

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 is capable of outcompeting both of them. Here, I apply
8 models from statistical physics, called Markov networks or Markov random fields, that can
9 predict the direct and indirect consequences of any possible species interaction matrix.

10 Interactions in these models can also be estimated from observed co-occurrence rates via
11 maximum likelihood, controlling for indirect effects. Using simulated landscapes with known
12 pairwise interaction strengths, I evaluated Markov networks and six existing approaches.

13 The Markov networks consistently outperformed other methods, correctly isolating direct
14 interactions between species pairs even when indirect interactions or abiotic factors largely
15 overpowered them. Two computationally efficient approximations, based on controlling for
16 indirect effects with linear or generalized linear models, also performed well. Indirect effects
17 reliably caused a common null modeling approach to produce incorrect inferences, however.

18 **Key words:** Ecological interactions; Occurrence data; Species associations; Markov network;
19 Markov random field; Ising model; Biogeography; Presence–absence matrix; Null model

20 **Introduction**

21 To the extent that nontrophic species interactions (such as competition) affect community
22 assembly, ecologists might expect to find signatures of these interactions in species
23 composition data (MacArthur 1958, Diamond 1975). Despite decades of work and several

Inferring species interactions

24 major controversies, however (Lewin 1983, Strong et al. 1984, Connor et al. 2013), existing
25 methods for detecting competition's effects on community structure are unreliable (Gotelli
26 and Ulrich 2009). In particular, species' effects on one another can become lost in a web of
27 indirect effects. For example, the competitive interaction between the two shrub species in
28 Figure 1A is obscured by their shared tendency to occur in unshaded areas (Figure 1B).
29 While ecologists have long known that indirect effects can overwhelm direct ones at the
30 landscape level (Levine 1976), the vast majority of our methods for drawing inferences from
31 observational data do not control for these effects (e.g. Diamond 1975, Strong et al. 1984,
32 Gotelli and Ulrich 2009, Veech 2013, Pollock et al. 2014). To the extent that indirect
33 interactions like those in Figure 1 are generally important, existing methods will not provide
34 much evidence regarding species interactions.

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). For
37 example, after partitioning the data set into shaded and unshaded sites, there will be
38 co-occurrence deficits in each subset that wouldn't otherwise be apparent. More generally,
39 controlling for other species in the network will often be important for obtaining reliable
40 estimates of direct (conditional, or all-else-equal) effects. This kind of precision is difficult to
41 obtain from null models, which only test the most extreme possible hypothesis: that *all*
42 direct and indirect interactions are exactly zero. Nevertheless, null models have dominated
43 this field for more than three decades (Strong et al. 1984, Gotelli and Ulrich 2009).

44 Following Azaele et al. (2010), this paper shows that Markov networks (undirected
45 graphical models also known as Markov random fields; Murphy 2012) can provide a
46 framework for understanding the landscape-level consequences of pairwise species
47 interactions, and for estimating them from observed presence-absence matrices. Markov

