



Duncan, A.J., Gunn, G.J., Umstatter, C. and Humphry, R.W.  
(2014) Replicating disease spread in empirical cattle networks by  
adjusting the probability of infection in random networks.  
*Theoretical Population Biology*, 98. pp. 11-18. ISSN 0040-5809.

**NOTICE: this is the author's version of a work that was accepted  
for publication in *Theoretical Population Biology*. Changes  
resulting from the publishing process, such as peer review, editing,  
corrections, structural formatting, and other quality control  
mechanisms may not be reflected in this document. Changes may  
have been made to this work since it was submitted for publication.  
A definitive version was subsequently published in *Theoretical  
Population Biology*, 98, Decemeber 2014, DOI:  
10.1016/j.tpb.2014.08.004**

Copyright © 2014 Elsevier B.V. All rights reserved.

<http://hdl.handle.net/11262/10508>

Deposited on: 28 October 2014

# Replicating disease spread in empirical cattle networks by adjusting the probability of infection in random networks.

A. J. Duncan<sup>a,\*</sup>, G. J. Gunn<sup>b</sup>, C. Umstatter<sup>c</sup>, R. W. Humphry<sup>b</sup>

<sup>a</sup>*Inverness College UHI, Longman Campus, 3 Longman Road, Longman South, Inverness, IV1 1SA*

<sup>b</sup>*Epidemiology Research Unit, SRUC (Scotland's Rural College), Drummondhill, Stratherrick Road, Inverness, IV2 4JZ*

<sup>c</sup>*Agroscope, Institute for Sustainability Sciences ISS, Tänikon 1, 8356 Ettenhausen, Switzerland*

---

\*Corresponding author

*Email address:* `andrew.duncan.ic@uhi.ac.uk` (A. J. Duncan)

---

**Abstract**

Comparisons between mass-action or “random” network models and empirical networks have produced mixed results. Here we seek to discover whether a simulated disease spread through randomly constructed networks can be coerced to model the spread in empirical networks by altering a single disease parameter – the probability of infection. A stochastic model for disease spread through herds of cattle is utilised to model the passage of an SEIR (susceptible–latent–infected–resistant) through five networks. The first network is an empirical network of recorded contacts, from four datasets available, and the other four networks are constructed from randomly distributed contacts based on increasing amounts of information from the recorded network. A numerical study on adjusting the value of the probability of infection was conducted for the four random network models. We found that relative percentage reductions in the probability of infection, between 5.6% and 39.4% in the random network models, produced results that most closely mirrored the results from the empirical contact networks. In all cases tested, to reduce the differences between the two models, required a reduction in the probability of infection in the random network.

*Keywords:* Network; Mass-action; Disease; Recorded contacts; SEIR  
simulation

---

## 1. Introduction

The assumption of random interactions, or mass-action mixing, is a method widely used in the modelling of disease (Anderson and May, 1991; Brauer et al., 2000; De Jong et al., 1995). With cheaper and easier methods of data capture now available to record contact networks (Craft and Caillaud, 2001) homogeneously mixed networks or “random networks” have been tested against the recorded contact networks with varying results (Duncan et al., 2012; Hamede et al., 2012; Kleinlützum et al., 2013; Salathé et al., 2010). In this publication we seek to discover whether a simple model of disease spread, based on the principles of homogeneous mixing, can approximate a recorded network if the probability of infection is suitably adjusted. If this is possible, we will also investigate: whether the simplicity of the model affects the closeness of fit to the recorded network; whether there is consistency in the adjustment of the probability of infection across a variety of random network models and whether there is a relationship between the network properties, through values of network metrics, and the adjustment to the probability of infection.

Results from comparisons of simulated disease spread on random and structured network, whether recorded, empirically derived (i.e. extrapolated from empirical data) or theoretically constructed, have been mixed. Some studies have found random networks to be a suitable substitute for structured network models (Bouma et al., 1995; Dobson and Meagher, 1996; Shirley and Rushton, 2005a) whilst others have found it inadequate (Barlow, 2000; D’ Amico et al., 1996; Hamede et al., 2012; Porphyre et al., 2008; Shirley and Rushton, 2005b). For

25 inter-herd contact networks, rather than the intra-herd networks discussed  
26 herein, it has been shown that models should be at least based on any movement  
27 data available (Vernon and Keeling, 2009). The modification of the transmission  
28 rate of disease on a random network model has been shown to provide a good  
29 representation of the results from theoretically constructed networks (Keeling,  
30 2005). Simplified models of a complete contact network which take account of  
31 rewiring or preferential mixing show closer agreement than a mean-field model  
32 (random/mass-action mixing) when modelling Tasmanian devil facial tumour  
33 disease (Hamede et al., 2012) and it was found that the networks had highly  
34 connected animals, which would not be found in random networks. When  
35 modelling spread of influenza in high school students (Salathé et al., 2010), it  
36 was found that a small-world network (Watts and Strogatz, 1998) with a high  
37 proportion of repeated contacts fitted the recorded data best, but a  
38 homogeneous (random/mass-action) mixing model might be sufficient.

39

40 In our previous work (Duncan et al., 2012) we presented two stochastic models  
41 of the passage of an SEIR (susceptible-latent-infected-resistant) disease  
42 through herds of cattle. One model was based on a contact network constructed  
43 via continuously recorded interaction data from two herds of cattle, the other, a  
44 matching network constructed using the assumption of random mixing. Four  
45 recorded contact datasets were produced by attaching proximity data loggers  
46 (Drewe et al., 2012; Swain and Bishop-Hurley, 2007) to two separate herds of  
47 cattle during two separate recording periods. For each dataset the network  
48 constructed using the principles of random mixing had the same number of

49 contacts as the recorded network but these contacts were distributed randomly  
50 amongst the animals. The differences shown between the two models were that  
51 a lower proportion of simulations of the recorded network produced any disease  
52 spread when compared to those simulations of the random network and, of  
53 those that did, fewer infected animals were predicted. In this publication we  
54 seek to estimate the optimal adjustment of the probability of infection of a  
55 susceptible animal given a contact with an infectious animal so as to minimise  
56 these differences.

57

58 We constructed four types of random networks, with increasing similarities to  
59 the recorded contact network, and by adjusting the probability of infection  
60 attempted to gain the best possible approximation for the recorded network.

