

Dargatz, Christiane

**Working Paper**

## A diffusion approximation for an epidemic model

Discussion Paper, No. 517

**Provided in Cooperation with:**

Collaborative Research Center (SFB) 386: Statistical Analysis of discrete structures - Applications in Biometrics and Econometrics, University of Munich (LMU)

*Suggested Citation:* Dargatz, Christiane (2007) : A diffusion approximation for an epidemic model, Discussion Paper, No. 517, Ludwig-Maximilians-Universität München, Sonderforschungsbereich 386 - Statistische Analyse diskreter Strukturen, München, <https://doi.org/10.5282/ubm/epub.1882>

This Version is available at:

<https://hdl.handle.net/10419/31042>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

# A Diffusion Approximation for an Epidemic Model

Christiane Dargatz  
Ludwig-Maximilian University Munich

January 26, 2007

**Abstract:** Influenza is one of the most common and severe diseases worldwide. Devastating epidemics actuated by a new subtype of the influenza A virus occur again and again with the most important example given by the Spanish Flu in 1918/19 with more than 27 million deaths. For the development of pandemic plans it is essential to understand the character of the dissemination of the disease. We employ an extended SIR model for a probabilistic analysis of the spatio-temporal spread of influenza in Germany. The inhomogeneous mixing of the population is taken into account by the introduction of a network of subregions, connected according to Germany's commuter and domestic air traffic. The infection dynamics is described by a multivariate diffusion process, the discussion of which is a major part of this report. We furthermore present likelihood-based estimates of the model parameters.

**Keywords:** General stochastic epidemic; Likelihood inference; Euler scheme; Influenza.

## 1 Introduction

The analysis of the spread of epidemics dates back to the 18th century. However, as transportation systems improve and people travel faster, further and more frequently than in former times, also the character of the geographical spread of epidemics changes. Baroyan, Rvachev, and Ivannikov (1977) were probably the first to model the spatial spread of influenza by considering the transportation network of a specific region, in this case the train system in the USSR. More recent studies such as Colizza, Barrat, Barthélemy, and Vespignani (2006a) and Colizza et al. (2006b) especially emphasize the role of the airline traffic; Brownstein, Wolfe, and Mandl (2006) provide empirical evidence for the importance of long-distance air travel.

In this report, we investigate the spatio-temporal spread of influenza in Germany. For this purpose we employ a so-called SIR model, which will be described in Section 2. In the SIR model, the two parameters of interest are the contact rate  $\alpha$ , which is the number of an individual's potentially infectious contacts per unit time, and the reciprocal average infectious period  $\beta$ . There are many approaches to empirically obtain these parameters from external data as for example from so-called contact diaries in which a representative part of the population accurately reports each contact the person has with other people (see e.g. Fu, 2005). Naturally, such approaches are always fraught with errors and uncertainty. Our main objective is thus to consider a probabilistic rather than deterministic model and to statistically estimate these parameters based on disease counts.

Since we consider a large population, we use a diffusion approximation to carry out estimation methods. Section 3 deals with the derivation of this diffusion process from the discrete Markov chain model description given in Section 2. The diffusion model is much more amenable to statistical analysis and also disposes the inconvenience of being computationally costly for sampling. Section 4 provides an introduction to simulation and estimation methods for stochastic differential equations with a special emphasis on the Euler approximation scheme, which will be applied in Section 5 for a simulation study

and in Section 6 for the application of the diffusion process to the modelling of the spread of influenza in Germany. The report is concluded in Section 7.

## 2 Model

As a basis, we use the widely adopted SIR model in which the population under consideration is classified into susceptible (S), infectious (I) and recovered (R) individuals. Transitions between these classes are



which means that each contact between a susceptible and an infectious individual will cause an infection with rate  $\alpha$ , resulting in two infected individuals, each of which will recover with rate  $\beta$ . The parameter  $\alpha$  is the contact rate of an infectious individual sufficient to spread the disease, and  $\beta$  is the reciprocal average infectious period. The infection dynamics in this model, which in the literature is also often referred to as the general stochastic epidemic, can be deterministically described by the set of ordinary differential equations

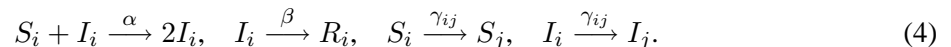
$$ds/dt = -\alpha sj, \quad dj/dt = \alpha sj - \beta j, \quad (2)$$

where  $s = S/N$  and  $j = I/N$  denote the fractions of susceptible and infective individuals of the total population of size  $N$ . Note that in this description the state space is considered to be continuous, which is an eligible assumption for large populations. The remaining fraction  $r = R/N$  can be calculated as  $r = 1 - s - j$  since we assume the population size to be constant during the time of consideration. The graphic on the left of Figure 1 shows the typical evolution of an epidemic following the deterministic description (2). The vertical line in this plot indicates the first time point at which the fraction of susceptibles falls below  $\rho^{-1} := \beta/\alpha$ , which is a crucial parameter called the basic reproduction number. Apparently, this mark agrees with the time point at which the epidemic reaches its maximum with respect to the number of infected individuals. However, since we are dealing with a process that is highly delicate to disturbances, we are not satisfied with a deterministic description as given by equations (2). Instead, we employ the stochastic differential equations (SDEs)

$$\begin{aligned} ds &= -\alpha sj dt + \sqrt{\frac{\alpha sj}{N}} dB_1(t) \\ dj &= (\alpha sj - \beta j) dt - \sqrt{\frac{\alpha sj}{N}} dB_1(t) + \sqrt{\frac{\beta j}{N}} dB_2(t), \end{aligned} \quad (3)$$

where  $B_1$  and  $B_2$  are independent Brownian motions, and  $dB_1$  and  $dB_2$  can hence be interpreted as Gaussian white noise forces accounting for fluctuations in transmission and recovery. This system is able to model the probabilistic character of the process. The graphs on the right of Figure 1 illustrate how the stochastic courses can fluctuate around the deterministic evolution of an epidemic. Section 3 deals with the formal derivation of these SDEs.

So far, a central assumption in our model is that the population under consideration mixes homogeneously. However, this situation is surely not given as soon as we regard the nationwide or even worldwide spread of a disease. We hence introduce a network of subregions  $i = 1, \dots, n$  of sizes  $N_i$  (Hufnagel, Brockmann, & Geisel, 2004), where in addition to the local infection dynamics, which again follows the standard SIR model, individuals travel between regions with rates which are summarized in a connectivity matrix  $\gamma = (\gamma_{ij})_{i,j=1,\dots,n}$ , where  $\gamma_{ii} = 0$  for all  $i$ . The transitions for this model are



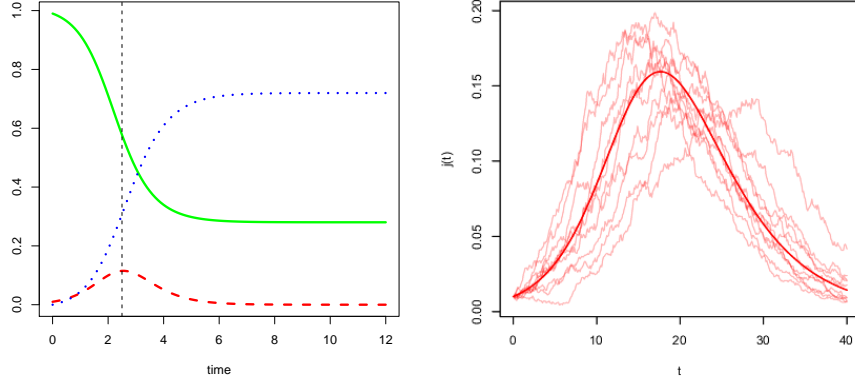


Figure 1: *Left*: Evolution in time of fractions of susceptible (solid), infected (dashed), and removed (dotted) individuals in the standard deterministic SIR model (2). The vertical line marks the point at which the fraction of susceptibles falls below  $\rho^{-1} = \beta/\alpha$ . *Right*: Deterministic (thick line) and stochastic (thin lines) courses of fractions of infectives during an epidemic according to (2) and (3). The stochastic simulations were performed with the Euler-Maruyama scheme, see Section 4.

