1 **Title:** Inferring species interactions from co-occurrence data with Markov networks

2 **Author:** David J. Harris: Population Biology; 1 Shields Avenue, Davis CA, 95616

3 **Abstract:** Inferring species interactions from co-occurrence data is one of the most
4 controversial tasks in community ecology. One difficulty is that a single pairwise interaction
5 can ripple through an ecological network and produce surprising indirect consequences. For
6 example, the negative correlation between two competing species can be reversed in the
7 presence of a third species that outcompetes both of them. Here, I apply models from
8 statistical physics, called Markov networks or Markov random fields, that can predict the
9 direct and indirect consequences of any possible species interaction matrix. Interactions in
10 these models can be estimated from observed co-occurrence rates via maximum likelihood,
11 controlling for indirect effects. Using simulated landscapes with known interactions, I
12 evaluated Markov networks and six existing approaches. Markov networks consistently
13 outperformed the other methods, correctly isolating direct interactions between species pairs
14 even when indirect interactions or abiotic factors largely overpowered them. Two
15 computationally efficient approximations, which controlled for indirect effects with partial
16 correlations or generalized linear models, also performed well. Null models showed no
17 evidence of being able to control for indirect effects, and reliably yielded incorrect inferences
18 when such effects were present.

21 **Introduction**

22 To the extent that nontrophic species interactions (such as competition) affect community

23 assembly, ecologists might expect to find signatures of these interactions in species

24 composition data (MacArthur 1958, Diamond 1975). Despite decades of work and several

1

₂₅ major controversies, however (Lewin 1983, Strong et al. 1984, Connor et al. 2013), existing

₂₆ methods for detecting competition's effects on community structure are unreliable (Gotelli

₂₇ and Ulrich 2009). In particular, species' effects on one another can become lost in a web of

₂₈ indirect effects. For example, the competitive interaction between the two shrub species in

₂₉ Figure 1A is obscured by their shared tendency to occur in unshaded areas (Figure 1B).

₃₀ While ecologists have long known that indirect effects can overwhelm direct ones (Levine

₃₁ 1976), most methods for drawing inferences from co-occurrence data do not control for these

₃₂ effects (e.g. Diamond 1975, Strong et al. 1984, Gotelli and Ulrich 2009, Veech 2013, Pollock

₃₃ et al. 2014). As a result, ecologists do not have tools that allow them to isolate direct

₃₄ interactions from indirect effects.

₃₅     While competition doesn't reliably reduce co-occurrence rates at the whole-landscape

₃₆ level (as most methods assume), it does still leave a signal in the data (Figure 1C). After

₃₇ controlling for the presence of the tree species (e.g. by splitting the data set into shaded and

₃₈ unshaded sites or by using a model that estimates conditional relationships amongs species),

₃₉ the two shrubs do have a negative association, with fewer co-occurrences in unshaded areas

₄₀ than would be expected if they occurred independently of one another.

₄₁     Following Azaele et al. (2010), this paper shows that Markov networks (undirected

₄₂ graphical models also known as Markov random fields; Murphy 2012) can provide a

₄₃ framework for understanding the landscape-level consequences of pairwise species

₄₄ interactions, and for estimating them from observed presence-absence matrices. Markov

₄₅ networks have been used in diverse scientific fields for decades, from physics (where nearby

₄₆ particles interact magnetically; Cipra 1987) to spatial statistics (where adjacent grid cells

₄₇ have correlated values; Harris 1974, Gelfand et al. 2005). While community ecologists

₄₈ explored some related approaches in the 1980's (Whittam and Siegel-Causey 1981), they used

49 severe approximations that led to unintelligible results (e.g. "probabilities" greater than one;

50 Gilpin and Diamond 1982).

51     Below, I demonstrate Markov networks' ability to produce exact predictions about the

52 direct and indirect consequences of an interaction matrix, and also to make inferences about

53 the interaction matrix based on co-occurrence rates. Using simulated data sets where the

54 "true" interactions are known, I compare this approach with several existing methods. Finally,

55 I discuss opportunities for extending the approach presented here to other problems in

56 community ecology, e.g. quantifying the overall effect of species interactions on occurrence

57 rates (Roughgarden 1983) and disentangling the effects of biotic versus abiotic interactions

58 on species composition (Harris 2016).