61 Alongside the simulation of disease, we examined the network properties via six  
62 network metrics: assortativity, average path length, closeness, clustering, degree  
63 distribution and our own metric – the number of repeated contacts. It has been  
64 shown that assortativity can be responsible for the lowering of the epidemic  
65 threshold (Molina and Stone, 2012) and clustering to lower the reproductive  
66 number  $R_0$  and increase the threshold of disease (Miller, 2009). We have  
67 already shown (Duncan et al., 2012) that the recorded networks had more  
68 repeated contacts, lower closeness and clustering but higher average path  
69 lengths. In this work we seek to relate any differences in these metrics to the  
70 adjustment in the probability of infection. Networks can now be constructed  
71 with algorithms, to have specific characteristics (Badham and Stocker, 2010a,b;  
72 Bansal et al., 2009; Håkansson et al., 2010). Therefore, if it were the case that a

73 metric value was linked to the optimal adjustment in the probability of  
74 infection, it would enable the use of specifically constructed theoretical networks  
75 in place of recorded contact networks where recording was not feasible.

## 76 **2. Materials and Methods**

### 77 *2.1. Disease*

78 The SEIR disease that is modelled through all of the network models can be  
79 described by the system of ordinary differential equations (ODEs) (Anderson  
80 and May, 1991),

$$\begin{aligned}\frac{dS}{dt} &= -\alpha\beta\frac{SI}{N}, \\ \frac{dE}{dt} &= \alpha\beta\frac{SI}{N} - \sigma E, \\ \frac{dI}{dt} &= \sigma E - \gamma I\end{aligned}\tag{1}$$

and

$$\frac{dR}{dt} = \gamma I,$$

81 with  $S + E + I + R = N$ , where  $N$  is the total (constant) population size. Each  
82 susceptible animal moves from the susceptible state (S) to the latent state (E)  
83 with rate  $\alpha\beta$  following a contact with an infectious animal, where  $\alpha$  is the  
84 probability of infection from a single contact with an infectious animal and  $\beta$  is  
85 the average number of daily contacts per animal. The parameter  $\sigma$  is the rate at  
86 which those in the latent class move to the infectious class and  $\gamma$  the rate at  
87 which animals move from the infectious class to the resistant class.

### 88 *2.2. Datasets*

89 Four datasets were available to us. These were recorded using two herds of  
90 cattle during two recording periods. The datasets are labelled 1A, 1B, 2A and

91 2B with the number denoting the recording period, first or second, and the  
92 letter representing the herd. Datasets 1A and 1B were recorded during July  
93 2009, both producing 30 complete days of usable data with both of the herds  
94 returning complete data for 29 animals. The final two datasets recorded 28  
95 complete days of data across August and September 2009 with 2A recording  
96 data for 21 animals whilst 2B returned data for 17 animals.

### 97 *2.3. Network Construction*

98 In order to answer the question about how close the approximation to our  
99 recorded network needed to be, we constructed four types of random network.  
100 Each type of network was constructed using increasing amounts of information  
101 taken from the recorded data. Details of how all the networks were constructed  
102 follows, including details on the construction of the recorded and  
103 matched-on-day network used in our previous publication (Duncan et al.,  
104 2012). The matched-on-day network was previously referred to as a  
105 mass-action or random network but for the purposes of this paper we are using  
106 the description “matched-on-day” to demonstrate its relationship to the other  
107 types of random network we present. The information required from the  
108 recorded network and the mathematical construction for each type of random  
109 network can be seen in table 1.

#### 110 *2.3.1. Recorded and Matched-On-Day Networks*

111 For each of the four datasets a contact network was established, with the nodes  
112 representing the animals, and the edges, the contacts. A contact was defined to  
113 be any recorded interaction that lasted longer than 4 minutes. Although the



114 term contact has been used, only close proximity of the animals can be assumed  
115 rather than actual physical contact. These networks were split into consecutive  
116 12 hour time steps to give a manageable number of edges for each step in the  
117 later disease simulation. An identical number of random networks were  
118 constructed by taking the total number of interactions recorded in the  
119 particular 12 hour period for a particular dataset, creating the same number of  
120 random contacts and randomly allocating each of these contacts to pairs of  
121 animals in the respective herd. For each dataset and 12 hour period this gave us  
122 two networks, a recorded contact network and a random (“matched-on-day”)  
123 network, with the same number of nodes and edges but with different edge  
124 distributions for each 12 hour period for each of the four datasets.

### 125 *2.3.2. Additional Random Networks*

126 For each dataset, in addition to the matched-on-day network, we constructed  
127 three other random networks: “constant-on-animal”, “constant-on-day” and  
128 “matched-on-animal”. For the constant-on-animal network all animals had the  
129 same number of contacts as one another for every 12 hour period. The contacts  
130 were randomly assigned amongst the animals whilst ensuring that each animal  
131 had the required number of contacts. The number of contacts per animal was  
132 calculated by averaging all the recorded contacts over the number of animals  
133 and the number of 12 hour time periods per dataset. Due to rounding, this  
134 meant that the total number of contacts for each of these networks was different  
135 from the total number of contacts in the recorded dataset they were derived  
136 from.

138 For the constant-on-day network, the same total number of contacts per 12  
139 hour time period as with the constant-on-animal network was used but the  
140 contacts were allocated randomly amongst all the animals. There were no other  
141 constraints on the number of contacts an individual animal could have. The  
142 structure of this network was seen as lying between that of the  
143 constant-on-animal network and the matched-on-day network. Very little  
144 information (see table 1) from the recorded network was used in the construction  
145 of either the constant-on-animal network or the constant-on-day networks.

146

147 In the matched-on-animal network each animal had exactly the same number  
148 of contacts as in the recorded network, for each 12 hour period, but those  
149 contacts were randomly distributed amongst the other animals subject to this  
150 condition i.e. that the number of contacts each animal had was the same as the  
151 recorded network. As with the other random network, matched-on-animal  
152 networks were constructed for all four datasets.

#### 153 *2.4. Network Metrics*

154 To investigate the differences between the five networks (constant-on-animal;  
155 constant-on-day; matched-on-day; matched-on-animal and recorded) six  
156 different network metrics were calculated. The first was our own metric, the  
157 number of repeated edges, chosen to quantify the observed difference in  
158 repeated contacts. The second was closeness, the inverse of the average length  
159 of the shortest paths to/from all the other vertices in the network (Csardi,

160 2013), and the third metric chosen was the clustering coefficient, a measure of  
161 the degree to which nodes in a network tend to cluster together (Newman,  
162 2003). The fourth metric that we used, average path length (Strogatz, 2001), is  
163 the average number of steps along the shortest path for all possible pairs of  
164 nodes. We also calculated the average degree distribution and finally the  
165 assortativity coefficient to establish whether assortative mixing, connections  
166 between nodes that are similar, was taking place (Molina and Stone, 2012).  
167 Each of these metrics were calculated for each network and for each dataset.