Again, we can express this dynamics in terms of SDEs:

$$\begin{aligned}
ds_i &= \left( -\alpha s_i j_i - \sum_{k=1}^n \gamma_{ik} s_i + \sum_{k=1}^n \gamma_{ki} s_k \right) dt \\
&\quad + \sqrt{\frac{\alpha s_i j_i}{N_i}} dB_1^{(i)}(t) + \sum_{k=1}^n \sqrt{\frac{\gamma_{ik} s_i}{N_k}} dB_3^{(i,k)}(t) - \sum_{k=1}^n \sqrt{\frac{\gamma_{ki} s_k}{N_i}} dB_3^{(k,i)}(t) \\
dj_i &= \left( \alpha s_i j_i - \beta j_i - \sum_{k=1}^n \gamma_{ik} j_i + \sum_{k=1}^n \gamma_{ki} j_k \right) dt - \sqrt{\frac{\alpha s_i j_i}{N_i}} dB_1^{(i)}(t) \\
&\quad + \sqrt{\frac{\beta j_i}{N_i}} dB_2^{(i)}(t) + \sum_{k=1}^n \sqrt{\frac{\gamma_{ik} j_i}{N_k}} dB_3^{(i,k)}(t) - \sum_{k=1}^n \sqrt{\frac{\gamma_{ki} j_k}{N_i}} dB_3^{(k,i)}(t) \\
dr_i &= \beta j_i dt - \sqrt{\frac{\beta j_i}{N_i}} dB_2^{(i)}(t)
\end{aligned} \tag{5}$$

for  $i = 1, \dots, n$ . The  $n$ -dimensional Brownian motions  $\mathbf{B}_1$  and  $\mathbf{B}_2$  and the collection of  $n \times n$  independent Brownian motions  $\mathbf{B}_3$  represent disturbances in transmission, recovery, and migration, respectively. See Section 3 for an in-depth analysis of this system. Figure 2 shows the evolution of the fractions of infectives during an epidemic in five regions which agree in all parameters but the initial numbers of infectives. In the graphic on the very left there is no migration between regions, while there is strong mixing on the right. Apparently, with increasing exchange of individuals between regions, the courses of the epidemics equalize. This fact is again illustrated in Figure 3, where the dotted vertical lines mark the instants at which the fractions of susceptibles in the deterministic courses fall below  $\rho^{-1}$ , while the dashed lines indicate the actual turning points of the deterministic courses of the epidemics. For regions with high fractions of infectives, the actual turning point lies before the one that is valid for the model without migration; for regions with relatively few cases, the opposite situation applies.

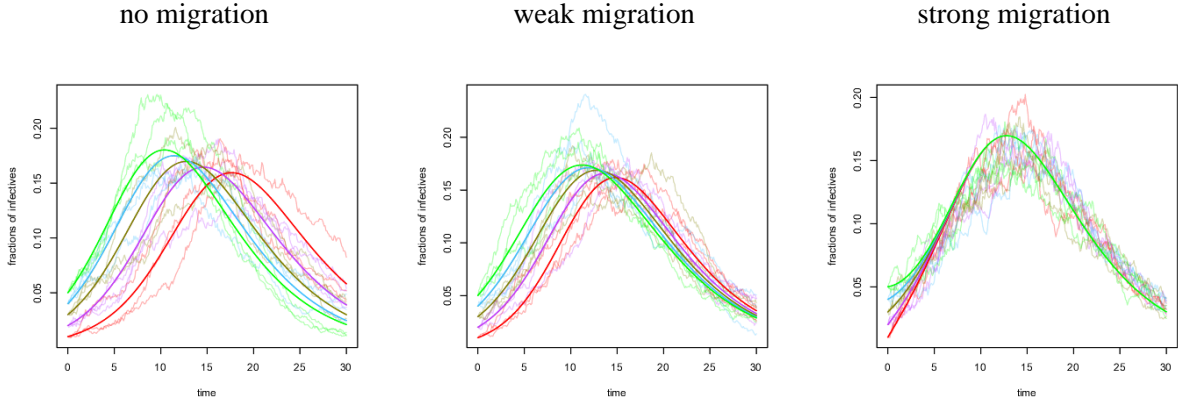


Figure 2: Evolution of the fractions of infectives in five regions which agree in all parameters but the initial numbers of infectives, which vary from one to five percent of the population. There is no traffic between regions in the graphic on the left, weak traffic in the middle, and strong traffic on the right. The thick lines show the deterministic evolution, the thin lines are stochastic simulations.

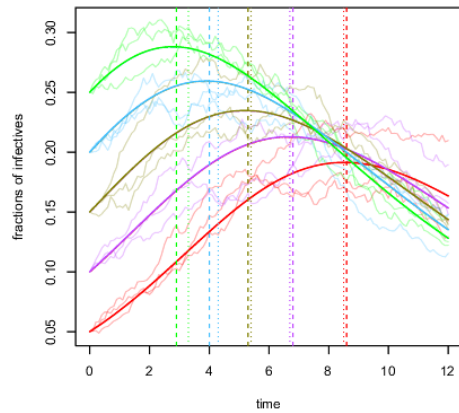


Figure 3: Evolution of the fractions of infectives in five regions between which people migrate and which agree in all parameters but the initial numbers of infectives, which vary from one to five percent of the population. The thick curve shows the deterministic evolution, the thin lines are stochastic simulations. The dotted vertical lines indicate the instants at which the fractions of susceptibles in the deterministic course fall below  $\rho^{-1}$ . The dashed vertical lines mark the actual turning points of the deterministic course of the epidemic. Without migration, these lines would agree within each region.

### 3 Diffusion Approximations of the Discrete State Space Processes

In this section we describe how to convert the model descriptions (1) and (4) into the systems of SDEs (3) and (5), respectively. There actually seems to be no general procedure for approximating discrete state space epidemics by diffusions; authors generally work through the specific examples which they cover in their papers, see for example Clancy and French (2001), Clancy, O’Neill, and Pollett (2001), Nasell (2002), and Pollett (2001). Bailey (1975) and Goel and Richter-Dyn (1974) treat approximations of the univariate birth and death process, but do not proceed to multivariate diffusion processes.

In the following we present two different approaches for the derivation of a diffusion process which follows the same law as the rather simple model description given by (1). The first method is the more descriptive one since it deals with the convergence of the process when changing to continuous time and continuous state space. This way, we obtain the *forward* Kolmogorov diffusion equation of a corresponding diffusion process. However, the second approach, in which we obtain the infinitesimal generator of the process as a *backward* Kolmogorov diffusion equation, is more straightforward and eligible for generalization. We hence apply the latter also for the transition from the discrete Markov chain description (4) of the spatial SIR model to the diffusion representation (5). Referring to the employed Kolmogorov diffusion equations, we call our methods the forward and backward approaches, respectively.

### 3.1 Forward Approach

The line of this procedure is as follows: We start by setting up the transition probabilities of the process in which we count the numbers  $S, I \in [0, N]$  of susceptible and infective individuals in the population. The state space of this process is discrete, and the evolution of the transition probabilities is described by differential-difference equations, i.e. first order differential equations in the (continuous) time variable and difference equations in the (discrete) space variable. These equations are called master equations. We then consider a sequence of discrete state space processes in which the state variables denote the fractions  $s$  and  $j$  of the susceptible and infective classes. For the population size tending to infinity, this sequence converges to a process with state variables changing continuously in space. The corresponding master equations of the discrete state space processes converge to a second order partial differential equation which is much more convenient for analytical analysis than the differential-difference equations. We can show that the Markovian limiting process is described by the limiting diffusion equation. The latter is of Fokker-Planck type and can be converted into the SDEs (3).

**Transition probabilities for the discrete state space process.** Assuming that at most one event can occur during a small time interval of length  $\Delta t$ , there are exactly three (disjoint) possibilities to obtain the state  $(S, I) \in [0, N - 1] \times [1, N - 1]$  with  $S + I \leq N$  at time  $t + \Delta t$ :

1. There were  $S + 1$  susceptibles and  $I - 1$  infectives at time  $t$ , and one infection occurred.
2. There were  $S$  susceptibles and  $I + 1$  infectives at time  $t$ , and one recovery occurred.
3. There were  $S$  susceptibles and  $I$  infectives at time  $t$ , and nothing happened.

The probability for the first event to occur is as follows: Each of the  $I - 1$  infectives at time  $t$  has  $\alpha$  potentially infectious contacts per time unit. On average,  $\alpha \cdot (S + 1)/N$  of these contacts actually cause an infection. The probability of the first event hence reads

1.  $(I - 1) \alpha \frac{S+1}{N} \Delta t + o(\Delta t)$ ,

where  $o(\Delta t)/\Delta t \rightarrow 0$  as  $\Delta t \rightarrow 0$ . Similarly, the probabilities of the second and third events are

2.  $\beta(I + 1)\Delta t + o(\Delta t)$ ,

3.  $1 - \alpha \frac{SI}{N} \Delta t - \beta I \Delta t + o(\Delta t)$ .

