



ELSEVIER

Computers in Biology and Medicine III (III) III-III

Computers in Biology  
and Medicine

www.intl.elsevierhealth.com/journals/cobm

# Simulations of infectious diseases on networks

G. Witten<sup>a, b, \*</sup>, G. Poulter<sup>a</sup>

<sup>a</sup>*Department of Mathematics and Applied Mathematics, University of Cape Town, Rondebosch 7701, Cape Town, South Africa*

<sup>b</sup>*Stellenbosch Institute for Advanced Study, 19 Jonkershoek Road, Mostertsdrift, Stellenbosch, 7600, South Africa*

Received 16 June 2005; accepted 21 December 2005

## Abstract

This paper examines the spread of diseases within populations in the context of networks of potentially disease-causing contacts. We examine the assumptions underlying classical mathematical models of epidemics and how more realistic assumptions can be made using contact networks. Several well-known kinds of contact networks are examined and simulated by evaluating their structural properties relevant to disease propagation. Algorithms used in the study of these networks are explained and numerical simulations of percolation and the epidemic process carried out to explore the effects that the network structure has on disease progression.

© 2006 Elsevier Ltd. All rights reserved.

*Keywords:* Algorithms; Networks; Diseases; Simulation; Percolation; Dynamic

## 1. Introduction

The purpose of modelling epidemics is to understand the processes by which they spread, and thus provide a rational basis for formulating more effective prevention programmes and combating outbreaks. To be able to decide beforehand whether a targeted vaccination programme is likely to work and to decide which individuals will contribute most to the transmission of the infection, is to make that much more effective use of limited resources in combatting disease spread.

The epidemic process is essentially a population growth model, the disease being represented by infected individuals and the remaining (limited) resources by those susceptible to infection. Models of epidemics are classified by their assumptions about the disease and population [1]. Disease assumptions include the mechanism of infection (such as by direct contact, carriers and vectors), and removal or recovery (fixed time, probabilistic). A model makes assumptions about the population structure, which may be homogeneous except for disease status, consist of heterogeneous

subgroups where there are different risk factors (such as age, gender, culture), or be generally heterogeneous on an individual basis, which can be included as a network representation. Regarding population dynamics, the model may assume a closed population, or an open one in which it may grow or shrink over time by including birth and deaths [2].

The primary assumption a model must make is in its choice of an exhaustive and mutually exclusive classification of disease status. The simplest possible model is the SI model in which individuals are either susceptible (S) or infective (I) and the progress of the epidemic is traced by transmission between infectives and susceptibles until it ends (at least mathematically) when the entire population is infected. The model explored here is the SIR model in which infectives become removed (R) probabilistically or after a time period. The R status may correspond to death, quarantine, or recovery with permanent immunity.

There are many variations of these models, such as the SEIR model, which includes an exposed (E) class in which the individual has the disease but does not pass it on to susceptibles. A long incubation period can cause new infectives to arrive in waves [3]. The SIS model includes infectives that become susceptible again, and the SIRS model represents diseases conferring temporary immunity—these

\* Corresponding author. Tel.: +27 21 650 3191.

E-mail address: [gareth@maths.uct.ac.za](mailto:gareth@maths.uct.ac.za) (G. Witten).

models are often used to model different epidemics such as HIV or SARS [2].

In the following sections, we first discuss the traditional fully mixed models and their success at describing the basic features of epidemics, and then how the contact network does away with the assumption of full-mixedness and opens many possibilities for more realistic models. We then investigate the structural measures of networks, the mapping of the final epidemic to bond percolation, and some of the more detailed epidemiological information that can be obtained by the SIR model on networks. We then examine common network models and compare their relative usefulness as models of real-world contact networks in terms of their structural measures. Algorithms for generating the important scale-free network models are described, as well as an efficient means of calculating the mean path length, clustering coefficient and percolation threshold, and carrying out the SIR model on a general network. We finish with a discussion of dynamic networks that do away with the important simplifying assumption of a static network structure and outline the challenges involved in doing so.

## 2. Fully mixed models

Traditional models assume that the population is fully mixed, that is, a susceptible has a fixed probability  $\beta$  per unit time of contracting the disease from any infective in the population. Combined with a fixed probability  $\gamma$  per unit time of any infective becoming removed, this allows the number  $s$ ,  $i$  and  $r$  of each class of individual in a closed population of size  $N$  to be modelled by a system of ordinary differential

equations first proposed by Kermack and McKendrick in 1927 [2]:

$$\frac{ds}{dt} = -\beta si, \quad \frac{di}{dt} = \beta si - \gamma i, \quad \frac{dr}{dt} = \gamma i, \quad (1)$$

where  $s$  is non-increasing,  $r$  is non-decreasing, and the last equation is redundant since  $(d/dt)(s + i + r) = 0$  implies that  $s + i + r = N$  at all times. This model has a number of useful features: it exhibits a threshold for epidemic outbreak, since we must have  $(di/dt)(0) > 0$  for an epidemic to occur, which implies that initially we must have  $s_0 > \rho$ , where  $\rho = \gamma/\beta$  is the relative removal rate.

Furthermore, if the population has no R-status individuals,  $s_0 + i_0 = N$ , then a non-zero portion of the susceptibles will remain uninfected as  $t \rightarrow \infty$ , the number  $r_\infty$  of individuals infected and removed being the unique root [1] of the equation

$$s_0 + i_0 - r_\infty = x_0 e^{-r_\infty/\rho}. \quad (2)$$

This survival property can be seen in the case of  $\beta = 0.02$ ,  $\gamma = 0.4$ ,  $s_0 = 50$  and  $i_0 = 1$  and plotted in Fig. 1.  $s$  is the decreasing curve,  $r$  the increasing one, and  $i$  peaks around  $t = 7$ . There is a 1 in 50 chance of contracting the disease from any infective and a 40% chance of recovery per day. The fully-mixed nature of the model might make this a realistic model of a mild influenza outbreak in a university residence with time measured in days. The model thus exhibits an epidemic threshold and finite survival rate. There are many variations, such as finite birth and death rates as in the SIS and SIRS models. Extra equations account for states such as in the SEIR model and stratified populations in which portions of susceptibles have different risk factors.

