Effects of Contact Network Models on Stochastic Epidemic Simulations

Rehan Ahmad and Kevin S. Xu

EECS Department, University of Toledo, Toledo, OH 43606, USA Rehan.Ahmad@rockets.utoledo.edu, Kevin.Xu@utoledo.edu

Abstract. The importance of modeling the spread of epidemics through a population has led to the development of mathematical models for infectious disease propagation. A number of empirical studies have collected and analyzed data on contacts between individuals using a variety of sensors. Typically one uses such data to fit a probabilistic model of network contacts over which a disease may propagate. In this paper, we investigate the effects of different contact network models with varying levels of complexity on the outcomes of simulated epidemics using a stochastic Susceptible-Infectious-Recovered (SIR) model. We evaluate these network models on six datasets of contacts between people in a variety of settings. Our results demonstrate that the choice of network model can have a significant effect on how closely the outcomes of an epidemic simulation on a simulated network match the outcomes on the actual network constructed from the sensor data. In particular, preserving degrees of nodes appears to be much more important than preserving cluster structure for accurate epidemic simulations.

Keywords: network model, stochastic epidemic model, contact network, degree-corrected stochastic block model

1 Introduction

The study of transmission dynamics of infectious diseases often involves simulations using stochastic epidemic models. In a compartmental stochastic epidemic model, transitions between compartments occur randomly with specified probabilities. For example, in a stochastic Susceptible-Infectious-Recovered (SIR) model [4,10], a person may transition from S to I with a certain probability upon contact with an infectious person, or a person may transition from I to R with a certain probability to simulate recovering from the disease.

The reason for the spread of infection is contact with the infectious individual. Hence, the contact network in a population is a major factor in the transmission dynamics. Collecting an actual contact network over a large population is difficult because of limitations in capturing all the contact information. This makes it necessary to represent the network with some level of abstraction, e.g. using a statistical model. A variety of statistical models for networks have been proposed [9]; such models can be used to simulate contact networks that resemble actual contact networks.



Fig. 1: For each of the susceptible (S), infectious (I), and recovered (R) compartments, the mean curve for simulations on the model (shown in blue) is compared to the mean curve for simulations on the actual network (shown in red). The closeness between the model and actual network is given by the sum of the shaded areas between the curves for each compartment (smaller is better).

Our aim in this paper is to evaluate different models for contact networks in order to find the best model to use to simulate contact networks that are close to an actual observed network. We do this by comparing the disease dynamics of a stochastic SIR model over the simulated networks with the disease dynamics over the actual network. One commonly used approach is to compare the epidemic size at the end of the simulation, i.e. what fraction of the population caught the disease [19,25]. A drawback of this approach is that it only considers the steady-state outcome and not the dynamics of the disease as it is spreading.

We propose to compare the dynamics at each time instant in the simulation by calculating the area between the mean SIR curves for the epidemic over the simulated and actual networks, shown in Fig. 1. A small area indicates that the dynamics of the epidemic over the simulated contact networks are close to those of the actual network. We use this approach to compare four contact network models (in increasing order of number of parameters): the Erdős-Rényi model, the degree model, the stochastic block model, and the degree-corrected stochastic block model. Our experiment results over six different real network datasets suggest that the degree-corrected stochastic block model provides the closest approximation to the dynamics of an epidemic on the actual contact networks. Additionally, we find that preserving node degrees appears to be more important than preserving community structure for accuracy of epidemic simulations.

2 Related Work

A significant amount of previous work deals with the duration [23], frequency [17], and type [6,24] of contacts in a contact network. These findings are often incorporated into simulations of epidemics over different types of contact models. The R package EpiModel [13] allows for simulation of a variety of epidemics over

	J					
	HVCCUPS	Friends $\&$	High	Infoctious	Primary	HOPE
	11100015	Family	School	mechous	School	
Number of nodes	43	123	126	201	242	1178
Sensor type	Wi-Fi	Blue to oth	RFID	RFID	RFID	RFID
Proximity range	N/A	5 m	11.5 m	11.5 m	11.5 m	Room
Graph density	0.326	0.228	0.217	0.0328	0.285	0.569
Clustering coefficient	0.604	0.496	0.522	0.459	0.480	0.748
Average degree	14.0	27.8	27.1	6.56	68.7	671
Maximum degree	28	73	55	21	134	1072

Table 1: Summary statistics from datasets used in this study.

temporal exponential random graph models for contact networks and has been used in studies of various different infectious diseases including HIV [14].