Let  $P(S, I; t)$  denote the probability that there are  $S$  susceptibles and  $I$  infectives at time  $t$ . Then, for  $(S, I) \in [0, N - 1] \times [1, N - 1]$ , where  $S + I \leq N$ ,

$$\begin{aligned} P(S, I; t + \Delta t) &= \left[ \alpha \frac{(S + 1)(I - 1)}{N} \Delta t + o(\Delta t) \right] P(S + 1, I - 1; t) \\ &+ \left[ \beta(I + 1) \Delta t + o(\Delta t) \right] P(S, I + 1; t) \\ &+ \left[ 1 - \left( \alpha \frac{SI}{N} + \beta I \right) \Delta t + o(\Delta t) \right] P(S, I; t). \end{aligned}$$

Subtract  $P(S, I; t)$  on both sides, divide by  $\Delta t$  and let  $\Delta t \rightarrow 0$ . We then get

$$\begin{aligned} \frac{\partial}{\partial t} P(S, I; t) &= \frac{\alpha}{N} (S + 1)(I - 1) P(S + 1, I - 1; t) \\ &+ \beta(I + 1) P(S, I + 1; t) \\ &- \left( \frac{\alpha}{N} SI + \beta I \right) P(S, I; t) \end{aligned} \quad (6)$$

as a description for the continuous time process with discrete state space. This is a differential-difference equation called the forward master equation. For the boundaries excluded above we obtain

$$\frac{\partial}{\partial t} P(S, 0; t) = \beta P(S, 1; t), \quad (7)$$

$$\frac{\partial}{\partial t} P(0, N; t) = \frac{\alpha}{N} (N - 1) P(1, N - 1; t) - \beta N P(0, N; t), \quad (8)$$

$$\frac{\partial}{\partial t} P(N, 0; t) = 0 \quad (9)$$

with  $S \in [0, N - 1]$  in the first formula. Equations (6) to (9) are subject to an initial condition  $(S_0, I_0)$  at time 0.

**Transition to continuous state space.** Instead of natural numbers  $S, I$  of susceptible and infectious individuals we now consider the respective fractions  $s = S/N$  and  $j = I/N \in (0, 1)$  of the total population. We consider a sequence of processes corresponding to a sequence of numbers  $N$  which tend to infinity. Define

$$\lambda(s, j) := N\alpha sj = \frac{\alpha}{N} SI \quad \text{and} \quad \mu(s, j) := N\beta j = \beta I.$$

The forward master equation (6) for each process then becomes for  $(s, j) \in (0, 1)^2$

$$\begin{aligned} \frac{\partial}{\partial t} P(s, j; t) &= \lambda(s + \varepsilon, j - \varepsilon) P(s + \varepsilon, j - \varepsilon; t) \\ &+ \mu(s, j + \varepsilon) P(s, j + \varepsilon; t) \\ &- (\lambda(s, j) + \mu(s, j)) P(s, j; t), \end{aligned}$$

where  $\varepsilon = N^{-1}$  fixed and  $P(s, j; t)$  now denotes the probability density for having fractions  $s$  and  $j$  of susceptible and infective individuals at time  $t$ . By simply adding and again subtracting some terms, we

can rewrite the right hand side as

$$\begin{aligned}
& \frac{1}{2} \left( \lambda(s + \varepsilon, j)P(s + \varepsilon, j; t) - \lambda(s, j)P(s, j; t) \right) \\
& + \frac{1}{2} \left( \lambda(s, j)P(s, j; t) - \lambda(s - \varepsilon, j)P(s - \varepsilon, j; t) \right) \\
& - \frac{1}{2} \left( \lambda(s, j + \varepsilon)P(s, j + \varepsilon; t) - \mu(s, j + \varepsilon)P(s, j + \varepsilon; t) - \lambda(s, j)P(s, j; t) + \mu(s, j)P(s, j; t) \right) \\
& - \frac{1}{2} \left( \lambda(s, j)P(s, j; t) - \mu(s, j)P(s, j; t) - \lambda(s, j - \varepsilon)P(s, j - \varepsilon; t) + \mu(s, j - \varepsilon)P(s, j - \varepsilon; t) \right) \\
& + \frac{1}{2} \left( \lambda(s + \varepsilon, j)P(s + \varepsilon, j; t) - 2\lambda(s, j)P(s, j; t) + \lambda(s - \varepsilon, j)P(s - \varepsilon, j; t) \right) \\
& + \frac{1}{2} \left( \lambda(s, j + \varepsilon)P(s, j + \varepsilon; t) + \mu(s, j + \varepsilon)P(s, j + \varepsilon; t) - 2\lambda(s, j)P(s, j; t) \right. \\
& \quad \left. - 2\mu(s, j)P(s, j; t) + \lambda(s, j - \varepsilon)P(s, j - \varepsilon; t) + \mu(s, j - \varepsilon)P(s, j - \varepsilon; t) \right) \\
& - \left( \lambda(s + \varepsilon, j)P(s + \varepsilon, j; t) - \lambda(s + \varepsilon, j - \varepsilon)P(s + \varepsilon, j - \varepsilon; t) - \lambda(s, j)P(s, j; t) \right. \\
& \quad \left. + \lambda(s, j - \varepsilon)P(s, j - \varepsilon; t) \right).
\end{aligned}$$

The first line (without the factor 1/2) becomes

$$\begin{aligned}
& \frac{1}{N} \lambda(s + \varepsilon, j)P(s + \varepsilon, j; t) - \frac{1}{N} \lambda(s, j)P(s, j; t) \\
& = \frac{\alpha(s + \varepsilon)j P(s + \varepsilon, j; t) - \alpha s j P(s, j; t)}{\varepsilon} \\
& \rightarrow \frac{\partial}{\partial s} \alpha s j P(s, j; t)
\end{aligned}$$

as  $\varepsilon \rightarrow 0$ . Proceed similarly with the remaining terms. Altogether, we can derive for  $\varepsilon$  tending to zero

$$\begin{aligned}
\frac{\partial}{\partial t} P(s, j; t) & = \frac{\partial}{\partial s} \alpha s j P(s, j; t) - \frac{\partial}{\partial j} (\alpha s j - \beta j) P(s, j; t) \\
& + \frac{1}{2} \frac{\partial^2}{\partial s^2} \frac{1}{N} \alpha s j P(s, j; t) + \frac{1}{2} \frac{\partial^2}{\partial j^2} \frac{1}{N} (\alpha s j + \beta j) P(s, j; t) \\
& - \frac{\partial^2}{\partial s \partial j} \frac{1}{N} \alpha s j P(s, j; t).
\end{aligned}$$

With  $\mathbf{x} = (s, j)'$ , this can be rewritten as

$$\frac{\partial}{\partial t} P(\mathbf{x}; t) = -\frac{\partial}{\partial \mathbf{x}} \left[ \mathbf{A}(\mathbf{x}) P(\mathbf{x}; t) \right] + \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \frac{\partial}{\partial \mathbf{x}} \left[ \boldsymbol{\Sigma}(\mathbf{x}) P(\mathbf{x}; t) \right], \quad (10)$$



where

$$\mathbf{A}(\mathbf{x}) = \begin{pmatrix} -\alpha sj \\ \alpha sj - \beta j \end{pmatrix} \quad \text{and} \quad \mathbf{\Sigma}(\mathbf{x}) = \frac{1}{N} \begin{pmatrix} \alpha sj & -\alpha sj \\ -\alpha sj & \alpha sj + \beta j \end{pmatrix}.$$

This is a so-called Fokker-Planck equation or forward (Kolmogorov) diffusion equation;  $\mathbf{A}$  is referred to as the drift term and  $\mathbf{\Sigma}$  the diffusion term. Again, equation (10) is to be solved with the initial condition that the process starts at  $\mathbf{x}_0 = (s_0, j_0)'$  at time 0. This condition will be included in the notation as  $P(\mathbf{x}|\mathbf{x}_0; t)$  in the following.