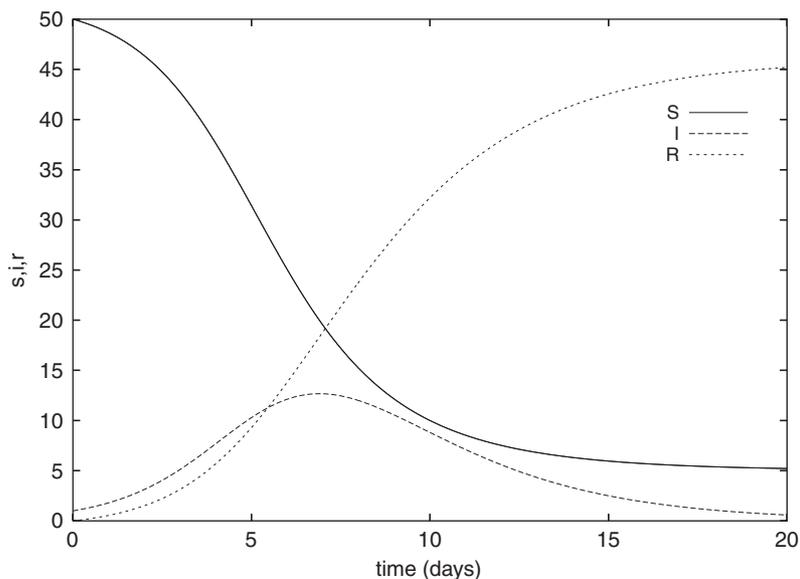


Fig. 1. Progress of a classical SIR epidemic model. Model parameters:  $\beta = 0.02$ ,  $\gamma = 0.4$ ,  $S_0 = 50$ , and  $i_0 = 1$ . The model exhibits an epidemic threshold and finite survival rate.

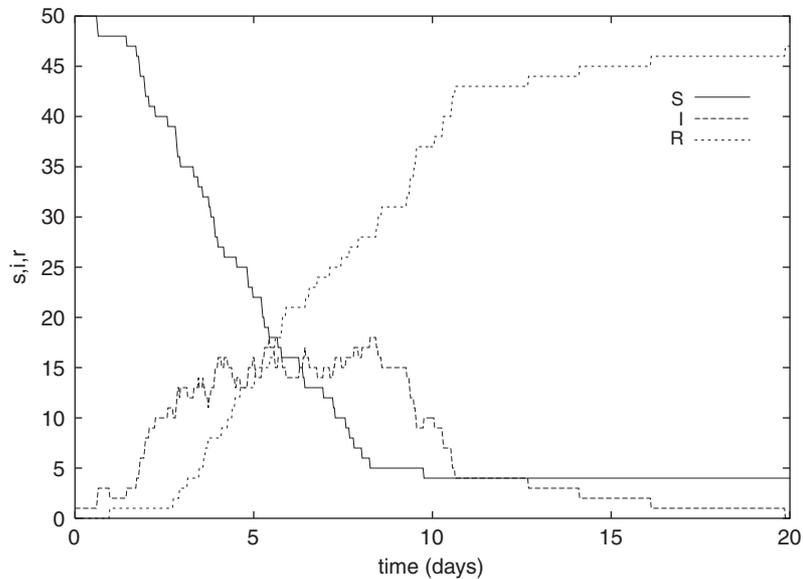


Fig. 2. Progress of a classical stochastic epidemic choosing a sufficiently small  $\Delta T = 1/2\gamma N$ .

There are also models for carriers as a different type of infective and separate populations for vector-borne diseases, such as modelling the mosquito population in the context of malaria.

$s$ ,  $i$  and  $r$  as continuous variables is a poor approximation in small populations, and so an alternative model is a Markov process with a discrete population and time, in which the probabilities of transitions from  $(s, i)$  to  $(s - 1, i + 1)$  and  $(s, i - 1)$  are given by  $\beta si\Delta T$  and  $\gamma i\Delta T$ , respectively. This can be simulated stochastically on computer by choosing sufficiently small  $\Delta T$ , so that the mean Poisson-distributed time of a state transition is larger than 1. The university residence example is a small population, and the stochastic simulation is shown in Fig. 2 for a time step of  $\Delta T = 1/(2\gamma N)$ . Note that the outcome of the simulation can be very different from the deterministic model, sometimes the epidemic does not occur at all or progresses more slowly or more quickly, as is the case in Fig. 2.

### 3. Network models

Network models do away with the assumption of a fully-mixed population. A population is represented as a collection of individuals, called vertices or nodes from the terminology of graph theory or sites in the context of percolation. Two nodes are connected by an edge if they are in regular contact and has the potential to transmit the disease if one of the nodes is infective. On a given day it is unlikely that there is a uniform probability of shaking hands with any given person in your city; you are far more likely to shake hands with one of the 20 or more people you see regularly, and this would be represented as 20 edges between yourself and 20 other nodes [4].

The structure of the network also provides many properties that influence the spread of an epidemic. With most diseases people you have contact with are likely to have contact with

each other as well or, in the context of HIV, a small number of high-degree nodes could be a significant influence in its spread [5]. We can also expect the structure to be different for a disease that is airborne, contact-transmitted, or sexually transmitted. The progression of an epidemic is modelled by initially putting one or more (usually one) of the nodes in a I state and the rest in a S state. At each successive time step, each neighbour of an infective node becomes infected with probability  $0 < \beta < 1$ ,  $\beta$  representing the rate of infective contacts, and each infective becomes removed with probability  $\gamma$ . Here  $\Delta T = 1$  time unit because, unlike the stochastic model, more than one infection is able to take place in each step. It can also be seen that the epidemic network is exactly the traditional stochastic model in the case of a complete graph, in which every node is connected to every other. More realistic assumptions can be added to these models. For example, an alternative to probabilistic removal is to put a fixed time limit  $\tau = 1/\gamma$  on how long a node remains infective. More complex models can be made by drawing  $\beta_{ij}$  and  $\tau_i$  for nodes  $i$  and  $j$  from probability distributions that might reflect risk factors amongst nodes of different types in a stratified population. Simulating networks have the disadvantage of being computationally intensive: an epidemic in a population of just a few hundred thousand nodes can take hours of computation to obtain averaged results over a range of model parameters. However, that is a once-off cost of doing away with the assumption of a fully-mixed population, and the range of possibilities it provides makes it a necessary trade-off.