59 **Methods**

60 **Markov networks.** Markov networks provide a framework for translating back and forth

61 between the conditional (all-else-equal) relationships among species (Figure 1C) and the

62 kinds of species assemblages that these relationships produce. Here, I show how a set of

63 conditional relationships can determine species composition. Methods for estimating

64 conditional relationships from data are discussed in the next section.

65     A Markov network defines the relative probability of observing a given vector of

66 species-level presences (1s) and absences (0s), $\vec{y}$ at a site, as

67
$$p(\vec{y}; \alpha, \beta) \propto exp(\sum_i \alpha_i y_i + \sum_{<ij>} \beta_{ij} y_i y_j),$$

68 where the second sum is over all $\frac{1}{2}n(n-1)$ pairs of $n$ species. In this model, $\alpha_i$ is an

69 intercept term determining the amount that the presence of species $i$ contributes to the

70 log-probability of $\vec{y}$; it directly controls the prevalence of species $i$. Similarly, $\beta_{ij}$ is the

71 amount that the co-occurrence of species $i$ and species $j$ contributes to the log-probability; it

72 determines the conditional relationship between two species, i.e. the probability that they

3

73  will be found together, after controlling for the other species in the network (Figure 2A,

74  Figure 2B). For example, if $\beta_{ij} = +2$, then each species' odds of occurrence would be $e^2$ times

75  higher when the other one is present (as compared with otherwise equivalent sites). The

76  relative probability of a presence-absence vector increases when positively-associated species

77  co-occur and decreases when negatively-associated species co-occur. As a result, the model

78  tends—all else equal—to produce assemblages where many positively-associated species pairs

79  co-occur and few negatively-associated pairs do (just as an ecologist might expect).

80  **Estimating $\alpha$ and $\beta$ coefficients from presence-absence data.** In the previous section,

81  the values of $\alpha$ and $\beta$ were known and the goal was to make predictions about possible species

82  assemblages. In practice, however, ecologists will often need to estimate the parameters from

83  an observed co-occurrence matrix (i.e. from a set of independent $\vec{y}$ vectors indicating which

84  species are present at each site on the landscape). When Equation 1 can be normalized

85  (Figure 2B), one can find exact maximum likelihood estimates for $\alpha$ and $\beta$ by numerically

86  optimizing $p(\vec{y}|\alpha, \beta)$. Fully-observed networks like the ones considered here have unimodal

87  likelihood surfaces (Murphy 2012), so optimizers will always find the global optimum.

88  When the number of species is larger than about 30, noralizing Equation 1 can become

89  intractable, and researchers will either need to approximate it (Lee and Hastie 2012) or

90  approximate its gradient (Harris 2016) if they want to fit a Markov network model. For the

91  analyses presented here, where the number of species did not exceed 20, approximations were

92  not necessary. Instead, I used the rosalia package (Harris 2015a) for the R programming

93  language (R Core Team 2015) to calculate $p(\vec{y}; \alpha, \beta)$ and its gradients exactly (Murphy 2012);

94  the package passes these functions to the "BFGS" method in R's general-purpose optimizer,

95  which finds values for $\alpha$ and $\beta$.

96  **Simulated landscapes.** I simulated several sets of landscapes using known parameters so

4

97 that model estimates could be compared with "true" values. The first set of landscapes

98 included the three competing species shown in Figure 1. For each of 1000 replicates, I

99 generated a landscape with 100 sites by sampling from a probability distribution defined by

100 the figure's interaction coefficients (Appendix 1). Each of the methods described below was

101 then evaluated on its ability to correctly infer that the two shrub species competed with one

102 another, despite their frequent co-occurrence in unshaded areas.

103      I then simulated landscapes with up to 20 interacting species at 25, 200, or 1600 sites

104 using three increasingly complex models (50 replicates for each combination of size and

105 model; Appendix 2). The simplest set of simulated landscapes were generated with Gibbs

106 samples from Equation 1. For each replicate, I randomly drew the "true" $\beta$ coefficient

107 magnitudes from an exponential distribution with rate 1 so that most species pairs interacted

108 negligibly but a few interactions were strong enough to propagate through the newtwork. I

109 randomly assigned 25% of the interactions to be positive; the remainder were negative.