The considered sequence of processes with discrete state space tends to a process with continuous state space, and the corresponding forward master equations converge to a forward diffusion equation. Goel and Richter-Dyn (1974) prove that the limiting process with continuous state space is described by the diffusion equation (10) if certain conditions are met. These are as follows: Define

$$\mathbf{M}_n(\mathbf{x}) := \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int_{\Omega} (z - \mathbf{x})^n P(z|\mathbf{x}; \tau) d\mathbf{z},$$

where  $\Omega$  denotes the state space of the continuous time process and  $n \in \mathbb{N}$ .  $\mathbf{M}_1(\mathbf{x})$  is the rate of growth of the mean of the vector  $(s, j)'$  when the process is at some state  $\mathbf{x}$ ;  $\mathbf{M}_2(\mathbf{x})$  is the rate of growth of the variance. If the growth rates of the higher moments vanish, i.e.  $\mathbf{M}_n(\mathbf{x}) = \mathbf{0}$  for all  $n \geq 3$ , and the process is Markovian, then the limit of the sequence of processes is described by the limit of the sequence of master equations. In our case, using  $P(z|\mathbf{x}; 0) = \delta(z - \mathbf{x})$  and integration by parts, we get

$$\begin{aligned} \mathbf{M}_1(\mathbf{x}) &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int_{\Omega} (z - \mathbf{x}) [P(z|\mathbf{x}; \tau) - P(z|\mathbf{x}; 0)] d\mathbf{z} \\ &= \int_{\Omega} (z - \mathbf{x}) \frac{\partial P}{\partial t}(z|\mathbf{x}; 0) d\mathbf{z} \\ &= \int_{\Omega} (z - \mathbf{x}) \left( -\frac{\partial}{\partial z} [\mathbf{A}(z) P(z|\mathbf{x}; 0)] + \frac{1}{2} \frac{\partial}{\partial z} \frac{\partial}{\partial z} [\mathbf{\Sigma}(z) P(z|\mathbf{x}; 0)] \right) d\mathbf{z} \\ &= \mathbf{A}(\mathbf{x}). \end{aligned}$$

Similarly,

$$\mathbf{M}_2(\mathbf{x}) = \mathbf{\Sigma}(\mathbf{x}) \quad \text{and} \quad \mathbf{M}_n(\mathbf{x}) = \mathbf{0} \quad \text{for } n \geq 3.$$

The conditions stated above are hence fulfilled and we can from now on consider the process  $\mathbf{x} = (s, j)'$  with continuous state space  $(0, 1)^2$  characterized by the forward diffusion equation (10).

**Transition to stochastic differential equations.** Eventually, according to Kloeden and Platen (1999) and Tory (2000), the Fokker-Planck equation (10) corresponds to a Markov process which is the solution of the multivariate SDE

$$d\mathbf{X}_t = \mathbf{A}(\mathbf{X}_t)dt + \mathbf{L}(\mathbf{X}_t)d\mathbf{B}_t,$$

where  $\mathbf{\Sigma} = \mathbf{L}\mathbf{L}'$ . In our case,

$$\mathbf{L} = \begin{pmatrix} \sqrt{\frac{\alpha sj}{N}} & 0 \\ -\sqrt{\frac{\alpha sj}{N}} & \sqrt{\frac{\beta j}{N}} \end{pmatrix} \quad (11)$$

is one solution. The corresponding system of SDEs is (3), as was to be shown. Since these SDEs agree with our model also for  $s = 0$  and  $j = 0$ , we extend the state space to  $[0, 1]^2$ .

### 3.2 Backward Approach

In this method, we aim to obtain the infinitesimal generator of the discrete Markov process described by (1), as the generator allows us to directly read out the drift vector and diffusion matrix of a corresponding diffusion process. We start by calculating the expected infinitesimal change in space for a function  $f$  over a small time interval of length  $\varepsilon$ :

$$\begin{aligned} \frac{1}{\varepsilon} \mathbb{E} \left( f(\mathbf{X}(t + \varepsilon)) - f(\mathbf{X}(t)) \mid \mathbf{X}(t) \right) &= \frac{1}{\varepsilon} \left( \frac{\alpha}{N} SI\varepsilon [f(S-1, I+1) - f(S, I)] \right. \\ &\quad \left. + \beta I\varepsilon [f(S, I-1) - f(S, I)] + \left( 1 - \frac{\alpha}{N} SI\varepsilon - \beta I\varepsilon \right) [f(S, I) - f(S, I)] \right). \end{aligned}$$

Changing the arguments of the function  $f$  to fractions instead of total numbers, and with  $\varepsilon = N^{-1}$ , the above terms become

$$\begin{aligned} &\alpha sj \varepsilon^{-1} [f(s - \varepsilon, j + \varepsilon) - f(s, j)] + \beta j \varepsilon^{-1} [f(s, j - \varepsilon) - f(s, j)] \\ = &\alpha sj \varepsilon^{-1} \left[ -\frac{1}{2} \left( f(s, j) - f(s - \varepsilon, j) \right) - \frac{1}{2} \left( f(s + \varepsilon, j) - f(s, j) \right) \right. \\ &\quad \left. + \frac{1}{2} \left( f(s, j) - f(s, j - \varepsilon) \right) + \frac{1}{2} \left( f(s, j + \varepsilon) - f(s, j) \right) \right. \\ &\quad \left. + \frac{\varepsilon^{-1}}{2N} \left( f(s + \varepsilon, j) - 2f(s, j) + f(s - \varepsilon, j) \right) \right. \\ &\quad \left. + \frac{\varepsilon^{-1}}{2N} \left( f(s, j + \varepsilon) - 2f(s, j) + f(s, j - \varepsilon) \right) \right. \\ &\quad \left. - \frac{\varepsilon^{-1}}{N} \left( f(s, j + \varepsilon) - f(s - \varepsilon, j + \varepsilon) - f(s, j) + f(s - \varepsilon, j) \right) \right] \\ + &\beta j \varepsilon^{-1} \left[ -\frac{1}{2} \left( f(s, j) - f(s, j - \varepsilon) \right) - \frac{1}{2} \left( f(s, j + \varepsilon) - f(s, j) \right) \right. \\ &\quad \left. + \frac{\varepsilon^{-1}}{2N} \left( f(s, j + \varepsilon) - 2f(s, j) + f(s, j - \varepsilon) \right) \right]. \end{aligned}$$

For  $\varepsilon \rightarrow 0$ , this tends to

$$\left( \alpha sj \left[ -\frac{\partial}{\partial s} + \frac{\partial}{\partial j} + \frac{1}{2N} \frac{\partial^2}{\partial s^2} + \frac{1}{2N} \frac{\partial^2}{\partial j^2} - \frac{1}{N} \frac{\partial^2}{\partial s \partial j} \right] + \beta j \left[ -\frac{\partial}{\partial j} + \frac{1}{2N} \frac{\partial^2}{\partial j^2} \right] \right) f(s, j).$$

Hence

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \mathbb{E} \left( f(\mathbf{X}(t + \varepsilon)) - f(\mathbf{X}(t)) \mid \mathbf{X}(t) \right) = \left( \mathbf{A} \nabla + \frac{1}{2} \nabla' \Sigma \nabla \right) f(\mathbf{X}(t)),$$

where

$$\mathbf{A} = \begin{pmatrix} -\alpha sj \\ \alpha sj - \beta j \end{pmatrix} \quad \text{and} \quad \Sigma = \frac{1}{N} \begin{pmatrix} \alpha sj & -\alpha sj \\ -\alpha sj & \alpha sj + \beta j \end{pmatrix}.$$

We can now (Øksendal, 2003) associate this generator to the Itô diffusion

$$d\mathbf{X}_t = \mathbf{A}(\mathbf{X}_t)dt + \mathbf{L}(\mathbf{X}_t)d\mathbf{B}_t,$$

where  $\mathbf{L}$  is the square root of  $\Sigma$ , compare with (11). Turning to the spatial SIR model (4), we again calculate

$$\begin{aligned} & \frac{1}{\varepsilon} \mathbb{E} \left( f(\mathbf{X}(t + \varepsilon)) - f(\mathbf{X}(t)) \middle| \mathbf{X}(t) \right) \\ &= \frac{1}{\varepsilon} \left[ \sum_{i=1}^n \left( \frac{\alpha}{N_i} S_i I_i \varepsilon \left[ f(S_i - 1, I_i + 1) - f(S_i, I_i) \right] + \beta I_i \varepsilon \left[ f(I_i - 1) - f(I_i) \right] \right. \right. \\ & \quad \left. \left. + \sum_{k=1}^n \gamma_{ik} S_i \varepsilon \left[ f(S_i - 1, S_k + 1) - f(S_i, S_k) \right] + \sum_{k=1}^n \gamma_{ik} I_i \varepsilon \left[ f(I_i - 1, I_k + 1) - f(I_i, I_k) \right] \right) \right], \end{aligned}$$

where we suppress the non-involved components of the process  $\mathbf{X} = (S_1, \dots, S_n, I_1, \dots, I_n)'$ . The first two summands in the sum over all  $i$  can be rewritten as above. We exemplarily expand the third summand, defining  $\varepsilon_i := N_i^{-1}$  and assuming  $(S_i - 1)/(N_i - 1) \approx (S_i - 1)/N_i$ :

$$\begin{aligned} & \sum_{k=1}^n \gamma_{ik} S_i \varepsilon_i^{-1} \left[ f(s_i - \varepsilon_i, s_k + \varepsilon_k) - f(s_i, s_k) \right] \\ &= \sum_{k=1}^n \gamma_{ik} S_i \varepsilon_i^{-1} \left[ -\frac{\varepsilon_k^{-1}}{N_k} \left( f(s_i, s_k + \varepsilon_k) - f(s_i - \varepsilon_i, s_k + \varepsilon_k) - f(s_i, s_k) + f(s_i - \varepsilon_i, s_k) \right) \right. \\ & \quad \left. - \left( f(s_i, s_k) - f(s_i - \varepsilon_i, s_k) \right) + \left( f(s_i, s_k + \varepsilon_k) - f(s_i, s_k) \right) \right], \end{aligned}$$