### 168 *2.5. Modelling Disease Spread*

169 All the models, using recorded or any of the four random network types, were  
170 implemented as stochastic due to the small numbers of animals in each of the  
171 datasets, and hence the increased influence of individual stochastic events on  
172 the overall disease transmission process (Brauer et al., 2000). Infection was  
173 always introduced by randomly infecting a single animal at the start of each  
174 model simulation, thus this animal began the simulation in the latent state.  
175 The periods of time each animal spends in the latent and infectious states were  
176 sampled from exponential distributions with means  $1/\sigma$  and  $1/\gamma$ . For simplicity,  
177 and because the largest dataset only contained 30 days of continuously recorded  
178 interactions, each infected animal had its length of resistance set to greater than  
179 30 days. Both models were simulated many times and it was found that the  
180 probability densities of the number of animals in each disease state at each time  
181 point, appeared to stabilise by 5000 simulations. All results presented were  
182 produced from 5000 simulations, where each simulation was run for the number

183 of days contained in the respective dataset with an initially infected animal  
184 randomly chosen for each simulation.

185

186 The value of  $\beta$ , the mean contact rate, used in the simulations was dependent  
187 on the dataset used, as each of the four datasets had a different average contact  
188 rate. Thus we had four values for  $\beta$  corresponding to our four datasets.

189

190 The disease spread through each model was a hypothetical disease with  
191 parameter values that allowed the peak of infection of an epidemic to occur  
192 within the 28 days of data available from the shortest dataset. Latent and  
193 infectious periods of six days were chosen. Using average values of  $\beta = 7.987$   
194 from our data and  $R_0 = 5$  (considered reasonable), a rounded value of  $\alpha = 0.1$   
195 was calculated from

$$R_0 = \frac{\alpha\beta}{\gamma}. \quad (2)$$

196 As each dataset has a different value of  $\beta$ , the contact rate, they will also have a  
197 different value of  $R_0$  but the characteristics specific to the disease ( $\alpha = 0.1$ ,  
198  $1/\sigma = 6$  days and  $1/\gamma = 6$  days) remain fixed across all datasets for the recorded  
199 network. For all random networks only the value of  $\alpha$  was altered. It was  
200 assumed that when an animal became infected its behaviour did not change  
201 such that its contacts continued as normal. This is not necessarily the case  
202 (Rush et al., 2008; Wilesmith, 1998) but until there exists actual contact data  
203 for a herd with spreading disease, it is parsimonious to use the actual data that  
204 we do have.

205 *2.6. Measuring the Differences in Disease Spread*

206 The results of our previous paper (Duncan et al., 2012) were divided into two  
207 parts: the proportion of 5000 simulations that produced no infection and  
208 percentiles of the number of infected animals predicted by those simulations  
209 that did produce infection. For all values of the disease parameters, the  
210 recorded network model had a higher proportion of simulations showing no  
211 infection and of those simulations that did show infection, fewer animals were  
212 modelled as infected. In an attempt to minimise the differences between the  
213 recorded and random network models the value of  $\alpha$  was altered in each type of  
214 random network model. The value of  $\alpha$  was chosen because the value of  $\beta$  was  
215 defined by the datasets and needed to be constant to maintain the continuity in  
216 number of contacts between the networks and  $\gamma$  has a basis in other diseases  
217 and was dependent on the amount of data available to us, a maximum of 30  
218 days. Additionally the large uncertainty in the estimates of the probability of  
219 infection for real diseases makes  $\alpha$  an attractive candidate for adjustment in  
220 random network models.

221

222 The standard value of  $\alpha = 0.1$  from our previous paper (Duncan et al., 2012)  
223 was used again for the recorded network model and a numerical study  
224 conducted on the value of  $\alpha$  for the various random network models. For each of  
225 the 40 equally spaced values of  $\alpha$  in the range  $0.025 \leq \alpha \leq 0.4$ , all random  
226 network models were run with 5000 simulations. The mean absolute difference  
227 in both the number of infected animals  $\text{M.A.D.}_{\text{No. Inf.}}$  and in the proportion of the  
228 5000 simulations showing no infection  $\text{M.A.D.}_{\text{Probn. Zero Sims.}}$  were calculated as shown in

229 equations (3) and (4). In these equations  $P_{rec}$  and  $P_{rand}$  represent the proportion  
 230 of the 5000 simulations that produced no infection for the recorded and random  
 231 network models respectively with  $\bar{I}_{rec}$  and  $\bar{I}_{rand}$  the mean number of infected  
 232 animals for each model from those simulations that did produce infection. The  
 233 *rand* refers to any of the four types of random network: constant-on-animal,  
 234 constant-on-day, matched-on-day and matched-on-animal. Each individual  
 235 time period is represented by  $t$  and  $T$  is the total number of time periods.

$$\text{M.A.D.}_{\text{No. Inf.}} = \frac{\sum_t |\bar{I}_{rec} - \bar{I}_{rand}|}{T}. \quad (3)$$

$$\text{M.A.D.}_{\text{Propn. Zero Sims.}} = \frac{\sum_t |P_{rec} - P_{rand}|}{T}, \quad (4)$$

236 This examination of  $\alpha$  gave an initial estimate of where the minima occurred for  
 237 each type of random network and dataset. To improve these estimates an  
 238 interval of length 0.05, including this first estimate, was examined in increments  
 239 of length 0.00125 for each type of network and each dataset. To get a single  
 240 value for the minima, splines were fitted to these data points for the mean  
 241 absolute difference in both number of infected animals and proportion of  
 242 simulations showing no infection, using the `smooth.spline` function of CRAN R  
 243 (CRAN-R, 2013) with a smoothing parameter of 0.7 which gave the closest  
 244 agreement with the visual minimum of the data points. This left two values of  $\alpha$   
 245 for each random network and dataset: one value minimising  $\text{M.A.D.}_{\text{No. Inf.}}$  and a  
 246 second minimising  $\text{M.A.D.}_{\text{Propn. Zero Sims.}}$ . The arithmetic mean of these two values was  
 247 calculated to leave one value  $\alpha_m$  to minimise the differences between the  
 248 recorded and random network models for each of the four random networks and

249 the four datasets. We conducted similar examinations to find  $\alpha_m$  for the  
250 matched-on-day network model when we set  $\alpha = 0.05$  and  $\alpha = 0.2$  in the  
251 recorded network model. This sensitivity analysis was carried out to establish  
252 whether the value of  $\alpha$  used in the recorded network model had any effect on  
253 the adjustment to find  $\alpha_m$ .

