

Utility of R_0 as a predictor of disease invasion in structured populations

PAUL C. CROSS^{1,2}, PHILIP L.F. JOHNSON³,
JAMES O. LLOYD-SMITH⁴, & WAYNE M. GETZ^{5,6}

¹ Northern Rocky Mountain Science Center, U.S. Geological Survey

² Department of Ecology, Montana State University

³ Biophysics Graduate Group, University of California at Berkeley

⁴ Center for Infectious Disease Dynamics, Pennsylvania State University

⁵ Department of Environmental Science, Policy and Management, University of California at Berkeley

⁶ Department of Zoology and Entomology, Mammal Research Institute, University of Pretoria, South Africa

Philip Johnson, 299 Life Sciences Addition, Berkeley CA 94720-3200

James Lloyd-Smith, 208 Mueller Lab, Pennsylvania State University, University Park PA 16802

Wayne Getz, 201 Wellman Hall, Berkeley CA 94720-3112

Author for correspondence:

Paul Cross, 229 AJM Johnson Hall, Montana State University, Bozeman MT 59717.

E-mail: pcross@usgs.gov

Word Count: 7741 including references and legends

1 **Abstract**

2 Early theoretical work on disease invasion typically assumed large and well-mixed host
3 populations. Many human and wildlife systems, however, have small groups with limited
4 movement among groups. In these situations, the basic reproductive number, R_0 , is likely to be a
5 poor predictor of a disease pandemic because it typically does not account for group structure
6 and movement of individuals among groups. We extend recent work by combining the
7 movement of hosts, transmission within groups, recovery from infection, and the recruitment of
8 new susceptibles into a stochastic model of disease in a host metapopulation. We focus on how
9 recruitment of susceptibles affects disease invasion and how population structure can affect the
10 frequency of superspreading events (SSEs). We show that the frequency of SSEs may decrease
11 with reduced movement and group sizes due to the limited number of susceptible individuals
12 available. Classification tree analysis of the model results illustrates the hierarchical nature of
13 disease invasion in host metapopulations. First the pathogen must effectively transmit within a
14 group ($R_0 > 1$), then the pathogen must persist within a group long enough to allow for movement
15 among groups. Therefore factors affecting disease persistence—such as infectious period, group
16 size, and recruitment of new susceptibles—are as important as local transmission rates in
17 predicting the spread of pathogens across a metapopulation.

18

19 Keywords: disease, invasion, metapopulation, SIR model, superspreader

1 1. INTRODUCTION

2 Early epidemiological models typically assumed that host populations were large and
3 well-mixed (e.g. Kermack & McKendrick 1927). Many human, wildlife, and livestock
4 populations, however, are structured into small groups with limited movement among groups
5 (Altizer et al. 2003; Kao et al. 2006). For example, communities of people that remain
6 unvaccinated for religious or philosophical reasons constitute isolated and weakly-linked patches
7 of susceptible hosts for diseases such as measles and pertussis (Feikin et al. 2000; Salmon et al.
8 1999). Similarly, the on-going spread of H5N1 influenza among wild birds underscores the
9 need to understand whether insights derived from the theory of epidemics in large human
10 populations can be applied accurately to diseases in wildlife. A number of studies have
11 considered the effects of spatial or social group structures on disease invasion and persistence
12 (e.g., Cross et al. 2004; Fulford et al. 2002; Hagenaars et al. 2004; Hess 1996b; Keeling 1999;
13 Keeling & Gilligan 2000a; Keeling & Gilligan 2000b; Keeling & Rohani 2002; Park et al. 2001;
14 Park et al. 2002; Swinton 1998; Thrall et al. 2000). Of particular importance is the research
15 investigating the effects of population structure in the form of households on disease invasion
16 and dynamics (e.g., Andersson 1997; Andersson & Britton 1998; Becker & Dietz 1995; Becker
17 & Starczak 1997; Schinazi 2002). In this study, we take a novel approach to investigating
18 disease invasion. Rather than analytically determining when a large outbreak is possible, we use
19 hierarchical statistical methods to determine what criteria predict successful disease invasion
20 most accurately. We then compare these results to more traditional thresholds to determine the
21 amount of prediction error arising from the different approaches.

22 The basic reproductive number, R_0 , is the expected number of infections caused by a
23 typical infectious individual in a completely susceptible population. $R_0 > 1$ is the threshold
24 condition traditionally applied for successful disease invasion (Anderson & May 1991;

1 Heesterbeek 2002; Heffernan et al. 2005). R_0 , as it is commonly used, assumes that the host
2 population size is sufficiently large that the depletion of susceptible individuals through death or
3 infection is negligible, and that the population is homogeneous or well-mixed (Anderson & May
4 1991; Keeling & Grenfell 2000). The R_0 metric has been widely studied and refined to address
5 more complex situations (e.g. multiple classes of host: Diekmann et al. 1990; spatial structure:
6 Keeling 1999; depletion of the susceptible pool: Keeling & Grenfell 2000). Although some
7 formulations of R_0 use a matrix-based approach to account for spatial or group structure (e.g.
8 Diekmann et al. 1990), R_0 is, by definition, an individual-based rather than group-based metric.
9 In this usage R_0 may be high, reflecting within-group transmission, while the probability of
10 between-group transmission remains low (Ball et al. 1997; Cross et al. 2005; Watts et al. 2005).
11 When social groups are small, understanding the processes affecting within-group invasion
12 becomes less important than understanding the processes regulating the spread of disease among
13 groups.

14 The natural invasion metric for disease in a metapopulation is R_* , defined as the number
15 of groups infected by individuals from the initially infected group (and hence the group-level
16 analogue of R_0 ; Ball et al. 1997). A similar metric, R_{H0} , was developed by Becker and Dietz
17 (1995) to assess the propagation of infection among households of variable sizes. In an idealized
18 metapopulation, analytic theory has proven R_* must be greater than one for a pandemic to occur
19 (Ball et al. 1997; Becker & Dietz 1995); under less restrictive assumptions, this same threshold
20 has been demonstrated by simulation (Cross et al. 2005). Unfortunately, R_* is difficult to
21 calculate analytically for any but the simplest metapopulation structures. Empirical estimation of
22 R_* from outbreak data would require contact tracing data at a group level, a formidable
23 challenge for wildlife or human diseases. Thus, while R_* brings conceptual clarity to the study
24 of disease in metapopulations, its immediate utility in applied settings is limited. Therefore, we

1 investigate the constituent parts of R_* to help focus field research on those parameters most
2 important to disease invasion in structured populations.

3 Many studies addressing R_0 in structured populations incorporate host movement via a
4 phenomenological mixing approach, whereby hosts do not move among groups but
5 simultaneously infect others locally and at a distance (Ball et al. 1997; Dobson & Foufopoulos
6 2001; Fulford et al. 2002; Keeling 1999; Park et al. 2001). Phenomenological mixing models are
7 often analytically tractable, but they overlook the fact that between-group movements are
8 discrete (and possibly rare) events, which can be crucial to understanding the stochastic
9 dynamics of disease invasion (Cross et al. 2005) and the role of superspreaders in fueling an
10 epidemic (Lloyd-Smith et al. 2005b). An alternative approach is to model host movement
11 mechanistically, explicitly tracking the movement of individuals between groups (eg. Cross et al.
12 2005; Hess 1996a; Keeling & Rohani 2002; Thrall et al. 2000).

13 Previously, we used mechanistic models to show that disease invasion across a
14 metapopulation depends crucially on the relative timescales of host movement and recovery from
15 disease (Cross et al. 2005). We showed that $R_0 > 1$ was insufficient for disease invasion when the
16 product of the average group size and the expected number of between-group movements made
17 by each individual while infectious (i.e. the ratio of movement rate to recovery rate) was less
18 than one (Cross et al. 2005). This previous study addressed settings where the rate host
19 population turnover was negligible relative to the rate of disease processes of infection and
20 recovery.

