

# Epidemic models on social networks—With inference

Tom Britton 

Stockholm University, Stockholm, Sweden

## Correspondence

Department of Mathematics, Stockholm University, 106 91 Stockholm, Sweden.  
Email: tom.britton@math.su.se

## Funding information

Vetenskapsrådet, 2015-05015

This article considers stochastic models for the spread of an infection in a structured community, where this structured community is itself described by a random network model. Some common network models and transmission models are defined and large population properties of them are presented. The focus is then shifted to statistical methodology: what can be estimated and how, depending on the underlying network, transmission model, and the available data? This survey article discusses several different scenarios, giving references to publications where more details can be found, and identifies important open problems.

## KEYWORDS

control measures, epidemic models, incidence data, random networks, sequence data, statistical inference

## 1 | INTRODUCTION

In the current article, we are concerned with stochastic models for how an infectious diseases spreads in a community, where the social structure of relevance for the disease spreading is described by a random network model. We will describe different situations when this can be applicable, give many references to where more can be found on the particular topic, and highlight some open problems on the way. A general reference containing more details and other perspectives is the recent book on the topic by Kiss, Miller, and Simon (2017).

In certain cases the underlying social structure is known, the presence of households being the prime example. Here we put more focus on the case where the underlying structure is not entirely known, which explains why a random network model is advocated.

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2020 The Authors. *Statistica Neerlandica* published by John Wiley & Sons Ltd on behalf of VVS.

Which network model to use will depend on the infectious disease under consideration and the community upon which it spreads. If considering diseases with airborne spreading like influenza and childhood diseases, the network should reflect pairs of individuals being in proximity of each other on a regular basis (preferably also adding additional random contacts). If spreading occurs through close physical contact such as Ebola, the network edges connect pairs of individuals having such contacts on regular basis, and if considering a sexually transmitted infection (STI), the underlying network will be that of sexual contacts.

The type of network model to be used hence depends on the disease and context. If we are considering short-term outbreaks, a static network may be sufficient, whereas if we are interested in longer time spans, a dynamic model, where individuals die and new are born but also where connections are dropped and new are created, may be preferred. Individuals are usually not all the same: some individuals have more social contacts than others, which should be reflected in the *degree distribution* (the degree of an individual is the number of social connections it has). Other characteristics of a network are *clustering*, reflecting how common triangles are in the network, and *degree correlation*: whereby positive (negative) degree correlation implies that individuals with high degree tend to mix with individuals of high degree. Given a set of such *egocentric* network properties, the network is then often defined by saying that, other than obeying to prespecified properties, the network is chosen randomly among all networks satisfying these properties. Of course, more complicated networks may also be considered, for example, by allowing different types of individuals, and having weighted edges affecting the transmission probabilities.

Our interest lies in how an infectious disease can spread on the (random) network. There are different infectious disease models, the base model being an Susceptible-Infectious-Recovered (SIR) epidemic model, where individuals are first susceptible, and if they get infected, they become infectious and later they recover and stay immune for the rest of the study period. Other versions are SEIR, where exposed individuals are first latent before they become infectious, SIRS where immunity eventually vanishes and the individual becomes susceptible again, and so on. In reality, infectivity usually builds up over time and after some further time starts dropping down to 0. Here we simplify the situation by assuming that an infected individual has constant infectivity during the whole infectious period  $I$ , and we focus on the two situations where  $I$  is either fixed and the same for all, or otherwise exponentially distributed. The two models are referred to as the Reed-Frost and the Markovian models, respectively. The Reed-Frost model is often studied in its discrete time version where infections happen sequentially in generations. The Markovian model assumes that infectious individuals infect each of their susceptible neighbors independently at rate  $\beta$  and recover at rate  $\gamma$ , whereas the Reed-Frost model assumes that an individual infected in generation  $k$  infects each of its susceptible neighbors in the next generation independently with probability  $p$  and then recovers. In both models, it is possible to also allow for transmission with randomly chosen individuals beside the neighbors in the network.

One reason for studying epidemic models on networks is to better understand what network features affect spreading the most, and in particular how it is possible to reduce spreading by means of public health measures such as: vaccination, (quicker) diagnosis and treatment, isolation, travel restrictions, and so on. This can be achieved by incorporating the relevant preventive measure(s) into the model and analyzing the outcome and then comparing with the outcome without prevention.

In order to draw conclusions about real-life situations, it is necessary to fit the models to the real-world situation, preferably by collecting network and/or disease data to infer model parameters from using proper statistical methods. If the entire underlying network is observed, this

is often straightforward. However, as mentioned above this situation is not a common situation. Instead some egocentric data may be available, perhaps together with some outbreak data, from which to perform inference, and then statistical methods are more complicated and many problems remain open.

In the current article, we will describe such network models, transmission models “on” the networks, models capturing control measures, and their inferential procedures. Needless to say, this whole area is bigger than that can be captured in one review article, so we will only touch upon most models and leave out several. Another focus of the article is to describe important unsolved problems, with the aim to stimulate more work in this important research area.

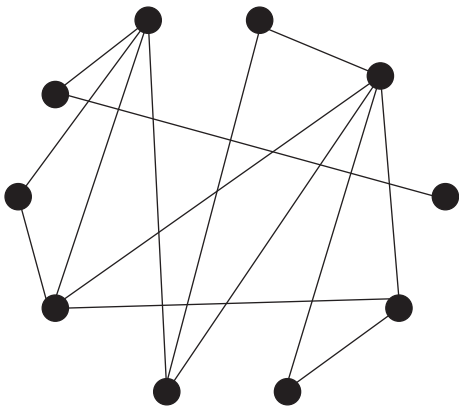
We start by describing a few different random network models (Section 2), then describe the two transmission models mentioned above in more detail, followed by models including prevention (Sections 3 and 4). In Section 5, we present known properties of the models with and without control measures, and in Section 6, we describe how to perform inference for several different models and data settings, also mentioning several open problems.

## 2 | SOCIAL NETWORK MODELS

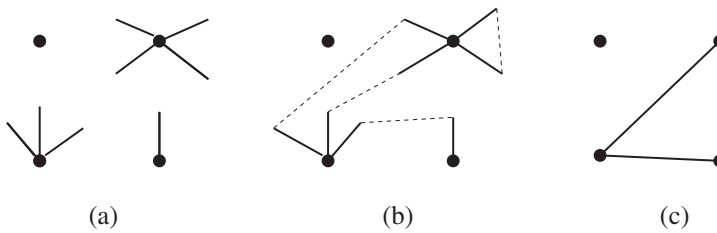
A network consists of nodes and edges. In our application, the nodes will be individuals and the edges, between pairs of individuals, reflect some type of social relationship. Unless otherwise mentioned, we consider static, unweighted, undirected edges, all being of the same type as shown in Figure 1. In what follows, we will assume that there are  $n$  nodes and that pairs of nodes are connected by an edge at random. We assume that the population size  $n$  is large and that the number of edges connecting pairs of individuals is of the same order  $O(n)$  as the number of nodes, implying that each individual has a mean degree  $E(D) = \mu$  ( $0 < \mu < \infty$ ) which remains fixed as  $n \rightarrow \infty$ , implying that the graph/network is sparse.

There exists several well-known random network models (the synonym “random graph models” is often used in more mathematically oriented articles).