110      The next set of landscapes provided a way to assess each method's ability to identify

111 direct interactions in the presence of environmental heterogeneity. Here, the $\beta$ coefficients

112 were calculated as above, but each species' $\alpha$ value depended linearly on two environmental

113 factors, which were drawn from independent Gaussians for each site. Once the local $\alpha$ values

114 were calculated, independent Gibbs samplers determined the species composition for each site

115 based on Equation 1.

116      In the final set of landscapes, I simulated each species' abundance (instead of just

117 presence/absence); furthermore, interactions between species occurred on a per-capita basis

118 in these simulations (i.e. each species' effect on the others is proportional to its abundance).

119 To prevent runaway mutualisms leading to infinite abundance, all interaction coefficients were

120 negative in these simulations. Good performance on these landscapes would indicate some

5

121 robustness to the mechanistic details of species interactions.

122 **Recovering species interactions from simulated data.** I compared seven techniques

123 for determining the sign and strength of the associations between pairs of species from

124 simulated data (Appendix 3). First, I used the rosalia package (Harris 2015a) to fit Markov

125 network models, as described above. For the analyses with 20 species, a weakly-informative

126 regularizer (equivalent to a logistic prior with location 0 and scale 2) ensured that the

127 estimates were always finite (Appendix 3).

128     I also evaluated six alternative methods: five from the existing literature, plus a novel

129 combination of two of these methods. The first alternative interaction metric was the sample

130 correlation between species' presence-absence vectors, which summarizes their marginal

131 association. Next, I used partial correlations, which summarize species' conditional

132 relationships. This approach, which is closely related to linear regression, is common in

133 molecular biology (Friedman et al. 2008), but is rare in ecology (see Albrecht and Gotelli

134 (2001) and Faisal et al. (2010) for two exceptions). In the context of non-Gaussian data, the

135 partial correlation (or partial covariance) can be thought of as a computationally efficient

136 approximation to the full Markov network model (Loh and Wainwright 2013). Partial

137 correlations are undefined for landscapes with perfectly-correlated species pairs, so I used the

138 regularized estimate provided by the corpcor package's `pcor.shrink` function with the

139 default settings (Schäfer et al. 2014).

140     The third alternative, generalized linear models (GLMs), also provide a computationally

141 efficient approximation to the Markov network (Lee and Hastie 2012). Following Faisal et al.

142 (2010), I fit regularized logistic regression models (Gelman et al. 2008) for each species, using

143 the other species as predictors. This produced two interaction estimates for each species pair

144 (one for the effect of species $i$ on species $j$ and one for the reverse). These two estimates were

6

<sub>145</sub> very tightly correlated (mean Pearson correlation of 0.95, Appendix 3); their arithmetic mean

<sub>146</sub> provided a consensus estimate of the overall interaction.

<sub>147</sub> The next method used the Pairs software described in Gotelli and Ulrich (2009). This

<sub>148</sub> program simulates new landscapes from a null model that retains the row and column sums

<sub>149</sub> of the original matrix (Strong et al. 1984) and calculates $Z$-scores to summarize a species

<sub>150</sub> pair's deviation from this null.

<sub>151</sub> The last two estimators used the latent correlation matrix estimated by the BayesComm

<sub>152</sub> package (Golding and Harris 2015) in order to evaluate the recent claim that the correlation

<sub>153</sub> coefficients estimated by "joint species distribution models" provide an accurate assessment

<sub>154</sub> of species' pairwise interactions (Pollock et al. 2014, see also Harris 2015b). In addition to

<sub>155</sub> using the posterior mean correlation (Pollock et al. 2014), I also used the posterior mean

<sub>156</sub> *partial* correlation, which should control better for indirect effects.

<sub>157</sub> **Evaluating model performance.** For the simulated landscapes based on Figure 1, I

<sub>158</sub> assessed whether each method's test statistic indicated a positive or negative relationship

<sub>159</sub> between the two shrubs (Appendix 1). For the null model (Pairs), I calculated statistical

<sub>160</sub> significance using its $Z$-score. For the Markov network, I used the Hessian matrix to generate

<sub>161</sub> approximate confidence intervals.

