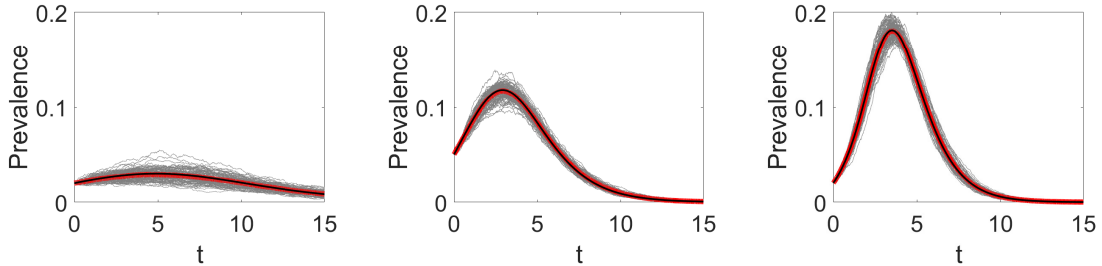
$$[\dot{S}] = -\frac{\gamma R_0}{(n-1) - R_0} [SI], \quad (9)$$

$$[\dot{I}] = +\frac{\gamma R_0}{(n-1) - R_0} [SI] - \gamma [I], \quad (10)$$

$$[\dot{R}] = +\gamma [I], \quad (11)$$

$$[\dot{SI}] = -\left(\frac{\gamma R_0}{(n-1) - R_0} + \gamma\right) [SI] + \kappa \frac{\gamma R_0}{(n-1) - R_0} \frac{[SI]([SS] - [SI])}{[S]}, \quad (12)$$

$$[\dot{SS}] = -2\kappa \frac{\gamma R_0}{(n-1) - R_0} \frac{[SS][SI]}{[S]}. \quad (13)$$



**Figure 1:** Prevalence based on Gillespie simulations. Thin lines/cloud in grey are the outcome of  $\sim 100$  individual realisations (10 networks with 10 realisations each) of an SIR stochastic epidemic on regular networks ( $n = 6$ ), with their average plotted in thick red lines. Epidemics are started with  $I_0 = 100$  (left panel) and  $I_0 = 250$  infectious nodes chosen at random (middle and right panels) and only epidemics that reach  $2I_0$  are kept and averaged over. The numerical solution of the corresponding pairwise model is plotted as a continuous black line. All networks have  $N = 10000$  nodes and the recovery rate is  $\gamma = 1$ . From left to right,  $\tau$  takes value 0.3, 0.4 and 0.5, respectively.

### 3 Data

Typically, real-world data for inference comes as daily counts of some quantity of interest (e.g., daily new cases or daily deaths) at discrete time steps, that is

$$(\mathbf{y}, \mathbf{t}) = \{(y_1, t_1), \dots, (y_{n_{obs}}, t_{n_{obs}})\}, \quad (14)$$

where  $(y_1, \dots, y_n) \in \{0, \dots, N\}$  and  $(t_1, \dots, t_{n_{obs}}) \in \{0, T\}$  with  $(0 \leq t_1 < t_2 < \dots < t_{n_{obs}} \leq T)$  are the counts and times respectively. However, we also consider data at the individual-level such as recovery and/or infection times, see equations (19)-(21) In this paper, data for inference is either simulated or taken directly from real-world epidemics. Full details are given below.

### 3.1 Data: PWM output with noise

Since the mean-field model is an approximation of the true stochastic process, we start by simulating data directly from the mean-field model and with varying levels of noise dispersion added in order to assess the ability of the inference schemes to recover the expected parameters, i.e., those used to generate the data (before noise). Since we mainly fit to daily reported cases, we first solve the PW model numerically with a given set of parameters and compute the daily new cases on day  $i$ ,  $([S](i+1) - [S](i))$ . Observations begin on the first day, at the earliest, and the initial conditions of the PWM are set at  $t = 0$ . Noise is introduced using draws from the Negative Binomial distribution. This is done such that the mean of the distribution is given by the model and the variance is controlled by the experimenter. For the Negative Binomial, and given a daily new cases count,  $y_d$ , from the true model without noise, we draw a sample from

$$X \sim \text{NegBin} \left( m(k) = \frac{1}{k}, p = \frac{1}{1 + ky_d} \right), \quad (15)$$

where the mean of this distribution is  $y_d$ , the variance is given by  $y_d + y_d^2 k$  with  $k$  the dispersion parameter, and the negative binomial distribution is interpreted as giving the probability of observing  $y_D$  failures given  $m$  successes, that is

$$\mathcal{P}(X = y_d) = \binom{y_d + m - 1}{y_d} p^m (1 - p)^{y_d}. \quad (16)$$

### 3.2 Data: stochastic simulations

Since the real challenge is to fit to stochastic data, in the first instance, we consider simulated data constructed by using the Gillespie algorithm [8] for a stochastic SIR epidemic on an explicit network of contacts. The idea behind the simulation is rather simple. Each node has its own rate, resulting in a rate vector  $(r_i)_{i=1,2,\dots,N}$ . A susceptible node with  $m$  infected neighbours will have rate  $\tau m$  and an infected node will have rate  $\gamma$ . Recovered or removed nodes have rate zero as they no longer play a role in the dynamics. The time to next event is chosen from an exponential with rate  $R = \sum_i r_i$ , and the event itself will be chosen at random from all possible  $N$ -events but proportionally to the values of the rate, e.g., event  $j$  will be chosen with probability  $r_j/R$ . Typical simulation plots are shown in Figures (1).

### 3.3 Data: real epidemic data

In addition to assessing the robustness of the inference schemes on synthetic data for which ground truth is known, we considered real-world outbreak data from three different data sets:

1. *The 2001 Foot-and-mouth (FMD) disease outbreak in the UK.* The 2001 FMD outbreak in the UK started towards the end of February in 2001 and ended in September 2001, impacting more than 2000 farms. Control efforts resulted in the culling of millions of livestock [5], see Figure 16.
2. *The A(H1N1) outbreak in Washington State University (WSU) campus at Pullman, Washington.* In April 2009, there was an outbreak of influenza virus in Veracruz, Mexico. After this initial outbreak, a new strain of the virus, A(H1N1)pdm09, started to spread around the world in the autumn. See [22, 12] for more details about this triple reassortment virus, which spread even among young, healthy adults. As a result, multiple outbreaks on college campuses were seen, one of which was on the Washington State University (WSU) campus in

Pullman, Washington in late August 2009. Within the space of three months, almost 2300 students came to the campus health centre with influenza-like illnesses that were treated as influenza A(H1N1) infections. Figure 16 shows the daily new cases starting on 22 August 2009.

3. *The third wave of COVID-19 in India* The COVID-19 pandemic has killed millions of people across the globe. Here, we consider the third wave in India. Similar to the other two datasets, we have daily incidence and prevalence of cases, recoveries and deaths from 15 February 2021 to 31 June 2021 (see Figure 16).

These datasets were chosen because they are thought to involve substantially different link densities, which is ideal since we are aiming to infer the number of links a node has. In addition, they were previously analysed in the literature (e.g., in [12, 6]), thus providing a good basis for comparison.

## 4 Inference methods

While most inference methods are based on the optimisation of a likelihood function, the likelihood function itself can be formulated based on different considerations of the underlying model and data. The most direct method typically focuses on matching model output and data as closely as possible, i.e., it is an error minimisation process. More sophisticated methods consider the underlying stochastic model in a more direct way and involve the timing of events, even if simplifying assumptions may be needed. To ensure that investigation into the possibility of inferring epidemic and network parameters using the pairwise model is not affected or biased by the data available or the inference scheme used, we consider two different methods: (i) the widely-used MLE-based approach when data comes in the form of daily new cases, and (ii) the dynamic survival analysis (DSA) inference method where individual-level times of infection and/or recovery (as opposed to counts) are used. These are detailed below.

### 4.1 Maximum-likelihood-based approach

In order to fit data produced by the PW model with the likelihood based on the PW model, we simply test how well the true parameters can be recovered. This scenario does not require any approximation. When fitting to stochastic data from an exact epidemic or a real epidemic, however, we are making the assumption that the exact forward model can be approximated by the PW model.

In this paper, we use the negative-binomial distribution as likelihood of choice, because of its flexibility. The distribution models the number of failures given a target number of successes,  $m$ , and the probability of each experiment's success  $p$ . Furthermore, vector  $\Theta$  denotes the parameters of the model, e.g.  $\Theta = (R_0, n, \gamma, k)$  or other combinations as required. Setting the parameters in  $\Theta$  to some concrete numerical values allows us to define the parameters of the Negative Binomial distribution at each time point  $t_i$ ,  $i = 1, 2, \dots, n_{obs}$ , where data was observed, these are:

$$m(k) = \frac{1}{k}, \quad p_{\Theta}(t_i) = \frac{1}{1 + ky_{\Theta}(t_i)}, \quad (17)$$

with  $k > 0$  being the dispersion parameter, which we also attempt to infer. For clarity,  $y_{\Theta}(t_i)$  is obtained by solving the PW model with concrete values of the parameters in  $\Theta$  and then finding quantities of interest at the desired times, such as daily new cases at time  $t_i$ . In this case, the

distribution at time  $t_i$  has mean  $y_{\Theta}(t_i)$  and variance  $y_{\Theta}(t_i) + y_{\Theta}(t_i)^2 k$ . This yields the following likelihood

$$\mathcal{L}_{NegBin}(\Theta | (\mathbf{t}, \mathbf{y})) = \prod_{i=0}^N \binom{y_i(t_i) + m - 1}{y_i(t_i)} p_{\Theta}(t_i)^m (1 - p_{\Theta}(t_i))^{y_i}, \quad (18)$$

where  $y_i(t_i)$  is the data at time point  $t_i$ , or simply  $y_i$  as introduced in Section 3, across all observations. Using  $\mathcal{L}_{NegBin}$  effectively decouples the mean and the variance of the distribution describing the data. This is expected to be sufficient to capture the variability of the data resulting from either natural stochasticity or variability due to how data was collected.

Parameter estimation was performed by minimising the negative log-likelihood (nLL thereafter) using the widely used direct search Nelder–Mead method. Because this technique can converge to non-stationary points (but see also Section 5.2 regarding the implications of unidentifiability), for each estimation process, multiple initial conditions (15) were used. To avoid biasing the search, initial conditions were drawn using Latin hypercube sampling, maximising the minimum distance between points. Because Latin hypercube sampling cannot prevent inappropriate parameter settings, initial conditions were only accepted if the ratio  $\tau/\gamma$  was not too large. Specifically, we enforced that the denominator in the expression of  $\tau$ , i.e.,  $n - 1 - R_0$ , was greater or equal than 1.5 (chosen empirically). On average, 10 out of 15 initial conditions survived.

## 4.2 Dynamical Survival Analysis

The statistical methodology Dynamical survival analysis (DSA) has recently been developed in a series of papers [12, 6, 25, 11] to address some of the shortcomings of traditional inference methods used in infectious diseases epidemiology. In essence, the method combines classical dynamical systems theory with tools from survival analysis. The crux of the methodology lies in interpreting the law of large numbers ODEs (representing population proportions in the continuous time Markov chain model) as describing probability distributions of transfer times, such as time to infection, time to recovery. Such a change in perspective allows one to use population-level mean-field ODEs to describe the dynamics of scaled compartment sizes as well as to write a likelihood function for individual-level trajectories based on transfer times, which may be censored, truncated or even aggregated.