The first and most well-studied random network model is the Erdős-Rényi random graph (Erdős & Rényi, 1959). This model has the least possible structure. It assumes that every pair of individuals is connected to each other, independently, with probability  $\lambda/n$ . This model contains a single parameter  $\lambda$  being the mean degree of individuals. An individual has  $n - 1$  possible connections, each being present with probability  $\lambda/n$ , so the number of neighbors any individual



**FIGURE 1** Illustration of a small random social network. In this network nodes have degrees between 1 and 4, and the mean degree equals  $E(D) = \mu = 2.8$



**FIGURE 2** Illustration of the configuration model for a very small network. (a) The degree of each vertex is drawn i.i.d. from the degree distribution. (b) Stubs are paired completely at random. (c) Multiple edges and self-loops are removed thus producing the final network

has is binomially distributed:  $D \sim \text{Bin}(n-1, \lambda/n)$  (so the exact mean is actually  $\lambda(1 - 1/n)$ ). As  $n \rightarrow \infty$ , it is well known that this distribution tends to the Poisson distribution with mean  $\lambda$ , so  $D \approx \text{Po}(\lambda)$ .

A second well-studied model is the Configuration model (e.g., Molloy & Reed, 1998 or Bollobás, 2001). This model is specified by an arbitrary degree distribution  $D \sim \{p_k\}_{k \geq 0}$ , where  $p_k = P(D = k)$ . The model is defined as follows (see Figure 2 for an illustration). Label the nodes  $1, \dots, n$ . Draw independent random variables  $d_1, \dots, d_n$  from  $D$  and create  $d_i$  stubs going out of node  $i$ ,  $i = 1, \dots, n$ , and put all these stubs into one long list of stubs. The network is then constructed by sequentially connecting the stubs pairwise at random. The first stub, of individual 1 say, is connected to any of the remaining stubs uniformly at random. This pair of stubs is then removed from the list of stubs, and the procedure is repeated until there are no stubs remaining on the list. This procedure may result in multiple edges between certain pair of nodes, self-loops, and if the number of stubs happens to be odd, there will be one remaining stub which cannot be paired. It has been proven that the *number* of such imperfections is bounded in probability as  $n$  tends to infinity, as long as the degree distribution has finite mean  $E(D)$ . As a consequence, removing multiple edges (keeping just one), self-loops and one possible odd stub will then have negligible effect on the network and its degree distribution. The configuration network is hence the network obtained after removing multiple edges, self-loops and the possible remaining odd stub.

Another popular network model is the Preferential attachment model due to Barabási and Albert (1999). This model contains, in its simplest form, one parameter  $r$  being a positive integer and is defined sequentially starting with a single node without any edges. At each time step  $k$ , one node equipped with  $r$  edges is added to the existing network. Each of  $r$  edges is connected independently to existing nodes, and each edge is connected to a specific node with a probability that is proportional to that node's current degree. As a consequence, there is a tendency to attach to nodes already having high degree: the “preferential” feature of the model. This procedure is continued until there are  $n$  nodes present in the network (and  $r(n-1)$  edges).

The last model we will describe is the Small world model by Watts and Strogatz (1998). In this model, all nodes are labeled and put on a line, which is made to a circle. The model has two parameters,  $k$  and  $\rho$ , the former being a positive integer and the latter a number between 0 and 1 (typically close to 0—sometimes it is scaled to  $\rho/n$ ). Each node is first connected to its  $k$  closest neighbor on both sides. For example, individual 1 is connected to  $2, \dots, 6$  and  $n, \dots, n-4$  if  $k = 5$ . Then each edge, independently and with probability  $\rho$ , rewires one of its end nodes to a node selected uniformly at random.

The models described above are all for static undirected networks with one type of nodes and no weight on the edges. There are many other random network models defined for particular

purposes. For example, Ball, Britton, and Sirl (2013) define a configuration type model, allowing also for arbitrary clustering and degree correlation.

Other models are dynamic in nature, allowing both individuals to die and new individuals are born, and/or for existing edges to disappear and new ones to appear. Often edges between unconnected pairs of individuals appear randomly in time, each potential edge with a fixed common rate, and existing edges disappear at constant rates, and the population size is either constant and equal to  $n$  (an individual who dies is replaced by a newborn individual without connections), fluctuates around  $n$ , or is growing according to a supercritical branching process (e.g., Britton, Lindholm, & Turova, 2011). There are also models where the network is affected by the ongoing epidemic, for example, individuals distancing themselves from infected individuals (e.g., Leung, Ball, Sirl, & Britton, 2018), but these are harder to analyze and will not be discussed further.

If there are different types of individuals, it is often natural to let the probability, of pairs of individuals to be connected, depend on the types of the two individuals. One such model, where the type space is continuous and reflecting “social activity,” is the Poissonian random graphs (Norros & Reittu, 2006). There each individual  $i$  is given a social activity weight  $W_i$  and the probability that individuals  $i$  and  $j$  are connected is given by  $\min(W_i W_j / n, 1)$  or a similar expression. A more specific model is where there are two types of individuals, and where an individual is only connected to individuals of the other type, heterosexual networks, and client-server-networks being two examples. Such networks are referred to as *bipartite networks* (Newman, Watts, & Strogatz, 2002).

Other random network models allow for directed as well as undirected edges, and edges having different weights, typically reflecting “closeness” of the social relationship (e.g., Barrat, Barthélemy, Pastor-Satorras, & Vespignani, 2004; Spricer & Britton, 2015;). A different class of network models are Random block models, where individuals can be grouped into a small number of (known or more often unknown) subcommunities, where the probability of two nodes being connected depends on the groups of the two involved individuals (e.g., Nowicki & Snijders, 2001).

Many of the models discussed above can be put under the general framework of inhomogeneous random graph models for which there is a rich theory of results available (e.g., Bollobás, Janson, & Riordan, 2007).

A slightly different very flexible class of random graph models are the Exponential random graph models, ERGMs in short (e.g., Snijders, Pattison, Robins, & Handcock, 2006). This class of models is inspired by statistical physics and allows for penalizing or favoring more or less any network feature. It could favor/penalize individual edges, but most often it does so for summary statistics  $S_1, \dots, S_k$ , such as mean degree, number of triangles, high/low degree correlation, and higher moments of the degree distribution. Given the set of chosen network statistics and corresponding model parameters  $\theta_1, \dots, \theta_k$ , the probability of a specific network/graph outcome  $G$  is defined by

$$P(G) \propto e^{\sum_j \theta_j s_j(G)},$$

where  $s_j(G)$  is the value of the summary statistic for the network  $G$ , and the proportionality constant is given by 1 divided by the corresponding sum over all  $2^{\binom{n}{2}}$  possible networks  $G_i$  of size  $n$ . There is no direct method for generating such networks, instead they are obtained by starting with an initial network and then adding/removing edges in an MCMC like manner until the chain is close to stationarity. Probabilistic properties of such graphs are less explored compared with the

models discussed above, and they have mainly been useful for inferring the importance of various network measures in small social networks. Their application for epidemics on networks is largely yet to be shown.

### 3 | INFECTIOUS DISEASE SPREADING MODELS

In the previous section, we described several random network models. In the current section, we assume the network is given to us, generated from a suitable network model or simply known. We now describe some epidemic models for such a network.

The two models we describe are so called SIR epidemic models, where individuals are first Susceptible, and if they get infected they become Infectious and after a while Recover and become immune.

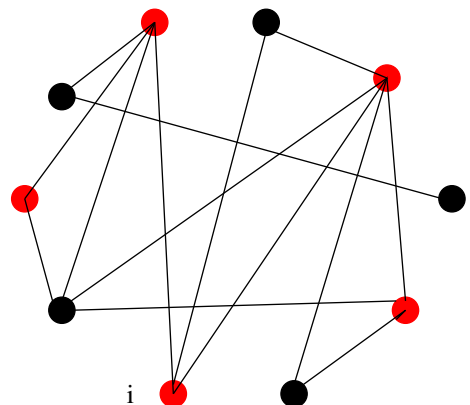
We describe first the discrete time epidemic model (Reed-Frost) and then a continuous time Markovian model. Both are defined on a static undirected network  $G$ .