Inferring species interactions

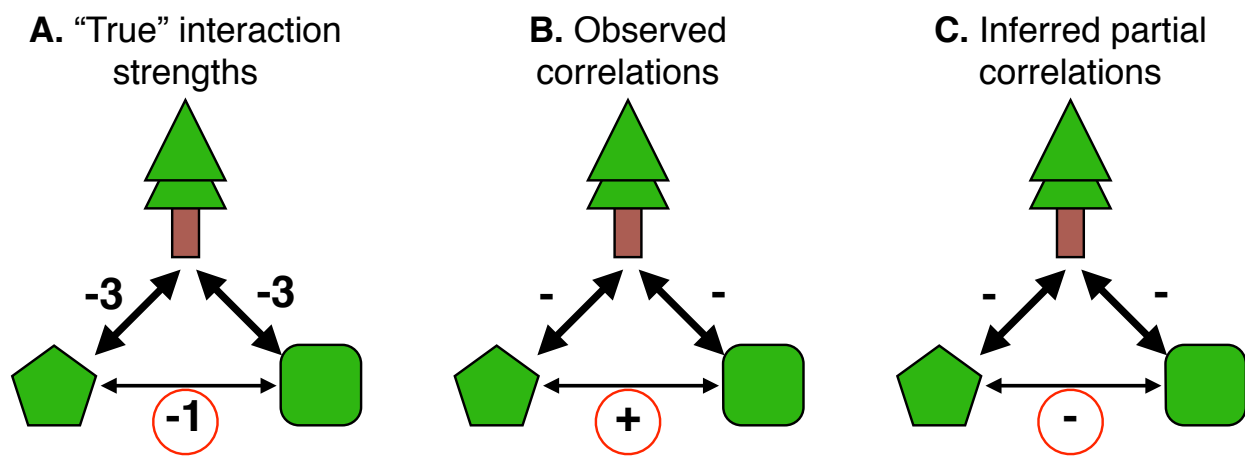


Figure 1: **A.** A small network of three competing species. The tree (top) tends not to co-occur with either of the two shrub species, as indicated by the strongly negative coefficient linking them. The two shrub species also compete with one another, but more weakly (circled coefficient). **B.** In spite of the competitive interactions between the two shrub species, their shared tendency to occur in locations without trees makes their occurrence vectors positively correlated (circled). **C.** Controlling for trees with a conditional (all-else-equal) approach such as a partial covariance or a Markov network leads to correct identification of the negative shrub-shrub interaction (circled).

Inferring species interactions

48 networks have been used in many scientific fields in similar contexts for decades, from
49 physics (where nearby particles interact magnetically; Cipra 1987) to spatial statistics (where
50 adjacent grid cells have correlated values; Harris 1974, Gelfand et al. 2005). While
51 community ecologists explored some related approaches in the 1980's (Whittam and
52 Siegel-Causey 1981), they used severe approximations that led to unintelligible results (e.g.
53 “probabilities” greater than one; Gilpin and Diamond 1982).

54 Below, I demonstrate Markov networks' ability to produce exact predictions about the
55 direct and indirect consequences of an interaction matrix, and also to make inferences about
56 the species interactions that contributed to an observed set of co-occurrences. Using
57 simulated data sets where the “true” interactions are known, I compare this approach with
58 several existing methods. Finally, I discuss opportunities for extending the approach
59 presented here to other problems in community ecology, e.g. quantifying the overall effect of
60 species interactions on occurrence rates (Roughgarden 1983) and disentangling the effects of
61 biotic versus abiotic interactions on species composition (Pollock et al. 2014).

62 **Methods**

63 **Markov networks.** Markov networks provide a framework for translating back and forth
64 between the conditional (all-else-equal) relationships among species (Figure 1C) and the
65 kinds of species assemblages that these relationships produce. Here, I show how a set of
66 conditional relationships can determine species composition. Methods for estimating
67 conditional relationships from data are discussed in the next section.

68 A Markov network defines the relative probability of observing a given vector of
69 species-level presences (1s) and absences (0s), \vec{y} at a site, as

$$70 \quad p(\vec{y}; \alpha, \beta) \propto \exp\left(\sum_i \alpha_i y_i + \sum_{\langle ij \rangle} \beta_{ij} y_i y_j\right),$$

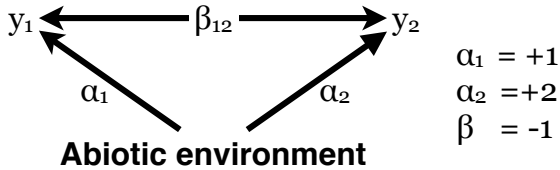
Inferring species interactions

71 where the second sum is over all $\frac{1}{2}n(n-1)$ pairs of n species. In this model, α_i is an
72 intercept term determining the amount that the presence of species i contributes to the
73 log-probability of \vec{y} ; it directly controls the prevalence of species i . Similarly, β_{ij} is the
74 amount that the co-occurrence of species i and species j contributes to the log-probability; it
75 controls the conditional relationship between two species, i.e. the probability that they will
76 be found together, after controlling for the other species in the network (Figure 2A, Figure
77 2B). For example, if $\beta_{ij} = +2$, then each species' odds of occurrence would be e^2 times higher
78 when the other one is present (as compared with otherwise equivalent sites). The relative
79 probability of a presence-absence vector increases when positively-associated species co-occur
80 and decreases when negatively-associated species co-occur. As a result, the model tends—all
81 else equal—to produce assemblages where many positively-associated species pairs co-occur
82 and few negatively-associated pairs do (just as an ecologist might expect). When all else is
83 *not* equal (e.g. Figure 1, where the presence of one competitor is associated with release from
84 another competitor), then predicting species' overall co-occurrence rates can be more
85 complicated, and may require summing over the different possible assemblages (Figure 2B).