### 254 **3. Results**

#### 255 *3.1. Network Metrics*

256 The 5000 simulations of the random contact networks, outlined above, were  
257 stored to calculate average values for the six metrics. For each dataset the  
258 contact networks were split into 12 hour periods and the metrics calculated on  
259 each of the 5000 simulations. The results were averaged across the simulations  
260 and then over the 12 hour periods. These were then compared to the equivalent  
261 metrics calculated for the recorded network which was split into 12 hour periods  
262 after the disease simulations.

263

264 Figure 1 shows the results of the metrics in six separate plots. Each plot shows  
265 results for all networks split by the four datasets. There is no clear result from  
266 the metrics as to which of the random networks provides the closest  
267 approximation to our recorded network. The recorded network had more  
268 repeated edges and lower closeness than any of the random networks and this  
269 was consistent across all the datasets. In all but one dataset the recorded  
270 network also had higher average path length than the random networks. The  
271 more information from the recorded network used to construct the random

272 network – the greater the number of repeated edges in the random networks  
273 and hence closer to that of the recorded network.

274

275 Each network shows disassortativity across all datasets. For three of the  
276 datasets the recorded network was more disassortative than all four random  
277 networks and, as with the repeated edges, the more information from the  
278 recorded network used by the random network, in general, the more  
279 disassortative they became. Generally speaking in, three metrics (average path  
280 length, average closeness and average repeated edges) increasing similarity with  
281 the recorded network was associated with the random model utilising increased  
282 information from the recorded network.

### 283 *3.2. Disease Spread*

284 A sample of the results for the mean absolute differences in both the number of  
285 infected animals and the proportion of 5000 simulations showing no infection  
286  $\left( \begin{matrix} \text{M.A.D.} & \text{M.A.D.} \\ \text{No. Inf.} & \text{Propn. Zero Sims.} \end{matrix} \right)$  can be seen in figure 2. These are the results for  
287 the matched-on-day network for all four datasets. The results for the other  
288 random networks can be seen in the supplementary information. The results for  
289  $\begin{matrix} \text{M.A.D.} \\ \text{No. Inf.} \end{matrix}$  are shown in the solid lines using the left hand axes with the results of  
290  $\begin{matrix} \text{M.A.D.} \\ \text{Propn. Zero Sims.} \end{matrix}$  plotted as dashed lines using the right hand axes.

291

292 For each of the datasets and across all the random networks the results were  
293 very similar with four points to note. First there is a single minimum value of  
294  $\alpha_m$  and the differences in  $\begin{matrix} \text{M.A.D.} \\ \text{No. Inf.} \end{matrix}$  and  $\begin{matrix} \text{M.A.D.} \\ \text{Propn. Zero Sims.} \end{matrix}$  at this value of  $\alpha_m$  are very



295 small. Secondly the value of  $\alpha_m$  is always less than the value of  $\alpha = 0.1$  used in  
 296 the recorded network. It is also consistent, across all networks and datasets,  
 297 that the value of  $\alpha$  that results in minimising the differences in the proportion  
 298 of the 5000 simulations showing no infection is larger than the respective value  
 299 of  $\alpha$  for the difference in the number of infected animals. Finally, there are clear  
 300 but not very large differences in the value of  $\alpha_m$  for each type of network across  
 301 the four datasets.

302

303 The results from the proportion of simulations with no infected animals and the  
 304 values of the 25<sup>th</sup>, 50<sup>th</sup> and 75<sup>th</sup> percentiles of the number of infected animals  
 305 from those simulations showing infection are plotted for both the recorded  
 306 network model ( $\alpha = 0.1$ ; black, solid lines) and the matched-on-day network  
 307 model ( $\alpha_m = 0.0696$ ; red, dashed lines) are plotted in figure 3 for dataset 1A.  
 308 Similar plots for the other random networks are shown in the supplementary  
 309 information. In all cases it is clear that by adjusting  $\alpha$  the results of simulated  
 310 disease spread through the random networks are extremely close to the results  
 311 from the recorded network. Using the single value of  $\alpha_m$  provides very close  
 312 agreement and it is not necessary to use both the value of  $\alpha$  that resulted in

313  $\frac{\text{M.A.D.}}{\text{No. Inf.}}$ , and the one that gives  $\frac{\text{M.A.D.}}{\text{Propn. Zero Sims.}}$ .

314

315 To compare the differences between the results for each of the four types of  
 316 random networks the minimum values of  $\frac{\text{M.A.D.}}{\text{No. Inf.}}$  and  $\frac{\text{M.A.D.}}{\text{Propn. Zero Sims.}}$  are shown  
 317 in figure 4. These were plotted for each dataset along with the relative  
 318 percentage decrease in  $\alpha$  needed to achieve  $\alpha_m$ . Figure 4 also shows the

319 differences  $\alpha_m$  for each type of network across the four datasets. It is clear from  
320 the plot that the mean differences in number of infected animals are much less  
321 than a single animal for each of the networks. The value is dependent on the  
322 network being used in the simulation as can be seen by the consistent order of  
323 results (constant-on-day, constant-on-animal, matched-on-day and  
324 matched-on-animal). It is worth noting that the network using the least  
325 information from the recorded network, constant-on-animal, is not the poorest  
326 performing. The relative percentage decrease needed to achieve  $\alpha_m$  is  
327 somewhere between 5.6% and 39.4% but this varies depending on the dataset  
328 and the random network used.

329

330 It is clear from the left-hand plot in figure 4 that the values of  $\frac{\text{M.A.D.}}{\text{No. Inf.}}$  are  
331 dependent on the simplicity of the model. The model using the most  
332 information, the matched-on-animal network, is closest to the recorded  
333 network. However the simplest network (constant-on-animal) was numerically  
334 closer to the recorded network than the second simplest network  
335 (constant-on-day). This was also the case for  $\frac{\text{M.A.D.}}{\text{Propn. Zero Sims.}}$  for all but dataset  
336 1B. The loss of representativeness that arises from choosing the simplest  
337 random network is not large.

338

339 The right-hand plot of figure 4 shows the relative percentage decrease of  $\alpha$   
340 needed to achieve  $\alpha_m$  for each the random networks and for each dataset. The  
341 patterns in the adjustment are not completely consistent either with regard to  
342 the datasets or networks. There appears by eye to be a dataset effect in the

343 right-hand plot of figure 4. General linear regression, included in the  
344 supplementary information, suggests there is evidence of both a dataset effect  
345 and random network effect. Each factor was fairly strongly significant after the  
346 addition of the other factor,  $p = 0.0007$  and  $p = 0.042$  for dataset and network  
347 respectively. The mean reduction in  $\alpha$  was 26.8% and the median reduction was  
348 30.0%.

