# Probabilistic predictions of SIS epidemics on networks based on population-level observations

T. Zerenner[a]*, F. Di Lauro[a,b], M. Dashti[a], L. Berthouze[c], I. Z. Kiss[a]†

[a]*Department of Mathematics, University of Sussex, Falmer, Brighton, BN1 9QH, UK*
[b]*Big Data Institute, Nuffield Department of Medicine, University of Oxford, Oxford, OX3 7FL, UK*
[c]*Department of Informatics, University of Sussex, Falmer, Brighton, BN1 9QH, UK*

November 11, 2021

## Abstract

We predict the future course of ongoing susceptible-infected-susceptible (SIS) epidemics on regular, Erdős-Rényi and Barabási-Albert networks. It is known that the contact network influences the spread of an epidemic within a population. Therefore, observations of an epidemic, in this case at the population-level, contain information about the underlying network. This information, in turn, is useful for predicting the future course of an ongoing epidemic. To exploit this in a prediction framework, the exact high-dimensional stochastic model of an SIS epidemic on a network is approximated by a lower-dimensional surrogate model. The surrogate model is based on a birth-and-death process; the effect of the underlying network is described by a parametric model for the birth rates. We demonstrate empirically that the surrogate model captures the intrinsic stochasticity of the epidemic once it reaches a point from which it will not die out. Bayesian parameter inference allows for uncertainty about the model parameters and the class of the underlying network to be incorporated directly into probabilistic predictions. An evaluation of a number of scenarios shows that in most cases the resulting prediction intervals adequately quantify the prediction uncertainty. As long as the population-level data is available over a long-enough period, even if not sampled frequently, the model leads to excellent predictions where the underlying network is correctly identified and prediction uncertainty mainly reflects the intrinsic stochasticity of the spreading epidemic. For predictions inferred from shorter observational periods, uncertainty about parameters and network class dominate prediction uncertainty. The proposed method relies on minimal data at population-level, which is always likely to be available. This, combined with its numerical efficiency, makes the proposed method attractive to be used either as a standalone inference and prediction scheme or in conjunction with other inference and/or predictive models.

## 1 Introduction

Mathematical models of the dynamics of directly transmitted infectious diseases can provide predictions about the future course of an ongoing epidemic and hence aid in decision making and epidemic control (e.g. Siettos and Russo, 2013). Statistical time series methods are utilised predominantly for epidemic nowcasting, i.e., shortest-term predictions and current state assessment under not yet complete data (e.g. Bastos et al., 2019; McGough

---

*t.zerenner@gmail.com

†i.z.kiss@sussex.ac.uk

et al., 2020) and epidemic surveillance, i.e., early identification of emerging epidemics (Unkel et al., 2012). Mechanistic/state-space models, which are based on a mathematical description of the spreading process, allow one to make predictions about the future course of an epidemic (Shaman and Karspeck, 2012; Tizzoni et al., 2012; Nsoesie et al., 2013) as well as simulating intervention strategies (Chao et al., 2011; Di Lauro et al., 2021b,a; Van Yperen et al., 2020). Many such models rely on a compartmentalisation of a population of $N$ individuals according to the individual's disease status. In diseases for which there is no immunity upon recovery, each individual is either susceptible (S) or infected/infectious (I) at any given time.

One assumption common to many compartmental models is that of random mixing of individuals, or of/within subgroups of a population (e.g., Kermack and McKendrick, 1927; Jacquez and Simon, 1993). While this assumption can be adequate in some instances (e.g., within households; Goeyvaerts et al., 2018), it is known that populations do not mix at random in general. Rather individuals have a finite set of contacts to whom they can pass on an infection. It is well-established that the contact network of a population significantly impacts epidemic dynamics (e.g., Shirley and Rushton, 2005; Yin et al., 2017). The importance of contact structure for epidemic dynamics has led to a close interaction between network science and mathematical epidemiology (Keeling and Eames, 2005; Danon et al., 2011; Pastor-Satorras et al., 2015; Kiss et al., 2017) whereby the spread of an epidemic within a population is understood and modelled as a stochastic process on a network. In such a model, each individual in the population corresponds to a node in the network, and a contact that represents a potential route for disease transmission between two individuals is a link in the network.

Drawbacks of network epidemiological models typically include their high dimensionality and the inaccessibility of the exact contact network of a population. Consider a SIS-epidemic on an undirected, unweighted network with $N$ nodes. At any given time, each node is either susceptible or infected/infectious. If the exact contact network is static and known, a complete description of the SIS dynamics is given by a continuous-time Markov-chain of dimension $2^N$ (one equation for each possible network state; e.g., Simon et al., 2011). Such a Markov-chain model is exact, but also high-dimensional even for modest values $N$. Hence, the numerical integration of the system of equations becomes unfeasible for most real-life networks. Consequently, analytical results based on the exact system are mostly out of reach, and existing results typically rely on mean-field approximations (e.g., Mata and Ferreira, 2013; Cota et al., 2018). Further, the exact contact network of a population is rarely accessible, but usually needs to be approximated either from limited observations and/or based on theoretical network models (e.g., Della Rossa et al., 2020; Xue et al., 2020).

In this study, we explore the suitability of a computationally inexpensive model to describe the stochastic process of an SIS epidemic spreading on Regular (Reg), Erdős–Rényi (ER) and Barábasi-Albert (BA) networks. The surrogate model utilised in this study was first introduced in Di Lauro et al. (2020b) and further expanded to include more network classes and consider the large $N$ limit in Di Lauro et al. (2020a). The core idea of the approach is a dimension reduction of the state space. In the surrogate model, the state of the epidemic at any given time is defined by the total number of infected nodes in the population. The effect of the contact structure on the spreading of the epidemic is accounted for by the model parameters. The continuous-time Markov-chain describing the SIS dynamics on the reduced state space takes the form of a Birth-and-Death (BD) process and is of dimension $N + 1$; that is, it is linear in $N$ and thus feasible also for large $N$. The parameters of the BD model correspond to recovery and infections rates. The

recovery rate is network-independent and here assumed to be known. The rate at which new infections occur depends on the network, but for particular network classes it can be well described by a three-parameter model, reducing the number of free parameters of the BD model from $2(N + 1)$ to only three. Di Lauro et al. (2020b) utilised the finding that different types of networks are associated with distinct regions in the space spanned by the three parameters to infer the type of network from population-level observations. To solve this inverse problem Di Lauro et al. (2020b) set up a Bayesian inference procedure and built network class specific prior distributions which then allow to identify the most likely network class from the posterior.

Here, we utilise the BD model to forecast the evolution of an on-going epidemic. We address the following questions:

- How well does the BD model capture the intrinsic stochasticity of an epidemic spreading on a network?

- How uncertain are model parameters when inferred from the kind of time-censored observations typically available in a realistic prediction scenario, and how does this uncertainty translate into prediction uncertainty?

- Can we use the BD model and Bayesian inference to provide epidemic forecasts with meaningful uncertainty information?

The manuscript is structured as follows: Section 2 introduces the BD model and Bayesian inference. We then outline the generation and evaluation of predictions using the BD model. We consider nine different combinations of networks and epidemic parameters: a small, a medium and a large epidemic on a network from each of the three aforementioned network classes (Section 3). An empirical validation of the BD model based on these nine cases is provided in Section 4.1. In Section 4.2, we evaluate the predictions obtained with the BD model for all nine cases. In particular, we study the sensitivity of network class and parameter inference on the number and timing of observations in realistic prediction scenarios and how uncertainty about network class and model parameters translates into prediction uncertainty. We conclude with a discussion including limitations of this work and future directions.

## 2 Methods

### 2.1 SIS epidemics on networks

We consider the standard SIS epidemic on a population of $N$ individuals whose contact structure is described by an undirected and unweighted network defined by its adjacency matrix $G = (g_{ij})$ with $i, j = 1, 2, \ldots N$ and $g_{ij} = 1$ if individuals (nodes) $i$ and $j$ are connected and $g_{ij} = 0$ otherwise. If two nodes $i$ and $j$ are connected, the disease can be transmitted from one to the other. In an SIS epidemic, each node is, at any given time, either susceptible (S) or infected/infectious (I). Thus, the epidemic state of the network at time $t$ is described by a Boolean vector $X(t) = (x_i(t))$ with $i = 1, 2, \ldots N$ where $x_i(t) = 0$ if node $i$ is susceptible and $x_i(t) = 1$ if node $i$ is infected at time $t$. Hence, there exists a total of $2^N$ distinct network states. The state of the network changes through two types of events: the recovery or the infection of a node. Infection and recovery are Markovian and act as homogeneous Poisson point processes with constant per-link infection rate $\tau$ and constant recovery rate $\gamma$. An infectious node can spread the infection only to neighbouring susceptible nodes. Infection dynamics thus depends on the network structure,

while recovery is network-independent. A complete description of the SIS-dynamics on a given network corresponds to a Markov-chain over a state-space of dimension $2^N$ for which numerical integration becomes intractable even for modest values of $N$. However, given network adjacency $G$, epidemic parameters $\tau$ and $\gamma$, and initial conditions $X_0 = X(t_0)$, realisations of the stochastic process can be readily obtained using the Gillespie algorithm (e.g., Gillespie, 1977; Kiss et al., 2017). This is computationally comparatively inexpensive and provides us with i.i.d. samples of the true stochastic process. Such samples serve as a reference for the validation of the BD model and for the evaluation of predictions. More precisely, we make use of the aggregated number of infected nodes in the network, i.e., $I(t) = \sum_{i=1}^{N} x_i(t)$, which describes the epidemic at population level.

## 2.2 BD model

Birth-and-death processes are intuitively linked to the population-level dynamics of SIS epidemics (e.g., Ganesh et al., 2005; Nagy et al., 2014; Devriendt and Van Mieghem, 2017). In this view, an increase in the number of infected individuals ($I \to I + 1$) corresponds to the 'birth of an infection'; a decrease ($I \to I - 1$) to the 'death of an infection'. Accordingly, the epidemic state is defined by the number of infected nodes $I \in \{0, \ldots, N\}$ in the network. The resulting model is, like the exact formulation (Sec. 2.1), a continuous-time Markov chain, but on a state space of dimension $N + 1$ only.