**Definition 1** (The discrete time Reed-Frost epidemic on a network). Initially, in generation  $k = 0$ , one randomly chosen index case is infectious, and the rest of the population/network is susceptible. Individuals who are infectious in generation  $k$ , infect each susceptible neighbor in  $G$ , independently with probability  $p$  and then recover (nothing happens with immune or infectious neighbors). Those who become infected by at least one infective become infectious in generation  $k + 1$ . The epidemic goes on until the first generation  $T$  at which no new infections arise and at this point the epidemic stops. Each individual is then either susceptible or recovered (i.e., has been infected). How the infected individuals are distributed in the network is characterized by the “final outbreak.” The number of individuals who get infected during the course of the epidemic (including the index case) is denoted by  $Z$  and called the “final size.”

Figure 3 illustrates the final outbreak in a small network, with final size  $Z = 5$ .

*Remark 1.* It is possible to have some other, random or nonrandom, set of index cases. The model is then the same except that it is started by more than one index case. This model can also be defined in continuous time as described in Remark 2.

**Definition 2** (The continuous time Markovian epidemic on a network). Initially, at time  $t = 0$ , one randomly chosen index case is infectious, and the rest of the population/network is susceptible. While an individual is infectious, it has infectious contacts with each susceptible neighbor in



**FIGURE 3** Illustration of the final outbreak in a small random social network (red nodes have been infected and black nodes have not). The final size of the outbreak is  $Z = 5$

$G$  randomly in time according to independent Poisson processes with rate  $\beta$ . Each infected individual remains infectious for a period  $I \sim \text{Exp}(\gamma)$  (exponentially distributed with mean  $1/\gamma$ ) after which it recovers and becomes immune. All infectious periods and contact processes are defined independently. The epidemic goes on until the first time  $T$  that there are no infectious individuals when the epidemic stops. The network then consists of susceptible and recovered (previously infectious) individuals. How they are distributed in the network is characterized by the “final outbreak.” The overall number of individuals that get infected during the course of the epidemic (including the index case) is denoted by  $Z$  and called the “final size.”

**Remark 2.** It is not hard to show that the Markovian network epidemic model can allow for a random latent period upon infection and before becoming infectious, without affecting who gets infected at the end (but of course affecting its time dynamics). The model can also be extended to let the infectious period  $I$  follow an arbitrary random distribution; however, then the model is no longer Markovian. One particular choice is when  $I \equiv \iota$  (nonrandom and equal for all individuals), denoted the “continuous time Reed-Frost epidemic.” It can be shown that the distribution of the final size for this choice of infectious period is identical to that of the discrete time Reed-Frost epidemic if the two models are calibrated by  $p = 1 - e^{-\beta \iota}$  (e.g., Diekmann, Heesterbeek, & Britton, 2013, section 3.2.1). The time dynamics of the two models are however different: in the discrete time version, all events happen at discrete generations, whereas in the continuous time Reed-Frost model, an individual can infect a neighbor any time during the infectious period and if it infects two neighbors, these events will happen at distinct points in time.

The two models described above are different in that the first model considers the disease outbreak to occur in discrete time referred to as generations and the latter in continuous time. In reality epidemic outbreaks of course take place in continuous time. However, if there is a fairly long latency period, followed by a short concentrated infectious period (such as measles and Ebola), then the description of generations makes sense, at least in the beginning of an outbreak. As described in Remark 2, there is also a continuous time version of the Reed-Frost network epidemic. The important mathematical difference between the Markovian network epidemic and either of the two Reed-Frost network epidemics lies in that the events of infecting different neighbors are *independent* in the Reed-Frost models, whereas they are positively correlated in the Markovian version. This is easy to show mathematically and comes from the fact that, when the infectious period has random length, the event to infect a given neighbor indicates a long infectious period with higher probability, which increases the risk for infecting also other friends.

In both models above it is only possible to infect neighbors in the network, where being neighbors reflects a social proximity of relevance for the disease under consideration (e.g., daily close contact). For many infectious diseases transmission also occurs from more random types of contacts, like sitting next to each other on a bus. The models defined above can add such random contacts as we now define.

**Definition 3** (Network epidemic models with global contacts). Start with either of the network epidemics defined in Definition 1 or 2 (or the extensions described in their remark). The discrete time network epidemic model with global contacts is defined from the discrete time network epidemic model as follows: at each time step and for each infective individual, beside potentially infecting susceptible neighbors in  $G$  with probability  $p$ , infective individuals can now also (globally) infect all susceptible individuals in the whole community (neighbor or not) independently with probability  $\beta_G/n$  ( $n$  is the population size). For the continuous time version where infectives infect each susceptible neighbor at rate  $\beta$ , we define the continuous time Markov epidemic



with global contacts by now also letting infectives have infectious (global) contacts with all other susceptible individuals (neighbor or not) independently at rate  $\beta_G/n$ .

*Remark 3.* As  $n$  grows, the rate of infecting a given nonneighbors is hence much smaller than infecting a given neighbor. The random global contacts are, for mathematical convenience, defined to happen also with neighbors. The total infection rate with a susceptible neighbor is hence  $\beta + \beta_G/n$  which can be approximated by  $\beta$  when the community size  $n$  is large. The contact rate/probability  $\beta_G/n$  with a specific individual outside the household is very small in a large community (as it should be). However, the *overall* rate that an infective makes global contacts equals  $n\beta_G/n = \beta_G$  which is not negligible.

There are numerous extensions to these models. Individuals may be categorized into different types, and the transmission rate/probability can then depend on the two types involved; there may be different types of edges, with each edge type having a specific infection rate/probability; the network may be dynamic, where infection only can take place along currently existing edges.

The models defined above assumed that all individuals beside the index case were susceptible. In reality, this may not be the case, since there might be immunity due to prior infection or vaccination in the community. If there are immune individuals in the community and assuming immunity remains over the time horizon of interest, such individuals can simply be neglected in the analysis, and edges connecting to them as well. The degree of susceptible individuals should hence reflect the number of *susceptible* neighbors, which should be kept in mind when, for example, estimating degree distributions from census data, so-called egocentric network data.

## 4 | MODELS FOR PREVENTION

In the previous section, it was mentioned that only initially susceptible individuals should be considered, whereas initially immune individuals and their connections could simply be ignored. In the current section, we assume all individuals to be initially susceptible, but now we consider what happens if some are immunized, for example, by vaccination (from now on we call them vaccinated). It is then natural to keep track also of the vaccinated individuals in order to study the effect of vaccination. Suppose that a vaccine giving complete immunity is available and that this can be distributed prior to the start of the outbreak. Mathematically, this can be modeled by labeling vaccinated individuals as recovered/immune and removing them and all their connections from the network. The effect is that the size of the network of susceptible (unvaccinated) becomes smaller, but more importantly, that the degrees of remaining individuals are reduced. The practical effect is hence that all vaccinated individuals are protected, but also that the nonvaccinated individuals profit in that they have fewer neighbors who can infect them. It is also possible to consider vaccines giving partial immunity (which reduces the risk of getting infected from neighbors), and possibly also reducing infectivity in case of still getting infected (see Halloran, Longini, & Struchiner, 2010), for a description of such a model and how the vaccine effects are estimated), but we will not consider such models further.