#### 3.1. Network measures

Network measures are statistical properties of network structure, which are useful for understanding the character of networks that might be far too large to visualise. We look qualitatively at how these measures might directly influence the

progress of a disease modelled on the network, and describe what we expect in an actual human contact network, later to be used to evaluate the realism of commonly used network structures.

The coordination number  $z$  is the average degree of a site in the network. In a human network of hand-shaking for modelling influenza, we would expect  $z$  to be quite high, and a disease will spread faster with higher coordination number. In a network of sexual contacts for modelling HIV,  $z$  would be low and does not reflect the small number of high-degree nodes, given by the degree distribution. The degree distribution,  $p(k)$ , is the probability of a node having degree  $k$ , and a number of social networks show highly skewed distributions of power law  $P(k) \sim k^{-\alpha}$  and power-law with exponential cutoff  $P(k) \sim k^{-\alpha}e^{-k/\kappa}$  forms.

The mean path length,  $\ell$ , is the average number of steps between any given node and any other node in the network, averaged over all nodes and all possible networks of the same type and parameter set. A closely related quantity is the network diameter,  $D$ , which is the average over all nodes in the network of the shortest path distance to the furthest node from it.

Related to  $\ell$  and  $D$  is the “small world” property of many real-world networks [4]. In particular, in contact networks as anyone who has come across surprising connections between arbitrary strangers and relatives, friends or friends of friends. A small-world network is quantified as one in which the mean path length and diameter increases as a logarithm of the number of nodes. The significance of the small world property of diseases is that it does not take long for an initial infection to reach remote parts of the contact network, since there are generally no remote parts.

The clustering coefficient  $C$  can be thought of as averaging over all nodes in the graph; “how many of my friends are friends with each other”. That is, each node’s contribution to the sum is the number of neighbours which are connected to each other, divided by the  $n(n-1)/2$  possible ways to pair  $n$  neighbours.

If the network is complete, then each node’s neighbours will form a complete network and the clustering coefficient will be  $C = 1$ . Thus,  $C$  gives a measure of how close the network is to the classical stochastic model. In a human contact network we expect a significant non-zero value of  $C$  and also expect that in disease modelling clusters will tend to get infected at once, as in the case of influenza in a university residence. Newman [5] also provides an alternative definition that is more difficult to compute but easier to derive analytically and does not weigh low-degree vertices as heavily.

The community structure of a network, in which there are denser connections within communities than between them, is also important but difficult to quantify. Newman [5] provides algorithms which can be used to estimate the number and strength of communities, and cites an Ohio study of the friendship network of children at a school which clearly shows divisions along lines of middle versus upper ages and by race. On the other hand, the artificial construction of community structure in existing network classes without disrupting degree distribution and other measures is a topic for further research that might

be approached by either thinning out a single network to produce community divisions, or by generating separate networks and appropriately rewiring connections to produce the desired inter-community connection density.

### 3.2. Percolation and epidemics

A property of networks that has direct relevance to epidemiology is the bond percolation threshold, which is directly equivalent to the epidemic threshold in traditional models.

In percolation, the nodes of a network are referred to as sites, and the edges or connections as bonds. If a site  $i$  is infected for time  $\tau_i$  and has infection-causing contacts with site  $j$  at a rate of  $0 < \beta_{ij} < 1$  per time unit, then the probability of  $j$  not being infected is  $(1 - \beta_{ij})^{\tau_i}$  for  $\tau_i$  time steps of  $\Delta t = 1$ . Therefore, the probability of infection across the edge  $ij$  is the transmissibility  $T_{ij}$ ,

$$T_{ij} = 1 - (1 - \beta_{ij})^{\tau_i}. \quad (3)$$

If an infection occurs, the bond is said to be occupied, otherwise, it is unoccupied. Thus, contact rate and infective time translates directly into a bond occupation probability. If  $\tau_i$  are equal across all nodes and  $\beta_{ij}$  across all edges, then Eq. (3) can be used to produce a uniform probability  $T$  of bonds occupied, otherwise, Newman [6] shows that the distributions of  $\beta_{ij}$  and  $\tau_i$  can be averaged to produce  $T$ .

Starting with a network of unoccupied bonds and randomly occupying them one by one corresponds to continuously increasing  $T$  from 0 to 1. Percolation is said to occur when occupied bonds connect enough sites together to create a percolating or spanning cluster, which corresponds to an epidemic with the size of the cluster. This occurs over a relatively narrow range of  $T$ , associated with a critical  $T_c$  called the percolation threshold, or in diseases, the epidemic threshold. On lattices a spanning cluster is defined as one which connects the top and bottom and/or left and right edges of the lattice, but for general networks it is described more loosely in terms of a giant component that covers a large portion of the graph. Taking  $T$  from just above  $T_c$  to just below will fragment the giant component into several much smaller ones. For  $T < T_c$  the smaller components correspond to non-epidemic outbreaks of disease, and their average size can be calculated as a function of  $T$ . Above  $T_c$  this size tends to infinity in the limit of infinitely large networks due to the giant component. Newman [6] found exact solutions for the percolation threshold on general networks, and also demonstrated by simulation that different values or distributions of  $\beta_{ij}$  and  $\tau_i$  that map to the same value of  $T$  show the same mean cluster and epidemic size, and thus that the result of simulating SIR on a network maps directly onto the problem of bond percolation.

### 3.3. Epidemiological measures

Beyond the epidemic size and distribution outbreak sizes obtained by percolation, there are a number of measures of an epidemic that require an actual SIR simulation to be carried out. For one, the distribution of infectives in the network is

important—are they concentrated, or spread out? Does the infection propagate outwards in a ring, leaving removed nodes behind, as in Grassberger’s lattice model [7]? In the context of HIV, it is important to know the proportion of high-degree nodes that are infected to determine their significance in the spread of the disease.