349

350 The exact values of  $\alpha_m$  are shown in table 2. For three of the datasets the  
351 highest value of  $\alpha_m$  occurred in the matched-on-animal network, the network  
352 using the most information from the recorded network. Nevertheless for dataset  
353 2A, the matched-on-animal had the second highest value of  $\alpha_m$ . For the first  
354 recording period (datasets 1A and 1B) the value of  $\alpha_m$  increases as the networks  
355 use more information from the recorded network and this trend is less clear for  
356 the second recording period.

357

358 Also included in the supplementary information are plots of the differences in  
359 the proportion of 5000 simulations that produced no infection and the median  
360 number of infected animals from those simulations that did produce infection.

#### 361 4. Discussion

362 It is clear from the simulations of disease spread that a simple homogeneous  
363 mixing model can approximate, very closely, a recorded network if the  
364 probability of infection,  $\alpha$ , is optimally adjusted. Each of our four types of  
365 random network can approximate the recorded network and can do so for each

366 of the four datasets. The adjustment was consistently a reduction in  $\alpha$ . The size  
367 of the adjustment was dependent on the dataset and random network used for  
368 the simulations. The relative percentage reduction in  $\alpha$  ranged from 5.6% to  
369 39.4%. The results of the sensitivity analysis shown in the supplementary  
370 information would suggest that the value of  $\alpha_m$  as a proportion of  $\alpha$  is  
371 negatively associated with the value of  $\alpha$  used in the recorded network, at least  
372 for the values of  $\alpha$  that we tested.

373

374 It has previously been shown that higher clustering tends to produce shorter  
375 path lengths within theoretical networks (Shirley and Rushton, 2005a), that  
376 clustering and assortativity can reduce epidemic size (Miller, 2009) and that  
377 increased clustering or increased assortativity can increase the likelihood of  
378 simulated disease spread occurring (Badham and Stocker, 2010a). There is  
379 however disagreement over whether clustering influences epidemics on  
380 undirected networks with regular (many repeated contacts) or random  
381 construction (Eames, 2008; Moslonka-Lefebvre et al., 2009).

382

383 Theoretical networks constructed with many repeated contacts show slower  
384 disease spread than random networks (Eames, 2008). This is also shown by  
385 both our earlier work (Duncan et al., 2012) and further demonstrated by  
386 random networks constructed here. In general, our random networks with lower  
387 repeated contacts, i.e. the simpler networks (contact-on-animal and  
388 contact-on-day) required smaller values of  $\alpha_m$  suggesting that disease spreads  
389 quicker through them.

390

391 As all the random networks are derived from the recorded network and the  
392 average degree distributions are either extremely close to one another or  
393 identical, we can gain little insight from degree distribution. However, degree  
394 distribution alone has been shown to not provide enough information for  
395 prediction of disease spread (Ames et al., 2011; Boily et al., 2007).

396

397 We found no clear relationship between the values of the metrics and the values  
398 of  $\alpha_m$  and formal inferential statistics are not possible given the sample size.  
399 Any inferential statistical relationship will, however, depend on a large number  
400 of herds being assessed in the same manner.

401

402 One of the largest differences between the recorded network and the random  
403 networks is the number of repeated edges. One possible reason for the high  
404 number of repeated edges in the recorded network was that the herds were  
405 constructed of cows with calves at foot. Of the repeated edges recorded, 15% to  
406 30%, depending on dataset, were between a cow and her calf. These repeated  
407 edges could also be a reason for the increased disassortativity found in the  
408 recorded network. Assortative mixing would normally entail cows contacting  
409 cows and calves contacting calves. With young calves present in the herd, the  
410 disassortative mixing, resulting from cow contacting calf, would seem probable.  
411 Assortativity has been shown to decrease epidemic size (Miller, 2009) and we  
412 have found that  $\alpha_m < 0.1$  for all networks and datasets, showing that the  
413 recorded network produces slower disease spread than the random networks.

414 The age of the calves may also explain why in the first recording period  
415 (datasets 1A and 1B) the value of  $\alpha_m$  increases as the random networks  
416 approach the recorded network. In the second recording period, where the  
417 calves were a little older, there is not such a clear pattern.

418

419 It has recently been shown that indirect, environmental or faecal, contact may  
420 aid the spread of disease in herds of cattle (Kleinlützum et al., 2013). These  
421 factors cannot be taken into account with the data available to us. Likewise we  
422 only have proximity data with which to construct our contact networks. We do  
423 not know the extent of the contacts and how likely each one is to spread disease.  
424 However, the only way to gather such data would be to film the animals at all  
425 times and to monitor real life spread of infection. Even those studies which  
426 attempt to take such things into account by observing animals and categorising  
427 the contacts by strength (Norton et al., 2012) are still summarising the contact  
428 networks as they extrapolate their networks from the observed data.

## 429 **5. Conclusion**

430 We have shown that it is possible to closely model disease spread through a  
431 network of recorded contacts with a network of randomly allocated contacts by  
432 adjusting the probability of infection. The adjustment in probability of infection  
433 is consistently a reduction and there appears to be a dataset effect in the value  
434 of the reduction. The exact values in adjustment varies between 5.6% and  
435 39.4% and as yet, with only four datasets, we have no clear relationship between  
436 the network properties and the adjustment in the probability of infection.

437 Recommended reductions in  $\alpha$  should not be made until further intra-herd  
438 contact data becomes available. Importantly, the simplest network, requiring  
439 least information to construct, performed reasonably well by giving a close  
440 match to disease spread in the recorded network. This is important because it  
441 suggests that in the absence of real contact data a good approximation to  
442 disease spread could be made if the correct adjustment in the probability were  
443 known.

## 444 **6. Acknowledgements**

445 The Scotland's Rural College (SRUC) receives financial support from the  
446 Scottish Government within WP6.1 of the RESAS 2011–2016 research  
447 programme. The authors would like to acknowledge the link between Inverness  
448 College UHI and the SRUC thus enabling this work to be carried out.

## **References**

- Ames, G. M., George, D. B., Hampson, C. P., Kanarek, A. R., McBee, C. D.,  
Lockwood, D. R., Achter, J. D., Webb, C. T., 2011. Using network properties  
to predict disease dynamics on human contact networks. *Proceedings of the  
Royal Society B: Biological Sciences* 278, 3544–3550.
- Anderson, R., May, R., 1991. *Infectious Diseases of Humans*. Oxford University  
Press, Oxford, UK.
- Badham, J., Stocker, R., 2010a. The impact of network clustering and  
assortativity on epidemic behaviour. *Theoretical Population Biology* 77,  
71–75.

