

1 Short communication:
2 Tools to study trends in community structure:
3 application to fish and livestock trading networks

4 Darren Michael Green^{a*}, Marleen Werkman^a,
 Lorna Ann Munro^b, Rowland Raymond Kao^c,
 István Zoltán Kiss^d & Leon Danon^{e,f}

5 January 11, 2011

6 ^a *Institute of Aquaculture, School of Natural Sciences, University of Stirling,*
7 *Stirling, Stirlingshire FK9 4LA, UK.*

8 ^b *Marine Scotland Marine Laboratory, PO Box 101, 375 Victoria Road, Aberdeen,*
9 *AB11 9DB, UK.*

10 ^c *Institute of Comparative Medicine, Faculty of Veterinary Medicine, University*
11 *of Glasgow, Glasgow, G61 1QH, UK*

12 ^d *Department of Mathematics, University of Sussex, Falmer, Brighton BN1 9RF,*
13 *UK.*

14 ^e *Mathematics Institute, University of Warwick, Gibbet Hill Road, Coventry*
15 *CV4 7AL, UK.*

16 ^f *Department of Epidemiology, Harvard School of Public Health, 677*
17 *Huntington Avenue, Boston 02115, USA.*

18 * *Author for correspondence. email: darren.green@stir.ac.uk*

Abstract

Partitioning of contact networks into communities allows groupings of epidemiologically related nodes to be derived, that could inform the design of disease surveillance and control strategies, e.g. contact tracing or design of ‘firebreaks’ for disease spread. However, these are only of merit if they persist longer than the timescale of interventions. Here, we apply different methods to identify concordance between network partitions across time for two animal trading networks, those of salmon in Scotland (2002-4) and livestock in Great Britain (2003-4). Both trading networks are similar in that they moderately agree over time in terms of their community structures, but this concordance is higher – and therefore community structure is more consistent – when only the ‘core’ network of nodes involved in trading over the whole time series is considered. In neither case was higher agreement found between partitions close together in time. These measures differ in their absolute values unless appropriate standardisation is applied. Once standardised, the measures gave similar values for both network types.

Keywords aquaculture; community; network; graph; movements

35 1. Introduction

36 Movement of farmed animals is an important route for disease spread in what are highly
37 structured industries. For example, sheep, cattle, and pigs were all involved in the UK
38 epidemic of foot-and-mouth disease in 2001 (Shirley & Rushton, 2005), and movements
39 of salmon were involved in the spread of infectious salmon anaemia (Murray et al. 2002).
40 A network representation, where farm sites are represented by ‘nodes’, and potentially
41 infectious contact by directed ‘arcs’ or undirected ‘edges’ is a powerful tool for studying
42 the potential for disease spread and control (for a review of networks in preventive
43 veterinary medicine, see Martínez-López et al., 2009).

44 Network communities represent partitions of nodes with a high level of
45 within-partition connectivity (for a review, see Fortunato, 2010). In a strongly
46 community-organised network, contact between communities may be relatively weak,
47 and community algorithms can provide us with natural groupings of epidemiologically
48 related nodes, derived from the network itself rather than artificially imposed.
49 Uncommon inter-community links might furthermore be considered as potential targets
50 for proactive targeted surveillance, or reactively in disease control (Kao et al. 2006; Green
51 et al. 2009; Salathé & Jones, 2010). That is to say, removing the disease transmission risk of
52 such contacts could reduce the size of potential epidemics by creating ‘firebreaks’,
53 particularly where these contacts are long distance. However, these analyses are only of
54 merit if partitions can be used predictively; that is, if community structure changes more
55 slowly than we collect data in order to inform surveillance or disease control strategy.