To track the progress of the disease, the rate of new infections is an important measure. It can be accumulated over each time step and plotted. The height of an epidemic can be thought of either as when the number of infectives is at a maximum, an important measure visible in the infective curve of Fig. 1, or when the infection rate is at a maximum. The proportion of untouched susceptibles is a measure of the impact of the epidemic, as a network with low degree nodes as well as high-degree ones may find a large number of low-degree nodes untouched in the epidemic.

#### 4. Common network models

In this section we describe a number of widely studied networks on which epidemic models are often based, and we provide some analytical properties and common variations. The models described are the random, Watts–Strogatz, lattice, Barabási–Albert, and power-law or “scale-free” networks, of which all but the lattice exhibit the small world property.

##### 4.1. Random networks

Random networks are simple models that can exhibit the small-world property—it consists of  $N$  sites and a probability  $p$  of any two sites being connected and the average degree is  $z = pN$ . An alternative approach is to specify that each site has degree  $z$ , in which case there are  $\frac{1}{2}zN$  connections. These are also called Poisson random graphs because their degree distribution is  $p(k) \approx z^k e^{-z}/k!$  [5].

Random networks are small-world, which Newman [8] explains by each site having  $z$  nearest neighbours, and  $z^2$  second-nearest neighbours and so on, such that the diameter  $D$  is given by  $N = z^D$  and hence the diameter and mean path length scale logarithmically with the number of nodes. However, random networks lack important properties of contact networks, in particular, it has a low clustering coefficient of  $C = z/N$ .

Also, random networks are said to be “well-mixed” (as opposed to fully mixed like classical models). They thus have a much lower epidemic threshold than expected in real populations (see Fig.3 for  $p = 0.05$  on a  $N = 1000$  graph (i.e.  $z = 50$  is reasonable for human social networks).

##### 4.2. Lattices

A  $d$ -dimensional lattice is a regular arrangement of sites in a space of dimension  $d$ , and in which each site is connected to nearest  $z(k)$  neighbours of distance  $k$  or less from itself. The number of sites  $N$  is found from the “side-length”  $L$ , such that  $N = L^d$ , with  $N = L^2$  for the usual case of a square lattice.

The simplest form is a line of sites,  $d = 1$ ,  $k = 1$ , which can be made into a circle using periodic boundary conditions (connecting the last site to the first). More useful for epidemics is the  $d = 2$  lattice in which points are laid out on a grid (or a torus if the boundaries wrap). The distance metric for  $k$  is usually either the Euclidian ( $\sqrt{(\Delta x)^2 + (\Delta y)^2}$ ) distance between the points, or the Manhattan ( $|\Delta x| + |\Delta y|$ ) distance that is obtained following only horizontal and vertical edges.

The advantage of a lattice is that it naturally incorporates spatial separation of sites and provides a simple analytical case of bond percolation with threshold  $T_c = 0.5$ , prompting Grassberger [7] to use the 2D lattice in his foundational research in dynamical percolation in the epidemic process. If  $T$  is only a little greater than 0.5, as in the epidemics shown in Fig. 3 simulated on an  $L = 100$ ,  $d = 2$ ,  $k = 1$  grid with  $\beta = 0.05$  and  $\tau = 15$  corresponding to  $T = 0.54$ , we find the epidemic to be long-term and always at a low level. For  $k > 1$ , we also find a non-zero clustering coefficient that is constant as  $N$  varies, as in social networks [8]:

$$C = \frac{3(z - 2d)}{4(z - d)}. \quad (4)$$

However, the lattice is inherently large since the mean path length  $\ell = \frac{1}{4}dL/k$  in the 2D lattice [9] scales linearly instead of logarithmically. All nodes also have the same degree, whereas social networks usually have skewed distributions.

##### 4.3. Watts–Strogatz networks

The Watts–Strogatz model adds the small-world property to a general lattice by creating shortcuts. The original model used a  $d = 1$  ring lattice and rewired each edge in the lattice with probability  $\phi$ . It therefore interpolates uniformly between a lattice at  $\phi = 0$  and a random graph at  $\phi = 1$ . Fig. 4 shows how the path length decreases with increasing  $\phi$ . Newman and Watts [9] modified the model to preserve the underlying lattice for each edge by adding a shortcut with probability  $\phi$ .

In the case of a  $d$ -dimensional lattice of side  $L$  in which  $2k$  nearest neighbours along each of the principle axes is connected, the Watts–Strogatz network averages  $\phi kdL^d$  shortcuts, and coordination number  $z = 2kd(1 + \phi)$ . At small  $\phi$ , the model is “large-world” with  $\ell \sim L$ , and at large  $\phi$  small-world with  $\ell \sim \log N = d \log L$  [9]. For small  $\phi$  the Watts–Strogatz network retains the clustering of the underlying lattice. Watts–Strogatz networks also show significantly higher percolation thresholds on account of the underlying lattice, where for well connected individuals ( $k = 4$ ) and many shortcuts ( $\phi = 0.25$ ) the percolation threshold is  $T_c \approx 0.12$ .

However, in networks where there is geographical separation, such as human social networks, distance should affect the probability of a connection. Newman [8] cites work by Kleinberg [10] in which he found that there exists a simple means of finding short paths between any two nodes in a small-world network on a 2D lattice using local information only if the

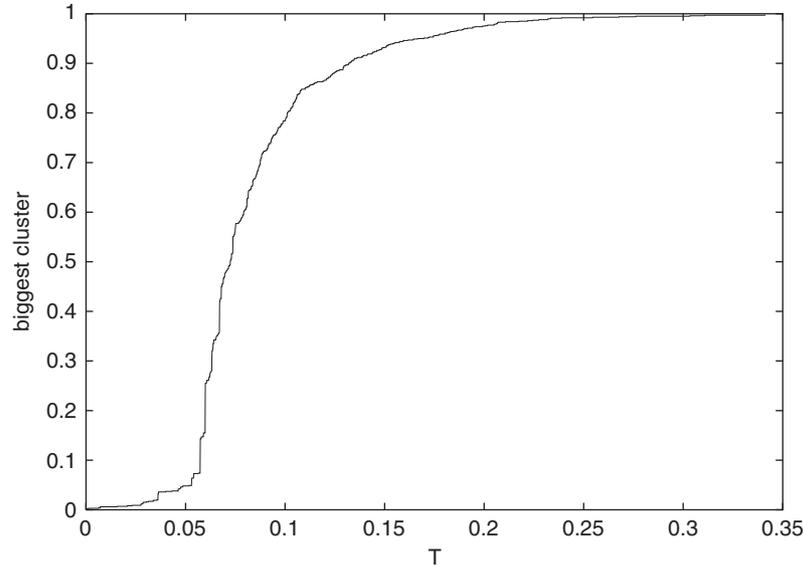


Fig. 3. Random network largest cluster with  $N = 1000$ ,  $p = 0.02$ . Random networks exhibit low epidemic thresholds.