21 Here we expand the earlier analysis to a much broader set of disease-host relationships,
22 exploring settings where the duration of immunity ranges from transient to lifelong, or where
23 demographic processes occur on comparable (or faster) timescales to disease processes. Rapid
24 replenishment of susceptibles allows qualitatively different dynamics compared to the earlier

1 study, including the possibility for diseases to remain endemic within a local group even if
2 movement is infrequent. Given $R_0 > 1$, we investigate additional factors that help explain the
3 remaining variation in whether or not a disease will become a pandemic. We also examine how
4 these additional factors alter the structure of epidemics through their effect on the frequency of
5 superspreading events (Lloyd-Smith et al. 2005b).

6

7 **2. METHODS**

8 (a) *Model structure*

9 We use two individual-based, stochastic, discrete-time SIR models that extend our
10 previous work (Cross et al. 2005). These models differ from each other and our previous
11 analyses only in the mechanism by which the susceptible pool is replenished. In the SIRS
12 model, immunity is transient so recovered individuals can return to the susceptible state; in the
13 SIR_BD model, immunity is permanent but births introduce new susceptibles, while deaths keep
14 the population size constant. In simulations of each model, we track each individual's spatial
15 position (group membership) and disease class (S-susceptible, I-infected, R-recovered)

16 In each model four processes occur: infection, recovery of infected hosts, creation of new
17 susceptibles, and movement among groups. We take disease transmission to be frequency-
18 dependent (Getz & Pickering 1983), whereby the instantaneous rate of infection for each
19 susceptible individual in group i is $\beta I_i/n_i$, where β is the transmission coefficient, I_i is the number
20 of infected individuals in group i , and n_i is the total number of individuals in group i . Because

21 our models operate in discrete time, the expression $1 - \exp\left(-\beta \frac{I_i}{n_i}\right)$ is used to depict the

22 saturating probability of infection per time step for each susceptible individual (implicitly

23 assuming that the force of infection is constant within each time step). All disease transmission

1 is assumed to occur within local groups, and contact among groups occurs only by movement of
2 individual hosts. We assume that infected individuals recover from infection to an immune class
3 with a constant probability γ per time step. We model movement among groups in a density-
4 independent fashion such that all individuals have a constant probability μ of leaving their
5 current group in each time step. In the SIRS model, recovered individuals lose their immunity
6 with probability ρ per time step, and births and deaths do not occur. In the SIR_BD model, all
7 individuals have probability δ of dying and being replaced by a susceptible individual in the
8 same group.

9 Groups are organized on a square lattice with periodic boundary conditions (i.e.
10 movement is on a torus), where individuals move to one of their four nearest-neighboring
11 groups, chosen at random. Each simulation starts with one infected individual and all groups
12 begin with the same number of individuals. Except where otherwise noted, we ran simulations
13 on an 11 x 11 array of groups. Since our spatial model was symmetric, group sizes remained
14 relatively constant during the course of each run. Thus, our assumption of frequency-dependent
15 transmission is approximately equivalent to a rescaling of density-dependent transmission.

16 In the continuous-time analogues of our models, $R_0 = \beta'/\gamma'$ for SIRS and
17 $R_0 = \beta' / (\gamma' + \delta')$ for SIR_BD (Anderson & May 1991; McCallum *et al.* 2001). The prime
18 indicates that, in continuous time, these variables are rates rather than probabilities. For the
19 discrete-time models used here, the ratio of β/γ is an approximation of R_0 that works well when
20 the timestep is small and group sizes are relatively large. These slight approximations do not
21 change our qualitative conclusions, so for succinctness we refer to these ratios as R_0 . Note that
22 in the SIR_BD model, increasing δ reduces R_0 because death removes individuals from the
23 infectious class. To allow full comparison of the SIRS and SIR_BD models while varying ρ or

1 δ , we present SIR_BD results for scenarios both where β is fixed (so R_0 changes with δ) and
2 where β is adjusted so that R_0 remains constant.

3 (b) *Simulations and analyses*

4 Using the models described above, we explore how different parameter interactions affect
5 the outcome of disease introductions. Past studies of this model structure indicate that, for the
6 parameter ranges we explore, most introductions result in extinction within the initial group or
7 relatively complete invasion of the entire metapopulation, i.e. a “pandemic” (Cross et al. 2005).
8 As a binary measure of invasion success we declare an invasion to be successful if >90% of
9 groups are ever infected following a single disease introduction. This definition of a pandemic
10 does not count disease persistence within a single patch as successful invasion, because we are
11 focused on disease spread at the broader metapopulation scale.

12 To capture the effect of a finite, diminishing pool of susceptibles, we calculate empirical
13 \hat{R}_0 and \hat{R}_* values during the simulations. In contrast to the theoretical R_0 values calculated from
14 model parameters, these estimates are based upon individual simulation results. For each
15 simulation we calculate the individual reproductive number, ν (Lloyd-Smith et al. 2005b), by
16 tracking the number of infections caused by the index case and then averaging ν over many
17 simulations to calculate \hat{R}_0 (Cross et al. 2005). Similarly, to calculate \hat{R}_* we take the average
18 over ν_* , which in turn is calculated by tracking the number of groups infected by individuals
19 from the index group. As estimates from model output, ν , ν_* , \hat{R}_0 and \hat{R}_* all incorporate the
20 effects of spatial structure, stochasticity, host movement, and depletion of the susceptible pool
21 within the infectious period of the index case (or group). We consider ν , ν_* , \hat{R}_0 and \hat{R}_* to be
22 ‘emergent’ quantities since they can only be estimated once the initial generations of a disease
23 invasion have occurred. Following Lloyd-Smith et al. (2005b), we assess the frequency of SSEs

1 in different population structures by constructing a histogram of infections caused by each index
2 case to calculate the proportion of the distribution beyond the point corresponding to the 99th
3 percentile of a Poisson distribution with the same mean. Since the distribution is not Poisson this
4 tail will not necessarily contain 1% of individuals, but rather $y\%$. The superspreading load (SSL)
5 is the observed number of SSEs divided by the expected based upon a Poisson distribution that,
6 when greater than one, predicts reduced invasion rates but more intense epidemics once invasion
7 occurs (Getz & Lloyd-Smith 2006; Lloyd-Smith et al. 2005b).

8 We used classification and regression tree analyses to explore which factors influence the
9 variation in disease invasion outcomes (Breiman et al. 1984). Classification tree analyses have
10 been used extensively in clinical risk assessments (e.g. Begg 1986; Steadman et al. 2000), and
11 are becoming more common in the ecological literature (e.g. Brose et al. 2005; De'ath &
12 Fabricius 2000; Karels et al. 2004; Usio et al. 2006). Classification trees divide data in a
13 hierarchical manner using binary rules based upon single predictor variables. Threshold criteria
14 are then chosen to partition the response variable into groups that are as homogeneous as
15 possible. We used the Gini index as the splitting criterion. Since larger trees will always predict
16 the learning dataset better, we used 10-fold cross-validation and the 1–SE rule to guide in the
17 choice of the ‘best’ tree size. This is a method to minimize the amount of prediction error on
18 testing data (not used in the construction of the tree) while also incorporating a penalty for
19 increasing tree size (Breiman et al. 1984). Since the classification analysis is intended to be
20 heuristic, for clarity of presentation we present trees that are slightly simpler than those trees
21 chosen according to the 1–SE rule, but resulted in only a minor increase in misclassification
22 (details on alternative trees are presented in the supplementary material). We explored three
23 different sets of explanatory variables for the classification analysis: 1) six raw model parameters
24 (β , γ , ρ , δ , μ and n), 2) five aggregate model parameters (β/γ , $\rho n/\gamma$, $\mu n/\gamma$, ρ/γ and ρn), and 3) the

1 five aggregate model parameters as well as ν and ν_* . Although we report results for all analyses
2 in Table 1, only the classification tree using the aggregate model parameters is shown in the
3 main text: the others are illustrated in the supplementary material.