which converges to

$$\sum_{k=1}^n \gamma_{ik} S_i \left[ -\frac{1}{N_k} \frac{\partial^2}{\partial s_i \partial s_k} - \frac{\partial}{\partial s_i} + \frac{\partial}{\partial s_k} \right] f(s_i, s_k)$$

as  $\varepsilon_i \rightarrow 0$  for all  $i$ . Altogether, we obtain the infinitesimal generator

$$\left( \tilde{\mathbf{A}} \nabla + \frac{1}{2} \nabla' \tilde{\Sigma} \nabla \right) f(\mathbf{X}(t)) \quad \text{where} \quad \tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{A}^s \\ \mathbf{A}^j \end{pmatrix} \quad \text{and} \quad \tilde{\Sigma} = \begin{pmatrix} \Sigma^{ss} & \Sigma^{sj} \\ \Sigma^{sj} & \Sigma^{jj} \end{pmatrix}$$

with components

$$\mathbf{A}_i^s = -\alpha s_i j_i - \sum_k \gamma_{ik} s_i + \sum_k \gamma_{ki} s_k, \quad \mathbf{A}_i^j = \alpha s_i j_i - \beta j_i - \sum_k \gamma_{ik} j_i + \sum_k \gamma_{ki} j_k$$

and

$$\begin{aligned} \Sigma_{ii}^{ss} &= \frac{\alpha s_i j_i}{N_i}, & \Sigma_{ik}^{ss} &= -\frac{\gamma_{ik} s_i}{N_k} - \frac{\gamma_{ki} s_k}{N_i} \quad (k \neq i) \\ \Sigma_{ii}^{jj} &= \frac{\alpha s_i j_i + \beta j_i}{N_i}, & \Sigma_{ik}^{jj} &= -\frac{\gamma_{ik} j_i}{N_k} - \frac{\gamma_{ki} j_k}{N_i} \\ \Sigma_{ii}^{sj} &= -\frac{\alpha s_i j_i}{N_i}, & \Sigma_{ik}^{sj} &= 0. \end{aligned}$$

These are exactly the drift and covariance terms that we obtain from the diffusion process (5). For example, computation of  $ds_i \cdot ds_k$  according to the rules  $dt \cdot dt = dt \cdot dB_t = dB_t \cdot dt = 0$ ,  $dB_t \cdot dB_t = dt$  and  $dB_t^1 \cdot dB_t^2$  for independent Brownian motions  $B^1$  and  $B^2$  yields

$$ds_i \cdot ds_k = -\frac{\gamma_{ik} s_i}{N_k} - \frac{\gamma_{ki} s_k}{N_i}.$$

We can therefore conclude that the diffusion process (5) is a valid approximation to the discrete Markov chain described by (4).

## 4 Simulation and Parameter Estimation for Diffusion Processes

In this section we give a brief introduction to modelling and estimation techniques for SDEs, all of which can be found in the literature. Our emphasis lies on the description of the so-called Euler-Maruyama approximation for non-explicitly solvable SDEs as this will be employed in Sections 5 and 6. Throughout this section, we consider the stochastic Itô differential equation

$$dX_t = a(X_t, \theta)dt + b(X_t, \theta)dB_t, \quad X_0 = x_0,$$

where  $X = (X_t)_{t \geq 0}$  is a stochastic process,  $t$  the time parameter,  $B = (B_t)_{t \geq 0}$  Brownian motion and  $a$  and  $b$  functions that fulfil the Lipschitz condition such that a unique solution of the differential equation exists.  $\theta$  is the possibly vector-valued parameter which is to be estimated.

### 4.1 Explicitly Solvable SDEs

A stochastic process satisfying an analytically explicitly solvable SDE can exactly be simulated for any value of  $\theta$ , and, conversely, for given data  $x_{t_0}, x_{t_1}, \dots, x_{t_u}$  the likelihood function for  $\theta$  can be calculated explicitly. Assume that a stochastic process can be written in the form

$$X_t = x_0 + d(t, \theta) + e(t, \theta)B_t$$

for  $t \geq 0$ , an initial condition  $X_0 = x_0$  and appropriate functions  $d$  and  $e$ . Then, for given  $\theta$ , a path of this process at time  $t$  can be sampled via

$$X_t = x_0 + d(t, \theta) + e(t, \theta) \cdot N(0, t)$$

or, in case the path is already given up to time  $s \geq 0$ , by

$$X_t | \mathcal{F}_s = x_0 + d(t, \theta) + e(t, \theta) \cdot (B_s + N(0, t - s)),$$

where  $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$  denotes the natural filtration. For the simulation of the process at intermediate time points, Brownian bridges can be employed. The log likelihood function of  $\theta$  in case of a time-homogeneous SDE and given data  $x_{t_0}, x_{t_1}, \dots, x_{t_u}$  reads

$$l(\theta) = \sum_{k=1}^u \log(p_{\Delta t_{k-1}}(x_{t_{k-1}}, x_{t_k}; \theta)),$$

where  $\Delta t_{k-1} := t_k - t_{k-1}$  and

$$p_t(v, w; \theta) = P(X_t \in dw | X_0 = v; \theta) / dw$$

is the transition probability from state  $v$  to  $w$  in the time interval  $[0, t]$  for  $v, w \in \mathbb{R}$  and  $t > 0$ . Because of

$$X_{\Delta t_{k-1}} - X_0 \sim N(d(\Delta t_{k-1}), e(\Delta t_{k-1})^2 \Delta t_{k-1})$$

and the time homogeneity of  $X$  we get

$$\begin{aligned} p_{\Delta t_{k-1}}(x_{t_{k-1}}, x_{t_k}; \theta) &= P\left(X_{\Delta t_{k-1}} \in dx_{t_k} \mid X_0 = x_{t_{k-1}}; \theta\right) / dx_{t_k} \\ &= P\left(X_{\Delta t_{k-1}} - X_0 \in dx_{t_k} - x_{t_{k-1}}; \theta\right) / dx_{t_k} \\ &= \phi\left(x_{t_k} - x_{t_{k-1}} \mid d(\Delta t_{k-1}, \theta), e(\Delta t_{k-1}, \theta)^2 \Delta t_{k-1}\right) \end{aligned}$$

for  $k = 1, \dots, u$ , where  $\phi(y | \mu, \sigma^2)$  denotes the probability density of the normal distribution with mean  $\mu$  and variance  $\sigma^2$  at point  $y$ .

## 4.2 Non-explicitly Solvable SDEs

However, only few SDEs are explicitly solvable, so in case the considered SDE is not, alternative methods need to be employed for simulation and parameter estimation. For sampling techniques, an extensive overview of established procedures is provided by Kloeden and Platen (1999). For inference on diffusions, which is a relatively young and highly developing research area, Sørensen (2004) reviews techniques including estimating functions, Bayesian analysis and Markov chain Monte Carlo (MCMC) methods, indirect inference, methods of moments and non-parametric approaches. This collection should be supplemented by an algorithm for exact estimation for discretely observed diffusion processes recently introduced in Beskos, Papaspiliopoulos, Roberts, and Farnhead (2006) and subsequent papers.

In the following, we will only deal with the most basic method for approximate simulation and estimation of SDEs, which is the Euler approximation (also called Euler-Maruyama approximation)

$$X_{t_k} = X_{t_{k-1}} + a(X_{t_{k-1}}, \theta) \cdot \Delta t_{k-1} + b(X_{t_{k-1}}, \theta) \cdot \Delta B_{k-1}, \quad (12)$$

where  $k = 1, \dots, u$  and  $\Delta B_{k-1} := B_{t_k} - B_{t_{k-1}}$ . With this, we can approximately sample the process  $X$  in the time interval  $[0, t]$  for discrete time points  $0 = t_0 < t_1 < \dots < t_u = t$  for given parameter  $\theta$  and initial value  $X_0 = x_0$ . Vice-versa, for the maximum distance between two consecutive instants tending to zero, the distribution of  $X_{t_k}$  conditional on  $X_{t_{k-1}}$  converges to a normal distribution. The conditional mean and variance can be obtained from (12):

$$\mathbb{E}(X_{t_k} | X_{t_{k-1}}) = X_{t_{k-1}} + a(X_{t_{k-1}}, \theta) \Delta t_{k-1} \quad \text{and} \quad \text{Var}(X_{t_k} | X_{t_{k-1}}) = b^2(X_{t_{k-1}}, \theta) \cdot \Delta t_{k-1}.$$

Hence, given data  $\mathbf{x} = (x_{t_0}, x_{t_1}, \dots, x_{t_u})$ , we can write down the approximate log likelihood function of  $\theta$  as

$$l(\theta) = \sum_{i=1}^u \log(\phi(x_{t_i} | x_{t_{i-1}} + a(x_{t_{i-1}}, \theta) \Delta t_{i-1}, b^2(x_{t_{i-1}}, \theta) \cdot \Delta t_{i-1}))$$

with  $\phi$  defined as above. In Sections 5 and 6, we will apply the Euler approximation scheme to estimate the transmission and recovery rates for our epidemic models.