To apply the DSA methodology, let us first define  $[D] = [SI]/[S]$ , which satisfies

$$[\dot{D}] = \tau(1 - \kappa)[D]^2 + \left( \kappa n \tau [S]^{(2\kappa-1)} - \tau - \gamma \right) [D],$$

with initial condition  $[D](0) = n\rho$  and  $[S](0) = 1$ , where, as before,  $\kappa = (n - 1)/n$  and  $[S]$  satisfies the pairwise mean-field equation with  $[S](0) = 1$  and  $[I](0) = \rho$ . The reason we normalize the system so that  $[S](0) = 1$  will be clear when we describe the DSA likelihood. Now, dividing the above equation by  $[S] = -\tau[S][D]$ , solving for  $[D]$  in terms of  $[S]$  with initial condition  $[S](0) = 1$  and then putting the solution back in  $[\dot{S}] = -\tau[S][D]$ , we get

$$-[\dot{S}] = n\tau(1 - [S]^{\kappa})[S]^{\kappa} + \frac{\gamma + \tau}{1 - \kappa}[S](1 - [S]^{\kappa-1}) + n\tau\rho[S]^{\kappa},$$

with initial condition  $[S](0) = 1$ . In essence, DSA interprets the susceptible curve as an improper survival function for the time to infection of a randomly chosen initially susceptible individual. That is,  $[S](t) = P(T_I > t)$ , where the random variable  $T_I$  describes the time to infection. This interpretation is exact (asymptotically in the limit of a large population of the continuous time

Markov chain model) and is justified by means of a Sellke construction argument, see [11, 12, 6]. In order to interpret  $[S](t)$  as a survival function following the Sellke construction [2], we set  $[S](0) = 1$ . This survival function is improper because  $\lim_{t \rightarrow \infty} [S](t) = \mathbf{P}(T_I = \infty) > 0$ . However, we can transform it into a proper survival function by conditioning it on a final observation time  $T \in (0, \infty)$ . We define the probability density function  $h_T$  on  $[0, T]$  as follows:

$$h_T(t) = -\frac{[\dot{S}](t)}{(1 - [S](T))}.$$

Given a random sample of infection times  $t_1, t_2, \dots, t_n$ , the likelihood contribution of the infection times is given by

$$\ell_I(\kappa, \tau, \gamma, \rho \mid t_1, t_2, \dots, t_n) = \prod_{i=1}^n h_T(t_i). \quad (19)$$

Note that DSA does not require knowledge of removal times. However, if individual recovery or removal times are known, they may be used to enhance the quality of inference. The likelihood contribution of a random sample of individual recovery times  $t'_1, t'_2, \dots, t'_m$  is given by

$$\ell_R(\kappa, \tau, \gamma, \rho \mid t'_1, t'_2, \dots, t'_m) = \prod_{i=1}^m r_T(t'_i), \quad (20)$$

where

$$r_T(t) = \frac{\int_0^t h_T(u) \gamma e^{-\gamma(t-u)} du}{\int_0^T \int_0^t h_T(u) \gamma e^{-\gamma(t-u)} du dt}$$

is the density of the individual recovery times. The density  $r_T$  is a convolution of two densities:  $h_T$  for the time of infection and the density of an exponential distribution with rate  $\gamma$  corresponding to the infectious period. In practice, it is convenient to differentiate the density  $r_T(t)$  with respect to  $t$  and then solve a system of ODEs.

Finally, the DSA likelihood function based on a random sample of infection times  $t_1, t_2, \dots, t_n$  and a random sample of recovery times  $t'_1, t'_2, \dots, t'_m$  is given by

$$\ell(\kappa, \tau, \gamma, \rho \mid t_1, t_2, \dots, t_n; t'_1, t'_2, \dots, t'_m) = \ell_I(\kappa, \tau, \gamma, \rho \mid t_1, t_2, \dots, t_n) \ell_R(\kappa, \tau, \gamma, \rho \mid t'_1, t'_2, \dots, t'_m). \quad (21)$$

Note that the likelihood function in (21) is exact when the underlying population size in the continuous time Markov chain model grows to infinity. For practical convenience (and as with the MLE-based approach), we work with the loglikelihood function, i.e., the logarithm of the likelihood function, rather than the likelihood function. It is, of course, possible to maximise the DSA likelihood function  $\ell$  in equation (21) to get point estimates of the parameter set  $(\kappa, \tau, \gamma, \rho)$ . Such a procedure would then be called a maximum likelihood approach and the difference between the two inference schemes discussed here would simply be that they maximise two different likelihood functions. An alternative way to perform parameter inference using DSA is to adopt a semi-Bayesian approach via a Laplace approximation to the posterior. In this paper, we adopted a fully Bayesian approach. Specifically, we draw posterior samples of  $(\kappa, \tau, \gamma, \rho)$  using a Hamiltonian Monte Carlo (HMC) scheme implemented in the *Stan* programming language [24, 7] interfaced with **R**.



Some of the datasets used in this paper (see relevant sections) provide daily new infection cases, rather than infection and/or recovery times. As mentioned earlier, the DSA methodology does not require knowledge of removal times. When these are not available, one can simply work with the likelihood function  $\ell_I$  (or the corresponding loglikelihood) in equation (19). Infection times, in turn, can be constructed from daily new cases as follows: If we observe 10 new cases on day  $t$ , then we simply draw 10 random samples from a uniform distribution over  $[t - 0.5, t + 0.5]$ . By repeating this procedure for all days for which daily new case counts are available and combining the individual infection times (samples from the uniform distributions), we can transform the original count data into data on infection times. A random sample of those infection times can then be fed into the likelihood function  $\ell_I$  in equation (19). In datasets in which daily recoveries are available, we can construct individual recovery times in a similar fashion: If we observe 5 recoveries on day  $t$ , we draw a random sample of size 5 from a uniform distribution over  $[t - 0.5, t + 0.5]$ . We repeat this procedure for all days for which we have daily number of recoveries available, and then combine the individual recovery times. A random sample of this data on individual recovery times is then fed into the likelihood function  $\ell_R$  in equation (20).

## 5 Results

In this section, we present numerical results, first using synthetic data and then real epidemic datasets. Whereas for the Gillespie data, we provide results using both the MLE and DSA approaches, for the data produced by the PW model with noise, we only provide results using the MLE approach. This is because the DSA method is based on the law of large numbers limit of the true continuous-time Markov chain model, which we simulate using the Gillespie’s algorithm [12, 6]. That is, the DSA likelihood function is the true likelihood function (for the infection and recovery times) in the limit. One could, if desired, still apply the DSA method on the PW model output with noise. However, it would be artificial as there is no natural survival perspective in the data generation process of the PW model with noise, unlike the stochastic model where the DSA likelihood function is justified by the Sellke construction.

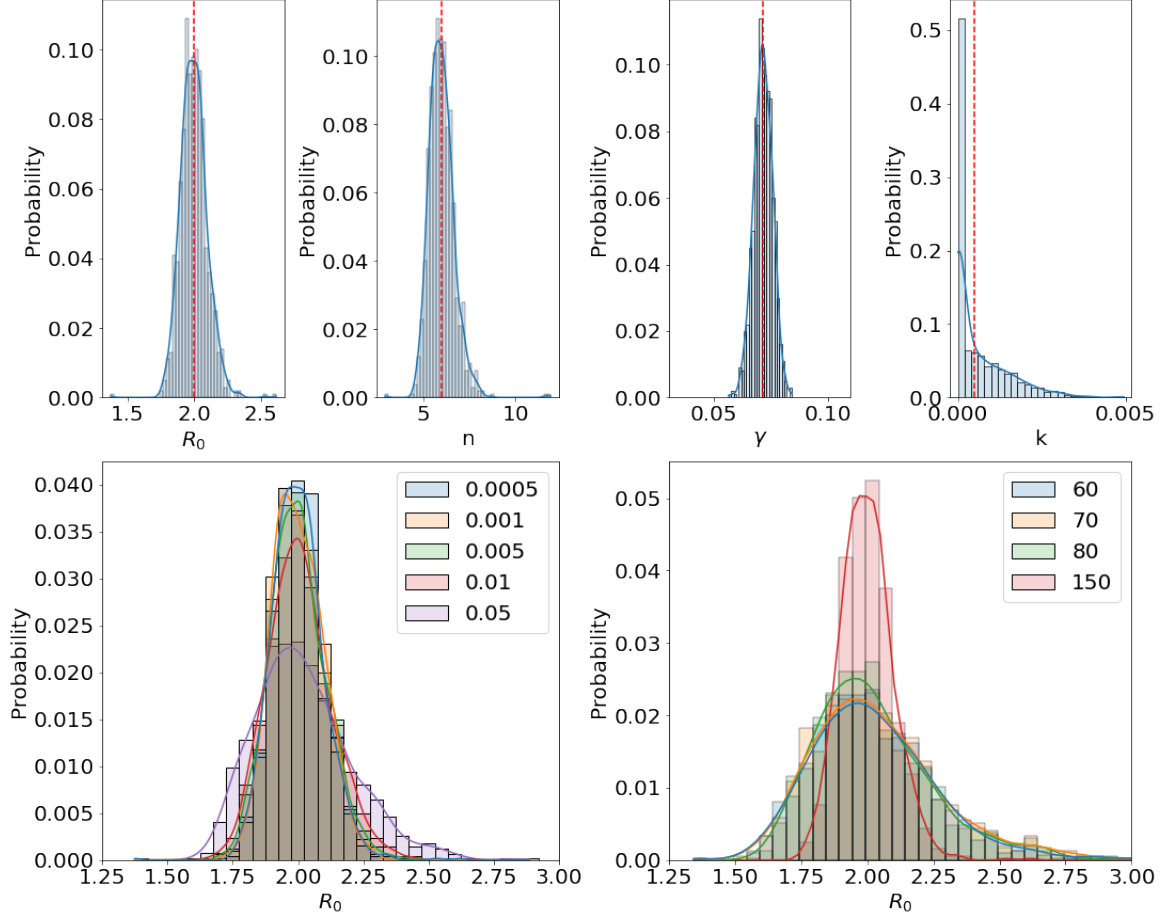
### 5.1 ML-based inference using data produced by the PW model

As a very first step toward assessing the ability of the inference scheme to recover the expected parameters, we first fitted the PW model (see Eqs. (9)-13) to daily cases data generated by the PW model and contaminated by some noise, whose dispersion was manipulated as will be described. Here, the values used to initialise the direct search Nelder-Mead method (see Section 4) for parameters  $R_0$ ,  $k$ ,  $n$  and  $\gamma$  were taken from  $[0.2, 10]$ ,  $[0.00001, 0.05]$ ,  $[3, 20]$  and  $[0.001, 0.1]$  respectively.

The top row of Figure 2 shows the histograms of parameters obtained when fitting  $M = 1000$  realisations, i.e. solving Eqs. (9)-13 with true  $[R_0, n, \gamma, I_0] = [2, 6, 1/14, 1]$  and  $N = 10000$ . Here, noise was simulated according to Eq. (15) using  $k = 0.0005$  (i.e., very low dispersion). These results confirm that the mean values are close to the true parameters, which is expected because the value of  $k$  is very small. For the avoidance of any confusion, we stress that in these histograms, each data point corresponds to the single-point estimate obtained for one of the  $M = 1000$  realisations and therefore these histograms should not be construed as posteriors.

To illustrate the sensitivity of the estimation process to the value of the dispersion parameter, we repeated the fitting process when considering 5 levels of dispersion, from 0.0005 to 0.01. As shown by the bottom left panel in Figure 2, as the dispersion level increases, so does the range of inferred  $R_0$  values. Nevertheless, the mean estimated value remains close to the true value in all cases.

Likewise, we found the inference process to be robust to the choice of time horizon (full epidemic  $t_{max} = 150$ , partial epidemic including the peak  $t_{max} = 80$ , epidemic up to the peak  $t_{max} = 70$ , partial epidemic not including peak  $t_{max} = 60$ ). As shown by the bottom right panel in Figure 2, as the time horizon reduces, the range of inferred  $R_0$  values increases but the average remains close to the true value. Importantly, whilst the inclusion of the peak does narrow the range of inferred values, it is not necessary for the inference process to correctly recover the expected value of  $R_0$ .

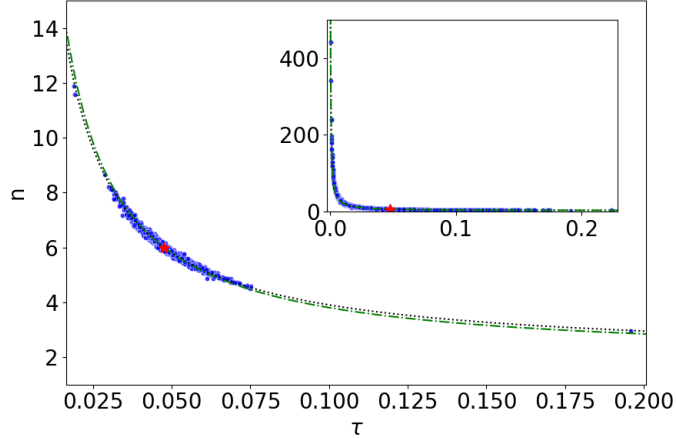


**Figure 2:** Inferring  $[R_0, n, \gamma, k]$  based on  $M = 10^3$  data realisations generated using  $[R_0, n, \gamma, k] = [2, 6, 1/14, 5 \times 10^{-4}]$  with  $N = 10^4$ ,  $I_0 = 1$ . Dashed lines indicate the true values of the parameters. Expected values were  $[1.999, 6.005, 0.0714, 0.00061]$ , respectively.

## 5.2 Identifiability

As Fig. 3 shows, the inferred values of  $\tau$  and  $n$  describe a hyperbola-like curve which indicates a clear identifiability problem; that is the values of  $\tau$  and  $n$  cannot be disentangled. However, we make two important remarks. First, it is possible to characterise this hyperbola analytically. Second, the values of  $\tau$  and  $n$  combine favourably into the expression of  $R_0$  whose inferred values are well behaved, see bottom panels in Fig. 2.