4 We compare the criteria for invasion from the classification tree analysis to more
5 traditional thresholds using a vocabulary taken from literature on diagnostics, where one assesses
6 the utility of a diagnostic tool according the proportion of times it yields false-positive and false-
7 negative results. In the case presented here, false-positives occur when the criteria for invasion
8 are met but the disease does not actually invade. False-negatives occur when the criteria for
9 invasion are not met and yet the disease does invade (recall that a successful invasion is defined
10 as the disease infecting individuals in over 90% of the groups of the metapopulation). Note that
11 $R_0 > 1$ and $R_* > 1$ are theoretical thresholds determining when disease invasions are possible; in
12 stochastic models (or a stochastic world), satisfying these criteria does not guarantee that
13 invasion will occur. The misclassification rate summarizes how well these thresholds work
14 when used to predict invasion in a single instance of the disease.

15 We generated simulation data for the classification tree analyses using a range of
16 parameter values chosen to reflect a diversity of disease/host systems. The length of the time
17 step in the model is arbitrary, but with a time step of one day in mind the average infectious
18 periods, $1/\gamma$, ranged from 10 days to 2.7 years ($\gamma = [0.001 - 0.1]$). Group sizes were relatively
19 small ($n = [3 - 300]$), and rates of movement between groups ranged from once every ten days to
20 once (or less) in a lifetime ($\mu = [0.0001 - 0.1]$). The theoretical R_0 (as described in Section 2a)
21 ranged from 0 to 19, while the probability of losing immunity (ρ) or dying (δ) ranged from
22 0.0001 to 0.1. All parameters were sampled on a log scale to emphasize low parameter values
23 where the disease is more likely to be near the invasion threshold. We simulated each model

1 with 6000 different parameter sets and ran each until the disease went extinct or every group of
2 the metapopulation had been infected.

3 Because the model was stochastic, we conducted many runs of each parameter set for
4 most analyses to determine average behaviour. For the classification tree analysis, however, we
5 conducted only one run of each parameter set. We chose this approach to highlight the binary
6 and stochastic nature of the invasion process: for real disease outbreaks, it is very rare to have
7 sufficient replicates of an invasion process to estimate the probability of success. Rather, we
8 were interested in the accuracy of different predictors in the stochastic context of single
9 outbreaks. This strategy also allowed us to sample the parameter space more intensively since
10 we ran each parameter set only once. Classification trees based on half as many runs were
11 identical in structure and similar in threshold values to those presented, so we feel confident that
12 this sampling approach was sufficient to yield robust results. All model simulations were run in
13 MATLAB 7.2 (Mathworks, Inc. 2006), which called spatial models written in C. Classification
14 tree analyses were conducted in R using the Rpart package (R Core Development Team 2005;
15 Therneau & Atkinson 2005).

16

17 **3. RESULTS**

18 Successful invasion of a disease into a host metapopulation is determined by many
19 factors in addition to the necessary, but not sufficient, threshold of $R_0 > 1$. As in our earlier study
20 (Cross et al 2005b), we find that the likelihood of a pandemic exhibits a clear threshold in the
21 ratio of movement rate to recovery rate (corresponding to the expected number of between-group
22 movements during each individual's infectious period). However, the location of this threshold
23 depends upon the recruitment of new susceptibles to the population (ρ/γ in the SIRS model and
24 δ/γ in the SIR_BD model), whereby faster recruitment of susceptibles results in lower movement

1 thresholds because the disease persists longer in each group (Fig. 1, top row). When β is fixed
2 for the SIR_BD model, the probability of a pandemic is influenced by δ via its effect upon R_0 ,
3 but δ does not alter the movement threshold (Fig. 1, second column). Results are generally
4 similar between the two model structures (SIRS and SIR_BD) when β is scaled so that R_0 values
5 are equal between the models (Fig. 1, first and third columns). The SIRS and SIR_BD models
6 also yield similar results for the classification tree analyses. Thus, we present only the SIRS
7 model results, but provide the SIR_BD model results in the supplementary material.

8 Inspection of Fig. 1 illustrates that \hat{R}_0 is not a reliable predictor of pandemics when group
9 sizes are small and movement between groups is limited, regardless of susceptible replenishment
10 rate. In many cases $\hat{R}_0 > 1$ but the disease invasion fails because movement among groups is too
11 infrequent compared to the infectious period of the disease (Cross et al. 2005). The quantity \hat{R}_* ,
12 on the other hand, is strongly associated with successful disease invasions across the
13 metapopulation, for all levels of susceptible recruitment (Fig. 1). Note in Fig. 1 that \hat{R}_0 is less
14 than R_0 (i.e. $\beta/(\gamma+\delta)$ or β/γ), primarily due to susceptible depletion effects that becomes
15 important in small groups. In the first and third columns of Fig. 1, R_0 predicts that the index case
16 will infect five others, on average, but the realized number of infections (\hat{R}_0) is lower owing to
17 competition among infectors for the limited pool of susceptibles. Depletion of the susceptible
18 pool also affects \hat{R}_* . When μ/γ is small, movement among groups is the limiting factor for \hat{R}_* ,
19 and \hat{R}_* increases with μ/γ (Fig. 1). As μ/γ approaches 10, however, \hat{R}_* declines due to
20 competition among groups to infect other groups.

21 Although R_* may not be analytically tractable, we can consider its constituent parts. The
22 probability that a disease propagates through a structured population depends upon at least two
23 factors: the frequency of between-group movements and the total duration that the disease

1 persists within a given group. The total infectious time (i.e. the sum of infectious host days in a
2 single isolated group) increases with group size and with susceptible recruitment (Fig. 2a). If
3 immune individuals are replaced by susceptibles sufficiently quickly, the disease can become
4 endemic even in small groups. In Figure 2, the average infectious period per individual ($1/\gamma$) is
5 100 time steps. When the per capita infectious time is 1000 time steps, each individual has been
6 infected 10 times on average, which we use as an indication that the disease is endemic within a
7 single group (though note that the choice of 10 infections is somewhat arbitrary; Fig. 2b).

8 The total infectious time within a group determines the threshold movement rate for a
9 pandemic. For example, when $n = 10$ and ρ/γ is low (say, 10^{-3}), the total infectious time is
10 roughly 800 time steps (Fig. 2a). In order for the expected number of between-group movements
11 of infectious individuals to exceed one, the movement probability per time step for each
12 individual (μ) must exceed $1/800$, or 0.00125. When the recovery rate (γ) is 0.01, a threshold of
13 $\mu/\gamma > 0.125$ is predicted, exactly as seen in Figure 1 for an SIRS model with low ρ/γ . Similarly,
14 when $n = 10$ and ρ/γ is high (say, 10), the total infectious time is $\sim 10^5$ time steps, so the predicted
15 threshold for μ/γ is $10^{-5}/0.01 = 10^{-3}$, again corroborated by Figure 1.

16 The classification tree analysis (Fig. 3a) indicates that disease-host combinations must
17 satisfy several criteria for a pandemic to be likely. First, the disease must be able to spread
18 successfully within the initially-infected group. Traditionally this is assessed using the
19 theoretical threshold $R_0 > 1$, above which invasion occurs with non-zero probability (Diekmann &
20 Heesterbeek 2000). In the statistical context, however, a higher threshold of $R_0 \geq 2$ minimizes the
21 amount of misclassification error, although it increases the probability of a false-negative result
22 where disease extinction is predicted but the disease actually invades (Fig. 3a, Table 1). If R_0 is
23 sufficiently high to favour within-group transmission, then the disease still needs to propagate
24 between groups, a process that depends upon group size, movement and the length of the

1 infectious period (yielding a threshold of $\mu n/\gamma \geq 2.7$). Similar to R_0 , the classification threshold
2 for $\mu n/\gamma$ exceeds the criterion $\mu n/\gamma > 1$ that we proposed in an earlier simulation study (Cross et
3 al. 2005). If the relative amount of movement between groups is low, then the disease may still
4 be able to invade the entire metapopulation if the recruitment of new susceptibles (ρn or δn)
5 scaled by the recovery rate (γ) is high. In the case we present, the classification threshold for
6 $\rho n/\gamma$ is ~ 7.2 ; this can be considered a loose statistical criterion for endemicity, above which the
7 disease persists long enough in each group that even infrequent between-group movements are
8 sufficient to maintain the disease.