## 5 Simulation Study: Local SIR Model

We now come back to the local SIR model introduced in Section 2, which was originally described as a discrete Markov chain with transitions (1) and which was in Section 3 proved to be approximated by the continuous state space process  $(s, j)'$  satisfying the system of stochastic differential equations (3). In this section we show simulation results for both these model descriptions and compare parameter estimates for  $\alpha$  and  $\beta$  based thereon.

### 5.1 Simulation

We repeatedly simulate the outbreak of an epidemic for different population sizes  $N$ , once based on the original discrete Markov chain model representation (1), and once as an Euler approximation to the stochastic differential equations (3).

For the discrete process, we fix the time step  $\Delta t$  and within one step transfer the discrete random vector  $(S, I) \in [0, N]^2$  to  $(S - 1, I + 1)$  with probability  $\alpha SI/N\Delta t$ , or to  $(S, I - 1)$  with probability  $\beta I\Delta t$ . Otherwise, the state of the Markov chain does not change. Note that  $\Delta t$  has to be sufficiently small to ensure that all involved probabilities are well-defined, i.e.  $\Delta t$  has to be less than  $(\max_{S,I \in [0,N]}(\alpha SI/N + \beta I))^{-1}$ . The initial state should satisfy  $S_0 + I_0 \leq N$ .

For the diffusion process, we employ the Euler scheme to subsequently simulate the state of the process  $(s, j)'$  at discrete equidistant time points in intervals of length  $\Delta t$ . In this case, there is no restriction to the size of the time step, though the approximation clearly improves as  $\Delta t$  tends to zero.

Figure 4 shows the fractions of infectives of ten independent samples for each model description and population sizes  $N = 10^3, 10^4, 10^5$ . For each sample, the step size was chosen to be  $\Delta t = 2.5 \cdot 10^{-5}$ , where time was measured in days. The model parameters were  $\alpha = 0.625$  and  $\beta = 0.25$ , and the initial fraction of infectives was one percent of the population. The stochastic samples twine around the deterministic course of the fractions of infectives, which was obtained with the standard Euler scheme from equations (2). However, the magnitude of the fluctuations decreases for larger population sizes, as was to be expected from (3).

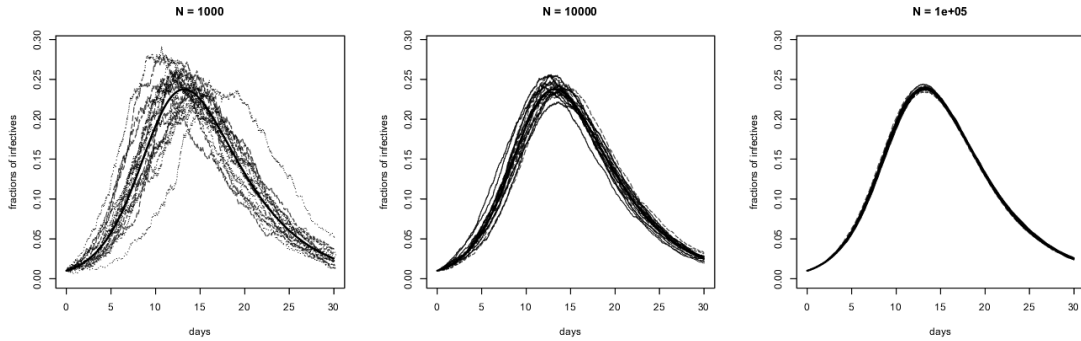


Figure 4: Independent samples from the discrete Markov chain model (dotted) and the discretized diffusion process (dashed) for population sizes  $N = 10^3$  (left),  $N = 10^4$  (middle), and  $N = 10^5$  (right). The solid lines indicate the deterministic course of the epidemic. The model parameters were  $\alpha = 0.625$  and  $\beta = 0.25$ , the fraction of infectives at time zero was 0.01, and the time step was chosen to be  $\Delta t = 2.5 \cdot 10^{-5}$ .

## 5.2 Likelihood Function

In this paragraph, we derive the log likelihood function of the process  $(s, j)'$  satisfying (3) according to Section 4. For the sake of clarity, we introduce functions  $a_p$  and  $b_{pk}$ ,  $p \in \{s, j\}$ ,  $k = 1, 2$ , and rewrite this as

$$d\mathbf{X}(t) = \begin{pmatrix} a_s(\mathbf{X}(t), \boldsymbol{\theta}) \\ a_j(\mathbf{X}(t), \boldsymbol{\theta}) \end{pmatrix} dt + \begin{pmatrix} b_{s1}(\mathbf{X}(t), \boldsymbol{\theta}) & b_{s2}(\mathbf{X}(t), \boldsymbol{\theta}) \\ b_{j1}(\mathbf{X}(t), \boldsymbol{\theta}) & b_{j2}(\mathbf{X}(t), \boldsymbol{\theta}) \end{pmatrix} \begin{pmatrix} dB_1(t) \\ dB_2(t) \end{pmatrix},$$

where  $\mathbf{X} = (s, j)'$  Lund  $\boldsymbol{\theta} = (\alpha, \beta)$ . The associated Euler scheme reads

$$\mathbf{X}(t_k) = \mathbf{X}(t_{k-1}) + \Delta t_{k-1} \underbrace{\begin{pmatrix} a_s(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) \\ a_j(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) \end{pmatrix}}_{=: \mathbf{A}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta})} + \underbrace{\begin{pmatrix} b_{s1}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) & b_{s2}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) \\ b_{j1}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) & b_{j2}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) \end{pmatrix}}_{=: \mathbf{M}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta})} \begin{pmatrix} \Delta B_1(t_{k-1}) \\ \Delta B_2(t_{k-1}) \end{pmatrix}$$

for  $k = 1, \dots, u$ , where  $t_u$  is the upper bound of the considered time interval. The conditional expectation and variance of the discretized process  $\mathbf{X}$  hence are

$$\mathbb{E}(\mathbf{X}(t_k) | \mathbf{X}(t_{k-1})) = \mathbf{X}(t_{k-1}) + \Delta t_{k-1} \cdot \mathbf{A}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta})$$

and

$$\text{Var}(\mathbf{X}(t_k) | \mathbf{X}(t_{k-1})) = \Delta t_{k-1} \cdot \mathbf{M}(\mathbf{X}(t_{k-1}), \boldsymbol{\theta}) \mathbf{M}'(\mathbf{X}(t_{k-1}), \boldsymbol{\theta})$$

for  $k = 1, \dots, u$ . Without any restrictions applying to  $\mathbf{X}$  we could now write down the associated log likelihood function for given data  $\mathbf{x} = (\mathbf{x}_{t_0}, \mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_u}) = ((s_{t_0}, j_{t_0})', (s_{t_1}, j_{t_1})', \dots, (s_{t_u}, j_{t_u})')$ ,

$$l(\boldsymbol{\theta}) = \sum_{k=1}^u \log \phi(\mathbf{x}_{t_k} | \mathbf{x}_{t_{k-1}} + \Delta t_{k-1} \mathbf{A}(\mathbf{x}_{t_{k-1}}, \boldsymbol{\theta}), \Delta t_{k-1} \mathbf{M}(\mathbf{x}(t_{k-1}), \boldsymbol{\theta}) \mathbf{M}'(\mathbf{x}(t_{k-1}))),$$

as introduced in Section 4.2. The function  $\phi(\mathbf{y} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  now denotes the probability density of the bivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  at  $\mathbf{y}$ . However, since the state variables  $s$  and  $j$  are defined as the fractions of susceptible and infectious individuals of the total population, respectively, they underlie the restriction of taking values in the interval  $[0, 1]$ . Hence, as soon as values less than zero or greater than one occur in the Euler simulation, these have to be corrected to zero or one, respectively. Vice-versa, when applying the likelihood function above, we assume the underlying data to be risen from an Euler approximation and hence to be potentially truncated. In order to determine the maximum likelihood estimator of  $\boldsymbol{\theta}$ , we hence maximise a combination of the density and the probability function of appropriately parameterized normal distributions; a justification for this procedure can be found in Klein and Moeschberger (1997). We therefore have to detect where data is possibly truncated. This is the case if either

$$s_{t_k} = 1 \quad \text{or} \quad j_{t_k} = 1 \quad \text{or} \quad s_{t_k} = 0 \wedge s_{t_{k-1}} \neq 0 \quad \text{or} \quad j_{t_k} = 0 \wedge j_{t_{k-1}} \neq 0$$