- Badham, J., Stocker, R., 2010b. A spatial approach to network generation for three properties: degree distribution, clustering coefficient and degree assortativity. *Journal of Artificial Societies and Social Simulation* 13 (1), 1–13.
- Bansal, S., Khandelwal, S., Meyers, L. A., 2009. Exploring biological network structure with clustered random networks. *BMC Bioinformatics* 10, 405.
- Barlow, N. D., 2000. Non-linear transmission and simple models for bovine tuberculosis. *J. Anim. Ecol.* 69 (4), 703–713.
- Boily, M., Asghar, Z., Garske, T., Ghani, A., Poulin, R., 2007. Influence of selected formation rules for finite population networks with fixed macrostructures: implications for individual-based model of infectious diseases. *Mathematical Population Studies* 14 (4), 237–267.
- Bouma, A., Dejong, M. C. M., Kimman, T. G., 1995. Transmission of pseudorabies virus within pig-populations is independent of the size of the population. *Prev. Vet. Med.* 23 (3-4), 163–172, doi:10.1016/0167-5877(94)00442-L.
- Brauer, F., van den Driessche, P., Wu, J. (Eds.), 2000. *Mathematical Epidemiology. Lecture Notes in Mathematics.* Springer.
- Craft, M., Caillaud, D., 2001. Network models: An underutilized tool in wildlife epidemiology? *Interdis. Perspec. Inf. Dis.ID* 676949, doi: 10.1155/2011/676949.
- CRAN-R, 2013. *Smooth.spline r documentation.* Last accessed 30/09/2013.



URL

<http://stat.ethz.ch/R-manual/R-patched/library/stats/html/smooth.spline.html>

Csardi, G., 2013. Package ‘igraph’. Last accessed 03/09/2013.

URL <http://cran.r-project.org/web/packages/igraph/igraph.pdf>

D’ Amico, V., Elkinton, J. S., Dwyer, G., Burand, J. P., Buonaccorsi, J. P.,  
1996. Virus transmission in gypsy moths is not a simple mass action process.  
*Ecol.* 77 (1), 201–206.

De Jong, M. C. M., Diekmann, O., Heesterbeek, H., 1995. How does  
transmission of infection depend on population size? Publications of the  
Newton Institute; Epidemic models: Their structure and relation to data,  
84–94.

Dobson, A., Meagher, M., 1996. The population dynamics of brucellosis in the  
yellowstone national park. *Ecol.* 77 (4), 1026–1036.

Drewe, J., Weber, N., Carter, S., Bearhop, S., Harrison, X., Dall, S., McDonald,  
R. a., 2012. Performance of proximity loggers in recording intra and  
inter–species interactions: A laboratory and field based validation study. *PLoS*  
*One* 7, 1–9.

Duncan, A., Gunn, G., Lewis, F., Umstatter, C., Humphry, R., 2012. The  
influence of empirical contact networks on modelling diseases in cattle.  
*Epidemics* 4, 117–123.

Eames, K., 2008. Modelling disease spread through random and regular contacts  
in clustered populations. *The Journal of Theoretical Biology* 73 (1), 104–111.

- Håkansson, N., Jonsson, A., Lennartsson, J., Lindström, T., Wennergren, U., 2010. Generating structure specific networks. *Advances in Computer Systems* 13 (2), 239–250.
- Hamede, R. K., Bashford, J., Jones, M., McCallum, H., 2012. Simulating devil facial tumour disease outbreaks across empirically derived contact networks. *Journal of Applied Ecology* 49, 447–456.
- Keeling, M., 2005. The implications of network structure for epidemic dynamics. *Theo. Pop. Biol.* 67 (1), 1–8.
- Kleinlützum, D., Weaver, G., Schley, D., 2013. Within-group contact of cattle in dairy barns and the implications for disease transmission. *Research in Veterinary Science* 95, 425–429.
- Miller, J. C., 2009. Percolation and epidemics in random clustered networks. *Physical Review E* 80, 020901(R).
- Molina, C., Stone, L., 2012. Modelling the spread of disease in clustered networks. *Journal of Theoretical Biology* 315, 110–118.
- Moslonka-Lefebvre, M., Pautasso, M., Jeger, M., 2009. Disease spread in small-size directed networks: epidemic threshold, correlation between links to and from nodes, and clustering. *Journal of Theoretical Biology* 260 (3), 402–411.
- Newman, M., 2003. The structure and function of complex networks. *SIAM Review* 45, 167–256.

- Norton, E., Benaben, S., Mbotha, D., Schley, D., 2012. Seasonal variations in physical contact amongst domestic sheep and the implications for disease transmission. *Livestock Science* 145 (1), 34–43.
- Porphyre, T., Stevenson, M., Jackson, R., McKenzie, J., 2008. Influence of contact heterogeneity on tb reproduction ratio  $r_0$  in a free-living brushtail possum *Trichosurus vulpecula* population. *Veterinary Research* 39, 31–43.
- Rush, J., de Passille, A., von Keyserlingk, M., Weary, D., 2008. *The Welfare of Cattle Volume 5*. Springer, Dordrecht, The Netherlands.
- Salathé, M., Kazandjieva, M., Lee, J., Levis, P., Feldman, M., Jones, J., 2010. A high-resolution human contact network for infectious disease transmission. *Proceedings of the National Academy of Sciences* 107 (51), 22020–22025.
- Shirley, M., Rushton, S., 2005a. The impacts of network topology on disease spread. *Ecological Complexity* 2, 287–299.
- Shirley, M., Rushton, S., 2005b. Where diseases and networks collide: lessons to be learnt from a study of the 2001 foot-and-mouth disease epidemic. *Epidemiology and Infection* 133, 1023–1032.
- Strogatz, S. H., 2001. Exploring complex networks. *Nature* 410 (6825), 268–276.
- Swain, D. L., Bishop-Hurley, G. J., 2007. Using contact logging devices to explore animal affiliations: Quantifying cow-calf interactions. *Applied Animal Behaviour Science* 102 (1-2), 1–11, doi:10.1016/j.applanim.2006.03.008.
- Vernon, M. C., Keeling, M. J., 2009. Representing the uk’s cattle herd as static

and dynamic networks. *Proceedings of the Royal Society B: Biological Sciences* 276, 469–476.

Watts, D., Strogatz, S., 1998. Collective dynamics of ‘small-world’ networks. *Nature* 393, 440–442.

Wilesmith, J., 1998. *Manual on bovine spongiform encephalopathy*, 34.