56 A key problem here is that objective measures of the rate of change of large-scale
57 network structure are not clearly defined, nor how large a change must be to heavily
58 compromise disease control strategies. In this short paper, we consider the first part of
59 this question, by comparing different methods for determining how network community
60 structures change, or not, over time. We apply these methods to two movement networks
61 of farmed animals, to investigate whether networks closer in time have more similar
62 network structure. The two networks are that of live Atlantic salmon *Salmo salar*
63 movements within Scotland 2002 to 2004, and that of livestock (pigs, sheep, cattle) in
64 Great Britain for 2003 to 2004.

65 2. Method

66 2.1. Data

67 The network of live fish movements in Scotland has been described for salmonid species
 68 (brown trout *Salmo trutta*, rainbow trout *Onchorhynchus mykiss*, Atlantic salmon *S. salar*)
 69 by Green et al. (2009) and Munro & Gregory (2009). Here, we extend and refine their
 70 analysis to a three-year dataset of Atlantic salmon alone for 2002-4. In brief, these data
 71 comprise movements of live fish (egg to adult) between registered sites in Scotland, where
 72 paper records of both off and on movements were legible and in agreement. Data are held
 73 by the Fish Health Inspectorate of Marine Scotland.

74 For the network of livestock movements, the partitions used here are derived from
 75 the data extract used by Kao et al. (2006). Their data set comprised data from January
 76 2003 to December 2004 for cattle (Cattle Tracing System) and sheep and pigs (Animal
 77 Movements Licence System, England and Wales; Scottish Agricultural Movements
 78 System, Scotland). A full description is given by Kao et al. (2006).

79 Both data sets provide source and destination premises, species and number
 80 moved, and date. Data were segregated into time periods (years for fish, four-week
 81 periods for livestock), with each network described by an adjacency matrix A . Here,
 82 $A_{ij} = 1$ implies movement of animals from node (site) i to node j (zero for no contact).
 83 The number of in and out connections for node i are given by k_i^{out} and k_i^{in} , the total
 84 number of nodes by n , and the total number of arcs by M .

85 2.2. Graph partitioning

86 Communitis were identified for the two datasets using related partitioning algorithms.
 87 For the fish network, the measure of community fit used is that defined by

$$Q = \frac{1}{M} \sum_{i,j} \left(A_{ij} - \frac{k_i^{\text{out}} k_j^{\text{in}}}{M} \right) [X_i = X_j]$$

88 where X_i is the community ‘label’ of node i . The Iverson bracket $[\cdot = \cdot]$ returns one if the
 89 condition inside is true, and zero otherwise. This formulation – as described by Kao *et al.*

(2006) and Leicht & Newman (2008) – accounts for the strong directed nature of the fish network. Higher Q indicates a larger fraction of arcs within communities. ‘Lone’ nodes in a network, with no movements during the period of interest, gain a unique label. However, it could be argued that without network activity, such nodes are not part of the network at all (further discussed in the Results section). The livestock network was treated similarly, except that the partition data available were based on undirected edges.

For both systems, we employ a ‘hill-climbing’ algorithm (Newman, 2004; Danon et al. 2005). This begins by assigning each node a unique community label $X_i = i$. Each possible merger of two communities is considered, with that providing the largest positive change in Q accepted. This step is repeated until a maximum Q is reached, for which the corresponding community assignments are taken as the ‘best fit’. Though other algorithms may find improved partitions, this one has the benefit of being practicable on very large networks such as that for livestock movements.

2.3. Entropy measures

Borrowing concepts from information theory, entropy-based measures can be used to compare multiple partitions of the same network (Strehl et al. 2002; Vinh et al. 2009). Beginning with two vectors X and Y containing community labels for two partitions, two vectors U and V are built containing the number of nodes present in each community in X and Y : $U_u = \sum_i [X_i = u]$; $V_u = \sum_i [Y_i = u]$. Also, an $n \times n$ matrix is defined containing the frequency combinations of communities in both X and Y : $W_{uv} = \sum_{i,j} [X_i = u][Y_j = v]$. For two networks with congruent partitions, this matrix contains only a single non-zero element in each row and column. The Shannon entropy (a measure of the information content of a dataset) is calculated for the partitions of each network ($H(U)$ and $H(V)$), and that of the matrix of community combinations, the ‘joint entropy’ $H(U,V)$.