To formally characterise the hyperbola, we rely on quantities that can be derived analytically from the PW model. These are the leading eigenvalue (or growth rate under some transformation) and



**Figure 3:** Scatter plots of the parameter estimates on the  $n, \tau$  plane with the two practical unidentifiability curves calculated as per Eqs. 22 (dotted line), and 23 (dashed line). The star denotes the true values, i.e., true  $n$  and calculated value of  $\tau$  given true values of  $R_0$  and  $n$ . Main panel: scatter plot when the full epidemic is used for inference. Inset: scatter plot when the time horizon does not include the peak, i.e.,  $t_{max} = 60$ . Note that an arbitrary cut-off of  $n < 500$  was used for clarity of the plot.

the final epidemic size. These are given below in terms of  $\tau$  as a function of  $n$ .

$$\tau = \frac{\lambda_L^* + \gamma^*}{n - 2}, \quad (22)$$

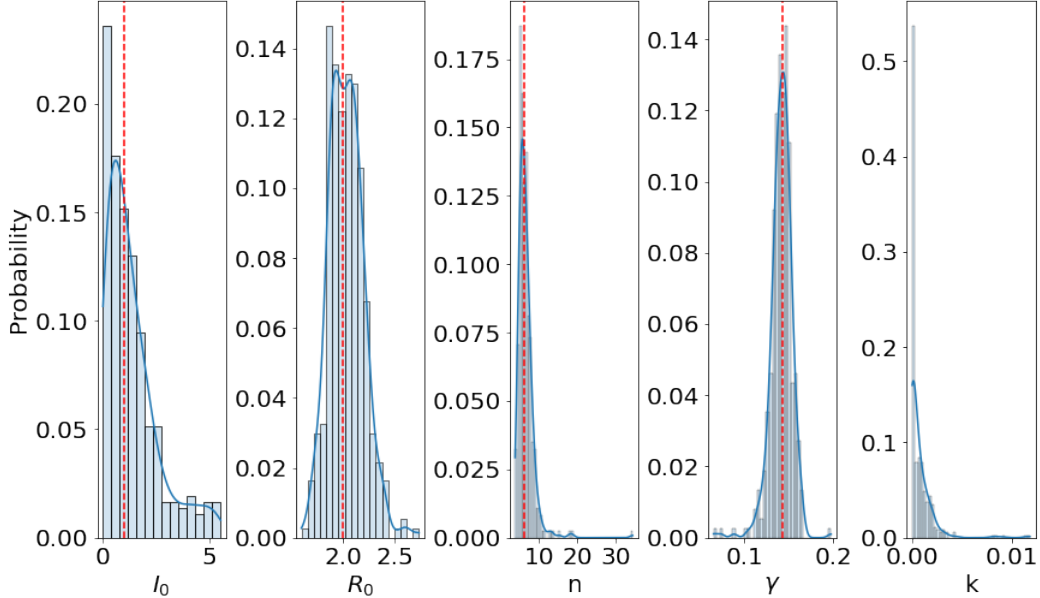
$$\tau = \gamma \frac{s_\infty^* 1/n - s_\infty^* 2/n}{s_\infty^* 2/n - s_\infty^*}, \quad (23)$$

where  $\lambda_L^*$  and  $s_\infty^* = S_\infty^*/N$  are obtained by setting all parameters to some desired values,  $(n, \tau, \gamma) = (n^*, \tau^*, \gamma^*)$ ; note that often  $R_0$  instead of  $\tau$  is given, with knowing the value of either being sufficient to have a well-defined system. The growth rate follows from the linear stability analysis of the pairwise model at the disease-free equilibrium, while the implicit formula for the final epidemic size can be found in [16] and is used here with initial conditions corresponding to the disease-free steady state.

### 5.3 ML-based inference using data from exact stochastic simulations

Five hundred Gillespie realisations were generated using parameters  $[R_0, n, \gamma, I_0] = [2, 6, 1/7, 1]$  and  $N = 10000$ . Of these 500 realisations,  $M = 370$  realisations did not die out. Figure 4 shows the histograms of the parameters estimated from fitting those realisations. Unlike with noisy realisations of the ODE, we also subjected  $I_0$  to the inference process. Results (not shown) obtained when assuming  $I_0 = 1$  during estimation revealed that the inclusion of  $I_0$  in the estimation process was key to being able to account for the stochasticity in the onset of the epidemic, or more precisely, the time elapsed before the growth becomes exponential. For the purpose of initialising Latin hypercube sampling, values were taken in  $[0.01, 10]$ . This particular choice has no bearing on our findings (results not shown). The mean of the estimated  $I_0$  was found to be 1.355, i.e., close to the expected 1; however, it showed a broad distribution of values, ranging from 0.012 to 5.534.

Comparing these histograms to those shown in Figure 2, we find that whilst the mean estimated values do not significantly differ, the variance in estimation is, not surprisingly, substantially larger. To quantify this more precisely, we calculated the mean (and standard deviation) of the confidence



**Figure 4:** Inferring  $[I_0, R_0, n, \gamma, k]$  based on  $M = 370$  data realisations generated using  $[I_0, R_0, n, \gamma] = [1, 2, 6, 1/7]$  with  $N = 10^4$ . Dashed lines indicate the true values of the parameters. Mean estimated values were  $[1.355, 2.029, 6.522, 0.141, 0.0071]$ , respectively.

intervals on  $R_0$  over all  $M = 370$  realisations. Specifically, we determined the nominal 99% profile likelihood confidence interval widths for  $R_0$  as described in [14]. Confidence intervals are  $0.534 \pm 0.203$  compared to  $0.498 \pm 0.071$  when fitting the ODE realisations with noise (dispersion level  $k = 0.0005$ ). These results are representative of those obtained when calculating confidence intervals for the other parameters (not shown).

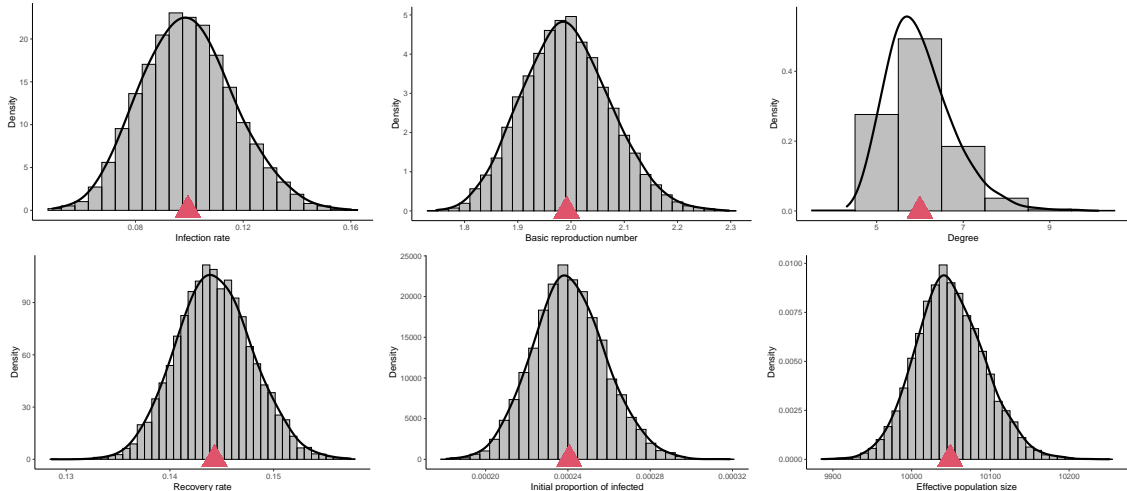
#### 5.4 Inference based on DSA

Before describing the results of DSA on the synthetic data, we highlight that, unlike the MLE-based approach which either assumes or infers both population size and initial number of infected individuals (see also Section 5.5.1), DSA inherently assumes an infinite population size (for both susceptible and infected individuals). Therefore, we do not infer the initial number of infected individuals. However, the ratio of initially infected to susceptible individuals, the parameter  $\rho$ , can be meaningfully inferred. In fact, having observed a finite number of infections in a given observation window  $[0, T]$ , DSA is also able to infer an *effective population size* using the discount estimator [12, 6]:

$$n_T = \frac{k_T}{1 - [S](T)}, \quad (24)$$

where  $k_T$  is the number of infections observed by time  $T > 0$ . It should be noted that estimates of the effective population size depend on the observation time  $T$ , and could be substantially different from the true population size when applying the method to a real epidemic. Nevertheless, as evidenced by the posterior distributions of the parameters  $(\tau, R_0, n, \gamma, \rho, n_T)$  shown in Figure 5, for this synthetic dataset, the method is able to infer the parameters well. The posterior distributions are unimodal, centred around the true values of the parameters. Here, at first random samples of individual infection and recovery times (of size 5000 each) were constructed from the count dataset

(one single trajectory of the Gillespie simulation) by drawing samples from appropriate uniform distributions (see Section 4.2). These random samples were then fed into the HMC scheme using four parallel Markov chains. Uninformative, flat priors were used except for domains  $(0.5, 10)$  for  $R_0$ ,  $(0.03, 0.3)$  for  $\gamma$ ,  $(0, 0.3)$  for  $\rho$  and an upper bound of 12 on  $n$ . The posterior distributions are not affected by the choice of those ranges on the prior distribution. For the sake of completeness, we have provided additional results in Appendix A when no ranges are imposed on the prior distribution.



**Figure 5:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using the DSA method on the synthetic data. The red triangles indicate the true values of the parameter. The means and the medians of the posterior distributions are  $(0.0994, 1.992, 5.997, 0.144, 0.0002, 10049)$  and  $(0.0989, 1.989, 5.891, 0.144, 0.0002, 10047)$ , respectively.

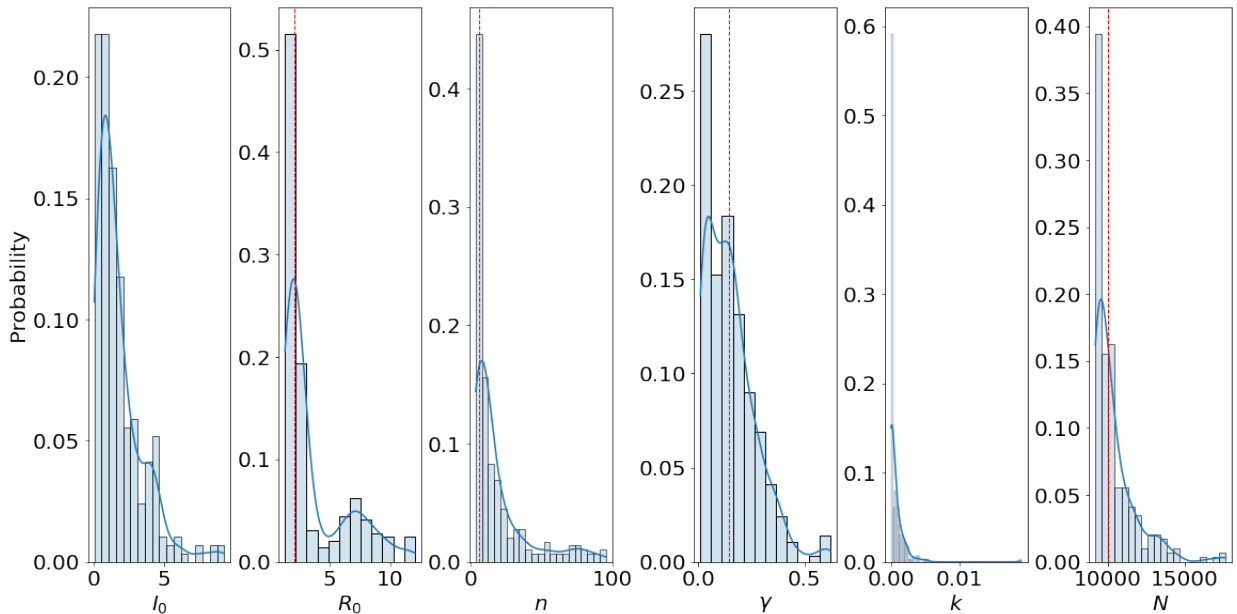
For this dataset, the parameter values estimated by both approaches are comparable. However, it is important to note that the two methods adopt two quite different likelihood constructions. Whilst the MLE-based approach relies on counts and the size of the population to construct the likelihood function, the DSA likelihood function only requires a random sample of infection times (and recovery times, if available). In other words, DSA identifies the probability laws of individual transfer times (infection and recovery times). These are often, even if censored, or truncated, more reliable and easily observed or derived statistical data than counts. For instance, even when we have partially observed count data on daily new infections, one can create a random sample of infection times (possibly censored/truncated). Even when the entire population is *not* monitored and only a set of randomly chosen individuals are followed through time and their transfer times are noted, the DSA methodology is still applicable. This advantage of DSA is particularly important when we fit the PW model to real epidemic data, which we do in the next section.