Table 1: Descriptions of how the four random networks relate to the recorded network and how much information from the recorded network was necessary to create them.

Information needed to construct random network	Random Network	Mathematical Comparison
Total number of contacts, number of animals, total number of time periods	constant-on-animal	$\sum_j x_{i,j,t} = k \quad \forall i, t$
Total number of contacts, number of animals, total number of time periods	constant-on-day	$\sum_{i,j;i>j} x_{i,j,t} = kN \quad \forall t$
Total number of contacts per time period, number of animals	matched-on-day	$\sum_{i,j;i>j} x_{i,j,t} = \sum_{i,j;i>j} r_{i,j,t} \forall t$
Total number of contacts per animal per time period, number of animals	matched-on-animal	$\sum_j x_{i,j,t} = \sum_j r_{i,j,t} \forall i, t$

where:

$x_{i,j,t}$  = a simulated contact between animals  $i$  and  $j$  during time period  $t$  with  $i \neq j$

$r_{i,j,t}$  = a recorded contact between animals  $i$  and  $j$  during time period  $t$  with  $i \neq j$

$$k = \text{round} \left( \frac{\sum_{i,j,t;i>j} r_{i,j,t}}{NT} \right)$$

$N$  = Total population size (Number of animals)

$T$  = Total number of timeperiods

Table 2: Values of  $\alpha_m$ , the value of the probability of infection  $\alpha$ , used to minimise the differences between the recorded and random network models for each of the four types of random networks - for each of the four datasets. A value of  $\alpha = 0.1$  was used for the recorded model across all simulations.

Network	$\alpha_m$ per dataset			
	1A	1B	2A	2B
constant-on-animal	0.0645	0.0684	0.0705	0.0944
constant-on-day	0.0649	0.0770	0.0606	0.0757
matched-on-day	0.0695	0.0830	0.0664	0.0844
matched-on-animal	0.0799	0.0915	0.0765	0.0886

Values of six metrics calculated for all networks – random and recorded

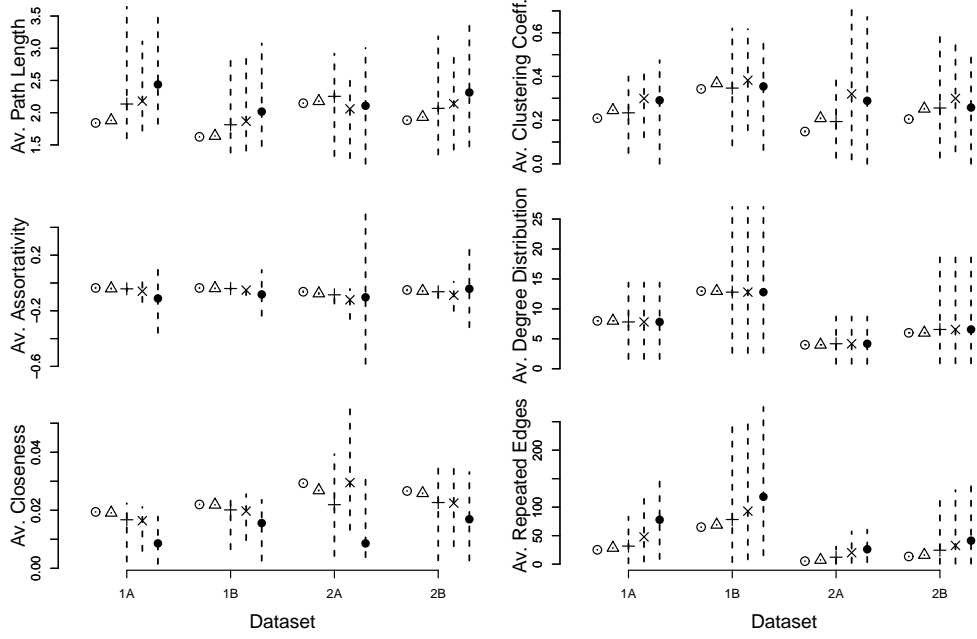


Figure 1: The average values of all six metrics calculated for each of the five networks. The symbols  $\circ$ ,  $\triangle$ ,  $+$ ,  $\times$  and  $\bullet$  denoting results from the constant-on-animal, constant-on-day, matched-on-day, matched-on-animal and recorded networks respectively. The vertical dashed lines represent the 95% percentiles for each metric.

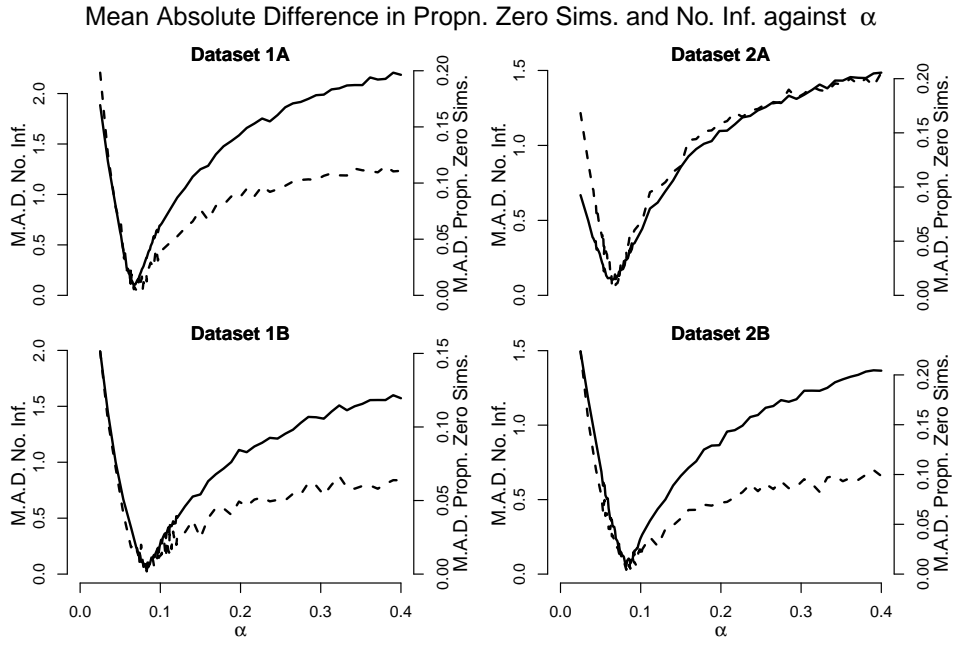


Figure 2: Plots of the mean absolute difference in the number of infected animals  $\left( \begin{matrix} \text{M.A.D.} \\ \text{No. Inf.} \end{matrix} \right)$  (left-hand axis, solid line) and mean absolute difference in the proportion of the 5000 simulations showing no infection  $\left( \begin{matrix} \text{M.A.D.} \\ \text{Propn. Zero Sims.} \end{matrix} \right)$  (right-hand axis, dashed line) against  $\alpha$  for all four datasets.  $\alpha = 0.1$  was used in the recorded network model.