There has also been prior work simulating the spread of disease over a variety of contact network models with the goal of finding a good approximation to the actual high resolution data in terms of the epidemic size, i.e. the final number of people infected [19,25]. Such work differs from our proposed area metric, which considers the dynamics as the disease is spreading and not just the steadystate outcome. In [3], the authors use the squared differences between the I curves (fraction of infectious individuals) of an epidemic model on simulated contact networks and on an actual contact network to calibrate parameters of the epidemic model when used on simulated contact networks. Although this metric does consider the dynamics of the epidemic, our proposed metric also involves the S and R curves for a more complete evaluation of population dynamics.

3 Datasets

We consider a variety of contact network datasets in this paper. Table 1 shows summary statistics for each dataset along with the sensor type. The HYCCUPS dataset was collected at the University Politehnica of Bucharest in 2012 using a background application for Android smartphones that captures a device's encounters with Wi-Fi access points [20]. The Friends & Family (F&F) dataset was collected from the members of a residential community nearby a major research university using Android phones loaded with an app that records many features including proximity to other Bluetooth devices [2]. The High School (HS) dataset was collected among students from 3 classes in a high school in Marseilles, France [7] using wearable sensors that capture face-to-face proximity for more than 20 seconds. The Infectious dataset was collected at a science gallery in Dublin using wearable electronic badges to sense sustained face-to-face proximity between visitors. [12]. We use data for one arbitrarily selected day (April 30) on which 201 people came to visit. The Primary School (PS) dataset was collected over 232 students and 10 teachers at a primary school in Lyon, France in a similar manner to the HS dataset [8]. Lastly, the HOPE dataset is collected from the Attendee Meta-Data project at the seventh Hackers on Planet Earth (HOPE)

conference [1]. We create a contact network where the attendees at each talk form a clique; that is, each person is assumed to be in contact with every other person in the same room, hence why this network is much denser.

4 Methods

We construct actual networks from the datasets by connecting the individuals (nodes) with an edge if they have a contact at any point of time. We evaluate the quality of a contact network model for simulations of epidemics by conducting the following steps for each dataset:

- 1. Simulate 5,000 epidemics over the actual network.
- 2. Fit contact network model to actual network.
- 3. Simulate 100 networks from contact network model. For each simulated network, simulate 50 epidemics over the network for 5,000 epidemics total.
- 4. Compare the results of the epidemic simulations over the actual network with those over the simulated networks.

These steps are repeated for each contact network model that we consider. We describe the stochastic epidemic model we use to simulate epidemics in Section 4.1 and the contact network models we use in Section 4.2. To get a fair evaluation of the dynamics of epidemics spreading over different contact network models, all of the parameters which are not related to the contact network model, e.g. probability of infection and probability of recovery are kept constant. Our aim is to single out the effect of using a particular contact network model while simulating an epidemic.

4.1 Stochastic Epidemic Model

An actual infection spread in a population experiences randomness in several factors which may aggravate or inhibit the spread. This is considered in stochastic epidemic models. The initial condition is, in general, to have a set of infectious individuals, while the rest of the population is considered susceptible. We consider a discrete-time process, where at each time step, the infectious individuals can spread the disease with some probability of infection to susceptible individuals they have been in contact with. Also, the infectious individuals can recover from the disease with some probability independent of the individuals' contacts with others. This model is known as the stochastic SIR model and is one of the standard models used in epidemiology [4,10].

We randomly choose 1 infectious individual from the population as the initial condition and simulate the epidemic over 30 time steps. We set the probability of infection for every interaction between people to be 0.025. The probability of recovery is also set to be 0.025. Note that the rate at which the disease spreads across the population is dependent not only on the infection probability but also the topology of the contact network; thus, by fixing these probabilities, we are exploring only the effects of the contact network.

4.2 Contact Network Models

In practice, it is extremely difficult to obtain accurate contact network data. An alternative is to simulate a contact network by using a statistical network model. We consider several such models, which we briefly describe in the following. We refer interested readers to the survey by Goldenberg et al. [9] for details.

Erdős-Rényi (E-R) Model In the E-R model, an edge between any two nodes is formed with probability p independent of all other edges. To fit the E-R model to a network, set the single parameter, the estimated edge probability $\hat{p} = M/{\binom{N}{2}}$, where N and M denote the number of nodes and edges in the actual network, respectively. By doing so, the expected number of edges in the E-R model will be $\binom{N}{2}\hat{p} = M$, the number of edges in the actual network.

Degree Model In several network models, including the configuration model and preferential attachment models, the edge probability depends upon the degrees of the nodes it connects [21]. We consider a model that preserves the expected rather than actual degree of each node, often referred to as the Chung-Lu model [5]. In this model, the probability of an edge between two nodes is proportional to the product of their node degrees, and all edges are formed independently. The model has N parameters, the expected degrees of each node.