If a fraction  $v$  of the community are vaccinated prior to the outbreak, it is of interest to see the effect of such a control program as compared to no prevention. In order to study this (which we briefly do in the next section), it is not enough to specify what fraction was vaccinated, and it is also necessary to specify the degrees of the vaccinated individuals. Needless to say, vaccinating individuals with many neighbors is better from a public health perspective than vaccinating individuals with no or few connections.



A vaccination policy which perhaps is practically easy to implement is where candidates are chosen at random, so that the group of vaccines is a uniformly-at-random chosen fraction  $v$  of the community. We call this strategy the *uniform vaccination strategy*. There are several other more efficient vaccination strategies, but these are often harder to implement practically as well as modelwise. One strategy is to choose the individuals with largest degree, so the fraction  $v$  being vaccinated are those with highest degree. This is often the best or close to best strategy, but on the other hand, it is rarely the case that the degree of individuals is known. Instead often other proxies are used to reach socially active individuals. For instance, when considering STIs, condoms are sometimes distributed freely at night clubs/discos and/or at STI clinics, thus reaching sexually active people. Another more mathematically formulated strategy is the so called *acquaintance vaccination strategy* on a configuration network in Cohen, Havlin, and Ben-Avraham (2003). In this strategy, individuals are chosen uniformly and then, for each selected individual, one of its neighbors are vaccinated, and this is done until a community fraction  $v$  has been vaccinated (in this procedure, it will happen that some selected neighbors have been vaccinated at an earlier occasion and then no vaccine is given, so the fraction of individuals selected at random will hence exceed  $v$ ). By vaccinating friends of randomly selected individuals, rather than the individuals themselves, individuals being vaccinated will tend to have higher degree. This follows from the somewhat sad network property that “your friends (typically) have more friends than you do.” In practice, the acquaintance vaccination strategy seems like an unethical vaccination strategy—selecting people and vaccinating a friend but not themselves, but a related strategy that has been implemented in practice is to give a vaccine (or other means of protection) to randomly selected individuals *and* their partners. A different form of vaccination strategy, often used for serious diseases and when not too many infectives are present, like smallpox and Ebola, is ring vaccination: whenever a new case is detected, all acquaintances (neighbors in the social network) are immediately vaccinated. Ring-vaccination may also be equipped with a second layer, meaning that also all acquaintances of the acquaintances to the case are vaccinated (cf. Henao-Restrepo et al., 2017).

There are many different forms of prevention. Mathematically, these can be classified into two different types. One aims at reducing the rate of contact between individuals, and the other aims at reducing the risk of transmission upon contact. Isolation (of infectious individuals) belongs to the first group, whereas vaccination to the second—so we have just seen that their mathematical effect may not differ, but often they do. When it comes to preparedness for a new pandemic influenza, school closure is often considered as one option for reducing disease spreading (e.g., Cauchemez, Valleron, Boëlle, Flahault, & Ferguson, 2008), but also vaccination once a vaccine has been developed for the new strain (e.g., Longini et al., 2005). When it comes to more serious diseases spreading locally, like SARS and Ebola, traveling bans are often discussed and their effects modeled (e.g., Poletto et al., 2014), and for STIs, increasing condom use is often recommended. School closure and traveling bans reduce contact rates, whereas condom use decreases the risk for transmission.

## 5 | MODEL PROPERTIES

In the current section, we state some results for the network epidemics defined above. We do this without complete rigor in order to avoid too many special cases and assumptions. Before doing this, we define the most important quantity of the network epidemic model.

**Definition 4** (Reproduction numbers). Consider a network epidemic taking place in a large community. The basic reproduction number is denoted by  $R_0$  and is (loosely) defined as the expected number of new infections caused by a typical infected individual during the early stage of the epidemic. The preventive reproduction number after a vaccination strategy  $S$  with vaccination coverage  $v$  has been implemented, is denoted by  $R_v^{(S)}$ , and is defined as the corresponding number after vaccination strategy  $S$  has been implemented.

The main reason why these reproduction numbers are important lies in their relation to the threshold value 1, which determines whether a major outbreak is possible or not:

**Result 1.** For the network epidemics defined in previous sections (and for a very wide class of epidemic models), the final *fraction* getting infected (during the entire outbreak),  $\tau_n = Z/n$ , where  $Z$  is the final size, satisfies  $\tau_n \rightarrow 0$  in probability if and only if  $R_0 \leq 1$ . Similarly, if  $R_0 > 1$ , a vaccination strategy  $S$  with vaccine coverage  $v$  has the effect that  $\tau_n \rightarrow 0$  in probability if and only if  $R_v^{(S)} \leq 1$ .

The implication of Result 1 is that a network epidemic having  $R_0 > 1$  may result in a positive community fraction getting infected (for many models, it is also possible to determine this fraction but we omit this type of results here). The aim for any vaccination (or other preventive) program is therefore to reduce the reproduction number to below 1, that is, to obtain  $R_v^{(S)} \leq 1$ , assuring that there will be no major outbreak infecting a positive fraction—instead only sporadic cases will appear. This state is known as “herd immunity,” which means that also unvaccinated individuals are protected from an outbreak.

It can be hard to determine the basic reproduction number for a network epidemic model, and for many network models with complicated structure, these are not available. However, for several simple networks,  $R_0$  is known:

**Result 2.** Consider the Erdős-Renyi network, the Configuration model network or the Preferential attachment network, having degree distribution  $D \sim \{p_k\}$  with mean  $\mu_D$  and variance  $\sigma_D^2$ . Let  $\tilde{D}$  have distribution  $\sim \{\tilde{p}_k\}$ , where  $\tilde{p}_k = kp_k/\mu_D$ . The basic reproduction number for the Reed-Frost epidemic (with transmission probability  $p$ ) on these networks is then given by

$$R_0^{(RF)} = p (E(\tilde{D}) - 1) = p \left( \frac{\sum_k k^2 p_k}{\mu_D} - 1 \right) = p \left( \mu_D + \frac{\sigma_D^2 - \mu_D}{\mu_D} \right).$$

The basic reproduction number for the Markovian epidemic on these networks is given by

$$R_0^{(M)} = \frac{\beta}{\beta + \gamma} (E(\tilde{D}) - 1) = \frac{\beta}{\beta + \gamma} \left( \frac{\sum_k k^2 p_k}{\mu_D} - 1 \right) = \frac{\beta}{\beta + \gamma} \left( \mu_D + \frac{\sigma_D^2 - \mu_D}{\mu_D} \right).$$

The basic reproduction number for the two models also having random global contacts is as above but adding the term  $\beta_G$  in the Reed-Frost model and the term  $\beta_G/\gamma$  in the Markov model. For the preferential attachment model, it is known that the degree distribution  $D$  has infinite variance, which hence implies that  $R_0 = +\infty$ .

We will not prove this result, but give a quick heuristic explanation of the first equality in each row—the latter equalities follow from simple algebra. The degree of infectives during the early stages of an epidemic is different from the degree of individuals in the community at large. Clearly, individuals with degree 0 (i.e., having no connections) will never get infected. More precisely, an

individual with degree  $k$  is  $k$  times more likely to get infected than an individual having degree 1. The consequence of this is that the degree of an infective at the early stage of an outbreak equals  $k$  with probability proportional to  $k p_k$ . This so called size-biased degree distribution is denoted  $\tilde{D}$ , with outcome probabilities  $\tilde{p}_k = k p_k / \mu_D$ . In the beginning of an outbreak, all neighbors except the infector will be susceptible, so there are  $\tilde{D} - 1$  possible individuals to infect, and the probability to infect any given neighbor is  $p$  in the Reed-Frost epidemic and  $\beta / (\beta + \gamma)$  in the Markov model. This explains the first equalities above: the expression to the right is the probability of infecting a susceptible neighbor multiplied by the expected number of susceptible neighbors of infected individuals during the early stage of an outbreak. The added term for the models also having random contacts is simply the mean number of such global infectious contacts (all will be with susceptible during the early stages of the outbreak).