<sub>162</sub> For the larger landscapes, I evaluated the relationship between each method's estimates

<sub>163</sub> and the "true" interaction strengths. To ensure that the different test statistics

<sub>164</sub> (e.g. correlations versus $Z$ scores) were on a common scale, I rescaled them using linear

<sub>165</sub> regression through the origin. I then calculated the proportion of variance explained for

<sub>166</sub> different combinations of model type and landscape size (compared with a baseline model

<sub>167</sub> that assumed all interaction strengths to be zero).

<sub>168</sub> For the null model and the Markov network, the probability of rejecting the null

7

169  hypothesis of zero interaction was estimated across a range of "true" interaction strengths

170  using a kernel smoother (Appendix 4). The probability of rejection when the "true" value of

171  $\beta$ was zero was defined as the Type I error rate. Because the coefficients' interpretation is

172  different for the abundance simulations than for the other two types, its error rates were not

173  analyzed in this way.

174  **Results**

175  **Three species.** As shown in Figure 1, the marginal relationship between the two shrub

176  species was positive—despite their competition for space at a mechanistic level—due to

177  indirect effects of the dominant tree species. As a result, the correlation between these

178  species was positive in 94% of replicates, and the randomization-based null model falsely

179  reported positive associations 100% of the time. Worse, more than 98% of these false

180  conclusions were statistically significant. The partial correlation and Markov network

181  estimates, on the other hand, each correctly isolated the direct negative interaction between

182  the shrubs from their positive indirect interaction 94% of the time (although the confidence

183  intervals overlapped zero in most replicates).

184  **Twenty species.** In general, each model's performance was highest for large landscapes

185  with simple assembly rules and no environmental heterogeneity (Figure 3). Despite some

186  variability across contexts, the rank ordering across methods was very consistent. In

187  particular, the four methods that controlled for indirect effects (the Markov network, the

188  generalized linear models, and the two partial correlation-based methods) always matched or

189  outperformed those that did not. The Markov network consistently performed best of all. As

190  anticipated by Lee and Hastie (2012), generalized linear models closely approximated the

191  Markov network estimates (Figure 4A), especially when the data sets were very large (Figure

192  3). As reviewed in Gotelli and Ulrich (2009), however, most analyses in this field of ecology

193 involve fewer than 50 sites; in this context, the gap between the methods was larger. As

194 shown in Appendix 4, the standard errors associated with the estimates in Figure 3 are small

195 (less than 0.01), so the differences among methods should not be attributed to sampling error.

196   Of the methods that did not control for indirect effects, Figure 3 shows that simple

197 correlation coefficients provided a more reliable indicator of species' true interaction strengths

198 than either the joint species distribution model (BayesComm) or the null model (Pairs). 95%

199 of the variance the Pairs test statistic was explained by correlation coefficients (controlling

200 for landscape size; Figure 4B); much of the remaining variance is due to sampling error.

201   Finally, we can evaluate the models' inferential statistics (focusing on the first two

202 simulation types, where the interaction coefficients are easiest to interpret). The Markov

203 network's Type I error rate was 0.02 for simulations that matched the model's assumptions,

204 and 0.14 for simulations that included environmental heterogeneity (see Appendix 4 for

205 confidence interval coverage across a range of $\beta_{ij}$ values). In contrast, the null model's Type I

206 error rates were 0.30 and 0.51, respectively—far higher than the nominal 0.05 rate. Figure

207 4C shows, across a range of true interaction strengths, the probability that the null model or

208 the Markov network will predict the wrong sign of the interaction with 95% confidence. The

209 null model makes such errors more than 8 times as often as the Markov network, even though

210 it only rejects the null hypothesis twice as often overall (Appendix 4). The Markov network's

211 errors were also more concentrated around 0, as it never misclassified strong interactions like

212 the null model did (Figure 4C).

### Discussion

214 The results presented above show that Markov networks can reliably recover species' pairwise

215 interactions from species composition data, even for cases where environmental heterogeneity

216 and indirect interactions cause ecologists' typical null modeling approaches to reliably fail.