86 **Estimating α and β coefficients from presence-absence data.** In the previous
87 section, the values of α and β were known and the goal was to make predictions about
88 possible species assemblages. In practice, however, ecologists will often need to estimate the
89 parameters from an observed co-occurrence matrix (i.e. from a set of independent \vec{y} vectors
90 indicating which species are present at each site on the landscape). When the number of
91 species is reasonably small (less than about 30), one can find exact maximum likelihood
92 estimates for all of the α and β coefficients given a presence-absence matrix by numerically
93 optimizing $p(\vec{y}; \alpha, \beta)$. Fully-observed Markov networks like the ones considered here have
94 unimodal likelihood surfaces (Murphy 2012), ensuring that this procedure will converge on

Inferring species interactions

A. Defining a small network



B. Solving the network

$$\begin{aligned}
 & \alpha_1 y_1 \quad \alpha_2 y_2 \quad \beta_{y_1 y_2} \\
 P[0 \ 0] &= e^{(+0 \ +0 \ +0)} / Z = e^{(0)} / Z \\
 P[y_1 \ 0] &= e^{(+1 \ +0 \ +0)} / Z = e^{(1)} / Z \\
 P[0 \ y_1] &= e^{(+0 \ +2 \ +0)} / Z = e^{(2)} / Z \\
 P[y_1 y_2] &= e^{(+1 \ +2 \ -1)} / Z = e^{(2)} / Z \\
 & (e^0 + e^1 + e^2 + e^2) / Z = 1
 \end{aligned}$$

C. Expected co-occurrence rates

		Species 1	
		Absent	Present
Species 2	Absent	5%	15%
	Present	40%	40%

D. Expected co-occurrence without competition ($\beta_{12}=0$)

		Species 1	
		Absent	Present
Species 2	Absent	3%	9%
	Present	24%	64%

Figure 2: **A.** A small Markov network, defined by its α and β values. The abiotic environment favors the occurrence of each species ($\alpha > 0$), particularly species 2 ($\alpha_2 > \alpha_1$). The negative β_{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. α_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 α_2 only when species 2 is present ($y_2 = 1$), and includes β_{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.

Inferring species interactions

95 the global maximum. I used the *rosalia* package (Harris 2015a) for the R programming
96 language (R Core Team 2015) to calculate $p(\vec{y}; \alpha, \beta)$ and its gradients (Murphy 2012); the
97 package passes these functions to the “BFGS” method in R’s general-purpose optimizer,
98 which then estimates the Markov network parameters.

99 **Simulated landscapes.** I simulated several sets of landscapes using known parameters to
100 evaluate different statistical methods’ performance. The first set of landscapes included the
101 three competing species shown in Figure 1. For each of 1000 replicates, I generated a
102 landscape with 100 sites by sampling from a probability distribution defined by the figure’s
103 interaction coefficients (Appendix 1). Each of the methods described below was then
104 evaluated on its ability to correctly infer that the two shrub species competed with one
105 another, despite their frequent co-occurrence.

106 I then generated landscapes with up to 20 interacting species at 25, 200, or 1600 sites
107 using three increasingly complex models (50 replicates for each combination of size and
108 model; see Appendix 2 for details). I randomly drew the “true” coefficient magnitudes for
109 each replicate landscape from exponential distributions so that most species pairs interacted
110 negligibly but a few pairs interacted strongly enough that their effects could propagate
111 indirectly to other species in the network.