The corresponding reproduction number after a vaccination program has been initiated is often more complicated to derive. However, the uniform vaccination strategy with vaccination coverage  $v$  has a simple form which we now show, together with its critical vaccination coverage. The critical vaccination coverage for a vaccination strategy  $S$  is denoted by  $v_c^{(S)}$  and is defined by the smallest  $v$  for which  $R_v^{(S)} \leq 1$ , thus avoiding a major outbreak.

**Result 3.** For the network epidemic models in Result 2 (and more or less all epidemic models), the reproduction number  $R_v^{(U)}$  after a uniformly chosen fraction  $v$  have been vaccinated, is related to  $R_0$  by:

$$R_v^{(U)} = R_0(1 - v).$$

The critical vaccination coverage for the uniform vaccination coverage is given by

$$v_c^{(U)} = 1 - 1/R_0.$$

In an unvaccinated community, a typical infective infects on average  $R_0$  individuals early in the epidemic. When a fraction  $v$  are initially vaccinated, infection only takes place when contacting unvaccinated individuals, so  $R_0$  is reduced to  $R_0(1 - v)$  explaining the first statement. The second statement follows immediately from the inequality  $R_v^{(U)} = R_0(1 - v) \leq 1$ , which is equivalent to  $v \geq 1 - 1/R_0$ .

For vaccination strategies that are more effective than the uniform strategy, typically by vaccinating individuals with higher degree, the corresponding reproduction number is smaller than that of the uniform strategy with the same coverage. Consequently, the critical vaccination coverage is smaller for such strategies compared with the uniform strategy. However, other strategies typically do not allow for explicit expressions for the reproduction number and critical vaccination coverage.

The formal proofs of the results above are based on branching process theory together with results showing that the initial stage of the outbreak is well approximated by a branching process when the community size  $n$  is large (e.g., section 12.6 in Diekmann et al., 2013, and references therein). Showing probabilistic results for the final size is much harder and remains open problem for many network epidemics. The two main reasons for the complexity are networks have a complex structure themselves, and to consider a random epidemic process taking place “on” the network clearly adds complexity. In particular, since all epidemic models induce dependent outcomes, that is, the two events that two individuals who are neighbors get infected are positively correlated, the final size will not be a sum of independent outcomes.

## 6 | STATISTICAL INFERENCE

We now move to the important area of making statistical inference for epidemics taking place on networks, an area which deserves more attention. One reason for the complications lies in the underlying probabilistic complexity of the network described in the previous section. However, a more important reason is the fact that often very little of the epidemic, and of the network in particular, is observed. Below, we describe a few different types of data, and give some important questions deserving attention for each of these data settings. The data could either only be the final size of an epidemic, or it may consist of disease incidence over time, or it could be disease incidence together with some information of the underlying network. In Section 6.3, we also briefly describe inference methods for outbreaks where virus sequence data of diagnosed individuals are available together with disease incidence, and methods which make use of the evolution of the virus by comparing virus sequences of diagnosed cases.

### 6.1 | Epidemic outbreak on known network

For diseases like influenza which spreads through airborne aerosols, the most important social structure of relevance for disease spreading is believed to be household structure. Since this is a known structure for which it is easy to collect information, a random model for the underlying network structure is not needed. Consider for example an epidemic outbreak taking place in a community built up of households. Suppose that the data consists of observing how many that got infected and how many that were not in each household at the end of the outbreak, so no temporal information is available. This data can be summarized by  $\{n_{h,i}; 0 \leq i \leq h \leq h_{\max}\}$ , where  $n_{h,i}$  denotes the number of households having  $h$  initially susceptible out of which  $i$  got infected during the outbreak. As mentioned earlier, it is important to only consider initially susceptible individuals—initially immunes are excluded and not counted in  $h$ . In practice, this can be achieved by testing for antibodies prior to the outbreak.

The simplest household epidemic model is equivalent to the Reed-Frost network epidemic model with random global contacts defined earlier, for the special case where the network consists of small, fully connected subgroups— the households.

This model can be approximated by assuming that all individuals are infected from outside the household *independently* with probability  $p_G$ , and to then assume that the outbreaks within each household progress independently. The approximation lies in assuming that the probability to get infected from outside is a fixed parameter  $p_G$  rather than depending on how many that get infected outside the household. For this approximation model, it is possible to express the probability that  $i$  out of  $h$  individuals in a household get infected  $\pi_{h,i} = \pi_{h,i}(p, p_G)$ , as a function of the two model parameters  $p$  and  $p_G$ . This is done by first conditioning on how many individuals in the household that get infected from outside the household (a binomial outcome), and then computing the probability that the rest get infected from household members. These latter probabilities are nontrivial and given by certain recursive formulas not presented here (see Longini & Koopman, 1982 or Addy, Longini, & Haber, 1991), who also consider different types of individual). The conclusion is that the likelihood is simple in terms of these nontrivial probabilities:

$$L(p, p_G) = \prod_{h=1}^{h_{\max}} \prod_{i=0}^h \pi_{h,i}(p, p_G)^{n_{h,i}}.$$

It is worth emphasizing that this is not the exact likelihood for the stochastic household model, but an approximation relying on a large community where households are close to independent, so it can be called a pseudolikelihood. Parameter estimates are obtained by maximizing the likelihood with respect to  $p$  and  $p_G$  (or when the original parametrization with individual global transmission probability  $\beta_G/n$  is used then  $p_G$  should be replaced by the community infection probability  $1 - e^{-\beta_G \tau}$ ,  $\tau$  being the overall fraction infected). For details on this derivation we refer to Longini and Koopman (1982). The pseudolikelihood gives a consistent estimate of the parameters  $p$  and  $\beta_G$ , but their uncertainty estimates are biased from neglecting dependencies—in Ball, Mollison, and Scalia-Tomba (1997), it is described how to correct for this error by using the proper stochastic epidemic model and its central limit theorem.

Suppose now that instead of fully connected subgraphs, such as households, we consider an arbitrary but known network, where we observe the final outbreak as illustrated on the small network in Figure 3. Then inference is much less straightforward, even when all transmissions take place along edges in the network so that random global infectious contacts are not considered. The main reason is that we do not observe who infected whom, or not even in which order individuals were infected. As an illustration, individual  $i$  in Figure 3 was infected during the outbreak and so were two out of three neighbors of  $i$ . How should this affect the likelihood? The answer will depend on *when*  $i$  was infected in relation to the neighbors. Clearly  $i$  was infected by one of the two infected neighbors (excluding the small probability that  $i$  was the index case), but was  $i$  infected by the neighbor who acquired infection first and did  $i$  then infect the other?, or was  $i$  not infected by the first of them? The former would indicate a larger value of  $p$  compared with the latter, but information on which of these chains of events (or other chains) is not contained in the data, thus making inference much harder. One approach is then to collect information also on time of infection (or diagnosis) making the likelihood manageable, or to augment the final size data with this information and use MCMC methods thus integrating over possible temporal outbreak data (see Britton & O'Neill, 2002). Another observation, illustrated in Figure 3, is that there is a higher tendency for individuals with high degree to get infected. There are however exceptions: the individual left of  $i$  has four neighbors, and all of them were infected, but still this individual escaped infection by chance.