₂₁₇ Partial correlations and generalized linear models can both provide computationally efficient

₂₁₈ approximations, but with somewhat lower accuracy (especially for typically-sized data sets

₂₁₉ with small numbers of sites; Gotelli and Ulrich 2009). The difference in accuracy may be

₂₂₀ even larger for real data sets than for the simulated landscapes in Figure 3; linear

₂₂₁ approximations to the Markov network make larger errors when the interaction matrix is

₂₂₂ structured (e.g. due to guilds or trophic levels; Loh and Wainwright 2013). Similarly, the

₂₂₃ separate generalized linear models for each species can severely overfit in some cases (Lee and

₂₂₄ Hastie 2012). The full Markov network should thus be preferred to the approximations when

₂₂₅ it is computationally tractable.

₂₂₆ Compositional data only contains enough degrees of freedom to estimate one interaction

₂₂₇ per species pair (Schmidt and Murphy 2012), so none of these methods can identify the exact

₂₂₈ nature of the pairwise interactions (e.g. which species in a positively-associated pair is

₂₂₉ facilitating the other). To estimate asymmetric interactions, such as commensalism or

₂₃₀ predation, ecologists could use time series, behavioral observations, manipulative

₂₃₁ experiments, or natural history. These other sources of information could also be used to

₂₃₂ augment the likelihood function with a more informative prior distribution, reducing

₂₃₃ ecologists' error and uncertainty relative to Figure 3's results.

₂₃₄ Markov networks have enormous potential to improve our understanding of species

₂₃₅ interactions. In particular, they make many fewer errors than existing approaches, and can

₂₃₆ make precise statements about the conditions where indirect interactions will overwhelm

₂₃₇ direct ones. They also provide a simple answer to the question of how competition should

₂₃₈ affect a species' overall prevalence, which has important implications for community-level

₂₃₉ modeling (Strong et al. 1984). Specifically, Equation 1 can be used to calculate the expected

₂₄₀ prevalence of a species in the absence of biotic influences as $e^{\alpha_i}/(e^0 + e^{\alpha_i})$. Competition's

10

effect on prevalence can then be estimated by comparing this value with the observed

prevalence (e.g. comparing Figure 2D with Figure 2C). This novel quantitative result

undermines most of our null models, which unreasonably assume that prevalence would be

the exactly same in the absence of competition as it is in the observed data (Roughgarden

1983).

Markov networks—particularly the Ising model for binary networks—are very well

understood, having been studied for nearly a century (Cipra 1987). Tapping into this

framework would thus allow ecologists to take advantage of into a vast set of existing

discoveries and techniques for dealing with indirect effects, stability, and alternative stable

states. Numerous extensions to the basic network are possible as well. For example, the

states of the interaction network can be modeled as a function of the local abiotic

environment (Lee and Hastie 2012, Harris 2016), which would lead to a better understanding

of the interplay between biotic and abiotic effects on community structure. Alternatively,

models could allow one species to alter the relationship between two other species (Whittam

and Siegel-Causey 1981, Tjelmeland and Besag 1998, cf Bruno et al. 2003).

Finally, the results presented here have important implications for ecologists' continued

use of null models for studying species interactions. When the non-null backdrop is not

controlled for, Type I error rates can skyrocket, the apparent sign of the interaction can

change, and null models can routinely produce misleading inferences (Figure 1, Figure 4C,

Gotelli and Ulrich (2009)). Null and neutral models can be useful for clarifying our thinking

(Harris et al. 2011, Xiao et al. 2015), but deviations from a given null model must be

interpreted with care (Roughgarden 1983). Even in small networks with three species, it may

simply not be possible to implicate specific ecological processes like competition by rejecting

a general-purpose null (Gotelli and Ulrich 2009), especially when the test statistic is

11

265  effectively just a correlation coefficient (Figure 4B).

266   Controlling for indirect effects via simultaneous estimation of multiple ecological

267  parameters seems like a much more promising approach: to the extent that the models'

268  relative performance on real data sets is similar to the range of results shown in Figure 3,

269  scientists in this field could often triple their explanatory power by switching from null

270  models to Markov networks (or increase it nearly as much with linear or generalized linear

271  approximations). Regardless of the methods ecologists ultimately choose, controlling for

272  indirect effects could clearly improve our understanding of species' direct effects on one

273  another and on community structure.