The Kolmogorov (or Master) equation of a standard BD process is given by

$$\forall k \in \{0, \ldots, N\}, \; \dot{p}_k(t) = a_{k-1} \, p_{k-1}(t) - (a_k + c_k) \, p_k(t) + c_{k+1} \, p_{k+1}(t), \tag{1}$$

where $p_k(t)$ denotes the probability of observing $k$ infected nodes at time $t$, and $a_k$ and $c_k$ denote population-level infection and recovery rate, respectively. We note that $a_{-1} = c_{N+1} = 0$. The population-level recovery rate $c_k$ can be directly obtained from the node recovery rate $\gamma$ as $c_k = \gamma k$. The population-level infection rate $a_k$ however depends on the number of links between susceptible and infected nodes (S-I links) present in the network in its current state and is thus a random variable depending on the precise network. Following Di Lauro et al. (2020b), $a_k$ is represented in the BD model by its expectation $\hat{a}_k = \tau \times$ *the time-averaged number of S-I links over the network states with $k$ infected nodes*. Di Lauro et al. (2020b) further demonstrated that for Regular, Erdős–Rényi and Barabási-Albert networks $\hat{a}_k$ can be well represented by a three-parameter model of the form

$$\forall k \in \{0, \ldots, N\}, \; a_k(C, \alpha, p) = C \, k^p \, (N - k)^p \left( \alpha \left( k - \frac{N}{2} \right) + N \right), \tag{2}$$

with $C$ serving as a general scaling parameter, $\alpha$ allowing to shift the peak of the curve with respect to $k = N/2$ and $p$ adjusting the flatness of the curve. The shape of the $a_k$ curves is network class-specific (Fig. 1). Whilst the peak is located near the centre ($k = N/2$) for Erdős–Rényi networks, it is shifted to the right for Regular networks, and to the left for Barabási-Albert networks. Accordingly, the $(C, \alpha, p)$-triplets for the different network types cluster in different regions in the three-dimensional parameter space. This observation is central to the network class inference of Di Lauro et al. (2020b).

For a given $\gamma$, a given $(C, \alpha, p)$-triplet, and initial conditions $p_k(t_0) = 1$ if $k = I(t_0)$ and $p_k(t_0) = 0$ otherwise, where $I(t_0) \in \{0, \ldots, N\}$ denotes the number of infected nodes at $t_0$, we can numerically integrate Eq. 1 to obtain predictions $p_k(t)$. In our experiments, we assume that the recovery rate $\gamma$ is known. Initial conditions are provided by the last
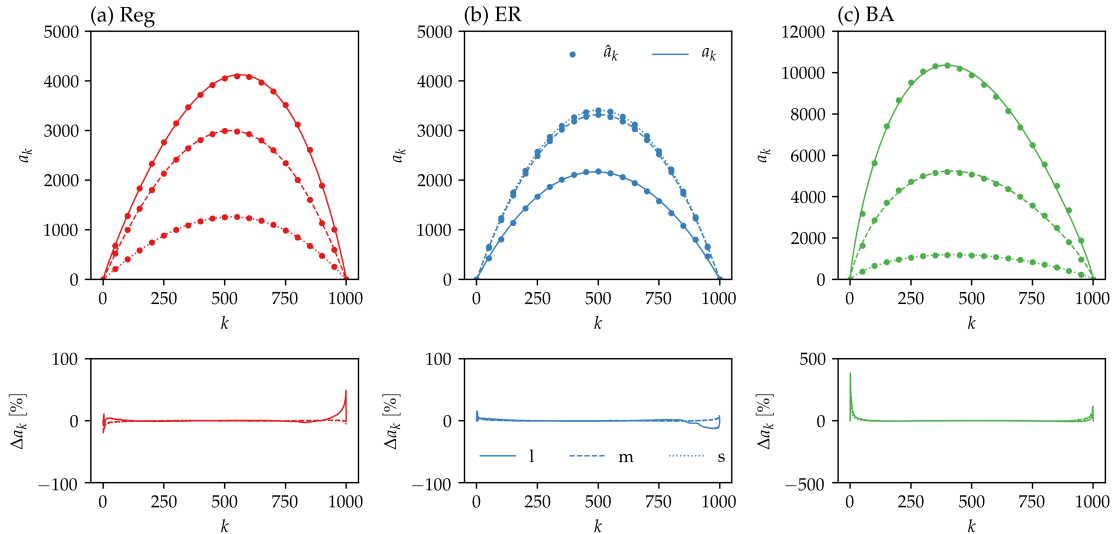
Figure 1: Parametric $a_k(C, \alpha, p)$ model fitted to $\hat{a}_k$ from Gillespie simulations of large (l), medium (m) and small (s) epidemics on the different networks. In the top panels, the dots indicate $\hat{a}_k$ and the lines correspond to $a_k(C, \alpha, p)$ with parameters $(C, \alpha, p)$ obtained from a least-squares fit of Eq. 2 to $\hat{a}_k$. The bottom panels show the relative error $\Delta a_k = (a_k(C, \alpha, p) - \hat{a}_k)/(\hat{a}_k + 1) \times 100$. (Adding one in the denominator allows to also include the values at and in the vicinity of $k = 0$ and $k = 1000$.)

observation available. Thus, the remaining task is the inference of $(C, \alpha, p)$ from the available observational data, here, the number of infected nodes at a set of discrete times. We use Bayesian inference to estimate the posterior distribution over the parameters $(C, \alpha, p)$ given observations. The required priors were derived from extensive Gillespie simulations on different networks and with different epidemic parameters. The particular challenge for making predictions lies in the limited observational period available for inference in a realistic prediction scenario, in which observations exist only up to the current state of the epidemic.

### 2.2.1 Bayesian inference

The detail of the inference framework can be found in Di Lauro et al. (2020b). Here, we only recall the main ideas and steps of the inference procedure. We denote the population level observations by $(y, s)$ where $y = (k_1, ...., k_n)$ with $k_j \in \{0, ..., N\}$ denotes the number of infected individuals at times $s = (t_1, ...., t_n)$. For brevity, we use $u$ to denote $(C, \alpha, p)$. We further denote the set of candidate network classes as $\Theta = \{\text{Reg, ER, BA}\}$.

In order to make predictions, we require the posterior over $u$ given the observations, i.e., $\pi(u|y, s)$, which we can write as

$$\pi(u|y, s) = \sum_{\Theta} \pi_\theta(u|y, s) \ \pi(\theta|y, s).$$

In Di Lauro et al. (2020b), the goal was network class inference, i.e., obtaining the posterior over $\Theta$ given observations $(y, s)$, $\pi(\theta|y, s)$. To this end, Di Lauro et al. (2020b) generated network class specific priors $\pi_{0,\theta}(u)$. Precisely, they carried out a large number of Gillespie simulations during which they kept track of the number of infected nodes $k$, the number of S-I links in the respective network states as well as the time spent in the various states.

The parametric $a_k(C, \alpha, p)$ model from Eq. 2 was then fitted to the $(k, \hat{a}_k)$ curves from the Gillespie simulations by a least-squares fit using a particle swarm algorithm (Kennedy and Eberhart, 1995). The resulting $(C, \alpha, p)$ triplets were used to infer Gaussian kernel density estimators (Pedregosa et al., 2011) for the priors $\pi_{0,\theta}(u)$. Assuming a non informative, uniform prior for network class $\theta$, the prior distribution over $\theta$ and $u$ is given by

$$\pi_0(u, \theta) = \frac{1}{3}\pi_{0,\theta}(u).$$

Employing Bayes' rule we obtain the network class specific posterior(s) over the parameter space as

$$\pi_\theta(u|y, s) \propto \mathcal{L}^u(y, s) \, \pi_{0,\theta}(u). \tag{3}$$

and the posterior over the network classes as

$$\pi(\theta|y, s) = \int \pi(u, \theta|y, s)du$$
$$\propto \int \mathcal{L}^u(y, s)\pi_{0,\theta}(u)du,$$

where $\mathcal{L}^u(y, s)$ denotes the likelihood of the observations under the forward model from Eq. 1 which is given by

$$\mathcal{L}^u(y, s) = \prod_{i=1}^{n-1} p_{k_i, k_{i+1}}^u (t_{i+1} - t_i). \tag{4}$$

Following Di Lauro et al. (2020b), the terms $p_{k_i, k_{i+1}}^u$ are computed using the algorithm from Crawford et al. (2014). The Python implementation routine estimating $\pi(\theta|y, s)$ is available at `https://github.com/BayIAnet/NetworkInferenceFromPopulationLevelData`. To estimate $\pi_\theta(u|y, s)$, we draw samples from $\pi_\theta(u|y, s)$ using the Metropolis–Hastings algorithm, making use of Eqs. 3 and 4.

### 2.2.2   Prediction and uncertainty

To obtain predictions one needs to integrate Eq. 1 with the parameters $u = (C, \alpha, p)$ obtained from the posterior $\pi_\theta(u|y, s)$. We generate and evaluate two different types of predictions. The first variant incorporates information on prediction uncertainty as encoded in $\pi(\theta)$ and $\pi_\theta(u)$. The second variant is based on a point estimate of $u$. The Python routine for generating the predictions will be made available at `https://github.com/tzerenner/EpidemicPredictionsFromPopulationLevelData`.