To fit the degree model to a network, we compute the degrees of all nodes to obtain the degree vector **d**. We then set the estimated edge probabilities $\hat{p}_{ij} = \alpha d_i d_j$, where the constant α is chosen so that the sum of all edge probabilities (number of expected edges) is equal to the number of edges in the actual network.

Stochastic Block Model (SBM) In the SBM [11], the network is divided into disjoint sets of individuals forming K communities. The probability of edge formation between two nodes depends only upon the communities to which they belong. This model takes as input a vector of community assignments **c** (length N) and a matrix of edge formation probabilities Φ (size $K \times K$), where ϕ_{ab} denotes the probability that a node in community a forms an edge with a node in community b, independent of all other edges. For an undirected graph, Φ is symmetric so the SBM has $N + {K+1 \choose 2}$ parameters in total.

To estimate community assignments, we use a regularized spectral clustering algorithm [22] that is asymptotically consistent and has been demonstrated to be very accurate in practice. We select the number of communities using the eigengap heuristic [18]. Once the community assignments $\hat{\mathbf{c}}$ are estimated, the edge probabilities can be estimated by $\hat{\phi}_{ab} = m_{ab}/n_{ab}$, where m_{ab} denotes the number of edges in the block formed by the communities a, b in the observed network, and n_{ab} denotes the number of possible edges in the block [16].

Degree-corrected Stochastic Block Model (DC-SBM) The DC-SBM is an extension to the SBM in a way that incorporates the concepts of the degree model within an SBM [16]. The parameters of the DC-SBM are the vector of community assignments **c** (length N), a node-level parameter vector $\boldsymbol{\theta}$ (length N), and a block-level parameter matrix Ω (size $K \times K$). In a DC-SBM, an edge between a node $i \in a$ (meaning node i is in community a) and node $j \in b$ is formed with probability $\theta_i \theta_j \omega_{ab}$ independent of all other edges. Ω is symmetric, so the DC-SBM has $2N + {K+1 \choose 2}$ parameters in total.

To fit the DC-SBM to an actual network, we first estimate the community assignments in the same manner as in the SBM using regularized spectral clustering. We then estimate the remaining parameters to be $\hat{\theta}_i = d_i / \sum_{j \in a} d_j$, for node $i \in a$, and $\hat{\omega}_{ab} = m_{ab}$ [16]. Using these estimates, we arrive at the estimated edge probabilities $\hat{p}_{ij} = \hat{\theta}_i \hat{\theta}_j \hat{\omega}_{ab}$.

5 Results

To evaluate the quality of a contact network model, we compare the mean SIR curves resulting from epidemic simulations on networks generated from that model to the mean SIR curves from epidemic simulations on the actual network. If the two curves are close, then the network model is providing an accurate representation of what is likely to happen on the actual network.

To measure the closeness of the two sets of mean SIR curves, we use the sum of the areas between each set of curves as shown in Fig. 1. By measuring the area between the curves rather than just the final outcome of the epidemic simulation (e.g. the fraction of recovered people after the disease dies out as in [19,25]), we capture the difference in transient dynamics (e.g. the rate at which the infection spreads) rather than just the difference in final outcomes.

The area between the SIR curves for each model over each dataset is shown in Fig. 2a. According to this quality measure, the DC-SBM is the most accurate model on F&F, HS, and PS; the degree model is the most accurate on HYCCUPS and HOPE; and the SBM is most accurate on Infectious. However, the SBM appears to be only slightly more accurate than the E-R model overall, despite having $N + \binom{K+1}{2}$ parameters compared to the single parameter E-R model. The contact network models were most accurate on the HOPE network, which is the densest, causing the epidemics to spread rapidly.

We compute also the log-likelihood for each contact network model on each dataset, shown in Fig. 2b. To normalize across the different sized networks, we compute the log-likelihood per node pair. Since all of the log-likelihoods are less than 0, we show the negative log-likelihood (i.e. lower is better) in Fig. 2b. Unsurprisingly, the DC-SBM, with the most parameters, also has the highest log-likelihood, whereas the relative ordering of the log-likelihoods of the degree model and SBM, both with roughly the same number of parameters, vary depending on the dataset.

Both the proposed area between SIR curves and the log-likelihood can be viewed as quality measures for a contact network model. A third quality measure is given by the number of parameters, which denotes the simplicity of the model. A simpler model is generally more desirable to avoid overfitting. These three



Fig. 2: Comparison of (a) area between SIR curves of each model with respect to actual network for each dataset and (b) negative log-likelihood per node pair for each model (lower is better for both measures). The DC-SBM model appears to be the best model according to both quality measures, but the two measures disagree on the quality of the degree model compared to the SBM.