## 5.5 Inference from real-world data

### 5.5.1 System size and the MLE approach

In deploying the MLE approach to the above data, we used our knowledge of the true value of  $N$ . With real-world datasets, however, such information is typically not available. Whilst this is not an issue for DSA since it can infer an effective system size, it is for the MLE-based approach particularly in light of the unidentifiability issue discussed in 5.2. In what follows, we infer the value of  $N$  along with the other parameters, accepting that the increase in dimensionality of the parameter space will likely exacerbates unidentifiability. Here, we investigate the robustness of the inference

process when inferring known parameters on the stochastic realisations presented in Section 5.3. The data presented in Figure 6 result from the 289 out of a possible 370 realisations who satisfied the following conditions: (a) good fit (as quantified by the ratio 1.2 to the smallest likelihood value 217.25 obtained over the 370 realisations) – this excluded 66 estimates, (b) reasonable  $n$  (i.e.,  $n < 500$  arbitrarily – this excluded 13 estimates) and (c) reasonable  $\gamma$  (i.e.,  $\gamma < 1$  – this excluded a further 2 estimates – interestingly those estimates had very large  $N$ , specifically 26038.74 and 30168.98 but still showed very low nLL (244.02 and 228.3 respectively). The median values for the 6 parameters were:  $I_0 = 1.256792$ ,  $R_0 = 2.11$ ,  $n = 8.84$ ,  $\gamma = 0.129$ ,  $k = 0.00002$  and  $N = 9877.83$ . These values are reasonably close to the theoretical values ( $I_0 = 1$ ,  $R_0 = 2$ ,  $n = 6$ ,  $\gamma = 0.14$  and  $N = 10000$ ) which is encouraging. In particular, the percentage error in  $N$  is under 1.5%' (For reference, the percentage error for DSA on a random sample of the same data is in the order of 0.01%). Nevertheless, as shown by Figure 6, there is substantial variance in the estimates including significantly higher values of both  $N$  and  $R_0$  (e.g., 70 estimates have  $R_0 > 4$ ) despite excellent fits.

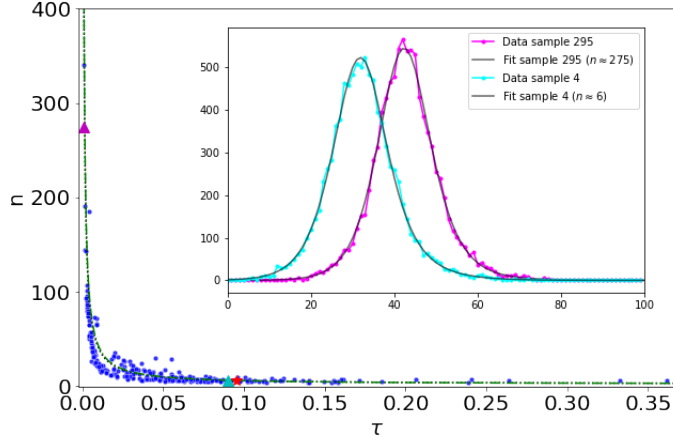


**Figure 6:** Inferring distributions for  $[I_0, R_0, n, \gamma, k, N]$  for the stochastic realisations. The ground truth parameter values ( $R_0$ ,  $n$ ,  $\gamma$  and  $N$ ) are denoted by vertical dashed lines. Data shown correspond to 289 out of the 370 stochastic realisations (see detail in text).

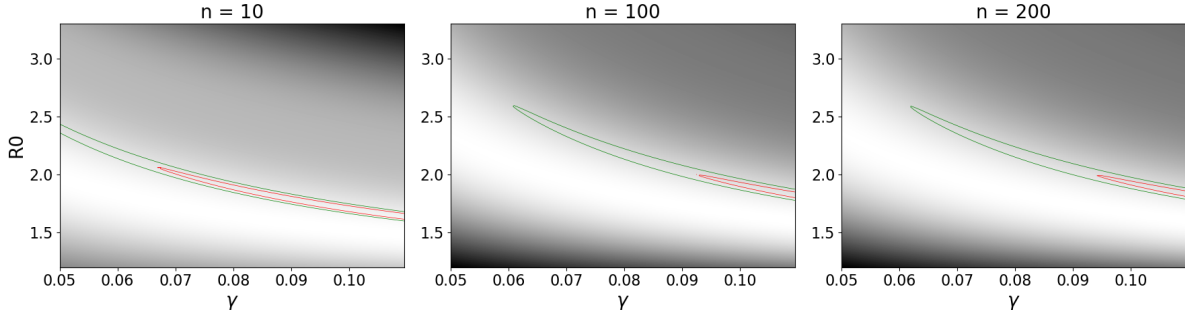
To illustrate this point, we plotted the estimates on the  $(\tau, n)$  plane (see Figure 7) and confirmed that they conform to the unidentifiability curve previously identified. The inset shows two stochastic realisations and the corresponding fits with one fit producing an estimate for the degree  $n$  close to the true value (6) and one producing an estimate magnitudes of order larger (275). As shown by the Figure (as well as the likelihood values), the fits are equally excellent. Inferred parameters for the data with the expected degree were:  $I_0 = 2.38$ ,  $R_0 = 2.18$ ,  $n = 6.05$ ,  $\gamma = 0.111$  and  $N = 9689.76$ , i.e., close to the ground truth data. In contrast, the inferred parameters for the data with the large degree were:  $I_0 = 0.24$ ,  $R_0 = 10.56$ ,  $n = 274.92$ ,  $\gamma = 0.024$  and  $N = 9281.53$ .

### 5.5.2 FMD data

Whereas DSA is Bayesian in nature and returns a posterior, the MLE approach only provides a single point estimate. When only one realisation of the process is available (as in the case of



**Figure 7:** Main panel: Scatter plot of the parameter estimates on the  $n, \tau$  plane with the two unidentifiability curves calculated as per Eqs. 22 (dotted line), and 23 (dashed line). The star denotes the true values, i.e., true  $n$  and calculated value of  $\tau$  given true values of  $R_0$  and  $n$ . Only those estimates who did not provide a good fit, as per the criterion above) were excluded, resulting in 304 surviving estimates. Inset: Empirical data and fit for two stochastic realisations corresponding to the triangles in the main panel with two significantly different inferred degree  $n$  (see detail in text).



**Figure 8:** Surface plots of the likelihood landscape for the FMD data (single realisation) obtained at 3 values of the degree parameter (10, 100 and 200) using a fine grid search (steps of 0.01 for  $R_0$  and 0.0005 for  $\gamma$ ). The nLL for the MLE (across all 100 initial conditions) was just under 129. The red and green contours denote levels 129 and 132 respectively and illustrate (a) the 'flatness' of the landscape around the unidentifiability curve for all values of  $n$  as well as (b) the fact that identical nLL values can be observed at vastly different values of the degree parameter (from 10 to 200).

empirical data), the identifiability issue discussed in Section 5.2 has significant practical implications for the interpretation of the results of the MLE process. Fig. 3 showed that when considering *multiple realisations* produced using the same parameters, the MLE estimates could be found widely distributed along the theoretical unidentifiability curves. To understand the numerical origin of this dispersion, we considered a *single realisation* and systematically investigated the likelihood landscape around the MLE estimate using a grid search. When using synthetic data with little noise (i.e., the data used to produce Fig. 3), we observe the presence of multiple local minima, relatively near to the known theoretical parameters and densely sampling the unidentifiability curve (not shown). In the presence of noise (e.g., shorter horizon) or when using empirical data, however, these minima manifest as isolated pockets over the full span of the unidentifiability curves in the high-dimensional parameter space. Fig. 8 illustrates this when considering the FMD dataset. Here,

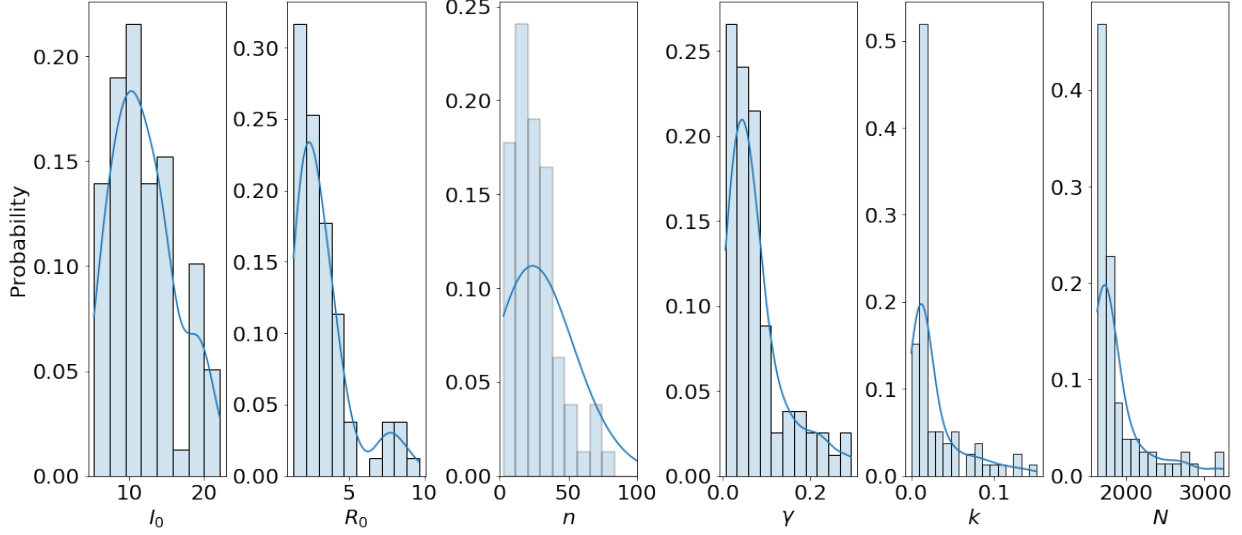
the same contour levels (set at, and very close to, the nLL corresponding to the MLE) are shown for 3 different degrees (10, 100 and 200) and demonstrate the presence of multiple local minima with the same nLL at vastly different degrees. In addition to forbidding the use of gradient-based methods, this type of landscape will trap most optimisation methods and will inevitably result in implausible parameters (e.g., extreme degree or  $R_0$  values).

For this reason, in the following 3 sections, we repeated the estimation process multiple times (100) using different initial conditions and provided histograms of the estimated values, after excluding those whose nLL was significantly different from the best nLL over all initial conditions (The number of estimates excluded for each dataset will be reported but will highlight the frequency with which the search algorithm can get stuck in very sub-optimal local minima). Although these histograms appear similar to the histograms summarising MLE estimates over multiple realisations when those are available, e.g., Fig. 2, they are fundamentally different in so far as they merely illustrate the diversity of parameter estimates that can be obtained within a single realisation. For the purpose of comparing MLE and DSA results, we always used the MLE, that is, the single point estimate with the smallest likelihood across all initial conditions, regardless of whether that estimate was away from the mean or mode of the 'distribution'. We also stress that, due to unidentifiability and rounding errors, the MLE parameters are not necessarily those closest to the theoretical parameters. Indeed, in many cases, some of the estimates were found to be closer to the (expected) parameters than the MLE estimate. In summary, these histograms should not be interpreted as posteriors. Although possibly confusing, we have included them because we believe that they are useful to highlight the numerical challenges posed by the likelihood profile (particularly when a large number of parameters are considered) and by this approach in general.