For brevity, in the following, we neglect the dependence on $t$ and instead consider some fixed point in time. We denote by $\nu_{\theta,k}$ the pushforward measure of $\pi_\theta$ under the forward solution operator $G_k : u \mapsto p_k$ defined by Eq. 1. The density of $\nu_{\theta,k}$ is then related to that of $\pi_\theta$ through $\nu_{\theta,k}(p_k) = \pi_\theta(G_k^{-1}(p_k))$. The conditional mean of $\nu_{\theta,k}$ is given by

$$m_k^\theta = \int p_k \, \nu_{\theta,k}(p_k) \, dp_k = \int G_k(u) \, \pi_\theta(u) \, du. \tag{5}$$

When additionally integrating over all network classes $\Theta$, we can further obtain the conditional mean of $\nu_k$, the pushforward measure of $\pi$, as
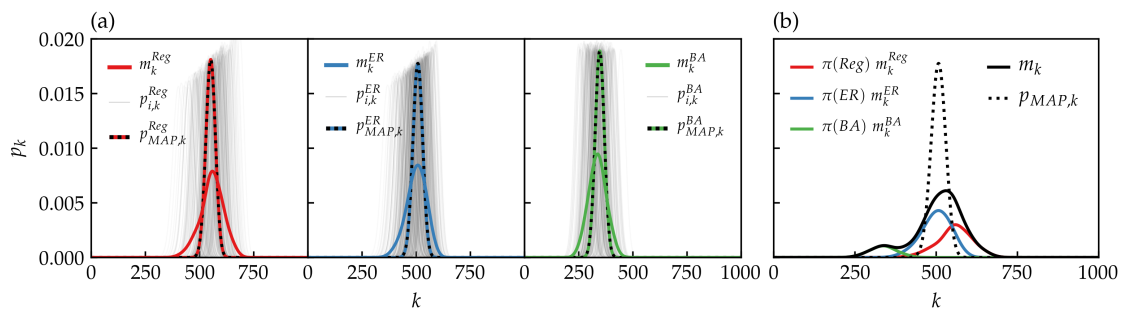
Figure 2: Illustration of the possible types of predictions for a single point in time during the growth phase of the epidemic. Panel (a) shows the conditional mean $m_k^\theta$ of the pushforward measures $\nu_{\theta,k}$ (Eq. 5) together with the pushforward of the mode of $\pi_\theta$ for $\theta \in \Theta$. The grey lines indicate the pushforward of the samples drawn from $\pi_\theta$. Panel (b) shows the resulting predictions $m_k$ (Eqs. 6,7) and $p_{MAP,k}$ (Eq. 12). The example shown here is a medium epidemic on an ER network (see Table 1) with inference based on 10 observations $y = (k_1, \ldots, k_{10})$ approximately equally spaced in time between $k_1 = 50$ and $k_{10} = 160$

$$m_k = \int p_k \left( \sum_{\theta \in \Theta} \pi(\theta)\, \nu_{\theta,k}(p_k) \right)\, dp_k = \int G_k(u) \left( \sum_{\theta \in \Theta} \pi(\theta)\pi_\theta(u) \right)\, du. \qquad (6)$$

We estimate $m_k$ from samples $(C, \alpha, p)_{i,1 \leq i \leq n}$ which we draw from the posterior distributions $\pi_\theta(C, \alpha, p)$, $\theta \in \Theta$, using the Metropolis–Hastings algorithm. Since integrating the Master equation with a large number of parameter combinations is computationally demanding, we thin the samples by including only every $i$-th draw, such that the autocorrelation between subsequent draws is $< 0.1$. To choose an appropriate size $n$ of the (thinned) sample, we consider its multivariate effective sample size (mESS), which is estimated by

$$\widehat{\mathrm{mESS}} = n \left( \frac{\det(\Lambda_n)}{\det(\Sigma_n)} \right)^{1/3},$$

where $\Lambda_n$ denotes the sample covariance and $\Sigma_n$ denotes the multivariate batch mean estimator of the covariance matrix in the Markov chain central limit theorem (Roy, 2020; Vats et al., 2020). The sample size $n$ is chosen to be the minimum $n$ such that $\widehat{\mathrm{mESS}} \geq 260$ which ensures a confidence level of $\delta = 0.1$ and a tolerance level of $\epsilon = 0.25$ for the expectation in the three parameter $(C, \alpha, p)$-space (Vats et al., 2019). To compute the mESS and the threshold for the desired confidence and tolerance levels, we used the Python implementation available at `https://github.com/Gabriel-p/multiESS`.
We then proceed to integrate the Master equation with each $(C, \alpha, p)$-triplet to obtain $p_{i,k}^\theta$ and finally approximate the conditional mean by

$$m_k = \sum_{\theta \in \Theta} \left( \pi(\theta)\, \sum_i \frac{p_{i,k}^\theta}{n} \right), \qquad (7)$$

as well as the respective cumulative density over $k$ by

$$M_k = \sum_{x \leq k} m_x. \qquad (8)$$

7

From the latter we obtain equal-tailed credible intervals for the predicted number of infected nodes. Equal-tailed intervals are defined such that the probability of being below the interval is as high as being above the interval and thus can be directly obtained from the quantiles of the cumulative density as

$$Q(x) = \inf\{k \in \{0, \ldots, 1000\} : x \le M_k\}. \tag{9}$$

The interval $[Q(0.05), Q(0.95)]$ for example corresponds to the 90% equal-tailed interval. When obtained from $m_k$ (Eq. 6), such intervals incorporate both prediction uncertainty arising from uncertainty about parameters and network class as encoded in $\pi_\theta(u)$ and $\pi(\theta)$, respectively, as well as prediction uncertainty arising from the intrinsic stochasticity of the epidemic spreading.

The second prediction variant is based on a point estimate, i.e., a single $u = (C, \alpha, p)$ inferred from the posterior $\pi$. We first identify the most likely network class, i.e, the mode of $\pi(\theta)$,

$$\hat{\theta}_{MAP} = \mathrm{argmax}_\theta\{\pi(\theta|y,s)\}, \tag{10}$$

where MAP stands for maximum a-posteriori, and then estimate the mode of $\pi_{\hat{\theta}}$,

$$\hat{u}_{\hat{\theta},MAP} = \mathrm{argmax}_u\{\pi_{\hat{\theta}}(u|y,s)\}, \tag{11}$$

using a combination of global and local optimisation routines (Di Lauro et al., 2020b). The predictions $p_{MAP,k}$ are obtained by integrating the Master equation with $\hat{u}_{\hat{\theta},MAP}$, i.e.,

$$p_{MAP,k} = G_k(\hat{u}_{\hat{\theta},MAP}). \tag{12}$$

Again, we compute the respective cumulative density over $k$ as

$$P_{MAP,k} = \sum_{x \le k} p_{MAP,x}, \tag{13}$$

from which we can obtain quantiles and equal-tailed prediction intervals. Such intervals are not credible intervals in the Bayesian sense, but solely represent the intrinsic stochasticity of the epidemic spreading. They are thus systematically narrower than the credible intervals discussed above. An illustration of the two prediction variants for one single point in time is provided in Fig. 2.

To compare uncertainty in the $p_k$-space for the two prediction variants we further consider the covariance of the pushforward $G : u \mapsto (p_0, p_1, \ldots p_k, \ldots p_N)^T$, $G(u) = (G_1(u), G_2(u), \ldots G_k(u), \ldots G_N(u)$ $\mathbb{R}^{(N+1)}$. We denote the mean of the pushforward of $\pi_\theta(u)$ under the forward solution operator $G$ by $m = (m_0, m_1, \ldots m_k, \ldots m_N)^T \in \mathbb{R}^{(N+1)}$ (Eq. 6). Its covariance $\mathcal{C} = (\mathcal{C}_{kl}) \in \mathbb{R}^{(N+1)\times(N+1)}$ is given by the outer product

$$\mathcal{C} = \int (G(u) - m)(G(u) - m)^T \, \pi_\theta(u) \, du.$$

We estimate $\mathcal{C}$ from the discrete samples (Eq. 7) as

$$\mathcal{C}_{kl} = \sum_{\theta \in \Theta} \left( \frac{1}{n-1} \sum_{i=1}^{n} (p_{k,i} - m_k)(p_{l,i} - m_l) \, \pi(\theta) \right), \quad k = 0, \ldots N, \ l = 0 \ldots N.$$

We can then evaluate

$$|\mathcal{C}| = \left| \int (G(u) - m)(G(u) - m)^T \pi_\theta(u) \, du \right|, \tag{14}$$

where $|\cdot|$ denotes the Euclidean norm in $\mathbb{R}^{(N+1)\times(N+1)}$. To accordingly evaluate the uncertainty of the predictions based on the point estimator from Eq. 11, we further evaluate

$$|\mathcal{C}_{MAP}| = \left| \int (G(u) - p_{MAP})(G(u) - p_{MAP})^T \pi_\theta(u) \, du \right|, \tag{15}$$

where $p_{MAP} = (p_{MAP,0}, p_{MAP,1}, \ldots m_{MAP,k}, \ldots m_{MAP,N})^T \in \mathbb{R}^{(N+1)}$ are the predictions from Eq. 12.

## 2.3 Performance assessment

We generate a reference for evaluating the BD model by carrying out a set of Gillespie simulations which provides us with an i.i.d. sample of the stochastic spreading of an SIS epidemic on a given network. We chose a sample size of 1000 and denote the set of 1000 epidemic trajectories obtained from the sample as $\{I_{r,i}(t)\}_{1 \le i \le 1000}$. To empirically validate the BD model, we compare the $p_k(t)$ obtained from the numerical integration of Eq. 1 against the reference. We evaluate the difference in expectation, i.e.,

$$\Delta \hat{I}(t) = \hat{I}_s(t) - \hat{I}_r(t) = \sum_{k=0}^{N} k \, p_k(t) - \sum_{i=1}^{1000} \frac{I_{r,i}(t)}{1000}, \tag{16}$$

where $\hat{I}_r(t)$ denotes the mean over the reference at time $t$ and $\hat{I}_s(t)$ denotes the mean number of infected nodes at time $t$ predicted by the BD model. We further evaluate the cumulative densities from the BD model in comparison to our reference using the integrated quadratic distance (IQD) which is given by

$$\text{IQD}(t) = \int_{-\infty}^{\infty} (F_{s,t}(x) - F_{r,t}(x))^2 \, dx. \tag{17}$$

Here $F_{s,t}(x)$ denotes the cumulative density at time $t$ predicted by the BD model,

$$F_{s,t}(x) = \sum_{k \le x} p_k(t),$$

and $F_{r,t}(x)$ denotes the reference empirical cumulative density,

$$F_{e,t}(x) = \frac{1}{1000} \sum_{i=1}^{1000} \mathbf{1}_{I_{r,i}(t) \le x},$$

where $\mathbf{1}$ is an indicator function equal to 1 if condition $I_{r,i}(t) \le x$ is true and 0 otherwise. A small IQD is obtained when not only the mean but also the intrinsic stochasticity of the epidemic spreading is adequately represented by the BD model. The goodness of fit between the cumulative densities is further illustrated by quantile-quantile plots. The quantile function is the inverse of the cumulative density, i.e., $Q_{s/r,t} = F_{s/r,t}^{-1}$, and hence given by

$$Q_{s/r,t}(P) = \inf\{x \in \{1, \ldots, 1000\} : P \le F_{s/r,t}(x)\}, \tag{18}$$

for the BD model (s), and the reference (r), respectively.