112 The first set of 20-species landscapes, like the landscapes with three species, were
113 generated directly from a Markov network to ensure that the model could recover the
114 parameters used to generate the “observed” co-occurrence data. Then, I added two
115 environmental factors that varied from location to location across the simulated landscapes,
116 and simulated a new set of co-occurrence data so that species’ α coefficients depended on the
117 local environment. The latter set of simulated landscapes provide an important test of the
118 methods’ ability to distinguish co-occurrence patterns that were generated from pairwise

Inferring species interactions

119 biotic interactions from those that were generated by external forces like abiotic
120 environmental filtering. This task was made especially difficult because—as with most
121 analyses of presence-absence data for co-occurrence patterns—the inference procedure did
122 not have access to any information about the environmental or spatial variables that helped
123 shape the landscape (cf Connor et al. 2013). I generated the final set of landscapes with an
124 abundance-based model that included per-capita interaction rates instead of per-species
125 interaction rates.

126 **Recovering species interactions from simulated data.** I compared seven techniques
127 for determining the sign and strength of the associations between pairs of species from
128 simulated data (Appendix 3). First, I used the *rosalia* package (Harris 2015a) to fit Markov
129 network models, as described above. For the analyses with 20 species, a weak regularizer
130 (equivalent to a logistic prior with scale 2) ensured that the estimates were always finite.

131 I also evaluated six alternative methods: five from the existing literature, plus a novel
132 combination of two of these methods. The first alternative interaction metric was the sample
133 correlation between species' presence-absence vectors, which summarizes their marginal
134 association. Next, I used partial correlations, which summarize species' conditional
135 relationships. This approach is common in molecular biology (Friedman et al. 2008), but is
136 rare in ecology (see Albrecht and Gotelli (2001) and Faisal et al. (2010) for two exceptions).
137 In the context of non-Gaussian data, the partial correlation can be thought of as a
138 computationally efficient approximation to the full Markov network model (Loh and
139 Wainwright 2013). Because partial correlations are undefined for landscapes with
140 perfectly-correlated species pairs, I used a regularized estimate provided by the *corpcor*
141 package's `pcor.shrink` function with the default settings (Schäfer et al. 2014).

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

Inferring species interactions

143 efficient approximation to the Markov network (Lee and Hastie 2012). Following Faisal et al.
144 (2010), I fit regularized logistic regression models (Gelman et al. 2008) for each species, using
145 the other species on the landscape as predictors. This produced two interaction estimates for
146 each species pair; one for the effect of species i on species j and one for the reverse. These
147 two coefficients are not identifiable from the data, however (Schmidt and Murphy 2012), so I
148 used their average as an overall measure of the pairwise relationship.

149 The next method, described in Gotelli and Ulrich (2009), involved simulating new
150 landscapes from a null model that retains the row and column sums of the original matrix
151 (Strong et al. 1984). I used the Z -scores computed by the Pairs software described in Gotelli
152 and Ulrich (2009) as my null model-based estimator of species interactions.

153 The last two estimators used the latent correlation matrix estimated by the BayesComm
154 package (Golding and Harris 2015) in order to evaluate the recent claim that the correlation
155 coefficients estimated by “joint species distribution models” provide an accurate assessment
156 of species’ pairwise interactions (Pollock et al. 2014, see also Harris 2015b). In addition to
157 using the posterior mean correlation (Pollock et al. 2014), I also used the posterior mean
158 *partial* correlation, which should control better for indirect effects.

159 **Evaluating model performance.** For the simulated landscapes based on Figure 1, I
160 assessed whether each method’s test statistic indicated a positive or negative relationship
161 between the two shrubs (Appendix 1). For the null model (Pairs), I calculated statistical
162 significance using its Z -score. For the Markov network, I used the Hessian matrix to
163 generate approximate confidence intervals.

164 For the larger landscapes, I evaluated the relationship between each method’s estimates
165 and the “true” interaction strengths. To ensure that the different test statistics
166 (e.g. correlations versus Z scores) were on a common scale, I rescaled them using linear

Inferring species interactions

167 regression through the origin. I then calculated the proportion of variance explained for
168 different combinations of model type and landscape size (compared with a null model that
169 assumed all interaction strengths to be zero).

170 **Results**

171 **Three species.** As shown in Figure 1, the marginal relationship between the two shrub
172 species was positive—despite their competition for space at a mechanistic level—due to
173 indirect effects of the dominant tree species. As a result, the correlation between these
174 species was positive in 94% of replicates, and the randomization-based null model falsely
175 reported positive associations 100% of the time. Worse, more than 98% of these false
176 conclusions were statistically significant. The partial correlation and Markov network
177 estimates, on the other hand, each correctly isolated the direct negative interaction between
178 the shrubs from their positive indirect interaction 94% of the time (although the confidence
179 intervals overlapped zero in most replicates).