$$H(U) = - \sum_u \frac{U_u}{n} \log \frac{U_u}{n}$$

$$H(U,V) = - \sum_{u,v} \frac{W_{uv}}{n} \log \frac{W_{uv}}{n}$$

Choice of logarithm base does not affect the end result below, and by definition, $0 \times \log 0 = 0$. The mutual information $I(U, V) = H(U) + H(V) - H(U, V)$ then measures the amount of information shared between the two partitions – and thus their similarity – with a lower bound of zero, but no upper bound. For comparison between networks, a normalised measure of similarity is required. A simple approach is to scale I by its maximum potential value (it cannot exceed the minimum of $H(U)$ and $H(V)$), giving the normalised mutual information $0 \leq NMI_1 \leq 1$:

$$NMI_1 = \frac{I(U, V)}{\min(H(U), H(V))}.$$

Alternatively, we can scale by the geometric mean of these two quantities (Strehl et al. 2002), $0 \leq NMI_2 \leq NMI_1$:

$$NMI_2 = \frac{I(U, V)}{\sqrt{H(U)H(V)}}.$$

For correlation coefficients such as Pearson's or Spearman's, a value of zero is obtained where there is no relationship, i.e. under the null hypothesis. However here, under a reasonable null hypothesis that communities are assigned randomly, the expectation of $I(U, V)$, $E_0(I(U, V))$ is not generally zero and depends upon the size distribution of communities (Vinh et al. 2009). A further approach is to normalise I against this expectation, providing the adjusted mutual information AMI (Vinh et al. 2009), with a maximum of one, zero under the null hypothesis, and negative where there is less agreement between network communities than would be expected by chance.

$$AMI = \frac{I(U, V) - E_0(I(U, V))}{\min(H(U), H(V)) - E_0(I(U, V))}$$

This definition of AMI is similar in form to that of Cohen's Kappa statistic, and has a lower value than NMI_1 except where $E_0(I(U, V))$ is vanishingly small. Vinh et al. (2009) suggest using $\max(\cdot, \cdot)$ not $\min(\cdot, \cdot)$, however the \min term has more in common with the formula for NMI_2 above. Unlike correlation coefficients, its minimum possible value is not defined to be -1 . A permutation test was employed to determine the mean and distribution of $E_0(I)$ allowing for calculation of AMI and its significance. One of the vectors X and Y is repeatedly shuffled, removing association between the node labels in X and Y . On each

permutation, $I(U, V)$ is recalculated. The original $I(U, V)$ can be compared with the distribution of these permuted versions.

2.4. Pair-based measures

Pairs of nodes can be examined with respect to whether or not they are in the same communities. Pairs of nodes that were in the same community in the two partitions were counted: $x = \sum_{i,j \neq i} [X_i = X_j]$ and $y = \sum_{i,j \neq i} [Y_i = Y_j]$, as well as pairs that were in the same community in both partitions: $a = \sum_{i,j \neq i} [X_i = X_j][Y_i = Y_j]$, or in different communities in both: $b = \sum_{i,j \neq i} [X_i \neq X_j][Y_i \neq Y_j]$.

From these values, the probability that a pair of nodes present in the same community in partition X are also in the same community in partition Y was calculated: $P(\text{pair in } Y | \text{pair in } X) = \frac{a}{x}$. However, this metric is not necessarily symmetric with respect to X and Y , unlike the earlier measures. Instead, the geometric mean of both possible probabilities was taken: $\bar{P} = \frac{a}{\sqrt{xy}}$ (Wallace, 1983; quoted in Meilă, 2007). These probabilities benefit from being easily interpretable. A further pair-based measure of clustering similarity, the Rand index $R = \frac{a+b}{n(n-1)}$ (Rand, 1971), was also calculated.