# 3    Data

We consider nine different network and epidemic parameter combinations. The network parameters, network class and mean degree $\langle k \rangle$, and the epidemic parameters, per link infection rate $\tau$ and recovery rate $\gamma$, are summarised in Table 1. Gillespie simulations were performed on a network of $N = 1000$ nodes and initialised with five infected nodes selected at random. The time $T$ is an approximate value of the time span between initialising the simulation and reaching quasi-steady state and in the following serves as a universal time scale which allows to plot the different cases onto the same time axis. Time is unit-free here. The simulated data can be re-scaled to physically-meaningful time scales by applying an appropriate multiplicative factor to the simulation time $t$. The parameters for the nine cases were chosen such that we obtained one large epidemic with $> 70\%$ of the population infected at quasi-steady state, one medium epidemic with $40\%$ to $60\%$ of the population infected at quasi-steady state, and one small epidemic with $< 40\%$ of the population infected at quasi-steady state for each network class.

In this study, we consider the epidemics at population-level, that is, we aim to predict the future course of the number of infected nodes. Network class and parameters of the BD model are inferred from population-level observations of the number of infected nodes at a set of discrete points in time.

| case | $\langle k \rangle$ | $\tau$ | $\gamma$ | $T$ |
|------|------|------|------|------|
| Reg l | 5 | 4.251 | 2.969 | 0.75 |
| Reg m | 10 | 1.265 | 5.773 | 1.5 |
| Reg s | 7 | 0.762 | 3.356 | 8 |
| ER l | 8.124 | 1.251 | 0.969 | 1.25 |
| ER m | 15.868 | 0.859 | 6.338 | 1.25 |
| ER s | 12.042 | 1.143 | 9.579 | 2 |
| BA l | 13.902 | 3.123 | 6.969 | 0.25 |
| BA m | 9.95 | 2.19 | 8.948 | 0.5 |
| BA s | 7.968 | 0.612 | 3.803 | 2.5 |

Table 1: Parameters of the simulated SIS epidemics on networks of $N = 1000$ nodes. Listed are the names of the different cases, which consist of the respective network class (Regular (Reg), Erdős–Rényi (ER) or Barabási-Albert (BA)) and epidemic size (large (l), medium (m), small (s)), mean node degree $\langle k \rangle$, per-link infection rate $\tau$, recovery rate $\gamma$ and the approximate time $T$ between initialisation and quasi-steady state in simulations initialised with five infected nodes selected at random.

# 4    Results

## 4.1    Validation of the BD model

We carry out a set Gillespie simulations during which we keep track of the number of infected nodes $k$, the number of S-I links over time and the time spent in the observed states. For each case, we carry out in total 200 simulations half of which are initialised with five infected nodes and half with 1000 infected nodes. With this choice of initial conditions, we obtain realisations of the random variable $a_k$ ($\tau \times$ #S-I links) for each $k = 0, \ldots, N$, from which we compute the expectations $\hat{a}_k$ following Di Lauro et al. (2020b) as

| case | $C \times 10^4$ | $\alpha$ | $p$ | RMSE |
|------|------|------|------|------|
| Reg l | 0.469 | 0.475 | 0.915 | 85.10 |
| Reg m | 0.110 | 0.182 | 1.007 | 6.90 |
| Reg s | 0.039 | 0.224 | 1.019 | 6.41 |
| ER l | 0.116 | -0.085 | 0.977 | 36.48 |
| ER m | 0.162 | 0.020 | 0.984 | 18.07 |
| ER s | 0.169 | 0.016 | 0.983 | 20.01 |
| BA l | 4.872 | -0.726 | 0.799 | 208.09 |
| BA m | 3.229 | -0.551 | 0.778 | 241.24 |
| BA s | 0.765 | -0.494 | 0.776 | 56.07 |

Table 2: $(C, \alpha, p)$-triples from a least-squares fit of the parametric $a_k$ model (Eq. 2) to $\hat{a}_k$ from Gillespie simulations (Eq. 19) for all nine cases (see Table 1). The right column shows the root mean square error between empirical $\hat{a}_k$ and parametric model $\mathrm{RMSE}(\hat{a}_k, a_k(C, \alpha, p))$.

$$\hat{a}_k = \tau \frac{\sum_i i\, t_{ik}}{\sum_i t_{ik}}, \quad k = 1, \ldots, N, \tag{19}$$

where $t_{ik}$ denotes the total lifetime of all network states with $k$ infected nodes and $i$ S-I links.

We then proceed to fit the parameters $(C, \alpha, p)$ of the parametric $a_k$ model from Eq. 2 to the $(k, \hat{a}_k)$ curves by a least-squares fit using a particle swarm algorithm. Table 2 lists the resulting $(C, \alpha, p)$-triples for each case along with the root mean square error (RMSE) between parametric $a_k$ curves and $\hat{a}_k$.

Figure 1 shows the empirical $\hat{a}_k$ curves as well as the fitted $a_k(C, \alpha, p)$ curves for the nine cases. The top panels illustrate the good agreement between the parametric model and $\hat{a}_k$. The bottom panels of Fig. 1 show relative errors. It is not surprising that the relative error is largest for $k$ close to zero or close to $N$, i.e., when the number of S-I links is small either because only very few nodes are infected or because almost the entire population is infected. The relative errors are lowest for the ER network class followed by the Regular networks. For the BA network class, we obtain larger errors.

We simulate each case by integrating the Master equation (Eq. 1) with the empirical $\hat{a}_k$ as well as the parametric $a_k(C, \alpha, p)$. We initialise simulations with $I(t_0) = 5$, 20 and 80 infected nodes at time $t_0 = 0$, the latter two values corresponding to two and four cycles of doubling from the initially five infected nodes. A set of 1000 Gillespie simulations of each case serves as a reference. Figure 3 shows the difference in the expected number of infected nodes between BD model and reference. Figure 4 shows the integrated quadratic distance (IQD) between the cumulative densities from BD model and reference. The errors remain small throughout the simulations with the empirical $\hat{a}_k$ (top panels), which confirms the suitability of the BD model to describe population-level infection rates in SIS epidemics on Reg, ER and BA networks. Not only is the expectation well captured by simulations with $\hat{a}_k$, but also the intrinsic stochasticity of the epidemic despite the mean-field approximation. We observe the largest errors for the BA network class and during the growth phase when simulations are initialised with $I(t_0) = 5$ (Fig. 4a). The BA network class exhibits a higher degree of heterogeneity than Regular and ER networks. Hence, a larger variance in the number of S-I links at a given $k$ is expected. Therefore, the mean-field approximation might be less well suited for that type of network than for Regular and ER networks.
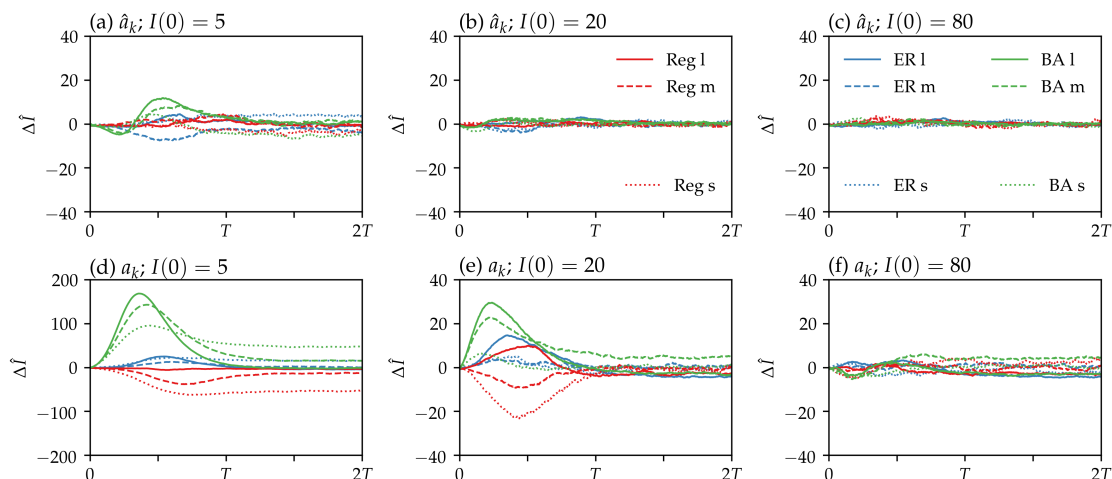
Figure 3: Difference in expected number of infected nodes between BD model and reference vs. time (Eq. 16). The top panels show the BD model with the empirical $\hat{a}_k$. The bottom panels show the BD model with parametric $a_k(C, \alpha, p)$. From left to right, three different initial conditions are shown: $I(t_0) = 5, 20, 80$ at $t_0 = 0$. The different colours and line styles indicate the nine different cases (see Table 1).

The simulations with the BD model with the parametric $a_k(C, \alpha, p)$ (bottom panels) exhibit larger errors compared to the simulations with $\hat{a}_k$. Ranking the different cases studied, the BD model achieves the lowest errors for the ER network class, followed by Regular networks and the BA networks. Again, errors are largest when the simulations are initialised with $I(0) = 5$ which appears to be caused by the larger relative errors of the parametric $a_k$-model for small $k$. The overestimation of $a_k$ at small $k$ by the parametric model for the BA network class (Fig. 1c) causes the number of infected nodes to increase too fast in the BD simulations, which leads to an over-estimation of the number of infected nodes during the growth phase (Figs. 3d, 5h). Conversely, the underestimation of $a_k$ for small $k$ for the Regular networks (Fig. 1b) causes the number of infected nodes to increase too slowly in the BD model simulations, and hence an under-estimation of the number of infected nodes during the growth phase (Figs. 3d, 5a). The errors peak during the growth phase and then decay until reaching an approximately constant value in the quasi-steady state.