9 The specific thresholds presented here are likely to depend upon the model structure and
10 parameter ranges used. Similar to previous work (Cross et al. 2005), we also simulated the
11 disease model using a ‘non-spatial’ array of groups where individuals could move to any other
12 group in one step (Supplementary Material). We found that the statistical threshold of $\mu n/\gamma$ in
13 the classification tree was lower (1.8 compared to 2.7) for the non-spatial array compared to
14 nearest-neighbor movement model, but the structure of the classification tree was the same
15 (compare Fig. 3a and Fig. E1). In addition, we simulated the SIRS model with only one group
16 and conducted a classification analysis on whether greater than 90% of that group was ever
17 infected. The best statistical threshold for disease invasion was $R_0 \geq 2.4$, which is similar to the
18 criteria for the multi-group metapopulation model.

19 To investigate the effect of different parameters on the classification tree analysis, we
20 constructed new classification trees using subsets of the data corresponding to particular ranges
21 of certain parameter values. The relative amount of error explained by different variables
22 depended upon the parameter space used, but the overall classification tree structure and
23 threshold values were very similar. For example, in all 6000 runs of the SIRS model the disease
24 invaded the metapopulation in 41.1% of the simulations. This percentage represents the total

1 amount of error associated with a classification tree with no nodes. Inclusion of the first node,
2 $R_0 \geq 2$, decreases the error rate to 25%, for a relative error rate of 0.62 (*i.e.* 0.25/0.41). Adding
3 the second node, $\mu n/\gamma \geq 2.7$, reduces the relative error rate to 0.38. The length of each branch of
4 the classification tree is proportional to the reduction in prediction error associated with that node
5 (Fig. 3). When we analyzed only the subset of the data where group sizes were greater than 100,
6 the first node alone, $R_0 \geq 1.9$, became a more important predictor, reducing the error rate from
7 0.46 to 0.16 (relative error = 0.34 compared to 0.616 with all group sizes) and the second node,
8 $\mu n/\gamma \geq 3.8$, only led to a marginal improvement (Fig. 3b). Thus, loosely stated, the predictive
9 ability of R_0 increased with larger group sizes while the importance of movement decreased.
10 Note, however, that the threshold values remained similar (Fig. 3a,b). When we analyzed the
11 subset of the dataset with shorter infectious period ($\gamma > 0.01$) the predictive power of R_0
12 decreased while the importance of $\mu n/\gamma$ increased (data not shown). Thus, for acute diseases
13 movement becomes a more important predictor of disease invasion (Cross et al. 2004; Cross et
14 al. 2005).

15 The theoretical threshold of $R_0 > 1$ determines when a disease invasion is possible in an
16 infinite population. In a large, but finite, population this threshold holds to close approximation
17 (Lloyd-Smith et al. 2005a), which makes it unsurprising that $R_0 > 1$ resulted in no false-negatives
18 in our simulations. However, at least for the parameter ranges we explored, the disease did not
19 invade in 35% of the simulations where R_0 was greater than one. These invasion failures
20 correspond to stochastic extinctions of the disease, but are counted as false-positive predictions
21 when $R_0 > 1$ is interpreted as a predictor. Our previous rule of thumb, $\mu n/\gamma > 1$, also resulted in
22 few false-negatives (2%) but many false-positive (42%). The false-positive rate is reduced when
23 using $R_0 > 1$ and $\mu n/\gamma > 1$ in combination, but these rules still do not account for the recruitment of
24 new susceptible individuals (Table 1).

1 All the classification trees we analyzed yielded lower misclassification rates on test data
2 (13-18%) than either $R_0 > 1$ or $\mu n / \gamma > 1$ (24-44%; Table 1). The ‘best’ classification tree, as
3 determined by the ‘1-SE rule’, was only marginally better at predicting disease invasion than the
4 reduced tree shown in Fig. 3a (13% vs. 14%, Table 1). The classification tree based upon the
5 raw model parameters β , γ , μ , and n did not perform quite as well as those based on aggregate
6 parameters β/γ , $\mu n/\gamma$ and $\rho n/\gamma$ (19% vs. 14%, Table 1). Threshold criteria based on the
7 emergent quantities ν and ν_* produced the lowest misclassification rate in the case of ν_* , which
8 was twice as good as that of ν (10% vs. 20%, Table 1). Our counting rules for ν_* did not account
9 for the possibility that the index group could lose the infection (all infected members moving
10 out) and then become re-infected (those same infected members moving back in, without having
11 transmitted in their new group) before finally going on to spread the infection. As a result, a few
12 simulations led to invasions when $\nu_* = 0$, which is at odds with the theoretical definition on ν_* ,
13 but this low probability event (33 out of 6000 simulations) does not change our overall
14 conclusions (Fig. 1, Table 1)

15 The analysis of individual reproductive numbers (Fig. 4) illustrates the strong influence
16 of population structure on SSEs. Owing to the constant recovery probability assumed in our
17 model, there is substantial individual variation in infectious periods. In a single large population,
18 this leads an overdispersed distribution of ν and numerous SSEs (31 SSEs out of 500
19 simulations). Compared to an expected 5 SSEs out of 500 individuals for a homogeneous
20 population, by our definition of an SSE, this yields a superspreading load of 31/5 or 6.2. In a
21 metapopulation of small populations ($n = 10$), the frequency of SSEs depends upon the
22 movement of hosts among groups. When movement rates are high ($\mu/\gamma = 10$), there were 56
23 SSEs for a superspreading load of 11, whereas when μ/γ equaled 0.001 there were 12 SSEs,
24 representing an SSL of just 2.4. The recruitment rate of new susceptibles did not have

1 significant impact upon SSEs (data not shown).

2

3 **4. DISCUSSION**

4 In socially or spatially structured host populations, $R_0 > 1$ is a necessary but not sufficient
5 condition for a pandemic. As R_0 increases beyond one the probability of disease invading the
6 initially-infected host group increases; but additional criteria are important to determining the
7 probability that the disease spreads to other groups. Disease transmission among groups depends
8 on the transmission rate among individuals (β), the frequency of individual host movement (μ)
9 and the duration of time (measured cumulatively over all infected hosts) the disease persists
10 within each group. Within-group persistence times increase due to longer individual infectious
11 periods ($1/\gamma$), greater group sizes (n), or faster replenishment of the susceptible pool (Bartlett
12 1957; Bjornstad et al. 2002; Grenfell et al. 2002; Lloyd-Smith et al. 2005a). To synthesize, the
13 disease is increasingly likely to invade the entire population for increasing $R_0 > 1$ and $\mu n/\gamma > 1$;
14 when movement is infrequent relative to host recovery ($\mu n/\gamma < 1$), a pandemic requires that the
15 recruitment of susceptible individuals is sufficiently fast to allow the disease to persist
16 endemically in infected groups (Figs. 1 and 3).

17 To our knowledge classification and regression tree analyses have not been used to
18 understand disease invasions, yet we found that the method was naturally suited to analyzing
19 simulation results and illustrating the hierarchical nature of disease invasion criteria. After
20 experimenting with many combinations of predictor variables (Supplementary Material), we
21 focused on a set of aggregate parameters that were most informative, hence resulting in small
22 trees, and corresponded to relevant biological processes: within-group transmission, R_0 (β/γ in
23 SIRS or $\beta/(\gamma+\delta)$ in SIR_BD); movement, $\mu n/\gamma$; and recruitment of new susceptibles, $\rho n/\gamma$ and
24 $\delta n/\gamma$. The classification tree analyses corroborated our previous rule of thumb (Cross et al. 2005)

1 that when transmission and recovery processes are fast relative to the recruitment of new
2 susceptibles, $\mu n/\gamma$ must exceed one for a pandemic to occur (Fig. 3a). Our expanded models,
3 however, revealed that the effects of low movement rates can be compensated for by faster
4 susceptible recruitment (e.g. $\rho n/\gamma > 7$, Fig. 3a).