Again, these measures do not equal zero under the null hypothesis that the two partitions are independent. The statistical significance of both was determined through a permutation test and – as with the mutual information – standardised according to $[\phi - E_0(\phi)]/[1 - E_0(\phi)]$, where ϕ is the measure of interest, giving an adjusted Rand index AR and an adjusted probability related to \bar{P} , $A\bar{P}$.

3. Results and Discussion

For the salmon movement network ($n = 502$), the unadjusted indices NMI_1 , \bar{P} and R gave numbers of different magnitude, despite their apparent normalisation (Table 1). This reflects their different values under their null models. A \bar{P} index of ~ 0.3 is easily interpretable as the proportion of same-community node pairs that persist across both partitions. Once ‘adjusted’, the range of values was narrower, with the pair-based indices giving almost coincidental values (Table 1). This coincidence was also evident for the

165 livestock network, thus in Figure 1 only the index R is shown. The null model for the
166 permutation test was amended for the fish network to account for variation in the activity
167 of nodes between years: Those nodes with no links were not considered during the
168 reshuffling process to prevent their single-node communities being spuriously reassigned
169 to other nodes.

170 For the much larger livestock network ($n = 141607$; see supplementary animation),
171 networks were built from four-week periods of data. As with the fish network, all
172 correlations were statistically significant ($P < 0.05$). These networks show a marked
173 seasonal pattern (Kao et al. 2006) with a higher density of arcs due to an autumn peak in
174 sheep trading. This seasonality was still noticeable despite normalisation as a peak in
175 AMI values for networks 13 four-week periods (i.e. one year) apart (Fig. 1). Though this
176 peak may represent a real similarity in the trading structure at particular times of year,
177 Meilä (2007) raises concerns over the use of adjusted indices for comparison purposes
178 where the baseline and actual values may vary non-linearly.

179 To explore this further, we accounted for seasonality in trading volume by
180 considering only a ‘core’ sub-network of nodes that were active in each of the 25 networks
181 examined ($n = 6424$). The AMI values together with the Rand index R are shown for this
182 core network in Figure 1, showing close agreement between the three statistics and much
183 reduced seasonality in community structure. Taking the ‘core’ network of $n = 208$ nodes
184 for the salmon network, a similar result is found as for the whole fish network, albeit with
185 higher values (Table 1).

186 Though both sets of networks show moderate agreement between partitions at
187 different time points, in neither case was a higher agreement between networks closer in
188 time apparent. One possible explanation of this is that there are no significant long-term
189 trends in community structure for either network, or that any such trends operate on
190 timescales either longer or shorter than examined in this study. There may also be other
191 trends and patterns within the data that remain observed. For example, the partitions
192 above are not absolutes: different measures and algorithms could produce different
193 groupings. Also, no allowance is made in this approach for the potential for sub- and
194 super-community network structure (Kao et al. 2006; Green et al. 2009).

195 The unadjusted indices give a wide selection of values for the same network,
196 however once adjusted they are more similar. Those for \bar{P} coincided with R . However,
197 whether this is in general the case or is network dependent remains to be established.
198 The computational efficiency of the measures varies: Despite their apparent complexity,
199 the entropy-based measures are relatively fast to compute, particularly for large networks,
200 since they do not rely on counting edges.

201 **4. Conclusions**

202 In conclusion, for both networks a significant and non-trivial level of concordance
203 between network partitions over time was seen. Dissimilarity in partitions, however,
204 appears to represent random variation rather than decay in partition similarity over time
205 for both networks. Characterising the way networks change over time remains a
206 challenging problem. Our results suggest that despite the fact that many features change,
207 a large part of the intermediate structure is conserved over time, particularly in the core
208 network. Nevertheless, the how stable a contact network must remain over time to be
209 epidemiologically useful for disease surveillance and control remains to be explored,
210 potentially through simulation of dynamic disease control measures on dynamic network
211 epidemic models.