For the majority of the cases, the quasi-steady state is well captured in both mean and variation around the mean, with the only exception being the small epidemics on Regular and BA networks. When the BD model is initialised at $I(0) = 5$, it starts from with a state from which some of the small epidemics will eventually die out and only some will eventually converge to the quasi-steady state. Due to the over-estimation of $a_k$ for small $k$ for the BA networks in the parametric model, the probability of an epidemic to proceed to the quasi-steady state from $I(0) = 5$ is over-estimated in the BD model. Hence, the expected number of infected nodes in the BD model is too large. For Regular networks, the opposite holds and the expected number of infected nodes is too small. When initialised with $I(0) = 20$ the errors are smaller, but the temporal pattern of the errors persists, i.e., errors peak during the growth phase and then decay until the quasi-steady state is reached (Figs. 3e, 4e). When initialised with $I(0) = 80$, we find that both the growth phase as well as the quasi-steady state is well captured by the BD model (Figs. 3f, 4f, 5c,g,j).
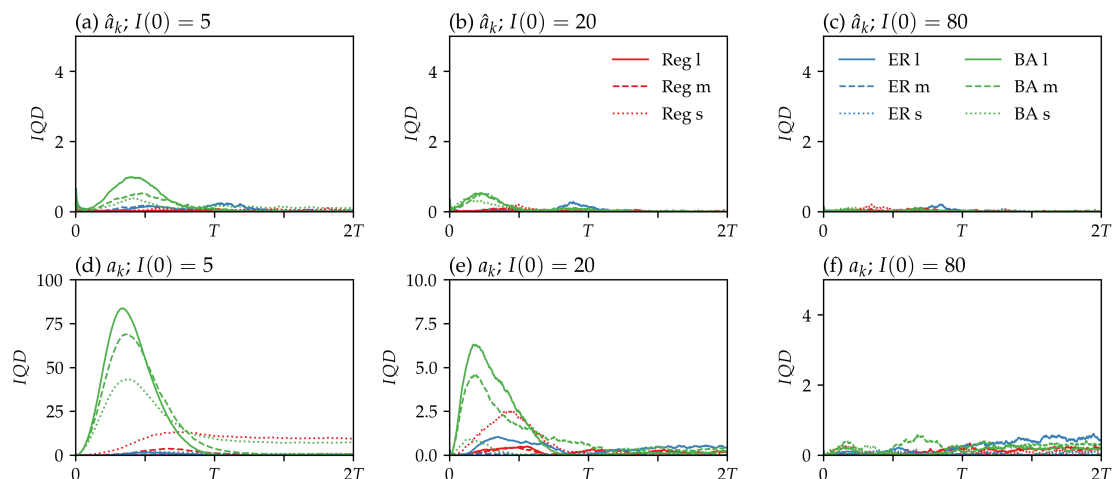
Figure 4: Integrated quadratic distance between between the cumulative density from the BD model and the reference vs. time (Eq. 17). The top panels show the BD model with the empirical $\hat{a}_k$. The bottom panels show the BD model with parametric $a_k(C, \alpha, p)$. From left to right, three different initial conditions are shown: $I(t_0) = 5, 20, 80$ at $t_0 = 0$. The different colours and line styles indicate the nine different cases (see Table 1).

## 4.2 Predictions

### 4.2.1 Network class inference

Di Lauro et al. (2020b) demonstrated that one can reliably recover the class of the underlying network from population-level observations, when observations of the full epidemic trajectory from an early stage up to quasi-steady state are available. When aiming to predict the future evolution of an ongoing epidemic, the inference of network class and parameters $(C, \alpha, p)$ can only utilise observations of the epidemic up to its current state. Therefore, the question arises as to when one has sufficient information during an epidemic to reliably predict its further course. We therefore carry out a sensitivity analysis by inferring the posterior distribution $\pi(\theta)$, $\theta \in \Theta = \{\text{Reg}, \text{ER}, \text{BA}\}$ from observation data sets covering different time windows during the evolution of the epidemic and incorporating different numbers of observations. For this analysis, we consider the medium epidemics (Table 1).

The results are summarised in Fig. 6. As expected, in general, the longer the observational period, the higher the (average) posterior probability of the true underlying network class is. When the observational period ranges from an early stage of the epidemic up to quasi-steady state (50 to quasi-steady state), ten out of ten realisations on the BA network and seven out of ten realisations on Regular and ER networks are classified correctly. While BA networks can be clearly separated when sufficient data is available, distinguishing between Regular and ER networks appears challenging. For some realisations, the respective trajectories largely overlap (Fig. 7a).

When the observation time span is shortened, the rate of correct classifications decreases (Fig. 6). As demonstrated in Fig. 7b, epidemics spreading on networks from the different classes may exhibit a similar shape during the earlier stage and only diverge later on. Thus, inferring the underlying network classes from population-level observations of a single realisation requires a sufficient observational time span. Increasing the number of observations from 10 to 100 does not have any visible effect on the classification. Ten
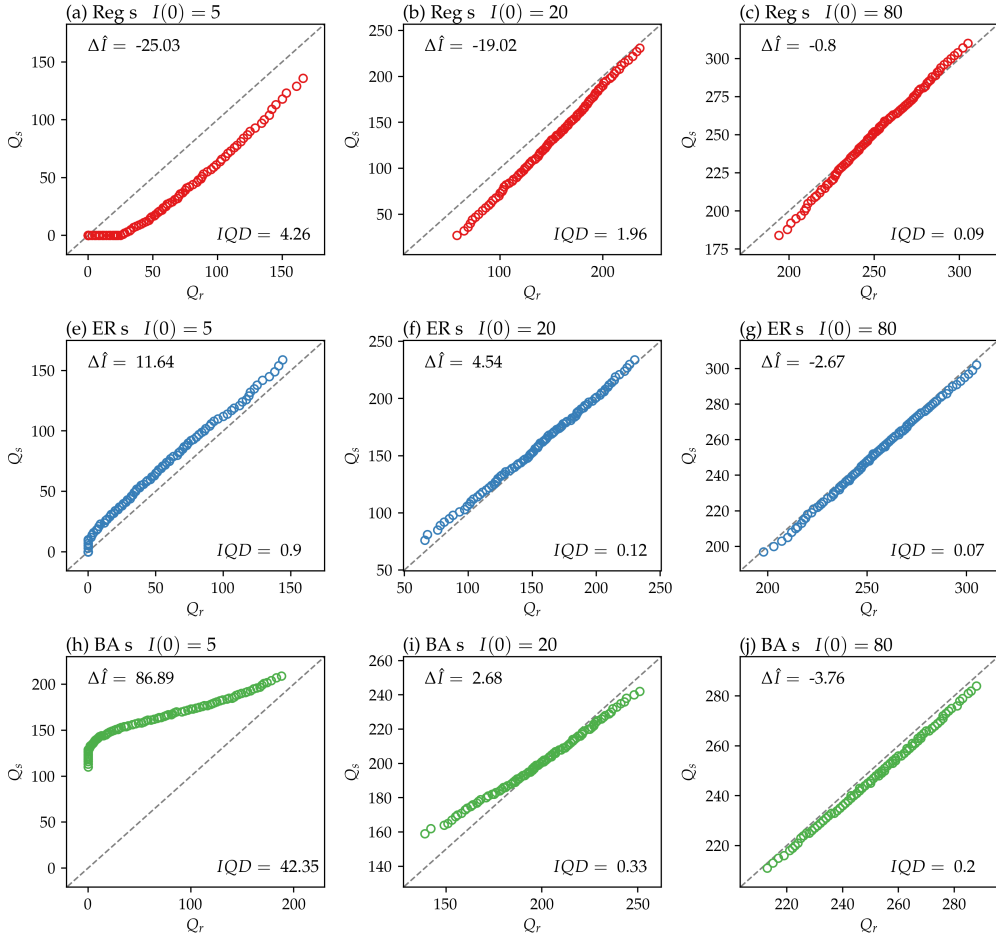
Figure 5: Quantile-quantile plots comparing BD model predictions (s) and reference (r). The circles indicate the quantiles $Q(0.05), Q(0.06), \ldots Q(0.95)$ at time $t = t_0 + \frac{1}{3}T$ after initialisation with $I(t_0) = 5, 20, 80$ infected nodes at time $t_0 = 0$.

observations provide a sufficient description of the epidemic trajectory.

We note that for the BA networks, classification accuracy increases when the very early stage of the epidemic up to $I \approx 50$ is excluded from the observational data set. In Fig. 6e,f, this is most obvious when comparing the posterior probabilities obtained for the observation intervals $[5, 400]$, $[5, 500]$ and $[5, qss]$ with those obtained for $[50, 400]$, $[50, 500]$ and $[50, qss]$, respectively. We believe this to be caused by the relatively large error of the parametric $a_k$-model for small $k$ for the BA network class (Fig. 1c). For small $k$, the average number of S-I links and hence $a_k$ are over-estimated by the model. Hence, the initial spreading of the epidemic on the network is expected to proceed significantly slower than the BD model with optimal $(C, \alpha, p)$-triplet would suggest. Further, we note that because Regular and ER networks are comparably close in the $(C, \alpha, p)$-space, confusion between Regular and ER networks is comparably likely, but predictions are also expected to be comparably robust to confusion between Regular and ER network classes.