We consider it an important challenge to come up with simpler pseudo-likelihood methods for the type of data illustrated in Figure 3. It is not clear if it is better to consider the likelihood for different *nodes*, or for different *edges* in the network. If considering the likelihood contribution from different individuals, such methods might make use of the fact that the likelihood contribution from an individual  $i$  escaping infection should be  $(1 - p)^{T_i}$ , where  $T_i$  denotes the total number of infected neighbors to  $i$ . If instead considering the likelihood in terms of edges, then from edges where both individuals escape infection there is no contribution to the likelihood, and for edges where one individual is infected and the other escapes infection the likelihood contribution is  $1 - p$ , since there was no transmission through that edge. The problem lies in edges where both individuals were infected. For such edges, it is not (always) clear if a transmission took place or if the second individual to get infected escaped infection from that neighbor and only later got infected from another neighbor.

It is also possible to consider other known network structures, such as schools together with households thus having overlapping complete subgraphs. The likelihood of course becomes more complicated and computer intensive statistical methods have to be used (e.g., Cauchemez et al., 2008).

If temporal data are available inference is often simplified, at least when the order of infection can be inferred from the temporal data. Suppose for simplicity that we observe when individuals

get infected and when they recover. We summarize this data by the following quantities for each individual  $i$ :  $(E_i, N_i, I_i)$ , where  $E_i$  is the accumulated infectious exposure time until individual  $i$  was infected, or to the end of the epidemic if  $i$  was not infected.  $E_i$  is hence equal to the sum of the infectious periods of all neighbors of  $i$  up until  $i$  was infected, or to the end of the epidemic if  $i$  did not get infected.  $N_i$  and  $I_i$  are only relevant if  $i$  got infected. In that case  $N_i$  denotes the number of infectious neighbors of  $i$  at the time when  $i$  gets infected, and  $I_i$  denotes the length of the infectious period of  $i$ . Clearly, this information is available when the network is known and we observe infection and recovery times of all infected individuals. The likelihood for this data is given by:

$$L(\beta, \gamma) = \prod_i e^{-\beta E_i} \prod_{i \in \text{Inf}} (\beta N_i) \prod_{i \in \text{Inf}} f(I_i) \propto \left( e^{-\beta \sum_i E_i} \right) \beta^Z \prod_{i \in \text{Inf}} f(I_i),$$

where we have left out a combinatorial factor not affecting estimation on the left-hand side,  $Z$  denotes the final size,  $i \in \text{Inf}$  means those individuals who became infected, and  $f(I_i)$  is the density of the infectious period—this likelihood is valid also for other distributions than the exponential distribution. From this we see that estimates of  $\beta$  and  $\gamma$  are more or less independent, and the inference for the infectious period distribution is just as when observing i.i.d. infectious periods. Furthermore, it is easy to show that the ML-estimate for  $\beta$  is given by

$$\hat{\beta} = Z / \sum_i E_i,$$

the number of infections divided by the sum of the exposure times of all individuals. To obtain standard errors for this estimate remains an open problem for most networks. However, using theory for counting processes (Andersen, Borgan, Gill, & Keiding, 1993) this should be possible to derive using similar techniques as when having temporal data for homogeneous epidemic models (see Diekmann et al., 2013, section 5.4.2).

## 6.2 | Epidemic outbreak on unobserved network (but known network model)

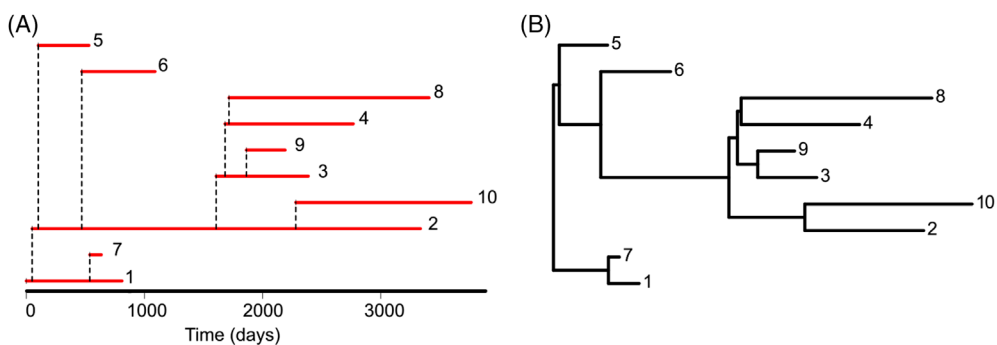
In contrast to Section 6.1, we now consider the more common situation that the underlying social network is not observed. Even when the underlying network is not observed, it is quite often the case that some global network features might be known. For example, it might be known how many neighbors individuals tend to have, in mean or the whole degree distribution. A common way to infer such network properties is then by collecting so-called egocentric network data using Respondent Driven Sampling (RDS). By egocentric network is meant a (typically anonymous) individual who gives information about how many connections (e.g., sex-partners) he or she has, together with other individual information. RDS is a sampling scheme for networks, with the aim to let respondents remain anonymous by letting the respondents select some of their peers anonymously to join the study, and so on (e.g., Gile & Handcock, 2010). As an illustration for STIs in a one-sex population: a study of sexual habits can be used to infer the degree distribution  $\{p_k\}$  for the number of sex-partners an individual has had in the last 12 months, and then the configuration model or the preferential attachment model, calibrated to the sexual habit study, might be used to describe the underlying sexual network even though it is not observed.

If an epidemic outbreak is observed for such a network, the inference methodology will depend on what is observed. If all that is observed is the final size (or equivalently the final *fraction* getting infected), a single data point, then very little can be done unless the underlying network model (including its parameter values) is known. If this is the case, it is in principle possible to estimate one parameter of the transmission model, at least for network epidemic models for which an expression for the mean final fraction infected have been derived. This (limiting) final fraction  $\tau$  will be a function of the network model (assumed known), and the transmission model parameters (e.g.,  $p$  in the Reed-Frost version), and the estimate  $\hat{p}$  is then the value of  $p$  for which  $\tau$  coincides with the observed final fraction infected  $\bar{\tau}$ . Such an estimate of disease spreading parameters relies on that the network model, including its parameter values, is fully known. If for example an incorrect value for the mean degree  $\mu$  is used, then  $\hat{p}$  will be biased.

If instead incidence is observed over time, and/or also network information of infected individuals is collected, then it should be possible to improve the inference procedures, but such methods are still to be developed.

### 6.3 | Epidemic outbreak on unknown network model using virus sequences

A related situation as considered in the previous section is given for epidemic data where pathogens of diagnosed individuals are sequenced, meaning that the DNA of the virus (or other disease agent) is sequenced. The underlying idea is that, for viruses that evolve (within infected individuals) at a similar time scale as the epidemic, individuals close to each other in the transmission tree, the tree describing the transmission events (cf. Figure 4), will have virus sequences that are more similar as compared to individuals that are further away from each other in the transmission tree. Hence, by observing the sequences of diagnosed individuals, it is possible to learn more about the underlying transmission tree and hence about the underlying social network upon which the disease spreads.



**FIGURE 4** Illustration of an transmission tree (a) and the corresponding virus genealogy (b). In the transmission tree, Individual 1 is the index case and this individual directly infects Individual 2, and later Individual 7 (infections are indicated by dashed upward lines and infectious periods of individuals are marked with red horizontal lines). In the virus genealogy, the information about who infected whom is lost, as are the starting times of individuals' infectious periods. The distance unit for the transmission tree is calendar time whereas it is evolutionary distance in the virus genealogy



This new research area of making inference using sequences, with or without traditional epidemic data, has exploded in the last 10–15 years. To describe it in detail would require many pages of modeling and statistical methodology—here we simply sketch it briefly. For a more thorough description we refer to other articles, for example, the survey chapter by Klinkenberg et al. in Held, Hens, O'Neill, and Wallinga (2019).