for any  $k = 1, \dots, u$ . In case of  $s_{t_k} = s_{t_{k-1}} = 0$  or  $j_{t_k} = j_{t_{k-1}} = 0$  we have no truncation since  $s_{t_{k-1}} = 0 \Rightarrow s_{t_k} = 0$ , and analogously for  $j$ . Define

$$\boldsymbol{\mu}(t_k) = \begin{pmatrix} \mu_1(t_k) \\ \mu_2(t_k) \end{pmatrix} = \mathbf{x}_{t_{k-1}} + \Delta t_{k-1} \mathbf{A}(\mathbf{x}_{t_{k-1}}, \boldsymbol{\theta})$$

$$\boldsymbol{\Sigma}(t_k) = \begin{pmatrix} \sigma_1^2(t_k) & \sigma_{12}(t_k) \\ \sigma_{12}(t_k) & \sigma_2^2(t_k) \end{pmatrix} = \Delta t_{k-1} \mathbf{M}(\mathbf{x}(t_{k-1}), \boldsymbol{\theta}) \mathbf{M}'(\mathbf{x}(t_{k-1}), \boldsymbol{\theta})$$

and

$$\mu_1^*(t_k) = \mu_1(t_k) + \frac{\sigma_{12}(t_k)}{\sigma_2^2(t_k)} (j_{t_{k-1}} - \mu_2(t_k)),$$

$$\mu_2^*(t_k) = \mu_2(t_k) + \frac{\sigma_{12}(t_k)}{\sigma_1^2(t_k)} (s_{t_{k-1}} - \mu_1(t_k)),$$

$$(\sigma_1^*(t_k))^2 = \sigma_1^2(t_k) - \frac{\sigma_{12}^2(t_k)}{\sigma_2^2(t_k)},$$

$$(\sigma_2^*(t_k))^2 = \sigma_2^2(t_k) - \frac{\sigma_{12}^2(t_k)}{\sigma_1^2(t_k)}.$$

Then, because of  $(s_{t_k}, j_{t_k}) | (s_{t_{k-1}}, j_{t_{k-1}}) \sim N(\boldsymbol{\mu}(t_k), \boldsymbol{\Sigma}(t_k))$ , we have

$$s_{t_k} | \{s_{t_{k-1}}, j_{t_{k-1}}\} \sim N(\mu_1(t_k), \sigma_1^2(t_k)), \quad j_{t_k} | \{s_{t_{k-1}}, j_{t_{k-1}}\} \sim N(\mu_2(t_k), \sigma_2^2(t_k))$$

and

$$s_{t_k} | \{j_{t_k}, s_{t_{k-1}}, j_{t_{k-1}}\} \sim N(\mu_1^*(t_k), (\sigma_1^*(t_k))^2), \quad j_{t_k} | \{s_{t_k}, s_{t_{k-1}}, j_{t_{k-1}}\} \sim N(\mu_2^*(t_k), (\sigma_2^*(t_k))^2).$$

The log likelihood function for given data  $\mathbf{x} = (\mathbf{x}_{t_0}, \mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_u})$  then reads

$$l(\boldsymbol{\theta}) = \log f(\mathbf{x} | \boldsymbol{\theta}) = \sum_{k=1}^u \log f(\mathbf{x}_{t_k} | \mathbf{x}_{t_{k-1}}, \boldsymbol{\theta}),$$

where  $f(\mathbf{x}_{t_k} | \mathbf{x}_{t_{k-1}}, \boldsymbol{\theta})$  equals

$$\begin{cases} \phi(\mathbf{x}_{t_k} | \boldsymbol{\mu}(t_k), \boldsymbol{\Sigma}(t_k)) & \text{if } s_{t_k}, j_{t_k} \text{ not truncated,} \\ \Phi(I(\mathbf{x}_{t_k}) | \boldsymbol{\mu}(t_k), \boldsymbol{\Sigma}(t_k)) & \text{if } s_{t_k}, j_{t_k} \text{ possibly truncated,} \\ \phi(s_{t_k} | \mu_1(t_k), \sigma_1^2(t_k)) \cdot \Phi(I(j_{t_k}) | \mu_2^*(t_k), (\sigma_2^*(t_k))^2) & \text{if } s_{t_k} \text{ not truncated, } j_{t_k} \text{ possibly truncated,} \\ \phi(j_{t_k} | \mu_2(t_k), \sigma_2^2(t_k)) \cdot \Phi(I(s_{t_k}) | \mu_1^*(t_k), (\sigma_1^*(t_k))^2) & \text{if } s_{t_k} \text{ possibly truncated, } j_{t_k} \text{ not truncated.} \end{cases}$$

Here,  $\Phi(I(z) | \mu, \sigma^2)$  denotes the integral of the univariate normal distribution function with mean  $\mu$  and variance  $\sigma^2$  over an interval  $I(z)$  depending on  $z$ .

### 5.3 Estimation Results

Finally, we apply the log likelihood function obtained above to the data sets sampled in Section 5.1. Recall that the value of the parameter was chosen to be  $(\alpha, \beta) = (0.625, 0.25)$ . Table 1 shows estimation results which are each based on a set of ten independent samples, obtained for the two different model descriptions and different population sizes, taking into account every 1000th sample point. Figure 5 displays the maximum likelihood estimates based on the single sample paths instead of on a set of ten. Maximization of the log likelihood function was performed by applying the R function `optim`. The results indicate that the diffusion process works well as an approximation to the original discrete model, since the estimation based on the data risen from (1) is of the same value as the estimation on the data from (3).

### 5.4 Remarks

The time step  $\Delta t = 2.5 \cdot 10^{-5}$  in Section 5.1 was chosen rather small for both model descriptions and all considered population sizes, which led to very detailed data and hence naturally to satisfying estimation results. However, this data situation does by no means comply with natural ones in epidemics. Still, even a data set of size ten (instead of 1200, as in our case) would have produced similar results. The reason for the chosen time step solely came from the upper bound restriction to  $\Delta t$  when sampling from the discrete Markov process (1) for  $N = 10^5$ . It however fulfilled the purpose of illustrating the resemblance of the sample paths drawn from the two models (1) and (3).

## 6 Application: Influenza in Germany

We now apply the spatial SIR model (5) in order to estimate the transmission and recovery rate for an influenza outbreak in Germany. As subregions we choose 438 rural and urban districts. A description of the data on influenza prevalence and of the connectivity matrix is given in the following.



discrete Markov chain	$N$	$\hat{\alpha}$	lower	upper	$\hat{\beta}$	lower	upper
	$10^3$	0.6198	0.6098	0.6299	0.2461	0.2421	0.2501
	$10^4$	0.6257	0.6217	0.6297	0.2505	0.2489	0.2521
	$10^5$	0.6256	0.6243	0.6269	0.2494	0.2489	0.2500
diffusion approximation	$N$	$\hat{\alpha}$	lower	upper	$\hat{\beta}$	lower	upper
	$10^3$	0.6284	0.6182	0.6385	0.2484	0.2444	0.2524
	$10^4$	0.6273	0.6233	0.6313	0.2485	0.2469	0.2501
	$10^5$	0.6246	0.6233	0.6259	0.2502	0.2497	0.2507

Table 1: Estimation results for  $(\alpha, \beta)$ , each based on a set of ten independent samples which were drawn from the discrete model (1) and the diffusion model (3) for population sizes  $N = 10^3, 10^4, 10^5$  in Section 5.1. From this data, every 1000th sample point was taken into account. The table displays the maximum likelihood estimates  $\hat{\alpha}$  and  $\hat{\beta}$  and the lower and upper bounds of the respective 95% Wald confidence intervals.

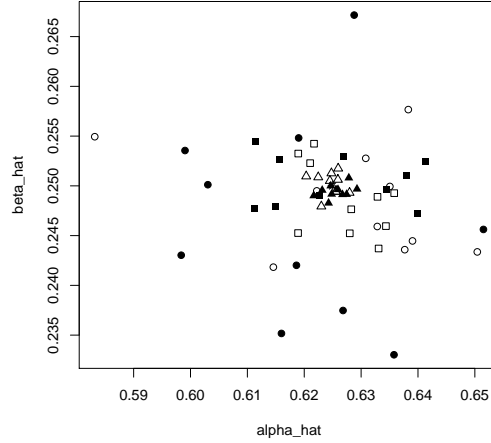


Figure 5: Maximum likelihood estimates for  $(\alpha, \beta)$ , each based on a single sample path obtained in Section 5.1. For the estimation, every 1000th sample point of each path was taken into account. The true parameter  $(\alpha, \beta) = (0.625, 0.25)$  is marked by an asterisk. Filled marks indicate estimates based on discrete model data, unfilled ones refer to those based on diffusion model data. Circles stand for the case  $N = 10^3$ , squares for  $N = 10^4$ , and triangles for  $N = 10^5$ .