180 **Twenty species.** In general, each model’s performance was highest for large landscapes
181 with simple assembly rules and no environmental heterogeneity (Figure 3). Despite some
182 variability across contexts, the rank ordering across methods was very consistent. In
183 particular, the four methods that controlled for indirect effects (the Markov network, the
184 generalized linear models, and the two partial correlation-based methods) always matched or
185 outperformed those that did not. The Markov network consistently performed best of all. As
186 anticipated by Lee and Hastie (2012), generalized linear models closely approximated the
187 Markov network estimates (Figure 4A), especially when the data sets were very large (Figure
188 3). As reviewed in Gotelli and Ulrich (2009), however, most analyses in this field of ecology
189 involve fewer than 50 sites, and the gap between the Markov network and GLMs is larger in
190 this context. As shown in Appendix 4, the standard errors associated with the estimates in

Inferring species interactions

191 Figure 3 are small (less than 0.01), so the differences among methods should not be
 192 attributed to sampling error.

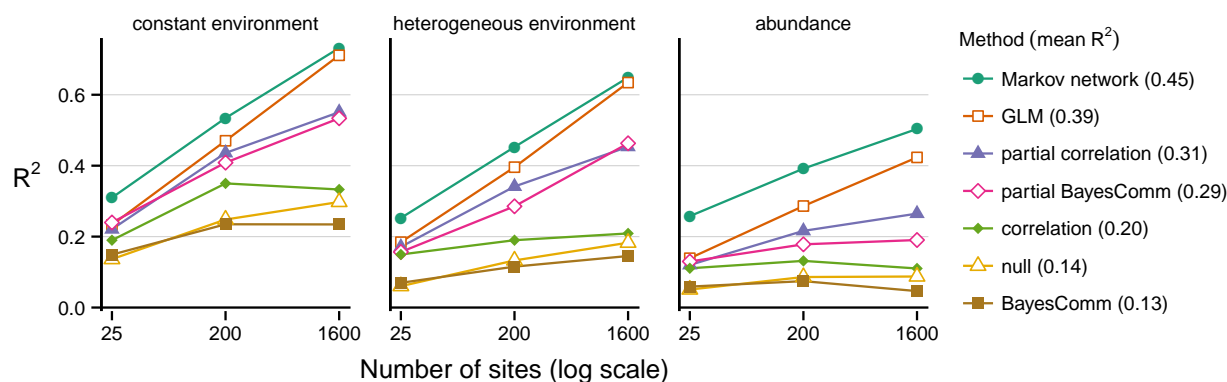


Figure 3: Proportion of variance in interaction coefficients explained by each method versus number of sampled locations across the three simulation types. For the null model (Pairs), two outliers with $|Z| > 1000$ were manually adjusted to $|Z| = 50$ to mitigate their detrimental influence on R^2 (Appendix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

193 Of the methods that did not control for indirect effects, Figure 3 shows that simple
 194 correlation coefficients provided a more reliable indicator of species' true interaction
 195 strengths than either the joint species distribution model (BayesComm) or the null model
 196 (Pairs). The estimates from these approaches were tightly correlated (after controlling for
 197 the size of the landscape) suggesting that the null model only contains a noisy version of the
 198 same information that could be obtained more easily and interpretably with simple
 199 correlation coefficients (Figure 4B).

200 Finally, we can evaluate the models' statistical inferences (focusing on the first two
 201 simulation types, for which the true interaction rates are easiest to interpret). The Markov
 202 network's approximate Type I error rate (defined here as the probability that 0 fell outside
 203 the 95% confidence interval for a pair of species where $|\beta_{ij}| < 0.1$) depended on the
 204 simulation type: 0.02 for simulations that matched the model's assumptions, versus 0.14 for
 205 simulations that included environmental heterogeneity (see Appendix 4 for confidence