5 Theoretical ecologists often search for thresholds or bifurcation points where system
6 behaviour qualitatively changes. The threshold $R_0 > 1$ demarcates when a disease outbreak is
7 possible, but as a predictor will lead to false-positive when the disease is predicted to invade but
8 goes extinct due to initial stochastic events. Thus $R_0 > 1$ is a conservative threshold for predicting
9 disease outbreaks and circumstances exist where more accurate (but less conservative)
10 predictions of invasion are useful. In 35% of simulations we conducted, the $R_0 > 1$ criterion was
11 satisfied but the disease failed to invade (Table 1). The combined threshold of $R_0 > 1$ and $\mu n/\gamma > 1$
12 resulted in fewer misclassifications (24%) but the classification tree criteria were more reliable,
13 misclassifying only $14 \pm 0.5\%$ (SD) of all simulations that were not used in the tree construction
14 (Table 1). We emphasize, though, that all the ‘thresholds’ we describe are necessarily fuzzy due
15 to the stochastic nature of disease invasion (Lloyd-Smith et al. 2005a).

16
17 All the criteria we applied, with the exception of $\nu_* > 1$, resulted in more false-positives
18 than false-negatives due to the high probability of stochastic extinction in the early generations
19 of disease invasion. The ν_* metric was the best predictor because it includes information on
20 initial stochastic events as well as the movement of infectious individuals among groups.
21 Predictions based on real empirical data are likely to suffer greater misclassification error rates
22 than the simulated data we present due to process-based variation and sampling error. Despite
23 these difficulties, our results emphasize the importance of understanding host movement and
24 those processes that allow diseases to persist for longer in spatially or socially structured host
25 populations.

1 Superspreading events (SSEs) result from heterogeneities in host, environment and
2 parasite factors (Lloyd-Smith et al. 2005b). Our analysis focuses on the interaction between
3 heterogeneity in the host factor of infectious period and in the environmental factor of contact
4 with susceptible individuals. In our simulations, all infectious individuals had constant and
5 identical probabilities per time step of recovering from disease, as well as moving between
6 groups, resulting in geometric distributions for the duration of infectiousness and the number of
7 groups visited while infectious. The heterogeneities embodied by these geometrically-
8 distributed quantities create the conditions necessary for SSEs; that is, they lead to distributions
9 of individual reproductive numbers that are overdispersed relative to the Poisson distribution
10 predicted when all infectious individuals (and their environments) are identical. Given these
11 individual heterogeneities, the frequency of SSEs may be constrained or facilitated by the
12 population structure where the individual resides. In a large or panmictic population,
13 transmission is not constrained by the supply of susceptible individuals. In contrast, when
14 groups are small and movement is infrequent, the number of potential contacts is limited and the
15 opportunity for SSEs is reduced even for individuals with extraordinarily long infectious periods.
16 The same qualitative effect would arise for individual heterogeneity in transmission rates, as
17 access to susceptibles is a prerequisite for transmission. The potential for superspreading in
18 structured populations would be amplified if positive correlations existed between movement
19 rates (and hence access to more susceptibles) and high transmissibility or slow recovery. Further
20 subtleties may arise if movement itself is linked to transmission (as in SSEs aboard airliners) or
21 increased risk of death (as in some wildlife systems).

22 The utility of simple, within-group calculations of R_0 as a predictive measure of disease
23 invasion is limited in systems where transmission between groups may be the primary factor
24 regulating the probability of a pandemic. Examples include many wildlife populations

1 (Woolhouse et al. 2001), livestock based on small holdings (Keeling et al. 2001; Woolhouse et
2 al. 2005), and human populations with small, weakly connected groups of susceptible individuals
3 (Feikin et al. 2000; Salmon et al. 1999). While further research should aim to advance analytic
4 theory, classification trees provide an effective means of connecting real-world, measurable
5 variables to the likelihood of invasion, particularly in structured populations where system
6 dynamics are governed by hierarchy of contributing factors. Our analyses have focused on a
7 relatively idealized system of equal group sizes and simplistic movement rules. Future work
8 should aim to extend our findings to more realistic, heterogeneous settings, and to link the ideas
9 presented here with empirical evidence from the field.

10

11 **ACKNOWLEDGEMENTS**

12 This research was funded by the NSF-NIH Ecology of Infectious Disease Grant DEB-0090323
13 to Wayne Getz, NIH-NIDA Grant R01-DA10135, the Center for Infectious Disease Dynamics
14 (CIDD) Postdoctoral Fellowship (JLS), a James S. McDonnell Foundation 21st Century Science
15 Initiative Award (WMG), and the U.S. Geological Survey (PCC). Many thanks to A. Shrag for
16 her help with R and classification trees. Previous versions of this manuscript were improved by
17 the comments of R. Plowright and several anonymous reviewers.

18

19 **REFERENCES**

20 Altizer, S., Nunn, C. L., Thrall, P. H., Gittleman, J. L., Antonovics, J., Cunningham, A. A.,
21 Dobson, A. P., Ezenwa, V., Jones, K. E., Pedersen, A. B., Poss, M. & Pulliam, J. R. C.
22 2003 Social organization and parasite risk in mammals: Integrating theory and empirical
23 studies. *Annual Review of Ecology Evolution and Systematics* **34**, 517-547.

- 1 Anderson, R. M. & May, R. M. 1991 *Infectious Diseases of Humans: Dynamics and Control*.
2 Oxford: Oxford University Press.
- 3 Andersson, H. 1997 Epidemics in a population with social structures. *Mathematical Biosciences*
4 **140**, 79-84.
- 5 Andersson, H. & Britton, T. 1998 Heterogeneity in epidemic models and its effect on the spread
6 of infection. *Journal of Applied Probability* **35**, 651-661.
- 7 Ball, F., Mollison, D. & Scalia-Tomba, G. 1997 Epidemics with two levels of mixing. *The*
8 *Annals of Applied Probability* **7**, 46-89.
- 9 Bartlett, M. S. 1957 Measles periodicity and community size. *Journal of the Royal Statistical*
10 *Society* **120**, 48-71.
- 11 Becker, N. G. & Dietz, K. 1995 The effect of household distribution on transmission and control
12 of highly infectious diseases. *Mathematical Biosciences* **127**, 207-219.
- 13 Becker, N. G. & Starczak, D. N. 1997 Optimal vaccination strategies for a community of
14 households. *Mathematical Biosciences* **139**, 117-132.
- 15 Begg, C. B. 1986 Statistical-Methods in Medical Diagnosis. *Crc Critical Reviews in Medical*
16 *Informatics* **1**, 1-22.
- 17 Bjornstad, O. N., Finkenstadt, B. F. & Grenfell, B. T. 2002 Dynamics of measles epidemics:
18 Estimating scaling of transmission rates using a Time series SIR model. *Ecological*
19 *Monographs* **72**, 169-184.
- 20 Breiman, L., Freidman, J. H., Olshen, R. A. & Stone, C. J. 1984 *Classification and Regression*
21 *Trees*. The Wadsworth statistics/probability series. New York: Chapman & Hall.
- 22 Brose, U., Berlow, E. L. & Martinez, N. D. 2005 Scaling up keystone effects from simple to
23 complex ecological networks. *Ecology Letters* **8**, 1317-1325.