### 4.2.2 Epidemic trajectories

Figure 8 shows the predictions incorporating the uncertainty encoded in the posterior distribution(s). Shown are three realisations of the medium epidemics on Regular, ER and
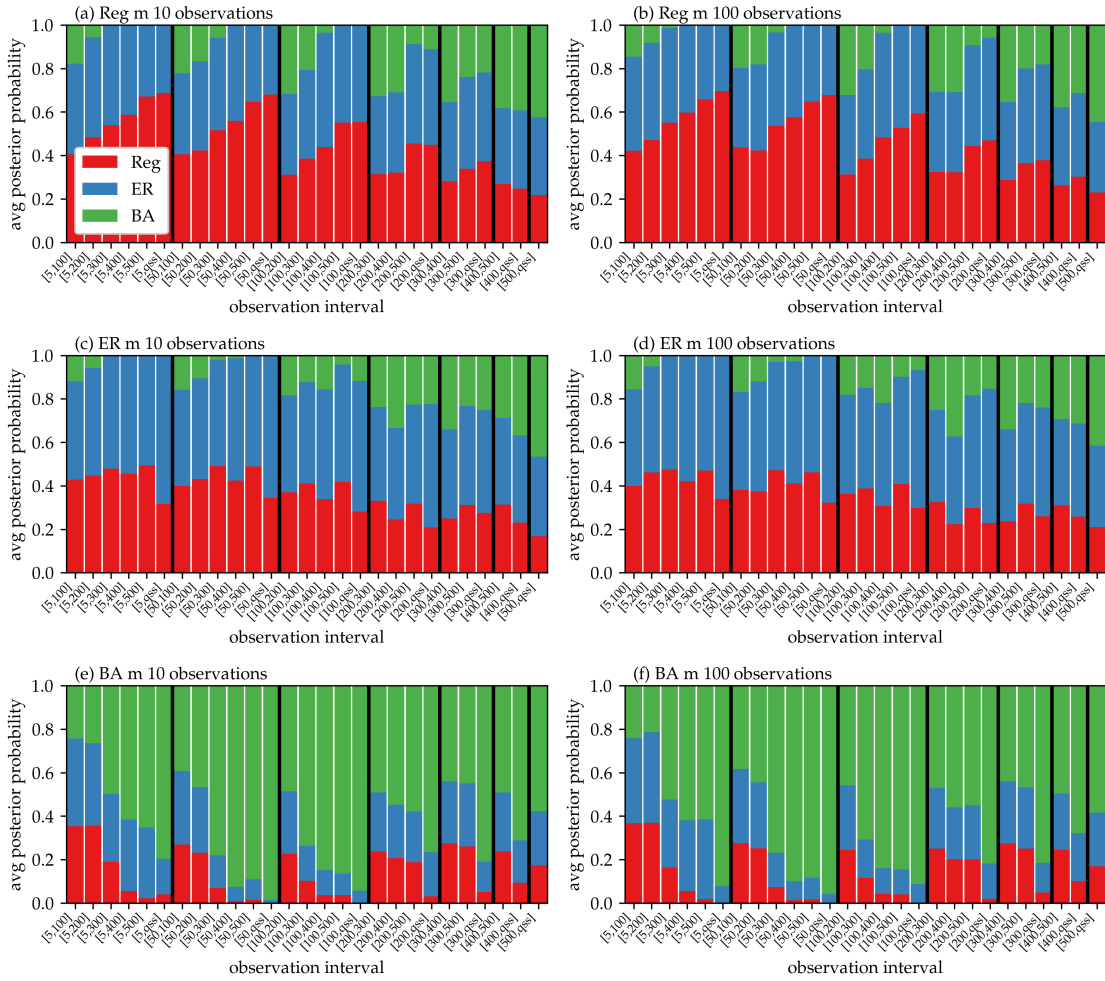
Figure 6: Sensitivity of network class inference on observational time span and number of observations. The bars indicate the average posterior probability $\pi(\theta)$, $\theta \in \Theta = \{\text{Reg}, \text{ER}, \text{BA}\}$ over ten realisations of the medium epidemic on Regular, ER or BA network (Table 1). For each case, 27 different observation intervals have been evaluated. The ten (or 100) observations are spaced approximately equidistantly in time throughout the observational period.

BA networks. Additional Figures for all ten realisations of the nine cases from Table 1 are provided in the Supplementary Material. The dots indicate the observations $(y, s)$. The grey-shaded areas indicate the predictions. Trajectories of 100 realisations of Gillespie simulations initialised at the network state associated with the last observational data point serve as reference.

For the majority of cases and realisations, the 90%-credible interval (CI) contains the reference. For some cases/realisations, the reference lies just outside 90%-CI (e.g., Fig. 8a realisation 9). For predictions of the medium epidemics initialised after five cycles of doubling ($k_{10} = 160$), we find the reference to lie just outside the 90%-CI for one out of ten realisations for the ER network, while for the Reg and BA networks, it lies partly outside of the 90%-CI for 3 out of 10 realisations each. The 90%-CI spans a range of up to $\approx 300$ and tends to be larger for the small and medium epidemics than for the large epidemics.
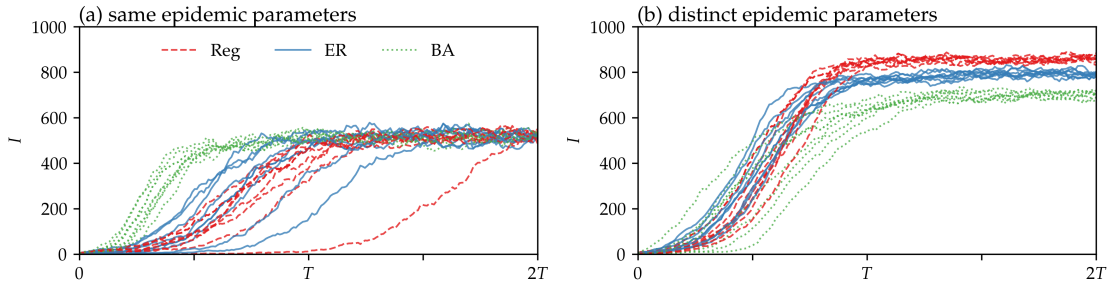
Figure 7: Examples of epidemic trajectories from Gillespie simulations on Regular, ER and BA networks. In panel (a), the epidemic parameters $\tau = 1.265$, $\gamma = 5.773$ and average node degree $k \approx 10$ (ER m, Table 1) are (approximately) the same for all three networks. In panel (b), epidemic parameters and average node degree are distinct for the three different networks and chosen such that trajectories exhibit a similar course during the early stage of the epidemic: $\tau_{ER} = 3.5$, $\gamma_{ER} = 2.969$, $k_{ER} = 5.046$; $\tau_{Reg} = 4.251$, $\gamma_{Reg} = 2.969$, $k_{Reg} = 5$ (Reg l, Table 1); $\tau_{BA} = 3.2$, $\gamma_{BA} = 2.969$, $k_{BA} = 3.992$. Eight realisations are shown for each network class and parameter combination.

When parameters and network class are inferred from observations up to five cycles of doubling ($k_{10} = 160$), prediction uncertainty is dominated by uncertainty about model parameters and network class. When inference is based on observations up to and including the quasi-steady state, uncertainty on parameters and network class is negligible and the uncertainty on the future course of the epidemic is dominated by the intrinsic stochasticity of the process (dotted brown lines in Fig. 8). Hence, the magnitude of the prediction uncertainty is sensitive to the observational time span available for inference. In any realistic setting, observations are of course not available beyond the point from which one aims to predict the future course of an epidemic. As illustrated in Fig. 8b,d and f, parameter and network class uncertainty is, as expected, reduced when a longer observational time span is available. For predictions initialised at 6 cycles of doubling ($k_{10} = 320$), prediction uncertainty due to intrinsic stochasticity and parameter/network class uncertainty is similar in magnitude and it depends on the particular case and realisation if prediction uncertainty is dominated by either one.

For the predictions of epidemics on BA networks, the reference tends to lie in the lower half of the 90%-CI whereas for Regular networks, it tends to lie in the upper half of the 90%-CI. Relatively large credible intervals are associated with relatively large uncertainty about the network class. As illustrated in Fig. 7b, epidemics on different networks with a similar trajectory during the early stage eventually diverge, with the epidemics on the BA networks converging to the lowest level of infection during quasi-steady state followed by ER and Regular networks. Thus, credible intervals obtained from observations of the beginning of one of the BA trajectories from Fig. 7b are expected to contain the true BA trajectories at their lower end. This behaviour can also be understood from the corresponding $a_k$-curves. Curves for networks from different classes that are similar for small $k$ will diverge for larger $k$, with the curve corresponding to the BA network having the smallest peak, typically followed by ER and finally the Regular network with the highest peak.

Figure 9 shows the point estimate-based predictions. The prediction intervals here do not account for network class and parameter uncertainty, but only represent the intrinsic stochasticity of the epidemic spreading. Accordingly, the prediction intervals are sys-
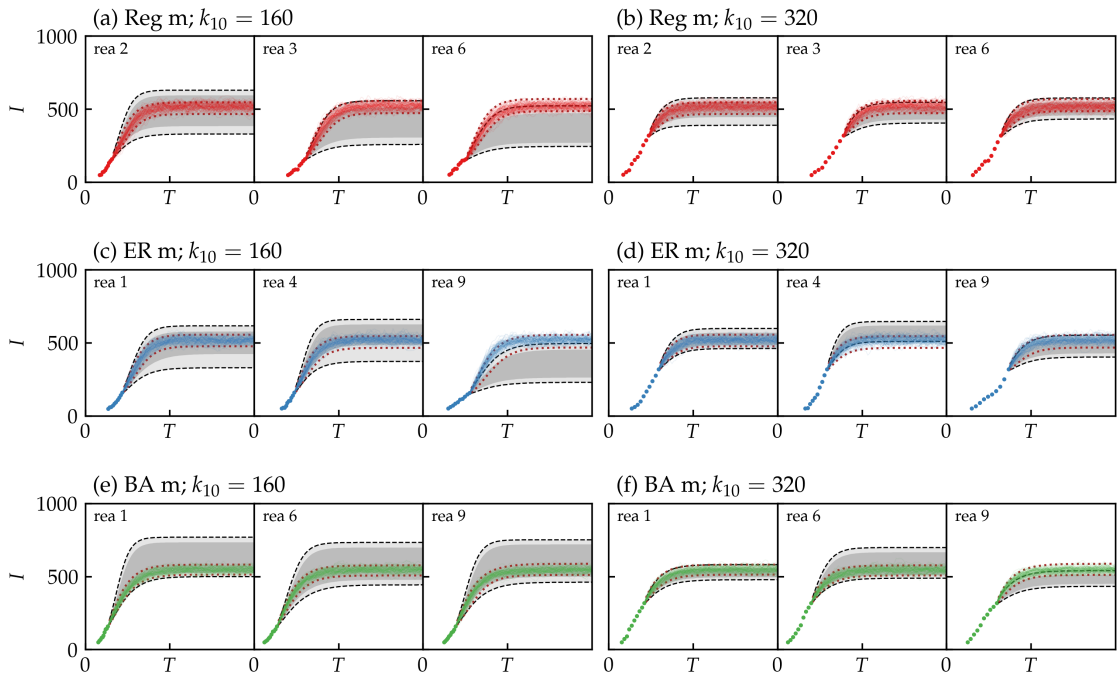
16