Plots of Proportion Zero Simulations and Percentiles of Infected Animals from Dataset 1A

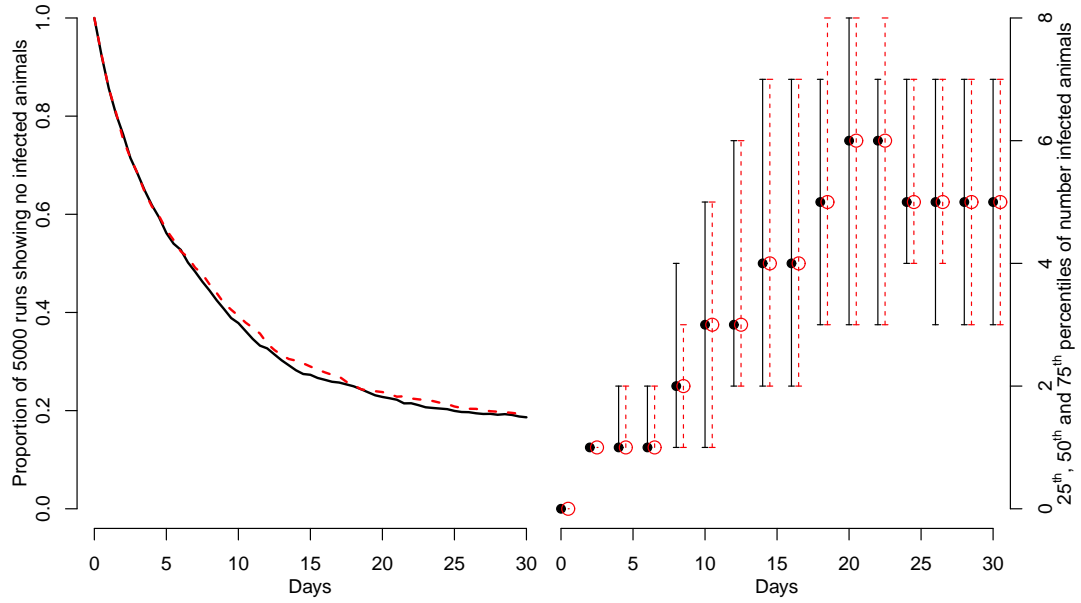


Figure 3: Left-hand plot: Proportion of 5000 simulations that produced no infection for the recorded network model with  $\alpha = 0.1$  (black, solid line) and the adjusted random network model with  $\alpha = \alpha_m$  (red, dashed line). Right-hand plot: The 25<sup>th</sup>, 50<sup>th</sup> and 75<sup>th</sup> percentiles of the number of infected animals from those simulations that did produce infection for the recorded network model with  $\alpha = 0.1$  (black, solid line) and the adjusted random network model with  $\alpha = \alpha_m$  (red, dashed line). Dataset 1A was used for both models.

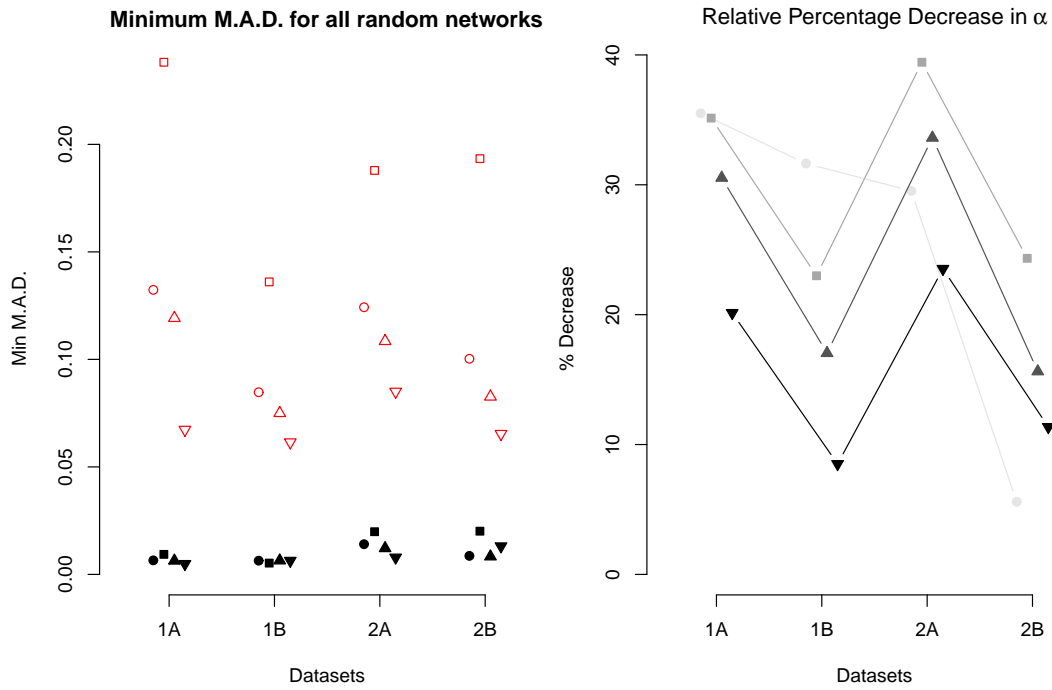


Figure 4: Left-hand plot: The values of the mean absolute difference in the number of infection animals  $\left( \begin{matrix} \text{M.A.D.} \\ \text{No. Inf.} \end{matrix} \right)$  (unfilled, red symbols) and the mean absolute difference in the proportion of 5000 simulations showing no infection  $\left( \begin{matrix} \text{M.A.D.} \\ \text{Propn. Zero Sims.} \end{matrix} \right)$  (filled, black symbols) for  $\alpha_m$  plotted for each of the four random networks. Right-hand plot: The relative percentage decrease in  $\alpha$  to achieve  $\alpha_m$  from the value of  $\alpha = 0.1$  used in the recorded network. The shading denotes the amount of information from the recorded needed to construct the random network, lightest representing the least information and the darkest representing the most information. In both plots the symbols  $\circ$ ,  $\square$ ,  $\triangle$  and  $\nabla$  represent the constant-on-animal, constant-on-day, matched-on-day and matched-on-animal networks.