Inferring species interactions

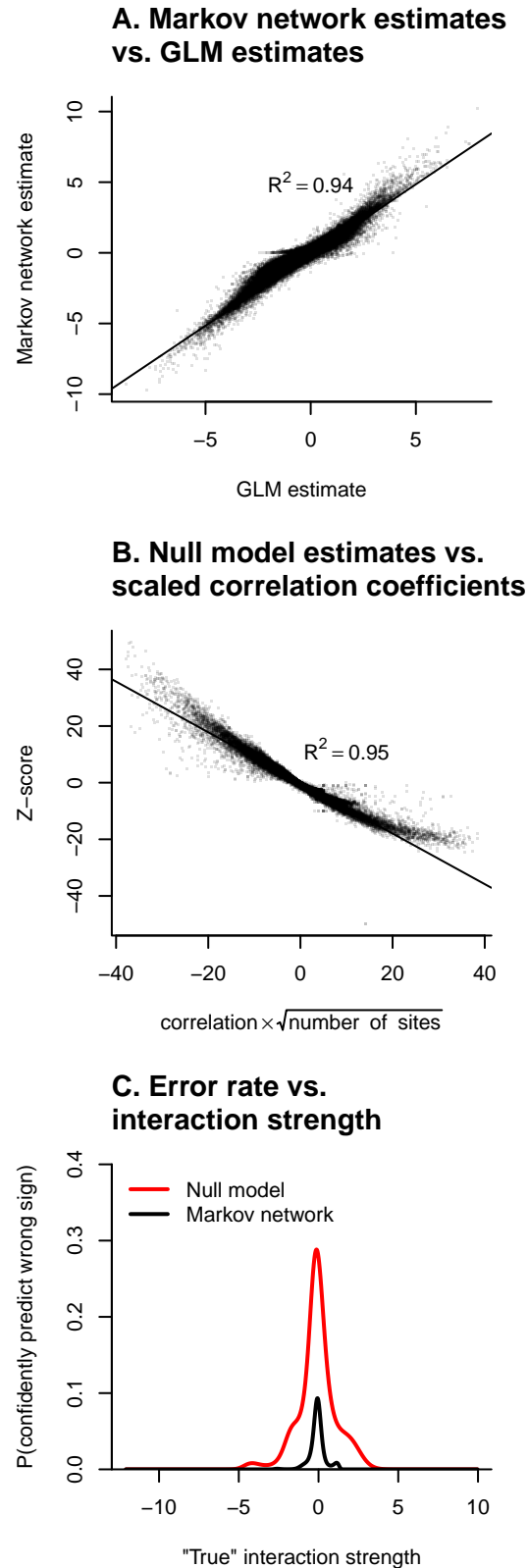


Figure 4: **A.** The Markov network's estimated interaction coefficients were generally very similar to the GLM estimates. **B.** The null model's estimates typically matched the (negative) correlation coefficient, after controlling for landscape size. **C.** For any given interaction strength, the null model was much more likely to misclassify its sign with 95% confidence than the Markov network was.

Inferring species interactions

206 interval coverage across a range of β_{ij} values). The null model's Type I error rates were 0.30
207 and 0.50 for the constant and heterogeneous landscapes, respectively—far higher than the
208 nominal 0.05 rate. Figure 4C shows, across a range of true interaction strengths, the
209 probability that the null model or the Markov network will predict the wrong sign of the
210 interaction with 95% confidence. The null model makes such errors more than 5 times as
211 often as the Markov network, even though it only reject the null hypothesis twice as often
212 overall (Appendix 4). The Markov network's errors were also more concentrated around 0, so
213 it never misclassified strong interactions like the null model did.

214 **Discussion**

215 The results presented above show that Markov networks can reliably recover species'
216 pairwise interactions from species composition data, even for cases where environmental
217 heterogeneity and indirect interactions cause ecologists' typical null modeling approaches to
218 reliably fail. Partial covariances and generalized linear models can both provide
219 computationally efficient approximations, but with somewhat lower accuracy (especially for
220 typically-sized data sets with small numbers of sites; Gotelli and Ulrich 2009). The difference
221 in accuracy may be larger for real data sets than for the simulated landscapes in Figure 3,
222 however; linear approximations to the Markov network make larger errors when the
223 interaction matrix is structured (e.g. due to guilds or trophic levels; Loh and Wainwright
224 2013). Similarly, the separate generalized linear models for each species can severely overfit
225 in some cases (Lee and Hastie 2012). The full Markov network should thus be preferred to
226 the approximations when it is computationally tractable.