## 6.1 Data

**Disease counts.** For the underlying data set on influenza prevalence, we employ a database of the Robert Koch Institute Berlin (RKI) as of 11 November 2006, which is available at <http://www3.rki.de/survstat>. This database contains reported incidences on a weekly basis. We assume the infectious period to usually not last longer than one week and hence data from subsequent weeks to be mutually independent. This way, we equate prevalence and incidence and also obtain the size of the susceptible group by subtracting the number of infectives of the previous week from the foregoing number of susceptibles. Table 2 and Figure 6 show the numbers of influenza A and A/B

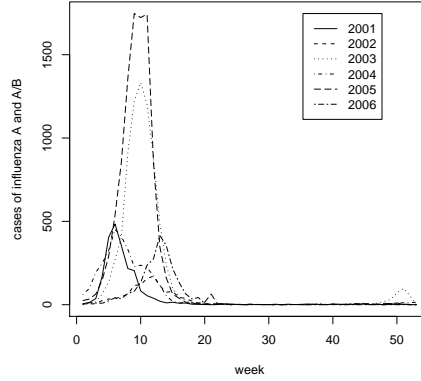


Figure 6: Numbers of reported cases of influenza A and A/B in Germany.

cases which were reported to the RKI between 2001 and 2006. Since the epidemics of subsequent years usually arise from different antigen mutants of the influenza virus, we consider the model parameters  $\alpha$  and  $\beta$  to depend only on the data of the respective season. We thus base the following estimation on the data of weeks 5 to 14 of the year 2005, which constitutes the heaviest reported influenza outbreak.

year	number
2001	2201
2002	1219
2003	7724
2004	3312
2005	10181
2006	2220

Table 2: Numbers of reported cases of influenza A and A/B in Germany.

**Connectivity matrix.** We assume the daily flow of commuters to be a sensible indicator for the migration between the rural and urban districts. We hence investigate data on commuter traffic, which is freely available for few parts of Germany from the Federal Agency for Work. From this data we construct a linear model, assuming that the outbound traffic from one region to another depends on the population densities  $x_{\text{dens1}}$  and  $x_{\text{dens2}}$  of both the origin and the target region (since a high population density can be taken as an indicator for an urban region with many working places and leisure amenities), on the distance  $x_{\text{dist}}$  between the two regions, and on the number  $x_{\text{num}}$  of neighbours of the starting region. In order to ensure that we obtain reasonable values for the rates of traffic  $\gamma_{ij}$ , we employ a logit function which transforms the real line to the interval  $[0, 0.7]$  and achieve

$$\text{logit}(\gamma_{ij}) = -3.207 \cdot 10^{-2} + 7.767 \cdot 10^{-4} x_{\text{dens1}} + 9.056 \cdot 10^{-4} x_{\text{dens2}} - 1.312 x_{\text{dist}} - 9.832 \cdot 10^{-2} x_{\text{num}}.$$

In addition, we take into account domestic air travel in order to cover also the long distance connections. Data on this can be obtained from the OAGflights database at <http://www.oagflights.com>. The overall resulting network is shown in Figure 7.

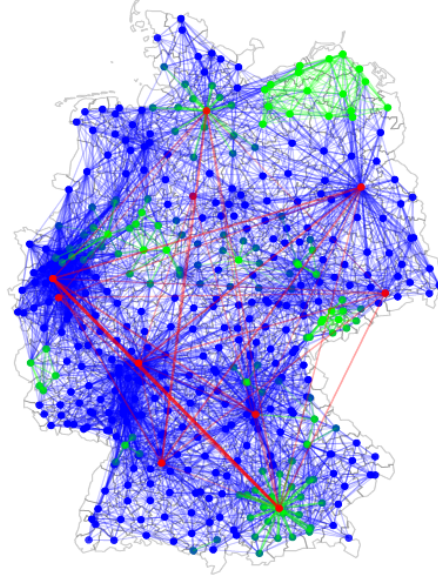


Figure 7: Estimated daily traffic network between the rural and urban districts of Germany. The thickness of each line represents the strength of migration between two regions. Green lines are given data, blue lines estimates based on the former; red lines are domestic air travel.

## 6.2 Results

Eventually, based of the data described above, we set up the log likelihood function for the parameter  $\theta = (\alpha, \beta)$  in the spatial SIR model (5) similarly to how it was done in Section 5.2. The maximum likelihood estimate  $\hat{\theta} = (0.257, 0.425)$ , which we again obtain by maximising the log likelihood function employing the R routine `optim`, is a somewhat chastening result since from real-world experience one would expect the fraction  $\rho = \alpha/\beta$  to be greater than one. However, a non-plausible result was to be foreseen from the fact that we deal with highly underreported disease counts (see Dargatz, Georgescu, & Held, 2006) and that with the Euler scheme we employ a method which approximates a diffusion only for sufficiently small time steps, which is surely not fulfilled for weekly reported cases. However, being aware of these difficulties, we have a strong motivation to further investigate improvements obtained by data augmentation and MCMC methods, which are part of our ongoing work.

## 7 Conclusion

In this report, we transformed a discrete state space epidemic Markov process to a continuous state space diffusion, which is much more convenient for sampling and estimation purposes. We performed simple Euler simulation and employed an approximate likelihood function for the estimation of the epidemic model parameters. Future work is clearly on further treatment of the statistical inference for the multivariate epidemic diffusion model (5), including data augmentation, Bayesian and non-Bayesian analysis and especially MCMC.

## References

- Bailey, N. (1964). *The Elements of Stochastic Processes*. New York: John Wiley & Sons.
- Bailey, N. (1975). *The Mathematical Theory of Infectious Diseases* (2nd ed.). London: Charles Griffin & Co. Ltd.
- Baroyan, O., Rvachev, V., & Ivannikov, Y. (1977). *Modelling and Forecasting of Influenza Epidemics in the Territory of the USSR (in Russian)*. Moscow.
- Beskos, A., Papaspiliopoulos, O., Roberts, G., & Farnhead, P. (2006). Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, 68, 333-382.
- Brownstein, J., Wolfe, C., & Mandl, K. (2006). Empirical evidence for the effect of airline travel on inter-regional influenza spread in the United States. *PLoS Medicine*, 3, 1826-1835.
- Clancy, D., & French, N. (2001). A stochastic model for disease transmission in a managed herd, motivated by *Neospora caninum* amongst dairy cattle. *Mathematical Biosciences*, 170, 113-132.
- Clancy, D., O'Neill, P., & Pollett, P. (2001). Approximations for the long-term behavior of an open-population epidemic model. *Methodology and Computing in Applied Probability*, 3, 75-95.
- Colizza, V., Barrat, A., Barthélemy, M., & Vespignani, A. (2006a). The modeling of global epidemics: Stochastic dynamics and predictability. *Bulletin of Mathematical Biology*, 68, 1893-1921.
- Colizza, V., Barrat, A., Barthélemy, M., & Vespignani, A. (2006b). The role of the airline transportation network in the prediction and predictability of global epidemics. *Proceedings of the National Academy of Sciences*, 103, 2015-2020.
- Dargatz, C., Georgescu, V., & Held, L. (2006). Stochastic modelling of the spatial spread of influenza in Germany. *Austrian Journal of Statistics*, 35, 5-20.
- Fu, Y. (2005). Measuring personal networks with daily contacts: A single-item survey question and the contact diary. *Social Networks*, 27, 169-186.
- Goel, N., & Richter-Dyn, N. (1974). *Stochastic Models in Biology*. New York: Academic Press.
- Hufnagel, L., Brockmann, D., & Geisel, T. (2004). Forecast and control of epidemics in a globalized world. *Proceedings of the National Academy of Sciences*, 101, 15124-15129.
- Klein, J., & Moeschberger, M. (1997). *Survival Analysis*. New York: Springer.
- Kloeden, P., & Platen, E. (1999). *Numerical Solution of Stochastic Differential Equations* (3rd ed.). Berlin, Heidelberg, New York: Springer.
- Nasell, I. (2002). Stochastic models of some endemic infections. *Mathematical Biosciences*, 179, 1-19.
- Øksendal, B. (2003). *Stochastic Differential Equations. An Introduction with Applications* (6th ed.). Berlin, Heidelberg: Springer.
- Pollett, P. (2001). Diffusion approximations for ecological models. *Proceedings of the International Congress of Modelling and Simulation*.

Risken, H. (1984). *The Fokker-Planck Equation*. Berlin: Springer.

Sørensen, H. (2004). Parametric inference for diffusion processes observed at discrete points in time: A survey. *International Statistical Review*, 72, 337-354.

Tory, E. (2000). Stochastic sedimentation and hydrodynamic diffusion. *Chemical Engineering Journal*, 80, 81-89.