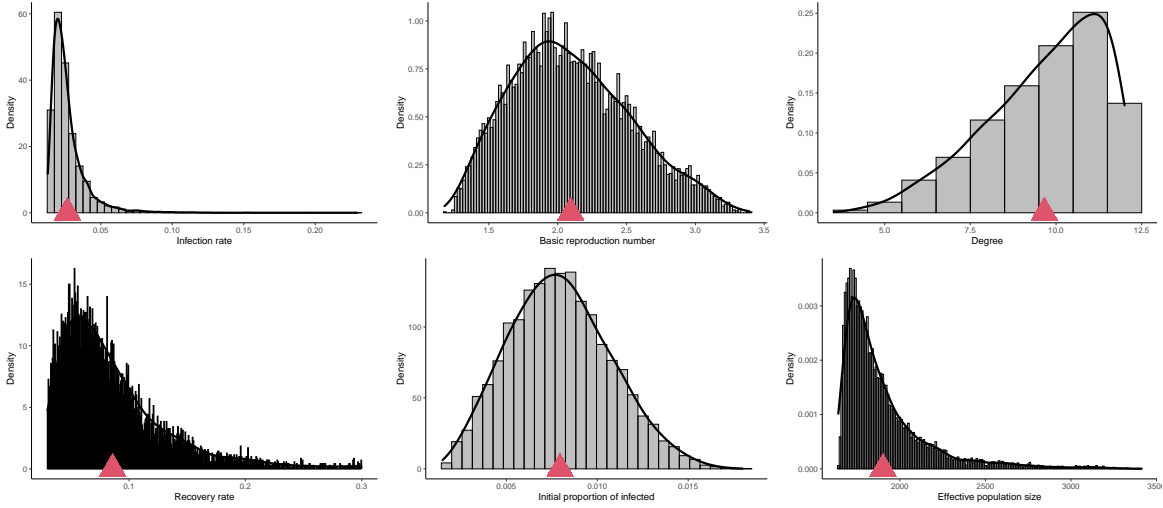
Histograms of inferred parameters for the FMD dataset using the MLE approach are shown in Figure 9. 11 out of 100 estimates were excluded because of an anomalous outcome of the inference process. The estimates with the lowest nLL are  $I_0 = 10.54$ ,  $R_0 = 2.58$ ,  $n = 153.67$ ,  $\gamma = 0.0723$ ,  $k = 0.010$ , and  $N = 1817.2$ . There is quite a bit of dispersion around the parameters, with fairly fat tails. For example, whilst the median for  $R_0$  (2.71, see Table 1) is relatively close to the best estimate, we also observe some fairly large values (in fact 10 out of 100 estimates were excluded because of  $R_0 > 10$ ). The best and median estimate for  $N$  was 1817 and 1747 respectively. This number is very likely implausible as many more than 2000 farms will have been involved in the epidemic, but see DSA results below. Likewise the inferred average degree seems far overinflated. The value of  $\gamma \approx 0.07$  implies 14 days for the infection period. Note that previous studies, see [6] for example, have reported a mean of 10.2 days. Importantly, the fits are good with all (accepted) estimates showing a very narrow range of nLL values (from 233.03 to 248.67 with a mean of 236.31 and a std of 4.06). This once again provides evidence of the fact that the MLE approach ascribes a likelihood to the trajectory produced by the inferred parameters rather than to the parameters themselves.

The posterior distributions obtained by DSA method on the FMD dataset are shown in Figure 10. It is important to note that, unlike with the MLE approach, these results were obtained when using an informative prior, an exponential distribution with mean 10.2 days, for the  $\gamma$  parameter following on the analysis in [6]. The posterior distributions are unimodal. The mean estimates are consistent with previously reported values, for example in [6]. Interestingly, and as with the MLE approach, the estimated effective population size is less than 2000. This is not to be confused with the number of farms, however (see brief explanation in Section 5.5.5). For the sake of completeness, we have also provided additional DSA results in Appendix A where informative priors are not used.





**Figure 9:** Inferred  $[I_0, R_0, n, \gamma, k, N]$  parameters when repeating the ML estimation process on the FMD data 100 times with different initial conditions for the parameter search algorithm. The MLE is provided in Table 1 and compared with the median from the DSA approach in Section 5.5.5. Five estimates for which  $n > 100$  (154, 156, 279, 294 and 368) were excluded from the figure (but not the statistics) for improved readability of the histogram.

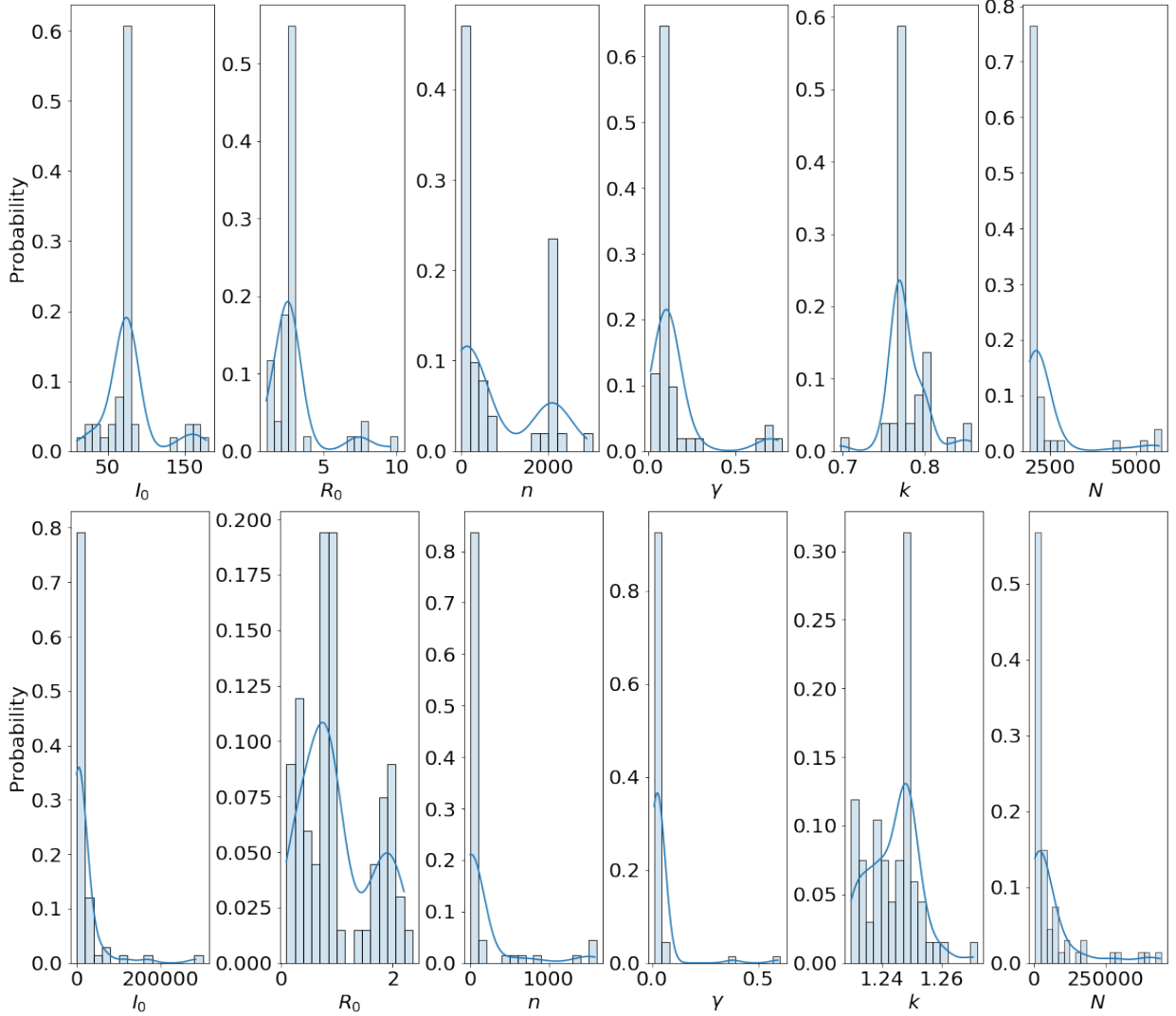


**Figure 10:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using DSA on the FMD dataset. The red triangles indicate the means of the posterior distributions. The means and medians of the posterior distributions are  $(0.0266, 2.095, 9.659, 0.0859, 0.0079, 1901)$  and  $(0.0233, 2.054, 9.982, 0.0737, 0.0078, 1819)$ , respectively.

### 5.5.3 H1N1-N18234

The A(H1N1) dataset presents an interesting challenge as it has a long persistent tail with visible stochastic effects. We therefore present two sets of results: one where we infer parameters on the full dataset (i.e., including the tail) and one when we restrict to  $T = 42$ . Figure 11 shows the results of the MLE-based approach for both scenarios. As clearly evidenced by the bottom right panel of Figure 16, when the full horizon is considered, the fits are poor, the noisy tail seemingly obfuscating the true trajectory of the epidemic. Not surprisingly, the parameter estimates appear meaningless

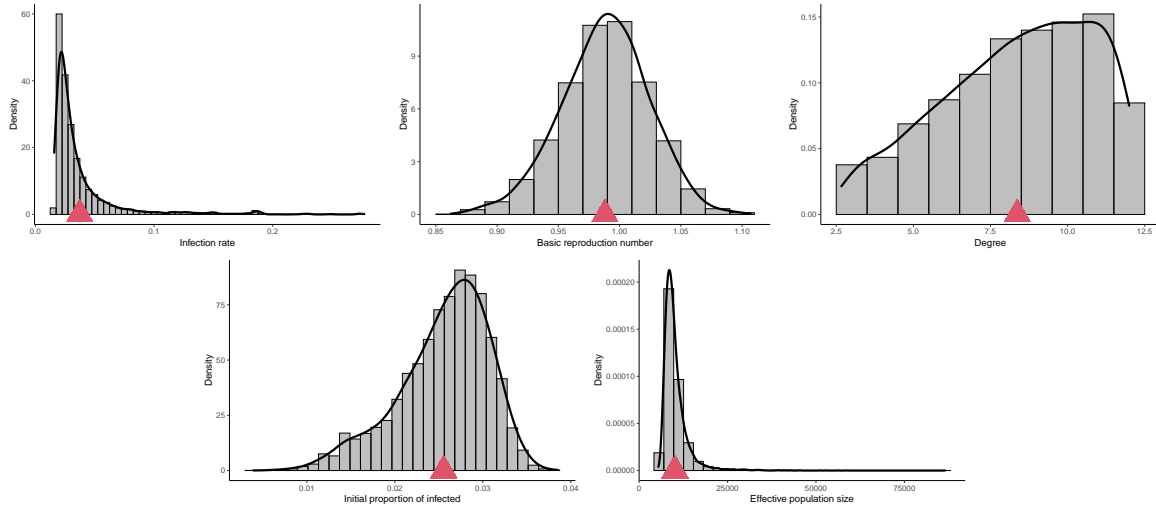
and highly variables from one round of inference to the other despite similar nLL. When restricting to  $T = 42$ , the fits are good and the parameter estimates are slightly better behaved albeit with not unimodal and with implausibly large  $n$  considering the inferred population size  $N$ . In fact, only 51 out of 100 parameter estimates survived once we excluded 3 estimates for being poor fits, 13 for excessive values of  $R_0$  ( $> 10$ ) and 33 estimates for excessive value of  $\gamma > 1$ . Interestingly, we note the high value of  $k$  inferred in both scenarios, with MLE correctly recognising the high dispersion of the counts.



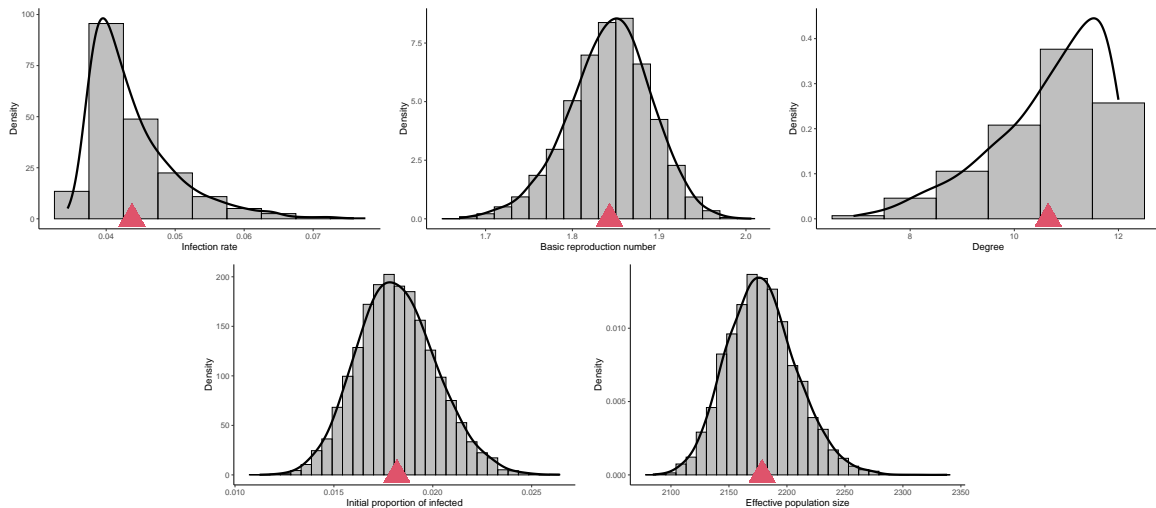
**Figure 11:** Inferred  $[I_0, R_0, n, \gamma, k, N]$  parameters when repeating the ML estimation process on the H1N1 data (with horizon restricted to 42, top panel and full data, bottom panel) 100 times with different initial conditions for the parameter search algorithm. The MLE is provided in Table 1.