Table 2: Quality measures (lower is better) averaged over all datasets for each model. Best model according to each measure is shown in bold.

Quality Measure	E-R	Degree	SBM	DC-SBM
Area between SIR curves	1.82	0.73	1.43	0.71
Negative log-likelihood per node pair	0.597	0.496	0.504	0.385
Number of parameters	1	319	328	647

quality measures for each model (averaged over all datasets) are shown in Table 2. The DC-SBM achieves the highest quality according to the area between SIR curves and the log-likelihood at the expense of having the most parameters. On the other hand, the E-R model has only a single parameter but is the worst in the other two quality metrics. Interestingly, the degree model and SBM appear to be roughly equal in terms of the number of parameters and log-likelihood, but the area between SIR curves for the two models differs significantly. This suggests that the degree model may be better than the SBM at reproducing features of contact networks that are relevant to disease propagation.

6 Discussion

The purpose of our study was to evaluate the effects of contact network models on the results of simulated epidemics over the contact network. While it is wellknown and expected that more complex models for contact network topology do a better job of reproducing features of the contact network such as degree distribution and community structure, we demonstrated that, in general, they also result in more accurate epidemic simulations. That is, the results of simulating an epidemic on a more complex network model are usually closer to the results obtained when simulating the epidemic on the actual network than if we had used a simpler network model. Moreover, models that preserve node degrees are shown to produce the most accurate epidemic simulations. Unlike most prior studies such as [19,25], we measure the quality of a network model by its area between SIR curves compared to the SIR curve of the actual network, which allows us to capture differences while the disease is still spreading rather than just the difference in the final outcome, i.e. how many people were infected.

Our findings suggest that the degree-corrected stochastic block model (DC-SBM) is the best choice of contact network model in epidemic simulations because it resulted in the minimum average area between SIR curves. Interestingly, using the degree model resulted in an average area between SIR curves to be only slightly larger than the DC-SBM despite having less than half as many parameters, as shown in Table 2. The SBM (without degree correction) also has half as many parameters as the DC-SBM, but has over twice the area between SIR curves. We note that the difference between the degree model and the SBM *cannot* be observed using log-likelihood as the quality measure, as both models are very close in log-likelihood. This leads us to believe that preserving degree has a greater effect on accuracy of epidemic simulations than preserving community structure. Furthermore, this finding demonstrates that one cannot simply evaluate the accuracy of a contact network model for epidemic simulations only by examining goodness-of-fit on the actual contact network!

In practice, one cannot often collect high-resolution contact data on a large scale, so having accurate contact network models is crucial to provide realistic network topologies on which we can simulate epidemics. In this paper, we estimated the parameters for each contact network model using the contact network itself, which we cannot do in practice because the contact network is often unknown. As a result, one would have to estimate the model parameters from prior knowledge or partial observation of the contact network, which introduces additional error that was not studied in this paper. It would be of great interest to perform this type of sensitivity analysis to identify whether the DC-SBM and degree model are still superior even when presented with less accurate parameter estimates. Also, there is a risk of overfitting in more complex models which should be examined in a future extension of this work. Both issues could potentially be addressed by considering hierarchical Bayesian variants of network models such as the degree-generated block model [27], which add an additional generative layer to the model with a smaller set of hyperparameters.

Another limitation of this study is our consideration of static unweighted networks. Prior work [15,19,23,25] has shown that it is important to consider the time duration of contacts between people, which can be reflected as weights in the contact network, as well as the times themselves, which can be accommodated by using models of dynamic rather than static networks, such as dynamic SBMs [26]. We plan to expand this work in the future by incorporating models of weighted and dynamic networks to provide a more thorough investigation.