212 **Conflict of interest**

213 None declared.

214 **Acknowledgements**

215 With thanks to the Fish Health Inspectorate for providing access to the movement
216 records, and to Malcolm Hall for comments on the manuscript. DMG and MW are
217 supported by Marine Scotland.

218 **References**

- 219 Danon, L., Duch, J., Díaz-Guilera, A., Arenas, A., 2005, Comparing community structure
220 identification, J. Stat. Mech. P09008
- 221 Fortunato S., 2010, Community detection in graphs. Physics Reports, 486, 75 – 174.
- 222 Green, D.M., Gregory, A., Munro, L.A. 2009. Small- and large-scale network structure of
223 live fish movements in Scotland. Prev. Vet. Med. 91, 261 – 269.
- 224 Kao, R.R., Danon, L., Green, D.M., Kiss, I.Z. 2006. Demographic structure and pathogen
225 dynamics on the network of livestock movements in Great Britain. Proc. R. Soc. B
226 273, 1999 – 2007.
- 227 Leicht, E.A., Newman, M.E.J. 2008. Community structure in directed networks. Phys.
228 Rev. Lett. 100, 118703.
- 229 Martínez-López, B., Perez, A.M., Sánchez-Vizcaíno, J.M., 2009. Social network analysis.
230 Review of general concepts and use in preventive veterinary medicine.
231 Transboundary and Emerging Diseases 56, 109 – 120.
- 232 Meilă, M., 2007. Comparing clusterings – an information based distance. Journal of
233 Multivariate Analysis 98, 873 – 895.
- 234 Munro, L., Gregory, A. 2009. Application of network analysis to fish movement data. J.
235 Fish Dis. 32, 641 – 644.
- 236 Murray, A.G., Smith, R.J., Stagg, R.R., 2002. Shipping and the spread of infectious
237 salmon anemia in Scottish aquaculture. Emerging Infectious Diseases 8, 1 – 5.
- 238 Newman, M.E.J., 2004. Fast algorithm for detecting community structure in networks,
239 Phys. Rev. E. 69, 066133.
- 240 Rand, W.M. 1971. Objective criteria for the evaluation of clustering methods. J. Amer.
241 Statist. Assoc. 66, 846 – 850.
- 242 Salathé M., Jones, J.H. 2010. Dynamics and control of diseases in networks with
243 community structure. PLoS Computational Biology 6, e10000736.

- 244 Shirley, M.D.F, Rushton, S.P. 2005. Where diseases and networks collide: lessons to be
245 learnt from a study of the 2001 foot-and-mouth disease epidemic. *Epidemiol.*
246 *Infect.* 133, 1023 – 1032.
- 247 Strehl, A., Ghosh, J., Cardie, C. 2002. Cluster ensembles – A knowledge reuse framework
248 for combining multiple partitions. *Journal of Machine Learning Research* 3, 583 –
249 617.
- 250 Vinh, N.X., Epps. J., Bailey, J. 2009. Information theoretic measures for clustering
251 comparison: Is a correction for chance necessary? *Proceedings of the 26th*
252 *international conference on machine learning*. Montreal, Canada, 2009.
- 253 Wallace, D.L. 1983. Comment. *J. Amer. Statist. Assoc.* 78, 569 – 576.

254 **Caption for figure**

Figure 1: Measures of agreement for network communities based on livestock movements of cattle, sheep and pigs in Great Britain (2003-4). Means and standard errors of measures for all possible combinations of four-week periods are shown, stratified by time difference in periods (1 to 24). Shown are the adjusted Rand index R (solid line) and entropy measure AMI (dashed line) for the ‘core’ network, with AMI for the entire network (dotted line). Probability \bar{P} coincided with R and is not shown.

255 **Caption for animation file [electronic]**

- 256 Animation of livestock movement network communities in Great Britain, 2003-4; nodes
257 (sites) sharing the same community label are indicated by the same colour.