When deploying DSA, once again, a prior was used for  $\gamma$  ( $\gamma^{-1} = 5.5$ ) based on published literature (see [22, 12]). Figures 12 and 13 show the posterior distributions of the parameters ( $\tau, R_0, n, \rho, n_T$ ) based on the full and partial data respectively. As with the MLE-based approach, when fitting to the full data, the DSA fit is poor, and in fact, very similar to that of the MLE approach (see bottom right panel of Figure 16). When removing the noisy tail of the data, the quality of inference improves significantly with both MLE and DSA producing near identical fits (bottom left panel

of Figure 16). However, unlike with FMD dataset, the inferred parameters are quite different although interestingly the ML-estimated population size and the DSA effective size are very similar (see Table 1). For the sake of completeness, we have provided additional DSA results in Appendix A where uninformative priors are used.



**Figure 12:** Posterior distributions of  $(\tau, R_0, n, \rho, n_T)$  using DSA on the full A(H1N1) outbreak data. The means and the medians of the posterior distributions are  $(0.0373, 0.9880, 8.369, 0.0255, 10146)$  and  $(0.0269, 0.9892, 8.665, 0.0264, 9286)$ , respectively.

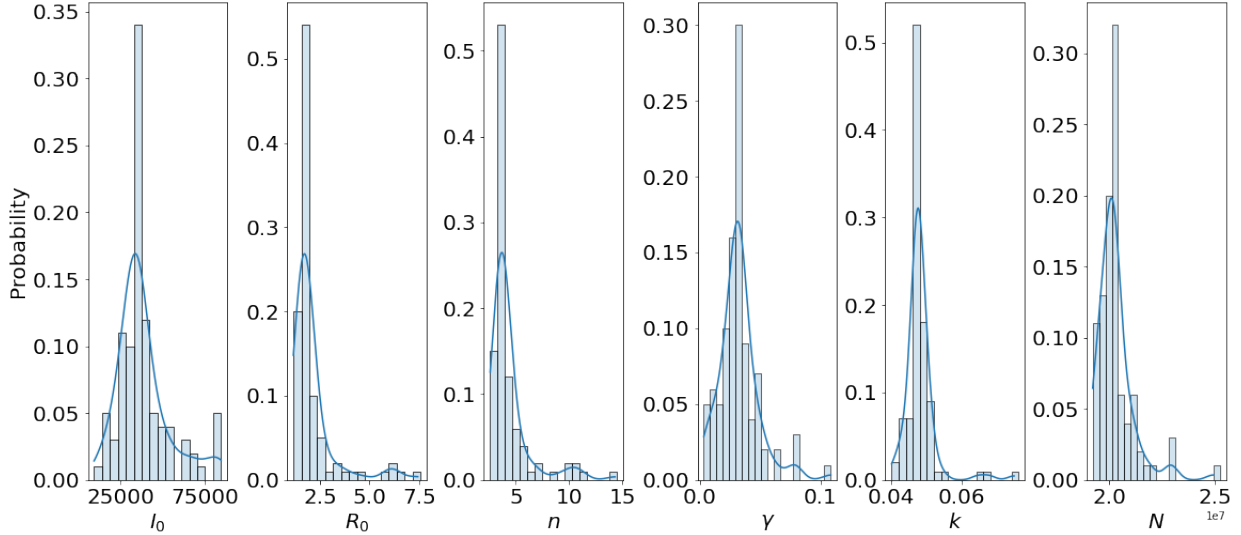


**Figure 13:** Posterior distributions of  $(\tau, R_0, n, \rho, n_T)$  using DSA on the A(H1N1) outbreak data restricted to time horizon  $T = 42$ . The means and medians of the posterior distributions are  $(0.0437, 1.843, 10.650, 0.0189, 2179)$  and  $(0.0418, 1.845, 10.908, 0.0189, 2177)$ , respectively.

### 5.5.4 COVID-19 in India

Figure 14 shows the histograms of the estimates obtained by the ML-based approach on the final dataset. Here, unlike with the previous dataset, there was high consistency between estimates over the 100 rounds with no exclusions needed. Curiously, this homogeneity of results is associated with an apparent mismatch between the fitted model and the data, as shown by the top right panel in

Figure16.



**Figure 14:** Inferred  $[I_0, R_0, n, \gamma, k, N]$  parameters when repeating the ML estimation process on the COVID-19 dataset 100 times with different initial conditions for the parameter search algorithm. The MLE is provided in Table 1.

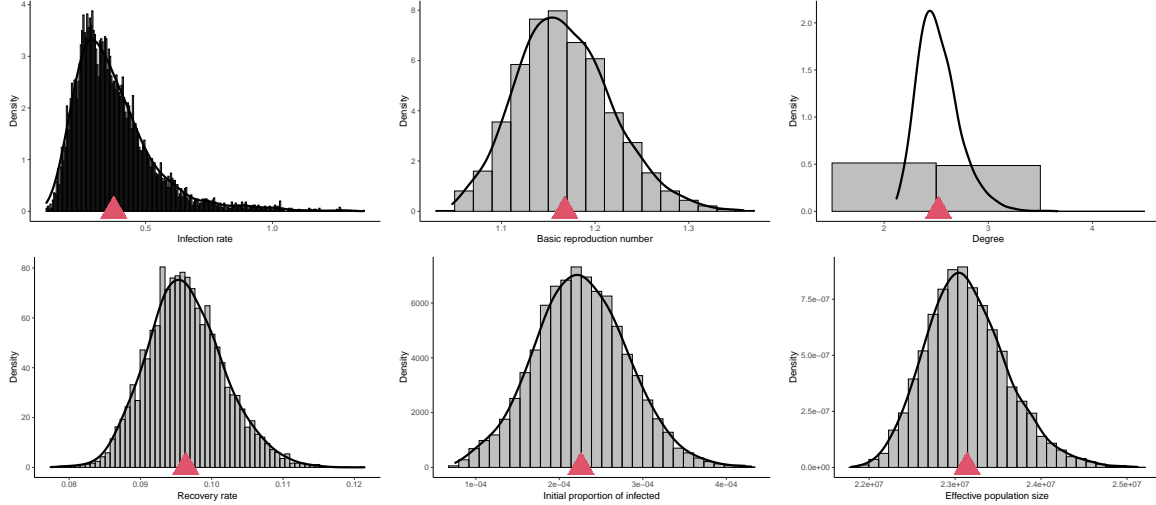
As in the synthetic data study, random samples of individual infection and recovery times (of size 5000 each) were constructed from the count dataset. These random samples were then fed into the HMC scheme using four parallel Markov chains. Uninformative, flat priors were used. The posterior distributions of the parameters  $(\tau, R_0, n, \gamma, \rho, n_T)$  using the DSA method are shown in Figure 15. The estimated parameters correspond to probability distributions that have similar measures of central tendency as those reported in an earlier analysis of the data in [6].

Interestingly, for both methods, the majority of the probability mass in the (posterior) distribution for the degree ( $n$ ) is concentrated around small values, indicating a low contact pattern. This is in agreement with various non-pharmaceutical interventions such as lockdowns that were put in place to reduce the spread of the virus. Finally, both ML-estimated population size and DSA effective size are in the same order of magnitude.

### 5.5.5 Comparison across real-world datasets

Figure 16 shows the data for all three real-world outbreaks together with fits produced when taking the best parameter estimates using the ML-based approach and the median values of the posteriors produced by DSA. Whilst our investigation of the COVID-19 dataset supports a like-for-like comparison between inference schemes, there are differences in the way the analyses of the FMD and the A(H1N1) datasets were carried out. Specifically, whereas no prior was involved in the MLE-based approach, informative priors (based on published literature) were used for the Hamiltonian Monte Carlo scheme for DSA. This reflects an important and fundamental difference between MLE-based approach and DSA methodology (here implemented via a Hamiltonian Monte Carlo scheme), namely that the latter follows a Bayesian route. It should be noted, however, that the effect of the choice of priors should vanish in the limit of a large number of data points, as suggested by the additional DSA results with uninformative priors provided in Appendix A.

With this in mind, we can make several observations:

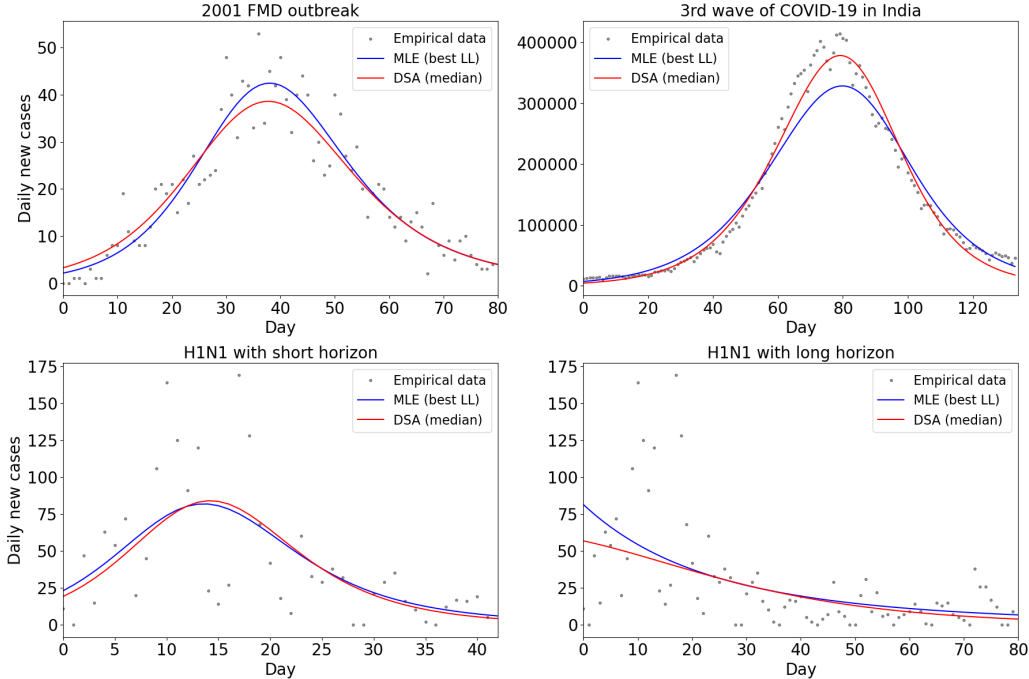


**Figure 15:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using the DSA method on the COVID-19 dataset. The means and the medians of the posterior distributions are  $(0.3745, 1.168, 2.522, 0.0963, 0.0002, 23139638)$  and  $(0.3401, 1.164, 2.493, 0.0961, 0.0002, 23101076)$ , respectively.

- In general, the fit to the real data is good except in two cases. In the COVID-19 data, despite relatively similar parameters between methods, the DSA fit appears to capture the trend of the data a lot better than the MLE fit where a clear mismatch is being observed. The scenario in which the full H1N1 epidemic is subjected to inference highlights the challenge of highly variable, potentially noisy, data, as well as the impact of the observation period. In particular, as shown by the bottom two panels of Figure 16, the longer observation window allows the long and noisy tail of the epidemic to dominate, with both approaches missing the rise and fall in the daily new cases.
- Table 1 only provides one single point estimate for the MLE approach even though the process was repeated a number of times to try to mitigate the impact of unidentifiability. Whilst this estimate is the 'true' MLE estimate (in the sense of being the one with maximum likelihood over all estimates of all rounds), it is worth remembering (as discussed in Section 5.5.2 and also shown by the histograms provided) that this estimate was not necessarily close to the median over the 100 rounds of estimation. In many cases, we observed a large difference between MLE and median. This is, once again, a manifestation of the unidentifiability problem whereby vastly different values of the mean-degree can result in likelihoods very close, or equal, to the best one (i.e., with the same quality of fit). Interestingly, we note that, in general (a few estimates were excluded as per the text), the impact of unidentifiability did not affect  $R_0$  as much as other parameters.
- The estimates for  $I_0$  and population size,  $N$ , are relatively similar across both inference approaches, except for A(H1N1) when the full dataset is considered and COVID-19. For the A(H1N1) outbreak, the MLE method appears to overestimate  $N$  by a large margin. Note that Washington State University campus is located in a relatively small town with a student population of size around 18000 and a resident population of size around 9000 [12]. For the COVID-19 wave in India, the DSA median estimate of 5204 for  $I_0$  appears smaller than the true count of 11592 new cases on 16 February 2021, whereas the MLE method seems to overestimate it (33130). It should be noted that the effective population size is a by-product of the DSA method (see Section 5.4). Strictly speaking, the parameters  $I_0$  and  $N$  are far less

meaningful in DSA than in MLE which requires them. However, keeping track of the DSA estimates  $n_T$  of the effective population sizes at times  $T$  is valuable in that it gives us a sense of the possible size of the epidemic and therefore, could be used for monitoring an ongoing epidemic [11].