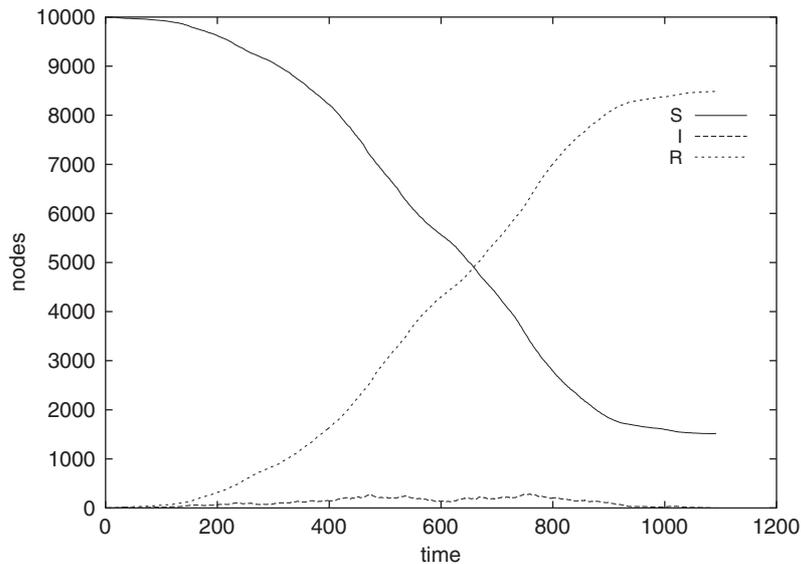


Fig. 4. SIR Lattice with,  $d = 2$ ,  $L = 100$ ,  $k = 1$ ,  $\beta = 0.05$ ,  $\tau = 15$ . The epidemic persists at a low level. For  $k > 1$ , there is a non-zero clustering coefficient that is constant as  $N$  varies.

probability of a shortcut fell as the inverse-square of the distance between nodes.

#### 4.4. Power-law networks

One feature lacking in Watts–Strogatz networks is a power-law degree distribution often found in human social networks, such as the network of sexual contacts studied by Liljeros [11]. These are of the form  $P(k) \sim k^{-\alpha}$ , and sometimes with an exponential cutoff  $P(k) \sim k^{-\alpha}e^{-k/\kappa}$ . Power-law networks are often called “scale-free” on account of the degree distribution being self-similar on different levels, although technically the network itself exhibits distinct scales.

The properties of power-law networks vary greatly with the exponent  $\alpha$ , described in detail by Newman [5]. For  $\alpha < \frac{7}{3}$  the clustering coefficient grows with graph size, is constant for  $\alpha = \frac{7}{3}$ , and shrinks as  $C \sim 1/N$ . For  $\alpha < 2$ , a random graph such as Fig. 5 will always have a giant component, and possibly a few small ones in finite graphs. For  $2 < \alpha < 3.4788 \dots$  the giant component fraction steadily decreases (and implies that we cannot reliably create a single-component power-law network with  $\alpha$  in that range), until for larger  $\alpha$  it ceases to exist. Also, increasing  $\alpha$  tends to increase path length and diameter due to fewer high-degree nodes.

Power-law networks are small world due to the small proportion of high-degree nodes, and are unrealistic primarily in the

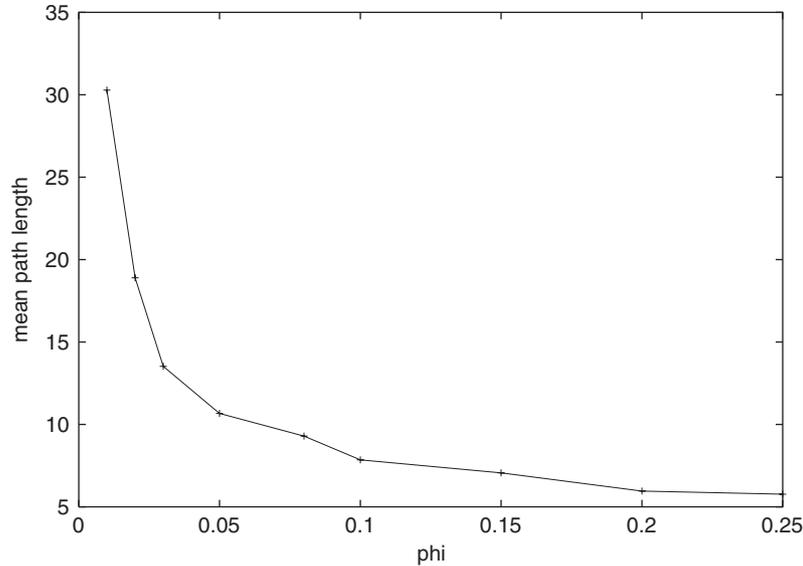


Fig. 5. Watts–Strogatz path length for  $N = 1000$  and  $k = 2$ .

random mixing (connections not being geographically correlated). Also, unlike the traditional epidemic threshold, May and Lloyd [12] demonstrated that there is a zero percolation threshold in transmissibility for  $2 < \alpha \leq 3$  in infinite graphs while finite graphs show small but non-zero thresholds. Further work by Newman [6] with a bipartite power-law graph does however exhibit the threshold  $T_c$  that rises from 0 to 1 over the range  $3 < \alpha < 3.4788$ , and since Liljeros’s [11] measured  $\alpha \approx 3.2$  for a network of sexual contacts it suggests that a random scale-free network still has the necessary properties for modelling STDs.