References

- aestetix, Petro, C.: CRAWDAD dataset hope/amd (v. 2008-08-07). Downloaded from http://crawdad.org/hope/amd/20080807 (2008)
- Aharony, N., Pan, W., Ip, C., Khayal, I., Pentland, A.: Social fMRI: Investigating and shaping social mechanisms in the real world. Pervasive and Mobile Computing 7(6), 643–659 (2011)
- Bioglio, L., Génois, M., Vestergaard, C.L., Poletto, C., Barrat, A., Colizza, V.: Recalibrating disease parameters for increasing realism in modeling epidemics in closed settings. BMC Infectious Diseases 16(1), 676 (2016)
- 4. Britton, T.: Stochastic epidemic models: a survey. Mathematical Biosciences 225(1), 24–35 (2010)
- Chung, F., Lu, L.: The average distances in random graphs with given expected degrees. Proceedings of the National Academy of Sciences 99(25), 15879–15882 (2002)
- Eames, K.: Modeling disease spread through random and regular contacts in clustered populations. Theoretical Population Biology 73(1), 104–111 (2008)
- Fournet, J., Barrat, A.: Contact patterns among high school students. PLoS ONE 9(9), e107878 (2014)
- Gemmetto, V., Barrat, A., Cattuto, C.: Mitigation of infectious disease at school: targeted class closure vs school closure. BMC Infectious Diseases 14(1), 695 (2014)
- Goldenberg, A., Zheng, A.X., Fienberg, S.E., Airoldi, E.M.: A survey of statistical network models. Foundations and Trends in Machine Learning 2(2), 129–233 (2010)
- Greenwood, P., Gordillo, L.: Stochastic epidemic modeling. In: Chowell, G., Hyman, J.M., Bettencourt, L.M.A., Castillo-Chavez, C. (eds.) Mathematical and statistical estimation approaches in epidemiology, pp. 31–52. Springer, Dordrecht (2009)
- Holland, P.W., Laskey, K.B., Leinhardt, S.: Stochastic blockmodels: First steps. Social Networks 5(2), 109–137 (1983)
- Isella, L., Stehl, J., Barrat, A., Cattuto, C., Pinton, J., Van den Broeck, W.: What's in a crowd? Analysis of face-to-face behavioral networks. Journal of Theoretical Biology 271(1), 166–180 (2011)
- Jenness, S., Goodreau, S.M., Morris, M.: EpiModel: Mathematical modeling of infectious disease (2017), http://epimodel.org/
- Jenness, S.M., Goodreau, S.M., Rosenberg, E., Beylerian, E.N., Hoover, K.W., Smith, D.K., Sullivan, P.: Impact of the Centers for Disease Control's HIV preexposure prophylaxis guidelines for men who have sex with men in the United States. The Journal of Infectious Diseases 214(12), 1800–1807 (2016)
- 15. Karimi, F., Holme, P.: Threshold model of cascades in empirical temporal networks. Physica A: Statistical Mechanics and its Applications 392(16), 3476–3483 (2013)
- Karrer, B., Newman, M.E.J.: Stochastic blockmodels and community structure in networks. Physical Review E 83, 016107 (2011)
- 17. Larson, R.C.: Simple models of influenza progression within a heterogeneous population. European Journal of Operational Research 55(3), 399–412 (2007)
- von Luxburg, U.: A tutorial on spectral clustering. Statistics and Computing 17(4), 395–416 (2007)
- Machens, A., Gesualdo, F., Rizzo, C., Tozzi, A.E., Barrat, A., Cattuto, C.: An infectious disease model on empirical networks of human contact: bridging the gap between dynamic network data and contact matrices. BMC Infectious Diseases 13(1), 185 (2013)

- Marin, R.C., Dobre, C., Xhafa, F.: Exploring predictability in mobile interaction. In: Proceedings of the 3rd International Conference on Emerging Intelligent Data and Web Technologies. pp. 133–139 (2012)
- Newman, M.: Networks: An Introduction. Oxford University Press, Inc., New York, NY, USA (2010)
- Qin, T., Rohe, K.: Regularized spectral clustering under the degree-corrected stochastic blockmodel. In: Advances in Neural Information Processing Systems 26. pp. 3120–3128 (2013)
- Smieszek, T.: A mechanistic model of infection: why duration and intensity of contacts should be included in models of disease spread. Theoretical Biology and Medical Modelling 6(1), 25 (2009)
- Smieszek, T., Fiebig, L., Scholz, R.W.: Models of epidemics: when contact repetition and clustering should be included. Theoretical Biology and Medical Modelling 6(1), 11 (2009)
- 25. Stehlé, J., Voirin, N., Barrat, A., Cattuto, C., Colizza, V., Isella, L., Régis, C., Pinton, J.F., Khanafer, N., Van den Broeck, W., Vanhems, P.: Simulation of an SEIR infectious disease model on the dynamic contact network of conference attendees. BMC Medicine 9(1), 87 (2011)
- Xu, K.S., Hero III, A.O.: Dynamic stochastic blockmodels for time-evolving social networks. IEEE Journal of Selected Topics in Signal Processing 8(4), 552–562 (2014)
- Zhu, Y., Yan, X., Moore, C.: Oriented and degree-generated block models: generating and inferring communities with inhomogeneous degree distributions. Journal of Complex Networks 2(1), 1–18 (2014)