- Although a strict comparison between the posteriors obtained by DSA and the histograms obtained by repeating the ML estimation process is meaningless (since these histograms are not posteriors), it is nevertheless interesting to note that when comparing them for the FMD data, we find the range of average degree obtained by DSA to be much better behaved than that obtained by MLE with mean and median being close and with a numerical value that seems more realistic. This observation holds for all datasets with DSA producing more realistic estimates. Unidentifiability aside, this is ultimately linked to the fundamental difference between how the likelihoods in the MLE and DSA approach are formulated. Whilst the MLE method simply minimises the mismatch between model trajectory and data, the DSA likelihood captures the underlying probability laws of individual infection and recovery times. More specifically, it models the underlying survival function through the  $[S](t)$  curve parameterized by  $(n, \tau, \gamma, \rho)$  (and implicitly, by the observation time  $T$ ).



**Figure 16:** Illustration of the real-world outbreak data (top-left - 2001 FMD outbreak in the UK, top-right - third wave of COVID-19 in India, bottom panels - H1N1 outbreak with short (left) and long horizon (right)) together with output from the pairwise model with point estimates from MLE (values with best likelihood) and DSA (median values). All parameter values are given in Table 1.

## 6 Discussion

In this paper, we have investigated the ability of a network-based mean-field model, i.e., the pairwise model, to infer not only disease parameters but also some of the underlying network. Outbreak data encapsulate the interplay between contact network and epidemic spreading. However, daily

	$I_0$	$R_0$	$n$	$\gamma$	$k$	$N$
FMD (MLE)	11	2.58	153.67	0.0723	0.0101	1,817
FMD (DSA median)	14	2.05	9.98	0.0737	-	1,819
H1N1-N18234 (MLE, 42 days)	76	2.70	2094.78	0.1073	0.7679	2,095
H1N1-N18234 (DSA median, 42 days)	39	1.85	10.91	0.1818	-	2,177
H1N1-N18234 (MLE, 80 days)	3463	0.85	2.61	0.0270	1.2468	20,256
H1N1-N18234 (DSA median, 80 days)	252	0.99	8.67	0.1818	-	9,286
covid (MLE)	33130	1.70	3.68	0.0333	0.0474	20,254,332
covid (DSA median)	5204	1.16	2.50	0.0961	-	23,101,076

**Table 1:** Summary statistics of the inferred parameters for the three empirical datasets considered in this study when using both MLE and DSA approaches. Estimates for  $I_0$  and  $N$  were rounded to the nearest integer for readability.

	$I_0$	$R_0$	$n$	$\gamma$	$k$	$N$
PW with $k=5e-4$ (MLE)	-	2.00	5.95	0.0714	0.0002	-
PW with $k=1e-3$ (MLE)	-	1.99	5.96	0.0718	0.0006	-
PW with $k=5e-3$ (MLE)	-	2.00	5.93	0.0713	0.0046	-
PW with $k=1e-2$ (MLE)	-	2.00	6.00	0.0718	0.0093	-
PW with $k=5e-2$ (MLE)	-	2.00	5.99	0.0715	0.0480	-
PW with Tmax=150 (MLE)	-	2.00	5.95	0.0714	0.0002	-
PW with Tmax=80 (MLE)	-	1.99	5.92	0.0712	0.0000	-
PW with Tmax=70 (MLE)	-	2.00	6.04	0.0707	0.0000	-
PW with Tmax=60 (MLE)	-	2.00	5.86	0.0697	0.0000	-
Gillespie with Tmax=100 (MLE)	1	2.02	6.10	0.1420	0.0002	-
Gillespie with Tmax=100 (MLE)	1	2.11	8.84	0.1290	0.0000	9878
Gillespie with Tmax=100 (DSA)	2	1.99	5.89	0.1442	-	10047

**Table 2:** Medians of the inferred parameters for the synthetic datasets considered in this study when using both the MLE approach and the DSA approach. Estimates for  $I_0$  and  $N$  (when estimated, shown as dash if not estimated) were rounded to the nearest integer for readability. Estimates of  $k$  below  $1e-4$  appear as 0.0000. The median estimates for full set of relevant parameters for the DSA method are also mentioned in the caption of Figure 5

new cases or other data incorporate network information only implicitly. Hence, it is interesting to investigate whether from such data one can learn about the underlying contact network. Several challenges arise; for example, an epidemic with a small transmission rate on a dense network may look very similar to an epidemic with a large transmission rate spreading on a sparser network. Hence, it is not a given that outbreak data hold a specific enough signature of the contact network. In fact, our investigation revealed an anti correlation between the value of the transmission rate and the density of the network. Regardless, the estimate of both parameters peaked at around the desired values, especially when ground truth was known.

While the pairwise model used in the paper assumes that the network is regular and only accounts for the number of links each node has, it is possible to relax this seemingly restrictive assumption. In [11], DSA was used for an SIR epidemic on a configuration model network with Poisson degree distribution. Recently, it has been shown [15] that the pairwise model remains exact for networks

with binomial, Poisson or negative binomial degree distribution; see also [13, Corollary 1, Section 5.2] where a similar result was derived for a susceptible-infected (SI) process on configuration model random graphs. The difference in the degree distributions manifests itself in the PW model via the type of closure one uses. For example, if the underlying network has a Poisson degree distribution, then  $\kappa$  is simply set to  $\kappa = 1$ , and the parameter of the Poisson distribution, and hence, the network enters the PW model via the initial conditions. A similar modification is possible for networks where the degree distribution is negative binomial thus separating mean from variance. These all offer extensions and improvements above and beyond what the PW model was able to capture about the network. Moreover, employing the edge-based compartmental model, another network-based mean-field model, which uses the probability generating function of corresponding to the degree distribution of the network makes it possible to aim for learning the degree distribution of the underlying network.

The crucial advantage of the DSA methodology is the change in perspective about the mean-field ordinary differential equations. In the DSA approach, we view the ODEs as descriptions of probability laws of individual times of infection and recovery, as opposed to their traditional interpretations as limiting proportions or scaled sizes of compartments. By doing so, we are able to directly model the underlying survival functions corresponding to the individual times of infection and recovery, and thereby, bring to bear the entire toolkit of survival analysis for the purpose of parameter inference. Even though the DSA methodology has now been applied to several compartmental models, both Markovian and non-Markovian, both under mass-action and network-based contact patterns, the law of large numbers-based DSA methodology needs further improvement to adjust for stochastic effects when applied to finite (often small) populations.

## A Additional DSA results

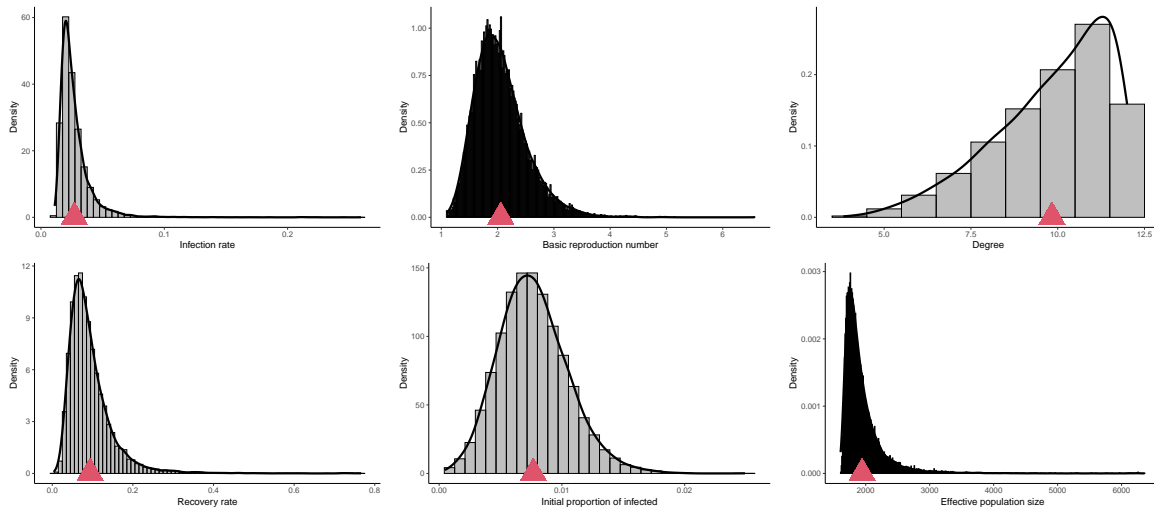
In the analysis of the FMD and the AH1N1 datasets using the dynamical survival analysis, we used informative priors so that the results could be directly compared against the analyses of the same datasets previously published in the literature. In this section, we present the results of dynamic survival analysis of the FMD and the AH1N1 datasets (interfacing with Stan using CmdStanR [7]) when informative priors are not used except for an upper bound of 12 on the parameter  $n$  (chosen arbitrarily). This is because when  $n$  is large,  $\kappa = (n - 1)/n \approx 1$  (and  $1 - \kappa \approx 0$  in the denominator of the differential equation for  $[S]$ , which can cause numerical instability and lead to slow mixing of the Hamiltonian Monte Carlo) suggesting the use of a Poisson distribution in the configuration model instead. If desired, the case of large  $n$  should be handled by taking a limit  $\kappa \rightarrow 1$  of the differential equation for  $[S]$  and then applying DSA with the limit, as done in [11]. In this paper, we do not pursue this additional complication.

The purpose of this presentation is to show the impact of the informative priors used in earlier analysis. The use of informative priors is not essential for the application of DSA. However, we do note that it is usually recommended to use informative priors when reliable information on the parameters are available.

### A.1 FMD data

The mean and the median estimates of  $(\tau, R_0, n, \gamma, \rho, n_T)$  obtained when uninformative priors are used are (0.027, 2.06, 9.38, 0.093, 0.0073, 1943), and (0.023, 1.99, 10.2, 0.080, 0.0074, 1848) respectively. They are comparable to the means (0.0266, 2.095, 9.659, 0.0859, 0.0079, 1901) and the medians (0.0233, 2.054, 9.982, 0.0737, 0.0078, 1819) of the posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$





**Figure 17:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using DSA with uninformative priors on the FMD dataset. The red triangles indicate the means of the posterior distributions.

when informative priors are used. Please see Figure 17 for the posterior distributions of the parameters, which are unimodal. The HMC chains reported convergence.

## A.2 AH1N1 data

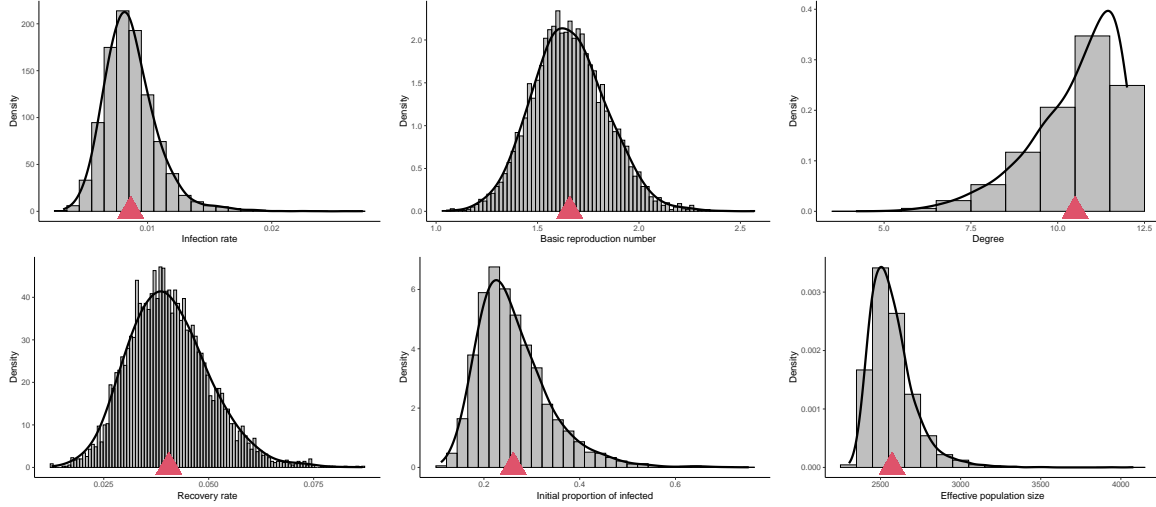
Next, we present the results of analysis of the AH1N1 dataset using DSA method with uninformative priors. As done in main body of the paper, we will present two sets of results. At first, we show results when the entire dataset is considered. See Figure 18. Next, we will present results when data for only the first 43 days are considered. See Figure 19. In case of the smaller dataset, the results are comparable to those obtained earlier when an informative prior was used. As before, when the full dataset is used, the method performs much worse than when only the first 43 days are considered. When the full dataset is used, the estimates of the effective population size are much smaller than the corresponding estimates under an informative prior, again suggesting poorer inference quality in the presence of tail noise.