- 1 Cross, P. C., Lloyd-Smith, J. O., Bowers, J., Hay, C., Hofmeyr, M. & Getz, W. M. 2004
2 Integrating association data and disease dynamics in a social ungulate: bovine
3 tuberculosis in African buffalo in the Kruger National Park. *Annales Zoologici Fennici*
4 **41**, 879-892.
- 5 Cross, P. C., Lloyd-Smith, J. O., Johnson, P. L. F. & Getz, W. M. 2005 Duelling timescales of
6 host mixing and disease recovery determine disease invasion in structured populations.
7 *Ecology Letters* **8**, 587-595.
- 8 De'ath, G. & Fabricius, K. E. 2000 Classification and regression trees: A powerful yet simple
9 technique for ecological data analysis. *Ecology* **81**, 3178-3192.
- 10 Diekmann, O. & Heesterbeek, J. A. P. 2000 *Mathematical Epidemiology of Infectious Diseases:*
11 *Model Building, Analysis and Interpretation*. Wiley Series in Mathematical and
12 Computational Biology. Chichester, England: John Wiley & Sons Ltd.
- 13 Diekmann, O., Heesterbeek, J. A. P. & Metz, J. A. J. 1990 On the Definition and the
14 Computation of the Basic Reproduction Ratio R_0 in Models for Infectious-Diseases in
15 Heterogeneous Populations. *Journal of Mathematical Biology* **28**, 365-382.
- 16 Dobson, A. & Foufopoulos, J. 2001 Emerging infectious pathogens of wildlife. *Philosophical*
17 *Transactions of the Royal Society of London B Biological Sciences* **356**, 1001-1012.
- 18 Feikin, D. R., Lezotte, D. C., Hamman, R. F., Salmon, D. A., Chen, R. T. & Hoffman, R. E.
19 2000 Individual and community risks of measles and pertussis associated with personal
20 exemptions to immunization. *Jama-Journal of the American Medical Association* **284**,
21 3145-3150.
- 22 Fulford, G. R., Roberts, M. G. & Heesterbeek, J. A. P. 2002 The metapopulation dynamics of an
23 infectious disease: Tuberculosis in possums. *Theoretical Population Biology* **61**, 15-29.

- 1 Getz, W. M. & Lloyd-Smith, J. O. 2006 Basic methods for modeling the invasion and spread of
2 contagious diseases. In *Disease Evolution: Models, Concepts, and Data Analysis*, AMS
3 vol. 71 (ed. Z. Feng, U. Dieckmann & S. A. Levin): AMS-DIMACS series.
- 4 Getz, W. M. & Pickering, J. 1983 Epidemic models: thresholds and population regulation.
5 *American Naturalist* **121**, 892-898.
- 6 Grenfell, B. T., Bjornstad, O. N. & Finkenstadt, B. F. 2002 Dynamics of measles epidemics:
7 Scaling noise, determinism, and predictability with the TSIR model. *Ecological*
8 *Monographs* **72**, 185-202.
- 9 Hagensars, T. J., Donnelly, C. A. & Ferguson, N. M. 2004 Spatial heterogeneity and the
10 persistence of infectious diseases. *Journal of Theoretical Biology* **229**, 349-359.
- 11 Heesterbeek, J. A. P. 2002 A brief history of R_0 and a recipe for its calculation. *Acta*
12 *Biotheoretica* **50**, 189-204.
- 13 Heffernan, J. M., Smith, R. J. & Wahl, L. M. 2005 Perspectives on the basic reproductive ratio.
14 *Journal of the Royal Society Interface* **2**, 281-293.
- 15 Hess, G. 1996a Disease in metapopulation models: Implications for conservation. *Ecology* **77**,
16 1617-1632.
- 17 Hess, G. R. 1996b Linking extinction to connectivity and habitat destruction in metapopulation
18 models. *American Naturalist* **148**, 226-236.
- 19 Kao, R. R., Danon, L., Green, D. M. & Kiss, I. Z. 2006 Demographic structure and pathogen
20 dynamics on the network of livestock movements in Great Britain. *Proceedings of the*
21 *Royal Society B-Biological Sciences* **273**, 1999-2007.
- 22 Karels, T. J., Bryant, A. A. & Hik, D. S. 2004 Comparison of discriminant function and
23 classification tree analyses for age classification of marmots. *Oikos* **105**, 575-587.

- 1 Keeling, M. J. 1999 The effects of local spatial structure on epidemiological invasions.
2 *Proceedings of the Royal Society of London Series B-Biological Sciences* **266**, 859-867.
- 3 Keeling, M. J. & Gilligan, C. A. 2000a Bubonic plague: a metapopulation model of a zoonosis.
4 *Proceedings of the Royal Society of London Series B-Biological Sciences* **267**, 2219-
5 2230.
- 6 Keeling, M. J. & Gilligan, C. A. 2000b Metapopulation dynamics of bubonic plague. *Nature*
7 **407**, 903-906.
- 8 Keeling, M. J. & Grenfell, B. T. 2000 Individual-based perspectives on R-0. *Journal of*
9 *Theoretical Biology* **203**, 51-61.
- 10 Keeling, M. J. & Rohani, P. 2002 Estimating spatial coupling in epidemiological systems: a
11 mechanistic approach. *Ecology Letters* **5**, 20-29.
- 12 Keeling, M. J., Woolhouse, M. E., Shaw, D. J., Matthews, L., Chase-Topping, M., Haydon, D.
13 T., Cornell, S. J., Kappey, J., Wilesmith, J. & Grenfell, B. T. 2001 Dynamics of the 2001
14 UK Foot and Mouth epidemic: stochastic dispersal in a heterogeneous landscape. *Science*
15 **294**, 813-817.
- 16 Kermack, W. O. & McKendrick, A. G. 1927 Contributions to the mathematical theory of
17 epidemics. *Proceedings of the Royal Society of Edinburgh* **115**, 700-721.
- 18 Lloyd-Smith, J. O., Cross, P. C., Briggs, C. J., Daugherty, M., Getz, W. M., Latto, J., Sanchez,
19 M. S., Smith, A. B. & Swei, A. 2005a Should we expect population thresholds for
20 wildlife disease? *Trends in Ecology & Evolution* **20**, 511-519.
- 21 Lloyd-Smith, J. O., Schreiber, S. J., Kopp, P. E. & Getz, W. M. 2005b Superspreading and the
22 effect of individual variation on disease emergence. *Nature* **438**, 355-359.
- 23 McCallum, H., Barlow, N. & Hone, J. 2001 How should pathogen transmission be modelled?
24 *Trends in Ecology and Evolution* **16**, 295-300.

- 1 Park, A. W., Gubbins, S. & Gilligan, C. A. 2001 Invasion and persistence of plant parasites in a
2 spatially structured host population. *Oikos* **94**, 162-174.
- 3 Park, A. W., Gubbins, S. & Gilligan, C. A. 2002 Extinction times for closed epidemics: the
4 effects of host spatial structure. *Ecology Letters* **5**, 747-755.
- 5 R Core Development Team. 2005 R: A Language and Environment for Statistical Computing.
6 Vienna, Austria: R Foundation for Statistical Computing.
- 7 Salmon, D. A., Haber, M., Gangarosa, E. J., Phillips, L., Smith, N. J. & Chen, R. T. 1999 Health
8 consequences of religious and philosophical exemptions from immunization laws -
9 Individual and societal risk of measles. *Jama-Journal of the American Medical*
10 *Association* **282**, 47-53.
- 11 Schinazi, R. B. 2002 On the role of social clusters in the transmission of infectious diseases.
12 *Theoretical Population Biology* **61**, 163-169.
- 13 Steadman, H. J., Silver, E., Monahan, J., Appelbaum, P. S., Robbins, P. C., Mulvey, E. P.,
14 Grisso, T., Roth, L. H. & Banks, S. 2000 A classification tree approach to the
15 development of actuarial violence risk assessment tools. *Law and Human Behavior* **24**,
16 83-100.
- 17 Swinton, J. 1998 Extinction times and phase transitions for spatially structured closed epidemics.
18 *Bulletin of Mathematical Biology* **60**, 215-230.
- 19 Therneau, T. M. & Atkinson, B. 2005 Rpart: Recursive Partitioning.
- 20 Thrall, P. H., Antonovics, J. & Dobson, A. P. 2000 Sexually transmitted diseases in polygynous
21 mating systems: prevalence and impact on reproductive success. *Proceedings of the*
22 *Royal Society of London Series B-Biological Sciences* **267**, 1555-1563.