Fig. 5 shows a simulated epidemic of  $\beta=0.02$ ,  $\tau=10$  ( $T=0.4$ ) on a complete power law network of  $\alpha = 2$ , created using the Monte Carlo Markov chain technique of the following section. The progress of the epidemic is markedly different from that of the classical model, in that even far above the epidemic threshold  $T_c \approx 0.07$  a large number of susceptibles remain uninfected.

One means of restoring a geographical metric is the lattice-based scale-free network of Warren et al. [13]. They scatter nodes randomly on a 2D lattice, with connections to lattice sites within a radius  $R$  drawn from a power-law distribution, and connecting nodes if their disks overlap. Their model also shows the interesting property of a non-zero percolation threshold in the density of the nodes and transmissibility.

#### 4.5. Barabási–Albert

The Barabási–Albert scale-free model is a network of the form  $P(k) \sim k^{-3}$ , which was designed to allow network growth, as opposed to the static nature of the networks above. The network starts with  $N_0$  nodes and no edges, and  $N - N_0$  nodes are added one by one. Each new node is connected to  $m$  existing nodes, each with a probability proportional to its degree. This is the same as  $m$  times randomly picking an edge, then randomly picking one of the nodes on it, and connecting to that node.

The result is preferential attachment instead of random mixing: high-degree nodes naturally accumulate more connections. Although simple, power-law and small-world, the model has a disadvantage over random mixing in that there are much shorter paths between high-degree nodes, so an epidemic can occur for any non-zero transmissibility even though it corresponds to  $\alpha = 3$ , as demonstrated in Figs. 6 and 7 for  $N = 5000$ , where finite-size effects produce  $T_c \approx 0.12$ . Note also the distinctly nonlinear curve in the size of the largest cluster. Like the ordinary power-law networks, the epidemic leaves many susceptibles behind in finite networks as in Fig. 6, despite the cluster being infinite in the limit of large network size.

#### 4.6. Generating scale-free networks

The algorithm, or sequence of steps, for generating a random network, lattice, Watts–Strogatz network or Barabási–Albert network follows trivially from the description of the network. However, to create a network with an arbitrary specified degree distribution, such as power-law, so as to maintain connectivity and a random choice from the ensemble of such networks, is rather more complicated.

Generating random mixing scale free networks is usually accomplished by first choosing a degree for each node from the distribution. An important result in the power-law network  $p(k) \sim k^{-\alpha}$  is the maximum value of  $k$  likely to occur [5], from which a discrete probability distribution for  $1, \dots, k_{\max}$  can be constructed:

$$k_{\max} \sim n^{1/(\alpha-1)}. \quad (5)$$

The approach then taken by Newman [5] is to treat the degree as the number of “spokes” on the node, which are connected randomly to other available spokes until none remain, thus selecting a random member of the ensemble of networks with that degree sequence. However, with  $\alpha > 2$  the graph will have

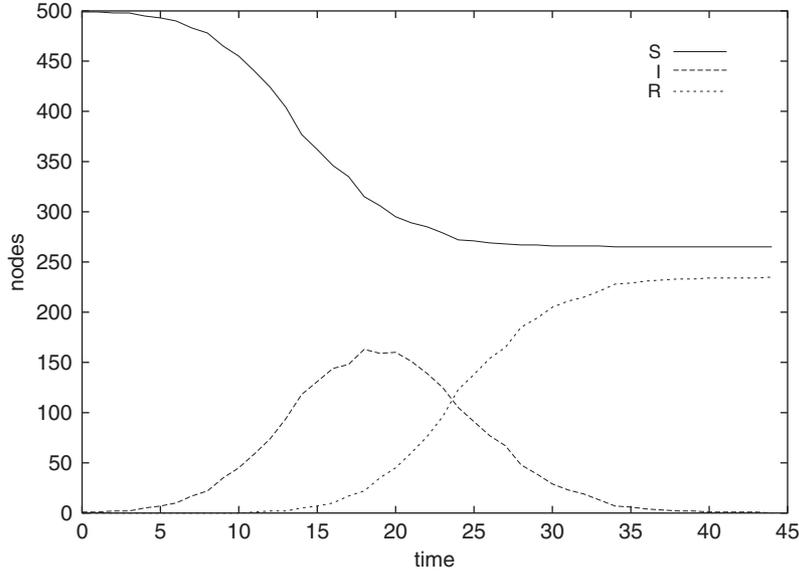


Fig. 6. Complete power-law SIR network with  $N = 500$  and  $\alpha = 2$ . Epidemic progress is different from the classical model.

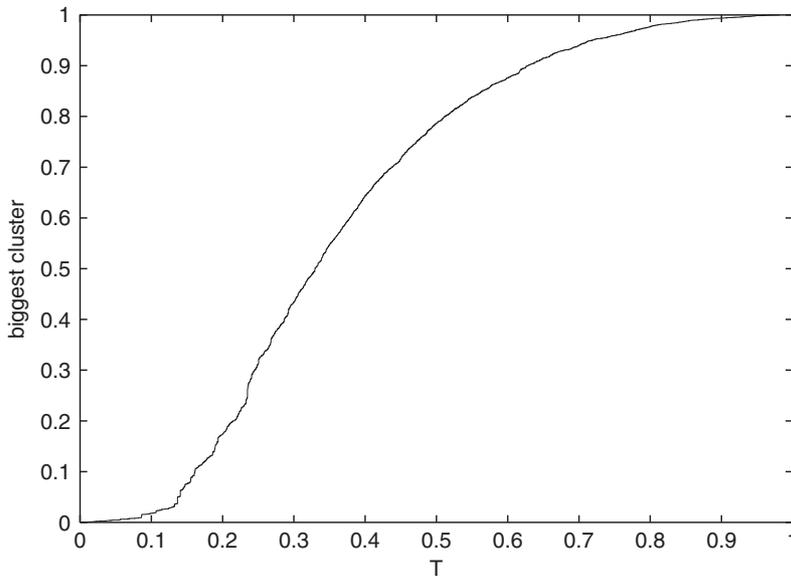


Fig. 7. Barbasar–Albert network percolation:  $N = 5000$ ,  $N_0 = 5$ ,  $m = 2$ .

a lot of small components. Interestingly, an attempt here to produce a bipartite network of accumulated sexual contacts by this algorithm with  $\alpha = 3.2$  [11] resulted mostly in pairs.