There are numerous models for how a virus population within a host evolves over time. A simple model assumes that an individual is infected by one single virus strain, and that this strain evolves over time within the individual, but that all copies of the virus are identical. The motivation is that mutations that are inferior quickly die out, whereas mutations that are superior quickly take over the entire virus population—this is called the “one dominant strain” assumption as opposed to allowing for within-host diversity. When applying the one-dominant-strain assumption, and a suitable evolutionary model for DNA mutations (see Felsenstein (2003)), it is possible to infer the virus genealogy of the virus sequences taken from the infected individuals. This virus genealogy, where distance is measured in evolutionary distance (e.g., number of mutations per 1,000 base pairs), is in turn related to the corresponding transmission tree, where distance is measured in calendar time. In case of one-dominant-strain and assuming a constant molecular clock (i.e., constant rate of mutation over time and across individuals), the two trees are identical, except that the virus genealogy does not contain information on who infected whom, which the transmission tree does (cf. Figure 4 for an illustration).

Hence by sequencing viruses isolated from infected individuals, it is possible to learn about the virus genealogy and hence also about the transmission tree. By comparing the inferred virus genealogy with typical virus genealogies from various network and transmission models (and their parameters), it is possible to learn what the underlying network structure might have been. These type of ideas are often used for outbreaks of HIV in order to learn more about spreading patterns (e.g., Giardina, Romero-Severson, Albert, Britton, & Leitner, 2017; Leventhal et al., 2012).

The statistical methodology is often quite involved and computationally intensive, for example, employing MCMC, Approximate Bayesian Computation (ABC), and Iterated filtering methods. As mentioned earlier, the main idea in MCMC for infectious disease data is to treat unknown features, such as the underlying network, as latent variables which are then integrated over in the MCMC chain. The main idea with ABC and iterated filtering methods is to simulate/generate output for different choices of models and parameters and to run additional simulations for models/parameters “close” to those of earlier simulations which resembled the observed data (importance sampling). The comparison of how well a simulation agrees with data can be performed using tree shape measures such as the Sackin index (cf. Leventhal et al., 2012).

A disadvantage with most work in this new area of infectious disease inference is that traditional epidemic data, in particular of exposed individuals avoiding infection, is rarely used. In fact, Li, Grassly, and Fraser (2017) even report that precision in inference worsens when the statistical analysis also makes use of incidence data beside the virus sequence data. It is my strong opinion that this should not be the case if a correctly performed statistical analysis is used on a suitable statistical model for community structure and disease transmission. It is an important area to develop statistical methodology for network epidemic data using both virus sequences, incidence and other epidemic and network data.

## 6.4 | Predicting effects of preventive measures

One of the main reasons for mathematical modeling and statistical analysis of (network) epidemics is to analyze the impact of prevention: next time there is a similar epidemic outbreak, or even during an ongoing outbreak, it is of interest to predict what would happen if various control measures are put in place. These control measures could be vaccination, increased condom use, isolation of infected cases, school closures, or traveling restrictions. The common way to proceed is to estimate parameters for a suitable model of the network and disease spreading, and then to mathematically analyze what would be the outcome if some control measure was put in place for the studied model and parameter values. As a very simple illustration, suppose the basic reproduction number has been estimated at  $\hat{R}_0 = 1.5$  with standard error  $\text{s.e.}(\hat{R}_0) = 0.1$ . If a vaccine giving 100% immunity is available, and a fraction  $v = 0.2$  of randomly selected individuals were vaccinated, an estimate of the new reproduction number would be  $\hat{R}_v^{(U)} = \hat{R}_0(1 - v) = 1.5 * 0.8 = 1.2$  (cf. Result 3). Similarly, the critical vaccination coverage is estimated at  $\hat{v}_c^{(U)} = 1 - 1/\hat{R}_0 = 0.33$ , meaning that the predicted fraction necessary to vaccinate in order to avoid a future outbreak is 33%. Using also the standard error, it is possible to construct an upper confidence bound on  $v_c^{(U)}$ .

In more complicated situations, like during an ongoing epidemic where vaccination is introduced, or when modeling effects of school closure in structured communities, the corresponding effects are most often studied by means of simulations (e.g., Cauchemez et al., 2008; Longini et al., 2005). The basic idea is however the same: to first use data to infer model parameters and then to study effects of intervention using the model and estimated parameters. The effectiveness of intervention, for example, how much susceptibility is reduced by the vaccine, or the effect of closing schools on transmission between school children, has to be known or estimated using some other data source.

A somewhat different type of prevention is contact tracing: when new individuals are diagnosed they are questioned about earlier contacts (i.e., neighbors in the network) who are then tested. The aim with contact tracing is to diagnose more infectious individuals thus preventing them from further spreading. Contact tracing is mainly used when the underlying contact is well defined, such as sexual intercourse when considering STIs and having had close physical contact with bleeding individuals for Ebola. Modeling effects of contact tracing in network epidemics remains largely an unexplored research area.

Statistical analysis of various effects of prevention is of course highly important and deserves more attention, both in terms of further mathematical analysis of effects of intervention in network models, but also to improve inference procedures for data from network epidemics.

## 7 | DISCUSSION AND EXTENSIONS

The ambition of this article was to describe some basic ideas behind stochastic modeling and statistical analysis of infectious disease outbreaks taking place in communities structured as social networks. The focus was on methodology rather than hands-on data analysis. Statistical modeling of network epidemics uses many different probabilistic and stochastic techniques, thus making it impossible to cover the area in any detail—instead we frequently referred to useful publications in different subareas for more details (see also the monograph by Kiss et al., 2017).

An inherent problem with statistical inference of network epidemics, beside its mathematical complexity stemming from randomness in both the underlying network and the disease spreading, is that very rarely the network structure is even partly known or observed. The area certainly

requires more research, in particular to make better use of virus sequences and other proxies (e.g., from egocentric networks) giving information about the underlying network, and to combine such information with incidence and exposure data for an improved combined statistical analysis.

Needless to say, we have left out, or only briefly mentioned, several issues, which are important to consider both in modeling and statistical analysis, for example, underreporting, asymptomatic cases, and partial immunity. Nearly all datasets from epidemic outbreaks miss some infected cases, perhaps because these were asymptomatic, but also for other reasons. Neglecting this will typically lead to under estimation of transmission parameters and necessary preventive measures. Similarly, neglecting partial immunity when an outbreak took place will make conclusions invalid when community immunity wanes.

We now mention some topics we have not touched upon: models and analysis of diseases that are endemic in the population, models and analysis for situations where individuals change behavior over time—perhaps as a result of the epidemic outbreak, statistical analysis of epidemics on dynamic networks, models for which the infectivity of an individual varies over time (e.g., the acute and chronic phase of HIV), and the important area of model fit.

As has been mentioned in several places, the statistical methods for models capturing both transmission and (often unobserved) network structure, are often too complicated for direct methods such as maximum likelihood estimation. Instead some numerically intensive method like MCMC, ABC, or Particle filtering can be adopted. Held et al. (2019) published a recent book describing such (and other) statistical methods for infectious disease data, but without focus on networks.