- 1 Usio, N., Nakajima, H., Kamiyama, R., Wakana, I., Hiruta, S. & Takamura, N. 2006 Predicting
2 the distribution of invasive crayfish (*Pacifastacus leniusculus*) in a Kusiro Moor marsh
3 (Japan) using classification and regression trees. *Ecological Research* **21**, 271-277.
- 4 Watts, D. J., Muhamad, R., Medina, D. C. & Dodds, P. S. 2005 Multiscale, resurgent epidemics
5 in a hierarchical metapopulation model. *Proceedings of the National Academy of*
6 *Sciences of the United States of America* **102**, 11157-11162.
- 7 Woolhouse, M. E. J., Shaw, D. J., Matthews, L., Liu, W. C., Mellor, D. J. & Thomas, M. R. 2005
8 Epidemiological implications of the contact network structure for cattle farms and the 20-
9 80 rule. *Biology Letters* **1**, 350-352.
- 10 Woolhouse, M. E. J., Taylor, L. H. & Haydon, D. T. 2001 Population biology of multihost
11 pathogens. *Science* **292**, 1109-1112.
- 12
- 13

Figure 1. Percentage of the metapopulation infected, \hat{R}_* , and \hat{R}_0 all depend upon host movement (μ), disease recovery (γ), and replenishment of the susceptible pool (indexed by ρ or δ for the SIRS and SIR_BD models, respectively). Each point shows the mean of 200 simulations with 10 individuals in each group and a recovery probability (γ) of 0.01. In the first and third columns $R_0=5$; in the second column R_0 varies from 0.45 to 5 depending on the value of δ .

Figure 2. The total infectious time (sum of infectious host days) and per capita infectious time in a single group of individuals. Infectious time increases due to the flow of new susceptibles, which is a function of group size (n) and the probability that a recovered individual returns to susceptibility (ρ). Above the dotted line individuals are infected more than 10 times, on average, indicating that the disease is endemic within the local group. Each point is the mean of 100 simulations of the SIRS model with a recovery probability (γ) of 0.01 and $R_0=5$. In the endemic range, simulations were stopped when infectious time was limited by the arbitrary maximum duration of the simulation.

Figure 3. Classification trees predicting the invasion or extinction of a disease introduced into a metapopulation using the SIRS model using all the simulation data (A) and only those runs with group sizes greater 100 (B). Threshold criteria are labeled above each node of the tree, and instances that satisfy the criteria are split off to the left. Labels underneath the terminal leaves indicate the number of simulations (out of 6000 for figure A and 1956 for figure B) resulting in invasions and extinctions, respectively, and in text the majority outcome for that set of classification rules.

Figure 4. Histograms of ν , the individual reproductive number (*i.e.* the number of individuals infected by the initial case), for different movement probabilities (μ) scaled by the probability of disease recovery ($\gamma = 0.01$) using the SIRS model. Mean values of ν are indicated by diamonds. Superspreaders are defined as those individuals beyond the 99th percentile of the Poisson distribution (vertical lines) with the same mean. Each parameter set was simulated 500 times with $\rho = 0.00001$ and $\beta = 0.05$ on an 11 x 11 toroidal array with 10 individuals in each group, with the exception of the top row which was one group of 1210 individuals.

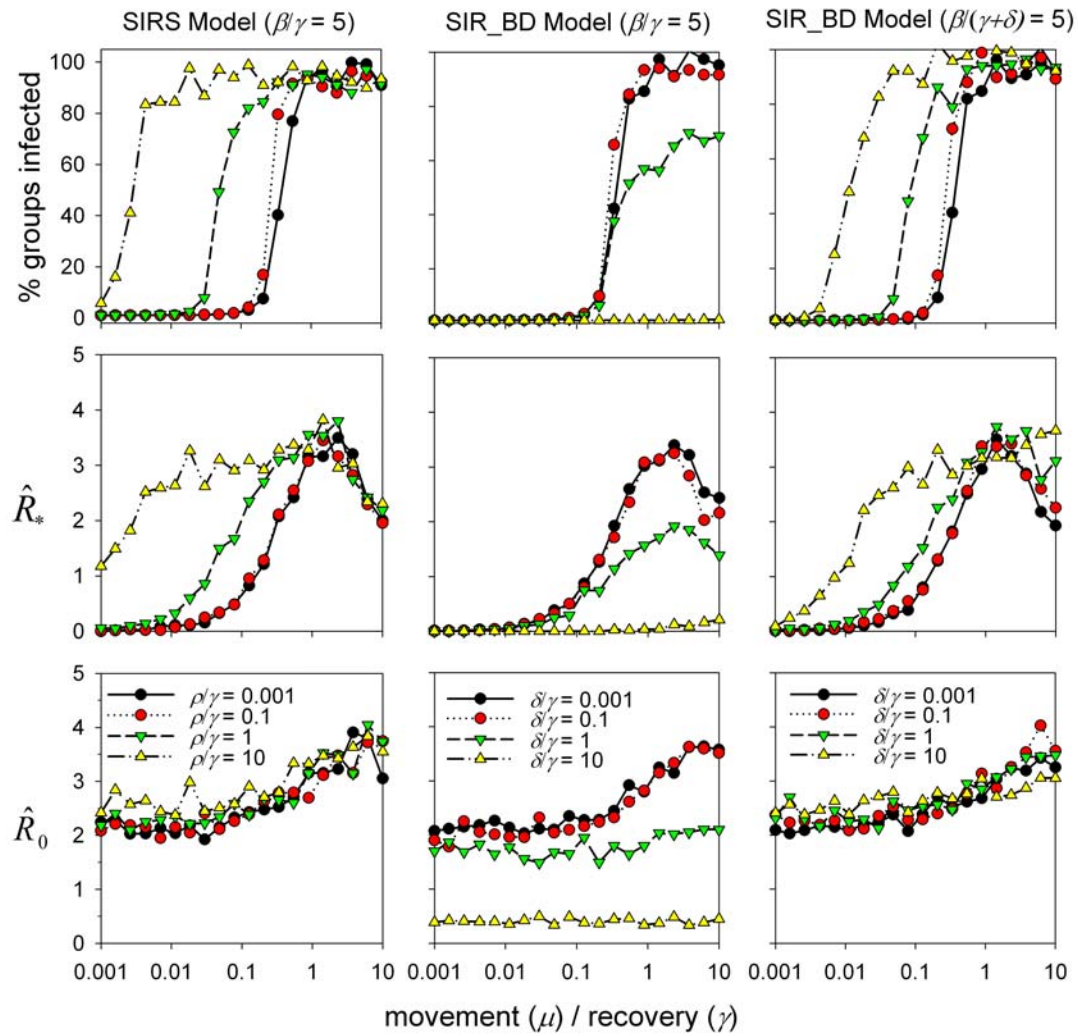
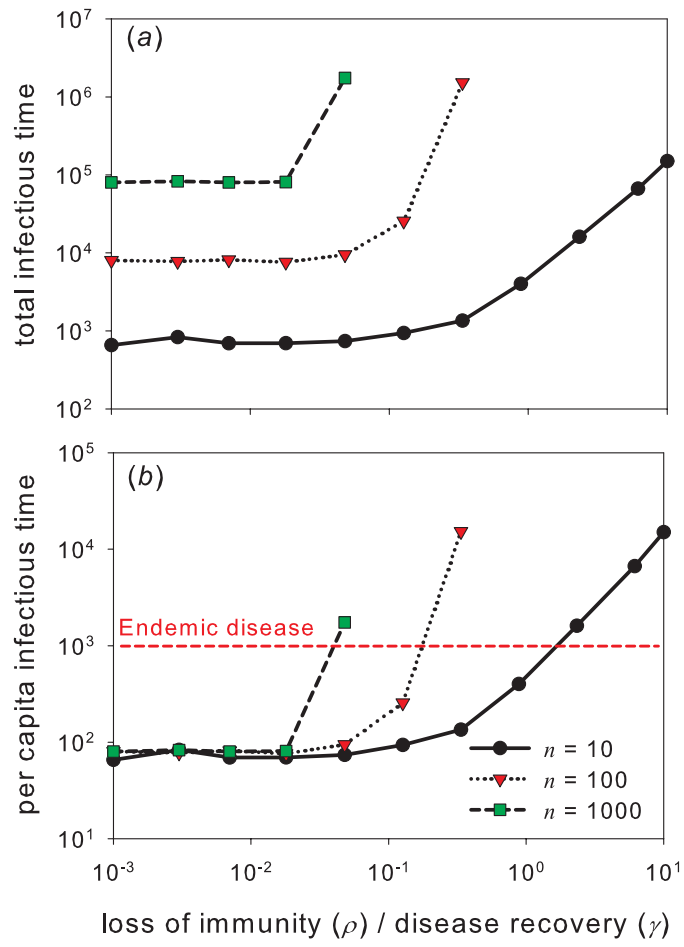
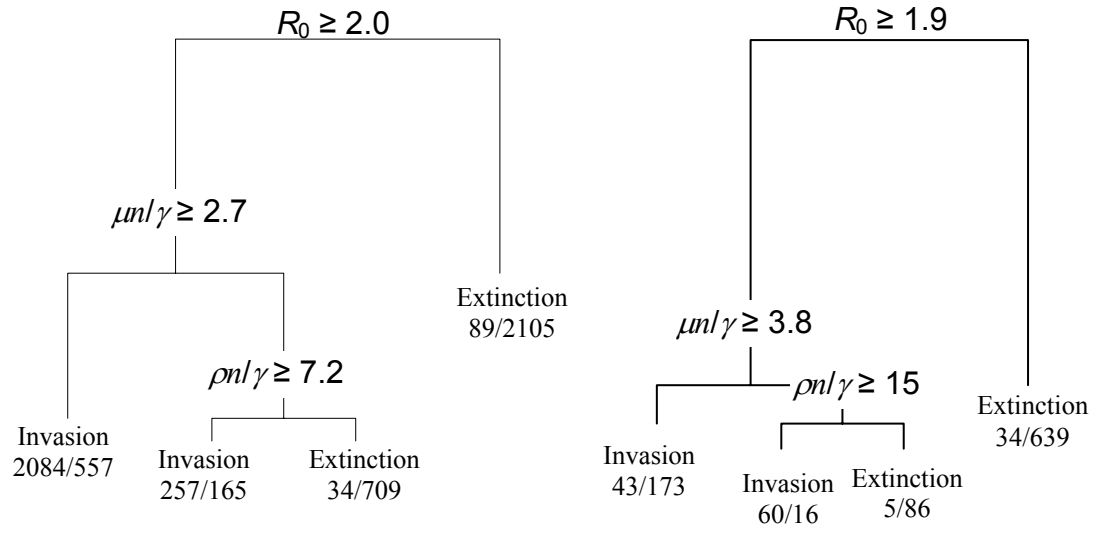