Multiple components is realistic in sexual contact networks that span only a few months or years, and possibly makes this kind of random network less useful for modelling HIV, where the long infective time necessitates a fragmentary network of regular contacts that changes on long timescales if an epidemic is to occur. Using networks accumulated over longer time periods to achieve connectivity would produce more high-degree nodes than is observed in any local time-period, and thus overestimate the spread of the disease.

One alternative taken here is to choose the network instead from the ensemble of *connected* graphs with the same power-

law degree distribution, which would then reflect a rapid spread when the disease reaches a relatively large component as might be found in a country-sized population. Unfortunately, a single-component form is only reliably possible for  $\alpha \leq 2$ .

The algorithm used in this paper is that of Gkantsidis et al. [14]. For this algorithm the degree distribution is first examined to see whether a connected network is possible for the degree sequence  $k_1 > k_2 > \dots > k_N$ . The first condition is that for each subset of the  $h$  highest degree nodes, the sum of their degrees can be absorbed by themselves and the outside nodes:

$$\sum_{i=1}^h k_i \leq h(h-1) + \sum_{i=h+1}^N \min\{h, k_i\}. \quad (6)$$

The second is that the graph contains a spanning tree,

$$\sum_{i=1}^N k_i \geq 2(N-1). \quad (7)$$

Together they are sufficient [14] to guarantee connectivity. The initial graph is constructed by picking random vertices and connecting them to a node with the highest residual degree. We found an efficient means of locating such a node by creating a table that lists the nodes of a given degree, and which is updated whenever two nodes are connected.

Deterministic attachment clearly does not produce a random instance of the graph, and the second step is to randomise it by a Markov chain process that preserves the degree distribution. The simplified approach used here is to pick any two edges  $(x, y)$  and  $(u, v)$  from the graph at step  $t$ . If  $(u, x)$  and  $(v, y)$  are not edges, and the graph obtained by replacing  $(x, y)$  and  $(u, v)$  with  $(u, x)$  and  $(v, y)$  is connected, then that is the graph at step  $t + 1$ . The mean path length and diameter of the graph converge from the low initial value to ensemble most probable values by  $t$ , roughly on the order of the number of edges.

It is possible to speed up the algorithm by performing multiple swaps between connectivity checks, and maintaining an undo list, where the window of swaps is incremented if the network is found to be connected, and halved if it is found to be disconnected, necessitating an undo.

## 5. Network algorithms

We now discuss the algorithms implemented here for numerically calculating structural and epidemiological measures on networks.

### 5.1. Mean path lengths

The mean path length is found by calculating the shortest path distance from a given node to each of the others, and averaging, and then averaging again over each possible starting node. The path distance can be found by doing a breadth first search in which for visiting a node at distance  $d$  from the start, its unvisited neighbours of distance  $d + 1$  are pushed onto a queue. The queue is first-in–first-out, so nodes are visited in increasing order of distance (if a node further away than the closest were visited out of order, it could put a longer than optimal path length on one of its neighbours). Formally, for shortest path distances  $d_{ij}$

$$\ell = \frac{1}{N} \sum_{i=1}^N \ell_i, \quad \ell_i = \frac{1}{N} \sum_{j=1}^N d_{ij}. \quad (8)$$

It was also found that a good estimate can be obtained by averaging over a random sample of starting nodes to save computation time on very large graphs.

### 5.2. Clustering coefficient

The clustering coefficient was found by averaging the clustering of each node in the graph. The clustering of a node  $i$  is found by, for each of the node's  $n$  neighbours  $j$ , counting those neighbours of  $j$  which are also neighbours of  $i$ , and then dividing by  $n(n-1)$ , since it is possible for each of  $n$  neighbours to be connected to at most  $n-1$  others. Again, a good estimate can be obtained using a random sample of nodes. Formally, for distances  $d_{ij}$ ,

$$C = \frac{1}{N} \sum_{i=1}^N C_i, \quad C_i = \frac{1}{n(n-1)} \sum_{j \in \text{adj}(i)} \delta(1, d_{ij}). \quad (9)$$

### 5.3. Percolation thresholds

Percolation thresholds were calculated using the union-find algorithm proposed by Ziff and Newman [15]. A continuum of  $T$  from 0 to 1 is obtained by randomly occupying bonds in the network one by one. At each step, the clusters formed are kept track of using a union-find structure, in which each site is labeled with a parent that is one of the other sites in the same cluster. The site at the root of the tree for that cluster is labelled with a negative number that is the number of sites in the cluster (lone sites are thus labeled “−1”).

Detecting whether a new bond joins two clusters is done by a “find” operation which follows the parents of the sites on either end of the new bond to their roots. If the roots are different, then the union of the clusters is carried out by setting the parent of root of the smaller cluster to be the root of the larger cluster, whose size is then updated.

Finds are made more efficient by path compression: after performing the find, the root of the tree is now known, so one more pass through the path to the root is carried out, setting the parent of each site along the way to point directly to the root.

Using the tree structure, the average cluster size can be found by dividing the size of the network by the number of clusters, and the percolation threshold occurs when several smaller clusters join to form one which covers a large part of the network. This can be detected by a sudden increase in the size of the largest cluster followed by a steady or decreasing positive gradient in that size until the graph is fully occupied.

### 5.4. SIR modelling

Other epidemiological measures require carrying out an SIR simulation as described in the section on “network models”. One means of avoiding passing through the entire network is to store the infectives in a separate list. By running through this list, only the neighbours of infectives are investigated and infectives whose time limit is reached can change state to removed. Also new infectives can be added to a separate list that is tagged onto the list of infectives at the end of the time step, which avoids making a copy of the network for updating purposes.

## 6. Dynamic networks

The networks discussed, thus, far have all been accumulated networks of potentially disease-transmitting contacts that occur between individuals on a regular basis. This assumption simplifies computation and is valid possibly on a timescale of several months. However, on longer timescales as in the case of HIV one might consider a regular contact network whose structure changes through addition and deletion of nodes and edges. Unfortunately, the network models here, except the random network, do not easily support modification since random changes will eventually lead to all the networks converging on a random graph. Constructing realistic random changes that maintain the network type, and then also the network parameters to within some tolerance, is a considerably more difficult task for all but Barabási–Albert network which was designed to support growth.