## A.3 Synthetic data based on Gillespie simulation

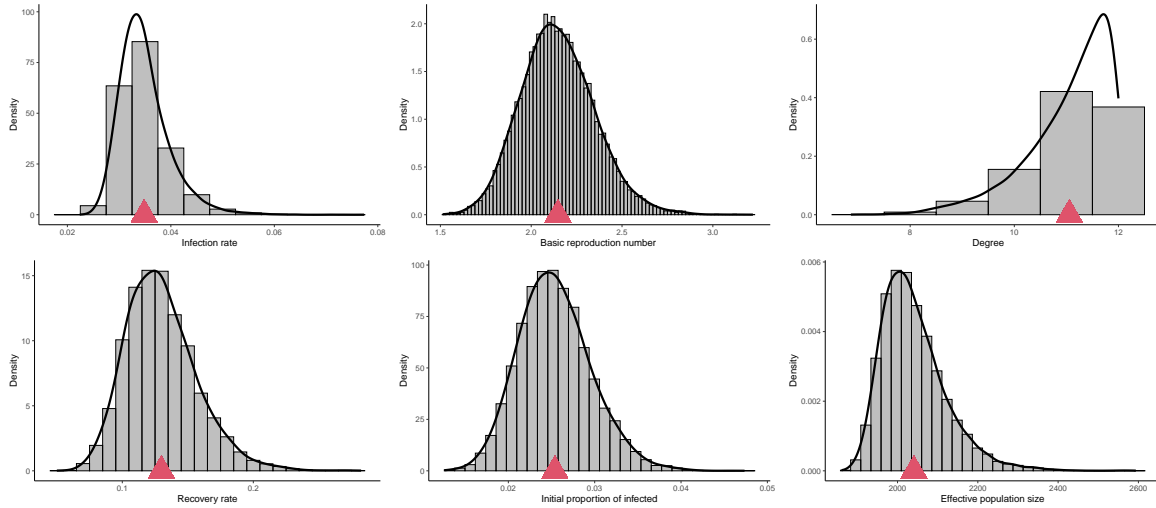
In the analysis of the synthetic data, we used flat, uninformative priors except for shorter domains of  $(0.5, 10)$  for  $R_0$ ,  $(0.03, 0.3)$  for  $\gamma$ ,  $(0, 0.3)$  for  $\rho$  and an upper bound of 12 on  $n$ . Here, we present DSA results based on flat priors on the whole region of validity for the parameters (i.e.,  $[0, \infty)$  for  $R_0$ ,  $[1, \infty)$  for  $n$ ,  $[0, \infty)$  for  $\gamma$ , and  $[0, 1]$  for  $\rho$ ) and show that such choices do not affect the quality of posterior inference. The means and the medians of the posterior distributions of the parameters  $(\tau, R_0, n, \gamma, \rho, n_T)$  are  $(0.109, 1.95, 5.62, 0.144, 0.0002, 10050)$  and  $(0.108, 1.94, 5.55, 0.144, 0.0002, 10047)$ , respectively. See Figure 20. Note that the posterior distributions are similar to those obtained earlier when flat priors over shorter intervals were used.

## B Data availability statement

H1N1 outbreak data is available at <https://github.com/cbskust/SDS.Epidemic>, data about the third COVID19 wave in India can be found at <https://data.covid19india.org/documentation/>

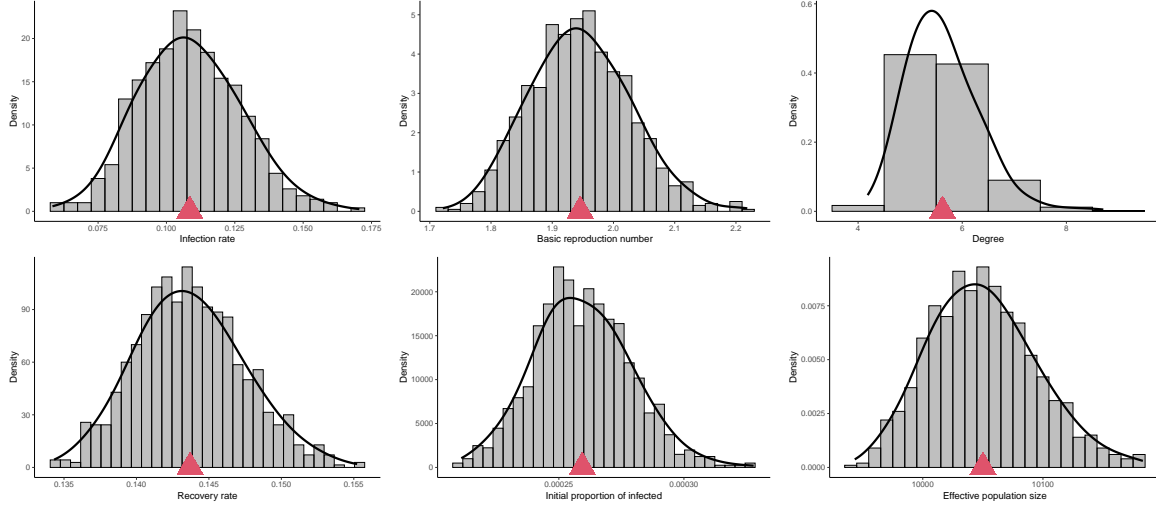


**Figure 18:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using DSA on the full A(H1N1) outbreak data. The means and the medians of the posterior distributions are  $(0.008, 1.66, 10.4, 0.04, 0.261, 2572)$  and  $(0.008, 1.65, 10.8, 0.039, 0.247, 2545)$ , respectively. There is noticeable difference between these estimates and the means  $(0.0373, 0.9880, 8.369, 0.0255, 10146)$  and medians  $(0.0269, 0.9892, 8.665, 0.0264, 9286)$  of  $(\tau, R_0, n, \rho, n_T)$  using DSA on the full A(H1N1) outbreak data with informative priors.



**Figure 19:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using DSA on the A(H1N1) outbreak data restricted to time horizon  $T = 42$ . The means and medians of the posterior distributions are  $(0.034, 2.15, 11.1, 0.13, 0.025, 2041)$  and  $(0.0341, 2.11, 10.1, 0.127, 0.025, 2027)$ , respectively. They are comparable to the means  $(0.0437, 1.843, 10.650, 0.0189, 2179)$  and medians  $(0.0418, 1.845, 10.908, 0.0189, 2177)$  of  $(\tau, R_0, n, \rho, n_T)$  obtained when informative priors are used.

csv/. All other datasets generated during and/or analysed during the current study are available from the corresponding author on reasonable request.



**Figure 20:** Posterior distributions of  $(\tau, R_0, n, \gamma, \rho, n_T)$  using the DSA method on the synthetic data. The red triangles indicate the true values of the parameter.

## C Code availability statement

The DSA scripts are available at <https://github.com/wasiur/PairwiseDSA.git>. The scripts for the MLE-based method are available at <https://github.com/berthouz/EpiPWMInf> for fitting the ODE realisations with negative binomial noise, <https://github.com/berthouz/EpiPWMInfwIO> for fitting the Gillespie stochastic realisations and <https://github.com/berthouz/EpiPWMInfwION> for fitting the real-world datasets.

## D Acknowledgements

L. Berthouze and I.Z. Kiss acknowledge support from the Leverhulme Trust for the Research Project Grant RPG-2017-370. The authors thank Prof Theodore Kypraios for useful discussions about approximate likelihoods.

## References

- [1] M Akian, L Ganassali, S Gaubert, and L Massoulié. Probabilistic and mean-field model of COVID-19 epidemics with user mobility and contact tracing. Technical report, 2020.
- [2] Håkan Andersson and Tom Britton. *Stochastic Epidemic Models and Their Statistical Analysis*. Springer New York, 2000.
- [3] Yi-Cheng Chen, Ping-En Lu, Cheng-Shang Chang, and Tzu-Hsuan Liu. A Time-dependent SIR model for COVID-19 with Undetectable Infected Persons. Technical report.
- [4] Kai Cui, Wasir R. KhudaBukhsh, and Heinz Koepl. Motif-based mean-field approximation of interacting particles on clustered networks. *Physical Review E*, 105(4), April 2022.
- [5] Gareth Davies. The foot and mouth disease (fmd) epidemic in the united kingdom 2001. *Comparative immunology, microbiology and infectious diseases*, 25(5-6):331–343, 2002.

- [6] Francesco Di Lauro, Wasiur R. KhudaBukhsh, István Z. Kiss, Eben Kenah, Max Jensen, and Grzegorz A. Rempała. Dynamic survival analysis for non-markovian epidemic models. *Journal of The Royal Society Interface*, 19(191):20220124, 2022.
- [7] Jonah Gabry and Rok Češnovar. CmdStanR, 2023. R package version 0.5.2.
- [8] Daniel T. Gillespie. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics*, 22(4):403–434, 1976.
- [9] Istvan Z. Kiss, Joel C. Miller, and Peter L. Simon. *Mathematics of epidemics on networks: from exact to approximate models*. 2016.
- [10] Matthew J Keeling. The effects of local spatial structure on epidemiological invasions. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 266(1421):859–867, 1999.
- [11] Wasiur R. KhudaBukhsh, Caleb Deen Bastian, Matthew Wascher, Colin Klaus, Saumya Yashmohini Sahai, Mark H. Weir, Eben Kenah, Elisabeth Root, Joseph H. Tien, and Grzegorz A. Rempała. Projecting COVID-19 cases and hospital burden in Ohio. *Journal of Theoretical Biology*, 561:111404, 2023.
- [12] Wasiur R. KhudaBukhsh, Boseung Choi, Eben Kenah, and Grzegorz A. Rempała. Survival dynamical systems: individual-level survival analysis from population-level epidemic models. *Interface Focus*, 10(1):20190048, 2020.
- [13] Wasiur R. Khudabukhsh, Casper Woroszylo, Grzegorz A. Rempała, and Heinz Koepl. A functional central limit theorem for SI processes on configuration model graphs. *Advances in Applied Probability*, 54(3):880–912, 2022.
- [14] Aaron A. King, Matthieu Domenech De Cellés, Felicia M.G. Magpantay, and Pejman Rohani. Avoidable errors in the modelling of outbreaks of emerging pathogens, with special reference to Ebola. *Proceedings of the Royal Society B: Biological Sciences*, 282(1806):0–6, 2015.
- [15] István Z. Kiss, Eben Kenah, and Grzegorz A. Rempała. Necessary and sufficient conditions for exact closures of epidemic equations on configuration model networks. *Journal of Mathematical Biology*, 87(2), August 2023.
- [16] István Z Kiss, Joel C Miller, Péter L Simon, et al. Mathematics of epidemics on networks. *Cham: Springer*, 598:31, 2017.
- [17] Philip E Paré and Carolyn L Beck. Modeling, estimation, and analysis of epidemics over networks: An overview. 2020.
- [18] Romualdo Pastor-Satorras, Claudio Castellano, Piet Van Mieghem, and Alessandro Vespignani. Epidemic processes in complex networks. *Reviews of modern physics*, 87(3):925, 2015.
- [19] Lorenzo Pellis, Frank Ball, Shweta Bansal, Ken Eames, Thomas House, Valerie Isham, and Pieter Trapman. Eight challenges for network epidemic models. *Epidemics*, 10:58–62, mar 2015.
- [20] Mason A. Porter and James P. Gleeson. Dynamical Systems on Networks. pages 49–51. 2016.
- [21] ReyerGerlagh. Closed-Form Solutions for Optimal Social Distancing in a SIR Model of COVID-19 Suppression Reyer Gerlagh. Technical report, 2020.

- [22] Elissa J Schwartz, Boseung Choi, and Grzegorz A Rempala. Estimating epidemic parameters. *Mathematical Biosciences*, 2014.
- [23] Neil Sherborne, Joel C. Miller, Konstantin B. Blyuss, and Istvan Z. Kiss. Mean-field models for non-Markovian epidemics on networks. *Journal of Mathematical Biology*, 76(3):755–778, feb 2018.
- [24] Stan Development Team. RStan: the R interface to Stan, 2023. R package version 2.21.8.
- [25] Harley Vossler, Pierre Akilimali, Yuhan Pan, Wasiur R. KhudaBukhsh, Eben Kenah, and Grzegorz A. Rempala. Analysis of individual-level data from 2018–2020 ebola outbreak in democratic republic of the congo. *Scientific Reports*, 12(1):5534, 2022.