Figure 1.



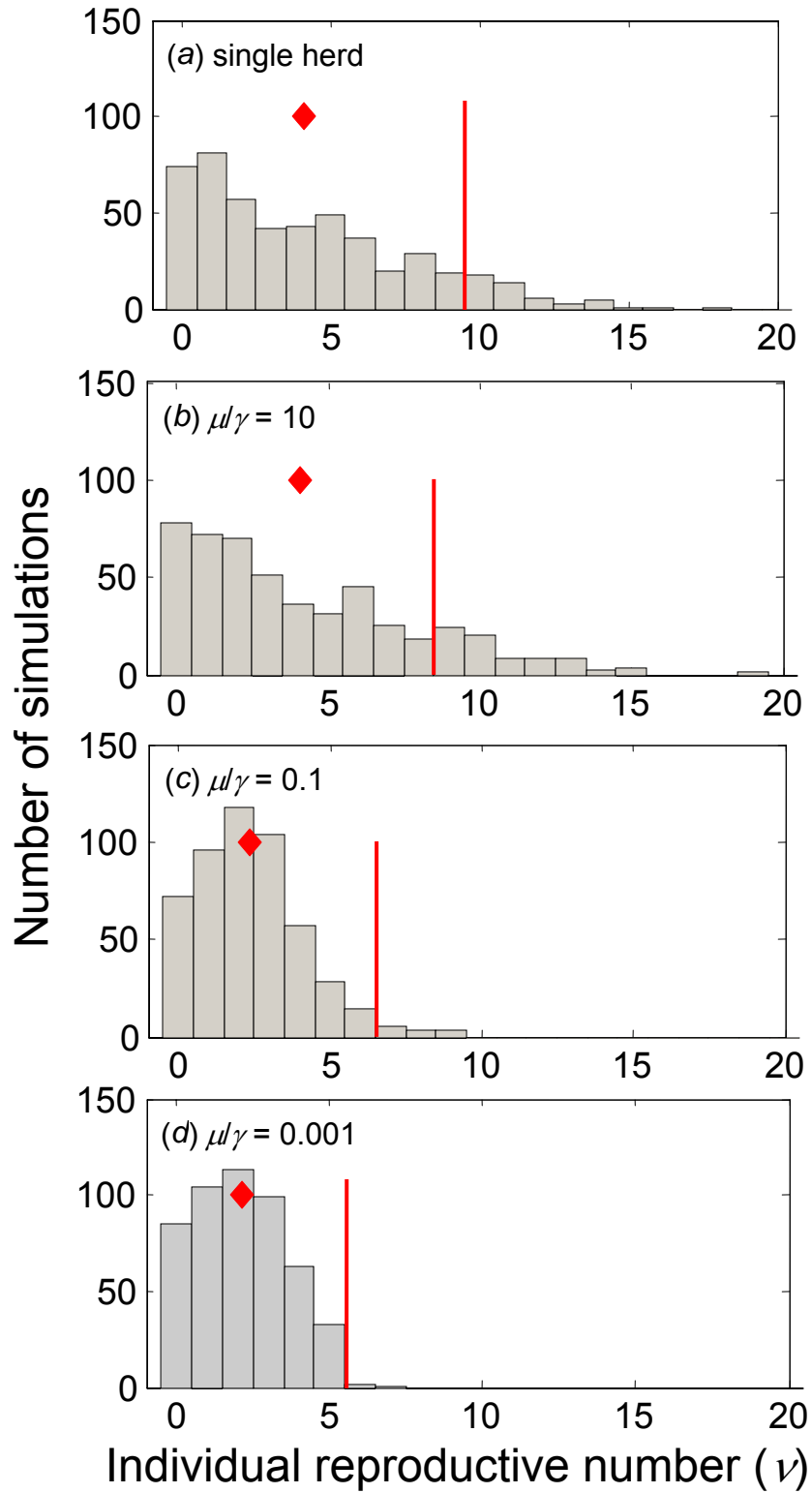
1
2
3
4 Figure. 2

1



2
3
4
5

Figure 3.



1
2
3 Figure 4.

1

Table 1. The proportion of SIRS model simulations where the disease invades the metapopulation and whether that invasion was predicted by theoretical thresholds or the classification tree analyses.

Rules for invasion	Correctly predicted invasions	Correctly predicted extinctions	False-positive ¹	False-negative ²	Total misclassified	Cross-validated misclassification ³	SD ³
$R_0 > 1$	0.411	0.240	0.353	0	0.353	--	--
$\mu n/\gamma > 1$	0.390	0.174	0.416	0.020	0.436	--	--
$R_0 > 1$ and $\mu n/\gamma > 1$	0.390	0.366	0.224	0.020	0.244	--	--
Best classification tree ⁴	0.383	0.485	0.104	0.028	0.132	0.141	0.0045
reduced classification tree ⁵	0.390	0.469	0.120	0.021	0.141	0.144	0.0045
raw parameter tree ⁶	0.327	0.485	0.105	0.084	0.188	0.205	0.0052
$\nu > 1$, <i>emergent</i> ⁷	0.355	0.444	0.145	0.056	0.201	--	--
$\nu_* > 1$, <i>emergent</i> ⁷	0.352	0.551	0.039	0.059	0.097	--	--

¹ Rules predicted invasions when the disease actually went extinct.

² Rules predicted extinctions when the disease actually invaded.

³ Average and standard deviation of error rates on test data not used in the construction of the classification tree using 10-fold cross-validation.

⁴ Using aggregate parameters not including ν and ν_* . The best tree had four nodes, further subdividing the 257/165 branch of the reduced tree (Fig. 3a), but this did little to improve accuracy. See Figure E2.

⁵ Using the aggregate parameters not including ν and ν_* . See Figure 3.

⁶ Using raw parameters not including ν and ν_* . See Figure E1.

⁷ ν and ν_* are considered emergent because they can only be estimated after the epidemic has begun and thus have an advantage over other metrics included in the table

2