The alternative is to consider a network of instantaneous contacts. The network is not stored directly, but at each time step the contacts between nodes is drawn from a probability distribution. This has the advantage of easily adding and removing nodes, since no edges need to be modified. However, the problem of persistence has merely been shifted from modifying edges to modifying a probability distribution, and measures such as path length and clustering become difficult if not impossible to quantify.

A conceptually simple approach to the dynamic network has been taken in a recent work by Verdasca et al. [16]. They use an underlying 2D lattice to provide spatial proximity. Contacts occur with probability  $p$  with one of the 12 nearest neighbours (8 around the node, 4 at a distance of 2 at the cardinal compass points), and small-world contacts at probability  $1 - p$  with any other node in the network. The primary motivation for the model was that it would easily accommodate birth by changing R nodes to S, and death by changing S or I nodes to R, allowing a network model of an SIR endemic.

Their simulation was carried out stochastically in a Monte Carlo fashion by performing  $N = L^2$  random site updates with relative rates of birth, death and infection of  $\mu$ ,  $\mu$  and  $\beta$ , but only carrying out birth if the chosen site is R, only carrying out death if it is S or I, and only carrying out infection if it is S and a second randomly chosen contact (neighbour with probability  $p$ , distant with  $p - 1$ ) is infective.

Their approach is not a fully general dynamic network and suffers from being computationally intensive: order of magnitude calculations show several weeks of computation behind their results carried out on 250 000 and 1 000 000-nodes graphs. Nevertheless, they succeeded in showing a percolation transition from outbreak to epidemic behaviour at  $p = 0.1$  in the SIR model and were able to fit an SEIR version of their model to produce three-year periodic oscillations in the rate of new infectives every two weeks, and average age at infection, using known parameters and data on measles. This much was accomplished without the difficult-to-measure parameters of seasonal forcing used in traditional models.

## 7. Conclusions

Using a contact network removes the assumption of a fully-mixed population inherent in traditional models. There are several common network models, each of which has similarities and differences from those expected of typical human contact networks. One important structure, the community, may be possible to incorporate as an additional step. Information can be gathered from a network's structural properties, the mapping of the removed nodes after an epidemic to a bond percolation problem, and actual simulations of the epidemic, and relatively efficient algorithms exist for all of these tasks. The next step in realism is to make the network dynamic, though the reduction in information available about the network structure and the complexity of computation makes this task difficult.

## References

- [1] D.J. Daley, J. Gani, Epidemic Modelling, Cambridge University Press, UK, 1999.
- [2] H.W. Hethcote, The mathematics of infectious diseases, *SIAM Rev.* 42 (4) (2000) 599–653.
- [3] B.T. Grenfell, O. Bjornstad, Sexually transmitted diseases—epidemic cycling and immunity, *Nature* 433 (2005) 366–367.
- [4] D. Watts, S. Strogatz, Collective dynamics of small world networks, *Nature* 433 (1998) 366–367.
- [5] M.E.J. Newman, The structure and function of complex networks, *SIAM Rev.* 45 (2) (2003) 167–257.
- [6] M.E.J. Newman, Spread of epidemic disease on networks, *Phys. Rev. E* 66 (2002) 016128.
- [7] P. Grassberger, On the critical behaviour of the general epidemic process and dynamical percolation, *Math. Biosci.* 63 (1983) 157–172.
- [8] M.E.J. Newman, Models of the small world, (<http://www.arxiv.org/cond-mat/000118>), 2000.
- [9] D.J. Watts, M.E.J. Newman, Scaling and percolation in the small-world network model, (<http://www.arxiv.org/cond-mat/9904419>), 1999.
- [10] J. Klienber, Navigation in a small world, *Nature* 406 (2001) 1849–1850.
- [11] F. Liljeros, C.R. Edling, L.A.N. Amaral, H.E. Stanley, Y. Aberg, The web of human sexual contacts, *Nature* 411 (2001) 907–908.
- [12] R.M. May, A.L. Lloyd, Infection dynamics on scale-free networks, *Phys. Rev. E* 64 (2001) 0066112.
- [13] C.P. Warren, L.M. Sander, I.M. Sokolov, Geography in a scale-free network model, (<http://www.arxiv.org/cond-mat/0207324> v3), 2002.
- [14] C. Gkantsidis, M. Mihail, E. Zegura, The Markov chain method for generating connected power law random graphs, (<http://www.siam.org/meetings/alnex03/Abstracts/CGkantsidis.pdf>), 2002.
- [15] R.M. Ziff, M.E.J. Newman, A fast Monte Carlo algorithm for site or bond percolation, *Phys. Rev. E* 64 (2001).
- [16] J.A. Verdasca, M.M. Telo da Gama, et al., Recurrent epidemics in small world networks, (<http://www.arxiv.org/cond-mat/0408002> v1), 2004.

**Gareth Witten** is a lecturer in the Department of Mathematics and Applied Mathematics at the University of Cape Town. His research is in the applications of quantitative methods to biology and medicine. He leads a vibrant research group under the South African Centre for Epidemiological Modelling and Analysis ([www.sacema.ac.za](http://www.sacema.ac.za)) that focuses on modelling diseases, in particular HIV, within the host. He is also a fellow of the Stellenbosch Institute for Advanced Study, a think-tank for complex challenges in South Africa. E-mail: [gareth@maths.uct.ac.za](mailto:gareth@maths.uct.ac.za), URL: <http://www.mth.uct.ac.za/~witten/>

**Graham Poulter** is in 2006 beginning research towards an M.Sc. in Computational Biology at the University of Cape Town. He was awarded a B.Sc.Hons with distinction in Applied Mathematics in 2005, and a B.Sc. with distinction in Physics, Applied Mathematics and Computer Science in 2004.

The core research for this paper, implementing and reviewing network models of epidemics, was carried out in 2004 during his B.Sc. final year project. E-mail: [graham@cbio.uct.ac.za](mailto:graham@cbio.uct.ac.za), URL: <http://mancala.cbio.uct.ac.za/~graham/>