Figure 8: Predictions incorporating uncertainty for three example realisations of the medium epidemics on Regular, ER and BA networks. The grey shaded areas indicate 70%- and 90%-equal tailed credible intervals of the predictions initialised at the last observation $I(t_{10}) = k_{10}$. The dots indicate the ten observations $(y, s)$ used for inference where $y = (k_1 \approx 50k_{10})$. The coloured lines show 100 realisations of Gillespie simulations initialised at the last observation. The dotted brown lines indicate the 90%-equal tailed credible intervals for predictions with inference from 10 observation up to (and including) quasi-steady state.

tematically narrower than the credible intervals. The width of the prediction intervals is consistent with the spread of the trajectories from the Gillespie simulations. For some cases and realisations, the point estimate-based predictions provide a near perfect fit to the reference (e.g., Fig. 9a realisation 1, e realisation 1). However for some cases/realisations prediction and reference differ by up to $\approx 300$ (e.g., Fig. 9c realisation 6). For epidemics on BA networks, the number of infected nodes is over-estimated in the point estimate-based predictions if the network is falsely identified as ER or Regular network. For epidemics on Regular networks, the number of infected nodes is under-estimated if the network is falsely identified as ER or BA. The reason for this is the same as for the tendencies of the reference to occur in different parts of the credible intervals for the different network classes discussed in the above paragraph. When inference is based on observations up to six cycles of doubling ($k_{10} = 320$), the errors of the point estimate-based prediction are visibly reduced (see also the Supplementary Material). Hence, when longer observation time spans are available also point estimate-based predictions are potentially useful.

Finally, in Fig. 10 we consider the uncertainty in the $p_k$-space as described by the covariance of the pushforward measure around the two different predictions $m_k$ and $p_{k,MAP}$. Shown is the medium epidemic on the ER network. As the predictions based on the conditional mean $m_k$ incorporate the uncertainty about the predicted $I(t)$ that stems from uncertainty about network class and model parameters it is systematically wider (Fig. 2). This width reflects the width of the pushforward and thus leads to lower values of $|\mathcal{C}|$
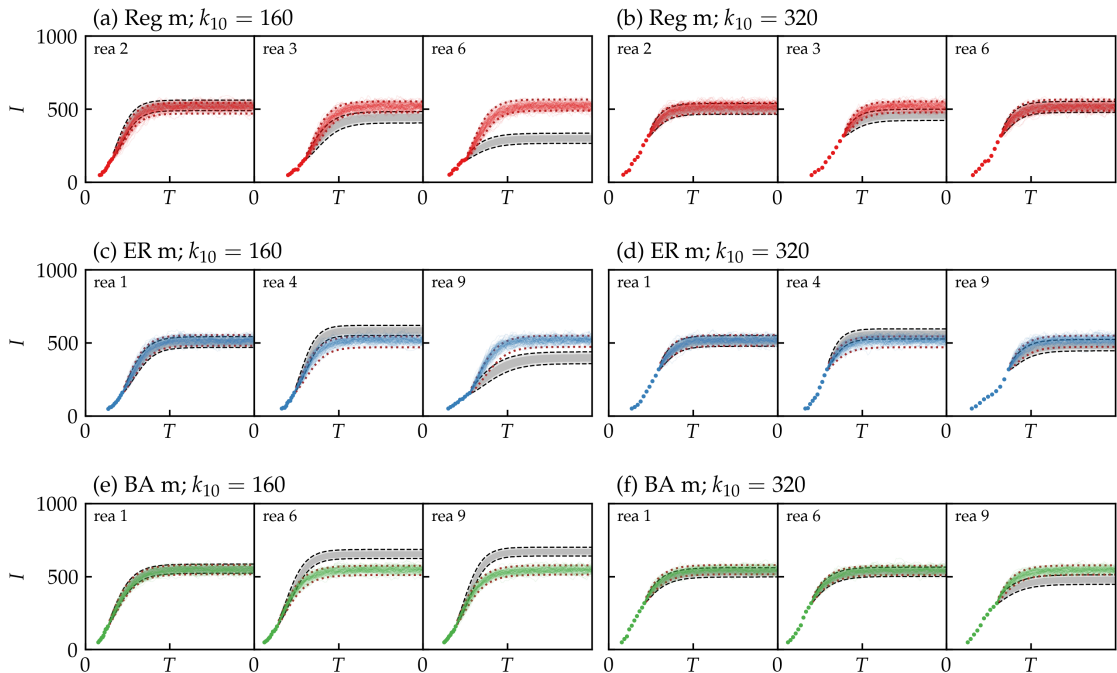
Figure 9: Point estimate-based predictions for three example realisations of medium epidemics on Regular, ER and BA networks. The grey shaded areas indicate 70%- and 90%-equal tailed prediction intervals of the predictions initialised at the last observation $I(t_{10}) = k_{10}$. The dots indicate the ten observations $(y, s)$ used for inference where $y = (k_1 \approx 50k_{10})$. The coloured lines show 100 realisations of Gillespie simulations initialised at the last observation. The dotted brown lines indicate the 90%-equal tailed credible intervals for predictions with inference from 10 observation up to (and including) quasi-steady state.

compared to the point-estimate based predictions. Further, the predictions incorporating uncertainty exhibit less variation of $|\mathcal{C}|$ among the different realisations than the point estimate-based predictions. As already discussed alongside Figs. 8 and 9, we find the uncertainty to be systematically lower when predictions are based on longer observational periods. The longer the available observation period, the narrower the posterior and the smaller the difference between the two types of predictions and their respective uncertainty in the $p_k$-space.

## 5  Discussion

We have explored a modelling and inference framework for forecasting SIS epidemics spreading on networks. The surrogate model is based on a BD process. The effect of the contact structure has been condensed into a birth-rate parameter, which is proportional to the average number of SI-links for a given number of infected nodes. Our empirical validation has confirmed that the BD model is well suited to describe the evolution of an SIS epidemic on a network (Figs. 3 d-f, 4 d-f). Both the expectation and the intrinsic stochasticity of the epidemic trajectories are well reproduced even though our model formulation contains a mean-field approximation. The parametric model for the number of SI-links, which has been introduced to enable the inference of network class and model parameters, is suitable for the range $50 \lesssim k \lesssim 950$ (Fig. 1). Hence, simulations with the
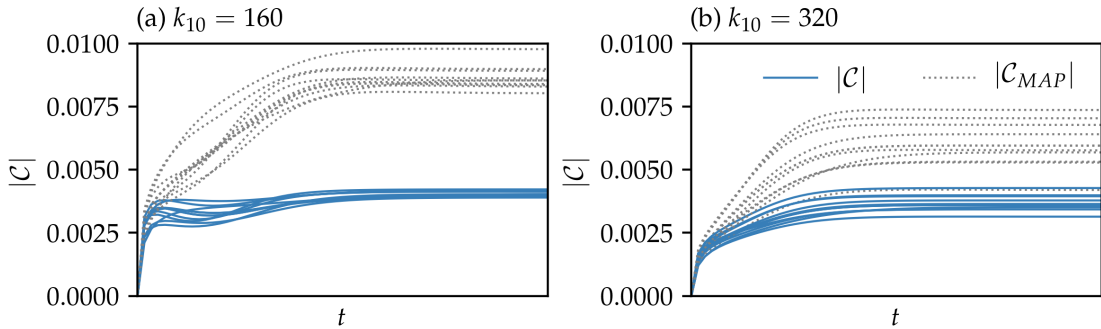
Figure 10: Uncertainty in the $p_k$-space. The solid blue lines show the Euclidean norm of the covariance $\mathcal{C}$ of the pushforward over time $t$ (Eq. 14). The grey dotted lines show the norm of the covariance around $p_{k,MAP}$ (Eq. 15). Shown are all ten realisations of the medium epidemic on an ER network with predictions initialised at and inferred from observations up to $I(t_{10}) = k_{10}$.

BD model with the parametric $a_k(C, \alpha, p)$ should not be initialised with fewer than 50 infected nodes (Fig. 3 a-c, 4 a-c).

Network class and epidemic parameters can be reliably inferred when observations are available from an early stage of the epidemic up to the quasi-steady state. However, in realistic prediction scenarios, observations are only available up to the current state of the epidemic. The accuracy of the network class inference is sensitive to the observational time span (Fig. 6). Uncertainty increases as observational time span is reduced. This is because epidemics, though spreading on networks from distinct classes, can exhibit very similar trajectories through their earlier stages and only diverge when approaching the quasi-steady state (Fig. 7). As discussed in Allen et al. (2021) for instance, the uncertainty of the future course of an emerging epidemic during its early stages is dominated by the intrinsic stochasticity of disease transmission. It is thus no surprise that observations from an early stage of the epidemic appear not to contain sufficient information about the network class.

In predictions based on observations up to and initialised at $I = 160$, the prediction uncertainty is dominated by the uncertainty of the parameters/network classes. In predictions based on observations up to and initialised at $I = 160$, the prediction uncertainty stems in about equal proportions from parameter/network class uncertainty and intrinsic stochasticity of the epidemic spreading (Fig. 8). Thus, and especially for shorter observational periods and hence predictions initialised early during the epidemic, considering parameter uncertainty is crucial for providing meaningful information about prediction uncertainty (see also Castro et al., 2020; Wilke and Bergstrom, 2020). The results suggest that for most cases the credible intervals obtained provide reliable uncertainty information for the epidemic forecasts (see also the Supplementary Material). If longer observational time spans are available, point estimate-based predictions are potentially useful as well (Fig. 9).