There are several network properties which have not received much attention in the current article. For example, clustering (quantifying the extent of short loops in the network) is only briefly discussed and illustrated with household epidemic models. Nearly all social networks exhibit clustering, and to study effects of clustering when estimating parameters and modeling prevention deserves more attention. Other network properties, such as centrality and betweenness, have been used when analyzing how social networks evolve in time either when left alone or when intervened in different ways (e.g., Valente, 2012). These ideas might be used for network epidemic prevention if aiming at changing risky behavior when knowledge of which individuals to specifically target is important. A disadvantage is however that if little is known about the network, then such targeted preventive measures are hard to implement.

Even though statistical analysis of network epidemics is quite hard, it is encouraging to see that preventive measures make use of ideas from network epidemic modeling. For example, for STI's vaccines and/or condoms are often distributed to sexually active individuals and their *partners* (mimicking the acquaintance vaccination strategy in Section 4), and Henao-Restrepo et al., 2017 studies effects of ring vaccination in Ebola outbreaks which clearly make use of social structures by targeting relatives and friends of all reported cases.

It is my strong belief that much progress can be made in this important research area by development of new statistical methodology in close collaboration with data collectors, thereby also giving suggestions to more informative data collection in future studies.

## ACKNOWLEDGEMENT

The author is grateful to the Swedish Research Council (grant 2015-05015) for financial support.

## ORCID

Tom Britton  <https://orcid.org/0000-0002-9228-7357>

## REFERENCES

- Addy, C. L., Longini, I. M., & Haber, M. (1991). A generalized stochastic model for the analysis of infectious disease final size data. *Biometrics*, 47, 961–974.
- Andersen, P. K., Borgan, Ø., Gill, R. D., & Keiding, N. (1993). *Statistical models based on counting processes*. New York, NY: Springer.
- Ball, F., Britton, T., & Sirl, D. (2013). A network with tunable clustering, degree correlation and degree distribution, and an epidemic thereon. *Journal of Mathematical Biology*, 66, 979–1019.
- Ball, F. G., Mollison, D., & Scalia-Tomba, G. (1997). Epidemics with two levels of mixing. *The Annals of Applied Probability*, 7, 46–89.
- Barabási, A.-L., & Albert, R. (1999). Emergence of scaling in random networks. *Science*, 286, 509–512.
- Barrat, A., Barthélemy, M., Pastor-Satorras, R., & Vespignani, A. (2004). The architecture of complex weighted networks. *PNAS*, 101, 3747–3752.
- Bollobás, B. (2001). *Random graphs* (2nd ed.). Cambridge, England: Cambridge University Press.
- Bollobás, B., Janson, S., & Riordan, O. (2007). The phase transition in inhomogeneous random graphs. *Random Structures & Algorithms*, 31, 3–122.
- Britton, T., Lindholm, T., & Turova, T. (2011). A dynamic network in a dynamic population: Asymptotic properties. *Journal of Applied Probability*, 48, 1163–1178.
- Britton, T., & O'Neill, P. D. (2002). Bayesian inference for stochastic epidemics in populations with random social structure. *Scandinavian Journal of Statistics*, 29, 375–390.
- Cauchemez, S., Valleron, A.-J., Boëlle, P.-Y., Flahault, A., & Ferguson, N. M. (2008). Estimating the impact of school closure on influenza transmission from sentinel data. *Nature*, 452, 750–754.
- Cohen, R., Havlin, S., & Ben-Avraham, D. (2003). Efficient immunization strategies for computer networks and populations. *Physical Review Letters*, 91, 247901.
- Diekmann, O., Heesterbeek, J. A. P., & Britton, T. (2013). *Mathematical tools for understanding infectious disease dynamics*. Princeton, New Jersey: Princeton University Press.
- Erdős, P., & Rényi, A. (1959). On random graphs. *Publicationes Mathematicae*, 6, 290–297.
- Felsenstein, J. (2003). *Inferring phylogenies* (2nd ed.). Sunderland, Massachusetts: Sinauer Associates Inc.
- Giardina, F., Romero-Severson, E. O., Albert, J., Britton, T., & Leitner, T. K. (2017). Inference of transmission network structure from HIV phylogenetic trees. *PLoS Computational Biology*, 13, e1005316.
- Gile, K. J., & Handcock, M. S. (2010). Respondent-driven sampling: An assessment of current methodology. *Sociological Methodology*, 40, 285–327.
- Halloran, M. E., Longini, I. M., & Struchiner, C. J. (2010). *Design and analysis of vaccine studies*. New York, NY: Springer.
- Held, L., Hens, N., O'Neill, P. D., & Wallinga, J. (Eds.). (2019). *Handbook of infectious disease data analysis*. New York, NY: CRC Press.
- Henao-Restrepo, A. M., Camacho, A., Longini, I. M., Watson, C. H., Edmunds, W. J., Egger, M., ... Kieny, M.-P. (2017). Efficacy and effectiveness of an RSV-vectored vaccine in preventing Ebola virus disease: Final results from the Guinea ring vaccination, open-label, cluster-randomised trial (Ebola Ça Suffit!). *Lancet*, 389, 505–518.
- Kiss, I. S., Miller, J. C., & Simon, P. L. (2017). *Mathematics of epidemics on networks Interdisciplinary applied mathematics* (Vol. 46). Cham: Springer.
- Leung, K., Ball, F., Sirl, D., & Britton, T. (2018). Individual preventive social distancing during an epidemic may have negative population-level outcomes. *Journal of The Royal Society Interface*, 15, 20180296.
- Leventhal, G. E., Kouyos, R., Stadler, T., Von Wyl, V., Yerly, S., Böni, J., ... Bonhoeffer, S. (2012). Inferring epidemic contact structure from phylogenetic trees. *PLoS Computational Biology*, 8, e1002413.
- Li, L. M., Grassly, N. C., & Fraser, C. (2017). Quantifying transmission heterogeneity using both pathogen phylogenies and incidence time series. *Molecular Biology and Evolution*, 34, 2982–2995.
- Longini, I. M., & Koopman, J. S. (1982). Household and community transmission parameters from final distributions of infections in households. *Biometrics*, 38, 115–126.
- Longini, I. M., Nizam, A., Xu, S., Ungchusak, K., Hanshaworakul, W., Cummings, D. A. T., & Halloran, M. E. (2005). Containing pandemic influenza at the source. *Science*, 309, 1083–1087.
- Molloy, M., & Reed, B. (1998). The size of the giant component of a random graphs with a given degree sequence. *Combinatorics, Probability and Computing*, 7, 295–305.

- Newman, M. E. J., Watts, D. J., & Strogatz, S. H. (2002). Random graph models of social networks. *PNAS*, 99, 2566–2572.
- Norros, I., & Reittu, H. (2006). On a conditionally Poissonian graph process. *Advances in Applied Probability*, 38, 59–75.
- Nowicki, K., & Snijders, T. A. B. (2001). Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96, 1077–1087.
- Poletto, C., Gomes, M. F. C., Piontti, A. P., Rossi, L., Bioglio, L., Chao, D. L., ... Vespignani, A. (2014). Assessing the impact of travel restrictions on international spread of the 2014 west African Ebola epidemic. *Eurosurveillance*, 19, 20936.
- Snijders, T. A. B., Pattison, P. E., Robins, G. L., & Handcock, M. S. (2006). New specifications for exponential random graph models. *Sociological Methodology*, 36, 99–153.
- Spricer, K., & Britton, T. (2015). The configuration model for partially directed graphs. *Journal of Statistical Physics*, 161, 965–985.
- Valente, T. W. (2012). Network interventions. *Science*, 337, 49–53.
- Watts, S. C., & Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393, 440–442.

**How to cite this article:** Britton T. Epidemic models on social networks—With inference. *Statistica Neerlandica*. 2020;74:222–241. <https://doi.org/10.1111/stan.12203>