227 Compositional data only contains enough degrees of freedom to estimate one interaction
228 per species pair (Schmidt and Murphy 2012), so none of these methods can identify the
229 exact nature of the pairwise interactions (e.g. which species in a positively-associated pair is

Inferring species interactions

230 facilitating the other). To estimate asymmetric interactions, such as commensalism or
231 predation, ecologists could use time series, behavioral observations, manipulative
232 experiments, or natural history. These other sources of information could also be used to
233 augment the likelihood function with a more informative prior distribution, reducing
234 ecologists' error and uncertainty relative to Figure 3's results.

235 Despite their limitations, Markov networks have enormous potential to improve
236 ecological understanding. In particular, they make many fewer errors than existing
237 approaches, and can make precise statements about the conditions where indirect
238 interactions will overwhelm direct ones. They also provide a simple answer to the question of
239 how competition should affect a species' overall prevalence, which has important implications
240 for community-level modeling (Strong et al. 1984). Specifically, Equation 1 can be used to
241 calculate the expected prevalence of a species in the absence of biotic influences as
242 $e^\alpha / (e^0 + e^\alpha)$. Competition's effect on prevalence can then be estimated by comparing this
243 value with the observed prevalence (e.g. comparing Figure 2D with Figure 2C). This novel
244 quantitative result conflicts with most of our null models, which unreasonably assume that
245 prevalence would be the exactly same in the absence of competition as it is in the observed
246 data (Roughgarden 1983).

247 Markov networks—particularly the Ising model for binary networks—are very well
248 understood, having been studied for nearly a century (Cipra 1987). Tapping into this
249 framework would thus allow ecologists to take advantage of into a vast set of existing
250 discoveries and techniques for dealing with indirect effects, stability, and alternative stable
251 states. Numerous extensions to the basic network are possible as well. For example, the
252 states of the interaction network can be modeled as a function of the local abiotic
253 environment (Lee and Hastie 2012), which would help incorporate networks of biotic

Inferring species interactions

254 interactions into species distribution models (Pollock et al. 2014) and lead to a better
255 understanding of the interplay between biotic and abiotic effects on community structure.
256 Alternatively, models could allow one species to alter the relationship between two other
257 species (Tjelmeland and Besag 1998, cf Bruno et al. 2003).

258 Finally, the results presented here have important implications for ecologists' continued
259 use of null models for studying species interactions. Null and neutral models can be useful
260 for clarifying our thinking (Harris et al. 2011, Xiao et al. 2015), but deviations from a
261 particular null model must be interpreted with care (Roughgarden 1983). Even in small
262 networks with three species, it may simply not be possible to implicate specific ecological
263 processes like competition by rejecting a general-purpose null (Gotelli and Ulrich 2009),
264 especially when the test statistic is effectively just a correlation coefficient (Figure 4B).
265 When the non-null backdrop is not controlled for, Type I error rates can skyrocket, the
266 apparent sign of the interaction can change, and null models can routinely produce
267 misleading inferences (Figure 1, Figure 4C, Gotelli and Ulrich (2009)).

268 Controlling for indirect effects via simultaneous estimation of multiple ecological
269 parameters seems like a much more promising approach: to the extent that the models'
270 relative performance on real data sets is similar to the range of results shown in Figure 3,
271 scientists in this field could often triple their explanatory power by switching from null
272 models to Markov networks (or increase it nearly as much with linear or generalized linear
273 approximations). Regardless of the methods ecologists ultimately choose, controlling for
274 indirect effects could clearly improve our understanding of species' direct effects on one
275 another and on community structure.

276 **Acknowledgements:** This work benefited greatly from discussions with A. Sih, M. L.
277 Baskett, R. McElreath, R. J. Hijmans, A. C. Perry, and C. S. Tysor. Additionally, A. K.