Our study differs from other approaches in network inference in so far as our aim here is not to infer the existence, or otherwise, of links but rather to infer the most likely network class that led to the observed population-level data resulting from an epidemic spreading on it. As a result, the data needed for inference does not contain node- or link-level information. There are both advantages and disadvantages to such an approach. On the one hand, the computation of the likelihood in our case is more straightforward and the

data needed for inference is modest. On the other hand, if more detailed data is available, the proposed model will not be able to capture it nor benefit from it. However, more complex models will need large quantities of detailed data (i.e., in the case of cascades, the data needs to contain cascades starting from, or involving, as many nodes as possible, Gomez-Rodriguez et al. (2012)) to produce acceptable results with large computational burden. The choice of model and inference will depend on the context.

There are many directions in which the current model and inference scheme can be developed. First, we only explored three network classes where the key difference was degree heterogeneity. However, networks displaying degree-degree correlations, clustering, spatial structure or some type of meso-scale structure, such as communities, may be of interest as they are more representative of real-world scenarios. Equally, from a theoretical viewpoint, lattices could be considered. This is a non-trivial task and depending on which network property or combination of properties we choose to model, it may turn out that the birth-rates of the BD process will no longer be parabola-like and the proposed parametric $a_k$-model may no longer provide a satisfactory fit. However, we expect that more complex models will be able to capture the birth-rates in the BD process resulting from such more exotic networks. Another natural extension would be to consider more complex epidemic models, such as SIR, where the corresponding BD model will now have $O(N^2)$ equations and the birth-rates of the BD process will define a surface rather than a curve. However, and perhaps more interestingly, the excellent agreement between the exact and surrogate model, leads us to believe that a rigorous proof that quantifies the error between the exact and BD models may be possible. For example, it is clear that as a Regular network becomes more densely connected, and in the limit of number of links going to $N-1$, the BD model becomes exact.

# Acknowledgements

# References

Allen, A. J., Boudreau, M. C., Roberts, N. J., Allard, A., and Hébert-Dufresne, L. (2021). Predicting the diversity of early epidemic spread on networks. *arXiv preprint arXiv:2107.03334*.

Bastos, L. S., Economou, T., Gomes, M. F., Villela, D. A., Coelho, F. C., Cruz, O. G., Stoner, O., Bailey, T., and Codeço, C. T. (2019). A modelling approach for correcting reporting delays in disease surveillance data. *Statistics in Medicine*, 38(22):4363–4377.

Castro, M., Ares, S., Cuesta, J. A., and Manrubia, S. (2020). The turning point and end of an expanding epidemic cannot be precisely forecast. *Proceedings of the National Academy of Sciences*, 117(42):26190–26196.

Chao, D. L., Matrajt, L., Basta, N. E., Sugimoto, J. D., Dean, B., Bagwell, D. A., Oiulfstad, B., Halloran, M. E., and Longini Jr, I. M. (2011). Planning for the control of pandemic influenza a (h1n1) in los angeles county and the united states. *American Journal of Epidemiology*, 173(10):1121–1130.

Cota, W., Mata, A. S., and Ferreira, S. C. (2018). Robustness and fragility of the susceptible-infected-susceptible epidemic models on complex networks. *Physical Review E*, 98(1):012310.

Crawford, F. W., Minin, V. N., and Suchard, M. A. (2014). Estimation for general birth-death processes. *Journal of the American Statistical Association*, 109(506):730–747.

Danon, L., Ford, A. P., House, T., Jewell, C. P., Keeling, M. J., Roberts, G. O., Ross, J. V., and Vernon, M. C. (2011). Networks and the epidemiology of infectious disease. *Interdisciplinary Perspectives on Infectious Diseases*, 2011.

Della Rossa, F., Salzano, D., Di Meglio, A., De Lellis, F., Coraggio, M., Calabrese, C., Guarino, A., Cardona-Rivera, R., De Lellis, P., Liuzza, D., et al. (2020). A network model of italy shows that intermittent regional strategies can alleviate the covid-19 epidemic. *Nature Communications*, 11(1):1–9.

Devriendt, K. and Van Mieghem, P. (2017). Unified mean-field framework for susceptible-infected-susceptible epidemics on networks, based on graph partitioning and the isoperimetric inequality. *Physical Review E*, 96(5):052314.

Di Lauro, F., Berthouze, L., Dorey, M. D., Miller, J. C., and Kiss, I. Z. (2021a). The impact of network properties and mixing on control measures and disease-induced herd immunity in epidemic models: a mean-field model perspective. *Bulletin of Mathematical Biology*, 83(117).

Di Lauro, F., Croix, J.-C., Berthouze, L., and Kiss, I. (2020a). Pde limits of stochastic sis epidemics on networks. *Journal of Complex Networks*, 8(4).

Di Lauro, F., Croix, J.-C., Dashti, M., Berthouze, L., and Kiss, I. (2020b). Network inference from population-level observation of epidemics. *Scientific Reports*, 10(1):1–14.

Di Lauro, F., Kiss, I. Z., and Miller, J. C. (2021b). Optimal timing of one-shot interventions for epidemic control. *PLoS Computational Biology*, 17(3):e1008763.

Ganesh, A., Massoulié, L., and Towsley, D. (2005). The effect of network topology on the spread of epidemics. In *Proceedings IEEE 24th Annual Joint Conference of the IEEE Computer and Communications Societies.*, volume 2, pages 1455–1466. IEEE.

Gillespie, D. T. (1977). Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry*, 81(25):2340–2361.

Goeyvaerts, N., Santermans, E., Potter, G., Torneri, A., Van Kerckhove, K., Willem, L., Aerts, M., Beutels, P., and Hens, N. (2018). Household members do not contact each other at random: implications for infectious disease modelling. *Proceedings of the Royal Society B*, 285(1893):20182201.

Gomez-Rodriguez, M., Leskovec, J., and Krause, A. (2012). Inferring networks of diffusion and influence. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 5(4):1–37.

Jacquez, J. A. and Simon, C. P. (1993). The stochastic SI model with recruitment and deaths i. comparison with the closed sis model. *Mathematical Biosciences*, 117(1-2):77–125.

Keeling, M. J. and Eames, K. T. (2005). Networks and epidemic models. *Journal of the Royal Society Interface*, 2(4):295–307.

Kennedy, J. and Eberhart, R. (1995). Particle swarm optimization. In *Proceedings of ICNN'95-International conference on Neural Networks*, volume 4, pages 1942–1948. IEEE.

Kermack, W. O. and McKendrick, A. G. (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, 115(772):700–721.

Kiss, I. Z., Miller, J. C., Simon, P. L., et al. (2017). Mathematics of epidemics on networks. *Cham: Springer*, 598.

Mata, A. S. and Ferreira, S. C. (2013). Pair quenched mean-field theory for the susceptible-infected-susceptible model on complex networks. *EPL (Europhysics Letters)*, 103(4):48003.

McGough, S. F., Johansson, M. A., Lipsitch, M., and Menzies, N. A. (2020). Nowcasting by bayesian smoothing: A flexible, generalizable model for real-time epidemic tracking. *PLoS Computational Biology*, 16(4):e1007735.

Nagy, N., Kiss, I. Z., and Simon, P. (2014). Approximate master equations for dynamical processes on graphs. *Mathematical Modelling of Natural Phenomena*, 9(2):43–57.

Nsoesie, E., Mararthe, M., and Brownstein, J. (2013). Forecasting peaks of seasonal influenza epidemics. *PLoS Currents*, 5.

Pastor-Satorras, R., Castellano, C., Van Mieghem, P., and Vespignani, A. (2015). Epidemic processes in complex networks. *Reviews of Modern Physics*, 87(3):925.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830.

Roy, V. (2020). Convergence diagnostics for Markov chain Monte Carlo. *Annual Review of Statistics and its Application*, 7:387–412.

Shaman, J. and Karspeck, A. (2012). Forecasting seasonal outbreaks of influenza. *Proceedings of the National Academy of Sciences*, 109(50):20425–20430.

Shirley, M. D. and Rushton, S. P. (2005). The impacts of network topology on disease spread. *Ecological Complexity*, 2(3):287–299.

Siettos, C. I. and Russo, L. (2013). Mathematical modeling of infectious disease dynamics. *Virulence*, 4(4):295–306.

Simon, P. L., Taylor, M., and Kiss, I. Z. (2011). Exact epidemic models on graphs using graph-automorphism driven lumping. *Journal of Mathematical Biology*, 62(4):479–508.

Tizzoni, M., Bajardi, P., Poletto, C., Ramasco, J. J., Balcan, D., Gonçalves, B., Perra, N., Colizza, V., and Vespignani, A. (2012). Real-time numerical forecast of global epidemic spreading: case study of 2009 a/h1n1pdm. *BMC medicine*, 10(1):1–31.

Unkel, S., Farrington, C. P., Garthwaite, P. H., Robertson, C., and Andrews, N. (2012). Statistical methods for the prospective detection of infectious disease outbreaks: a review. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 175(1):49–82.

Van Yperen, J., Campillo-Funollet, E., and Madzvamuse, A. (2020). Covid-19: measuring the impact on healthcare demand and capacity and exploring intervention scenarios. *arXiv preprint arXiv:2012.15392*.

Vats, D., Flegal, J. M., and Jones, G. L. (2019). Multivariate output analysis for Markov chain Monte Carlo. *Biometrika*, 106(2):321–337.

Vats, D., Robertson, N., Flegal, J. M., and Jones, G. L. (2020). Analyzing Markov chain Monte Carlo output. *Wiley Interdisciplinary Reviews: Computational Statistics*, 12(4):e1501.

Wilke, C. O. and Bergstrom, C. T. (2020). Predicting an epidemic trajectory is difficult. *Proceedings of the National Academy of Sciences*, 117(46):28549–28551.

Xue, L., Jing, S., Miller, J. C., Sun, W., Li, H., Estrada-Franco, J. G., Hyman, J. M., and Zhu, H. (2020). A data-driven network model for the emerging covid-19 epidemics in Wuhan, Toronto and Italy. *Mathematical Biosciences*, 326:108391.

Yin, Q., Shi, T., Dong, C., and Yan, Z. (2017). The impact of contact patterns on epidemic dynamics. *PLoS One*, 12(3):e0173411.