281  **References:**

282  Albrecht, M., and N. J. Gotelli. 2001. Spatial and temporal niche partitioning in grassland

283   ants. Oecologia 126:134–141.

284  Azaele, S., R. Muneepeerakul, A. Rinaldo, and I. Rodriguez-Iturbe. 2010. Inferring plant

285   ecosystem organization from species occurrences. Journal of theoretical biology

286   262:323–329.

287  Bruno, J. F., J. J. Stachowicz, and M. D. Bertness. 2003. Inclusion of facilitation into

288   ecological theory. Trends in Ecology & Evolution 18:119–125.

Cipra, B. A. 1987. An introduction to the ising model. American Mathematical Monthly 94:937–959.

Connor, E. F., M. D. Collins, and D. Simberloff. 2013. The checkered history of checkerboard distributions. Ecology 94:2403–2414.

Diamond, J. M. 1975. The island dilemma: Lessons of modern biogeographic studies for the design of natural reserves. Biological conservation 7:129–146.

Faisal, A., F. Dondelinger, D. Husmeier, and C. M. Beale. 2010. Inferring species interaction networks from species abundance data: A comparative evaluation of various statistical and machine learning methods. Ecological Informatics 5:451–464.

Friedman, J., T. Hastie, and R. Tibshirani. 2008. Sparse inverse covariance estimation with the graphical lasso. Biostatistics 9:432–441.

Gelfand, A. E., A. M. Schmidt, S. Wu, J. A. Silander, A. Latimer, and A. G. Rebelo. 2005. Modelling species diversity through species level hierarchical modelling. Journal of the Royal Statistical Society: Series C (Applied Statistics) 54:1–20.

Gelman, A., A. Jakulin, M. G. Pittau, and Y.-S. Su. 2008. A weakly informative default prior distribution for logistic and other regression models. The Annals of Applied Statistics 2:1360–1383.

Gilpin, M. E., and J. M. Diamond. 1982. Factors contributing to non-randomness in species co-occurrences on islands. Oecologia 52:75–84.

Golding, N., and D. J. Harris. 2015. BayesComm: Bayesian community ecology analysis.

Gotelli, N. J., and W. Ulrich. 2009. The empirical bayes approach as a tool to identify non-random species associations. Oecologia 162:463–477.

Harris, D. J. 2015a. Rosalia: Exact inference for small binary markov networks. r package version 0.1.0. Zenodo. http://dx.doi.org/10.5281/zenodo.17808.

313  Harris, D. J. 2015b. Generating realistic assemblages with a joint species distribution model.

314  Methods in Ecology and Evolution.

315  Harris, D. J. 2016. Estimating species interactions in large, abiotically structured

316  communities with markov networks and stochastic approximation. *in* Multi-process

317  statistical modeling of species' joint distributions. Figshare.

318  http://doi.org/10.6084/m9.figshare.3114226.v1.

319  Harris, D. J., K. G. Smith, and P. J. Hanly. 2011. Occupancy is nine-tenths of the law:

320  Occupancy rates determine the homogenizing and differentiating effects of exotic species.

321  The American naturalist 177:535.

322  Harris, T. E. 1974. Contact interactions on a lattice. The Annals of Probability 2:969–988.

323  Lee, J. D., and T. J. Hastie. 2012. Learning mixed graphical models.

324  Levine, S. H. 1976. Competitive interactions in ecosystems. The American Naturalist

325  110:903–910.

326  Lewin, R. 1983. Santa rosalia was a goat. Science 221:636–639.

327  Loh, P.-L., and M. J. Wainwright. 2013. Structure estimation for discrete graphical models:

328  Generalized covariance matrices and their inverses. The Annals of Statistics 41:3022–3049.

329  MacArthur, R. H. 1958. Population ecology of some warblers of northeastern coniferous

330  forests. Ecology 39:599–619.

331  Murphy, K. P. 2012. Machine learning: A probabilistic perspective. The MIT Press.

332  Pollock, L. J., R. Tingley, W. K. Morris, N. Golding, R. B. O'Hara, K. M. Parris, P. A. Vesk,

333  and M. A. McCarthy. 2014. Understanding co-occurrence by modelling species

334  simultaneously with a joint species distribution model (JSDM). Methods in Ecology and

335  Evolution:n/a–n/a.

336  R Core Team. 2015. R: A language and environment for statistical computing. R Foundation

337  for Statistical Computing, Vienna, Austria.

338  Roughgarden, J. 1983. Competition and theory in community ecology. The American

339   Naturalist 122:583–601.

340  Schäfer, J., R. Opgen-Rhein, V. Zuber, M. Ahdesmäki, A. P. D. Silva, and K. Strimmer.

341   2014. Corpcor: Efficient estimation of covariance and (partial) correlation.

342  Schmidt, M., and K. Murphy. 2012. Modeling discrete interventional data using directed

343   cyclic graphical models. arXiv preprint arXiv:1205.2617.

344  Strong, D. R., D. Simberloff, L. G. Abele, and A. B. Thistle. 1984. Ecological communities:

345   Conceptual issues and the evidence. Princeton University Press.

346  Tjelmeland, H., and J. Besag. 1998. Markov random fields with higher-order interactions.

347   Scandinavian Journal of Statistics 25:415–433.

348  Veech, J. A. 2013. A probabilistic model for analysing species co-occurrence. Global Ecology

349   and Biogeography 22:252–260.

350  Whittam, T. S., and D. Siegel-Causey. 1981. Species interactions and community structure

351   in alaskan seabird colonies. Ecology 62:1515–1524.

352  Xiao, X., D. J. McGlinn, and E. P. White. 2015. A strong test of the maximum entropy

353   theory of ecology. The American Naturalist 185:E70–E80.

354  **Figure Captions**

355  **Figure 1: A.** A small network of three competing species. The tree (top) tends not to

356  co-occur with either of the two shrub species, as indicated by the strongly negative coefficient

357  linking them. The two shrub species also compete with one another, but more weakly

358  (circled coefficient). **B.** In spite of the competitive interactions between the two shrub species,

359  their shared tendency to occur in locations without trees makes their occurrence vectors

360  positively correlated (circled). **C.** Controlling for trees with a conditional (all-else-equal)

15

approach such as a partial correlation or a Markov network leads to correct identification of the negative shrub-shrub interaction (circled). See Appendix 1 and the results for "three species" for more details.

**Figure 2: A.** A small Markov network, defined by its $\alpha$ and $\beta$ values. The abiotic environment favors the occurrence of each species ($\alpha > 0$), particularly species 2 ($\alpha_2 > \alpha_1$). The negative $\beta_{12}$ coefficient is consistent with competition between the two species. **B.** The coefficients determine the probabilities of all four possible presence-absence combinations for Species 1 and Species 2. $\alpha_1$ is added to the exponent whenever Species 1 is present ($y_1 = 1$), but not when it is absent ($y_1 = 0$). Similarly, the exponent includes $\alpha_2$ only when species 2 is present ($y_2 = 1$), and includes $\beta_{12}$ only when both are present ($y_1 y_2 = 1$). The normalizing constant $Z$, ensures that the four probabilities sum to 1. In this case, $Z$ is about 18.5. **C.** The expected frequencies of all possible co-occurrence patterns between the two species of interest, as calculated in the previous panel. **D.** Without competition (i.e. with $\beta_{12} = 0$, each species would occur more often.

**Figure 3:** Proportion of variance in interaction coefficients explained by each method versus number of sampled locations across the three simulation types. For the null model (Pairs), two outliers with $|Z| > 1000$ were manually adjusted to $|Z| = 50$ to mitigate their detrimental influence on $R^2$ (Appendix 5).

**Figure 4: A.** The Markov network's estimated interaction coefficients were generally very similar to the GLM estimates. **B.** The null model's estimates typically matched the (negative) correlation coefficient, after controlling for landscape size. **C.** For any given interaction strength, the null model was much more likely to misclassify its sign with 95% confidence than the Markov network was. As with the other analyses based on inferential statistics, this panel only shows data from the first two simulation types.

16