Inferring species interactions

278 Barner, E. Baldrige, E. P. White, D. Li, D. L. Miller, N. Golding, N. J. Gotelli, C. F.
279 Dormann, and two anonymous reviewers provided very helpful feedback on the text. This
280 research was partially supported by a Graduate Research Fellowship from the US National
281 Science Foundation and by the Gordon and Betty Moore Foundation's Data-Driven
282 Discovery Initiative through Grant GBMF4563 to E. P. White.

283 **References:**

- 284 Albrecht, M., and N. J. Gotelli. 2001. Spatial and temporal niche partitioning in grassland
285 ants. *Oecologia* 126:134–141.
- 286 Azaele, S., R. Munepeerakul, A. Rinaldo, and I. Rodriguez-Iturbe. 2010. Inferring plant
287 ecosystem organization from species occurrences. *Journal of theoretical biology*
288 262:323–329.
- 289 Bruno, J. F., J. J. Stachowicz, and M. D. Bertness. 2003. Inclusion of facilitation into
290 ecological theory. *Trends in Ecology & Evolution* 18:119–125.
- 291 Cipra, B. A. 1987. An introduction to the Ising model. *American Mathematical Monthly*
292 94:937–959.
- 293 Connor, E. F., M. D. Collins, and D. Simberloff. 2013. The checkered history of
294 checkerboard distributions. *Ecology* 94:2403–2414.
- 295 Diamond, J. M. 1975. The island dilemma: Lessons of modern biogeographic studies for the
296 design of natural reserves. *Biological conservation* 7:129–146.
- 297 Faisal, A., F. Dondelinger, D. Husmeier, and C. M. Beale. 2010. Inferring species interaction
298 networks from species abundance data: A comparative evaluation of various statistical and
299 machine learning methods. *Ecological Informatics* 5:451–464.
- 300 Friedman, J., T. Hastie, and R. Tibshirani. 2008. Sparse inverse covariance estimation with
301 the graphical lasso. *Biostatistics* 9:432–441.

Inferring species interactions

- 302 Gelfand, A. E., A. M. Schmidt, S. Wu, J. A. Silander, A. Latimer, and A. G. Rebelo. 2005.
303 Modelling species diversity through species level hierarchical modelling. *Journal of the*
304 *Royal Statistical Society: Series C (Applied Statistics)* 54:1–20.
- 305 Gelman, A., A. Jakulin, M. G. Pittau, and Y.-S. Su. 2008. A Weakly Informative Default
306 Prior Distribution for Logistic and Other Regression Models. *The Annals of Applied*
307 *Statistics* 2:1360–1383.
- 308 Gilpin, M. E., and J. M. Diamond. 1982. Factors contributing to non-randomness in species
309 Co-occurrences on Islands. *Oecologia* 52:75–84.
- 310 Golding, N., and D. J. Harris. 2015. BayesComm: Bayesian Community Ecology Analysis.
311 Gotelli, N. J., and W. Ulrich. 2009. The empirical Bayes approach as a tool to identify
312 non-random species associations. *Oecologia* 162:463–477.
- 313 Harris, D. J. 2015a. Rosalia: Exact inference for small binary Markov networks. R package
314 version 0.1.0. Zenodo. <http://dx.doi.org/10.5281/zenodo.17808>.
- 315 Harris, D. J. 2015b. Generating realistic assemblages with a Joint Species Distribution
316 Model. *Methods in Ecology and Evolution*.
- 317 Harris, D. J., K. G. Smith, and P. J. Hanly. 2011. Occupancy is nine-tenths of the law:
318 Occupancy rates determine the homogenizing and differentiating effects of exotic species.
319 *The American naturalist* 177:535.
- 320 Harris, T. E. 1974. Contact Interactions on a Lattice. *The Annals of Probability* 2:969–988.
- 321 Lee, J. D., and T. J. Hastie. 2012, May. Learning Mixed Graphical Models.
- 322 Levine, S. H. 1976. Competitive Interactions in Ecosystems. *The American Naturalist*
323 110:903–910.
- 324 Lewin, R. 1983. Santa Rosalia Was a Goat. *Science* 221:636–639.
- 325 Loh, P.-L., and M. J. Wainwright. 2013. Structure estimation for discrete graphical models:

Inferring species interactions

- 326 Generalized covariance matrices and their inverses. *The Annals of Statistics* 41:3022–3049.
- 327 MacArthur, R. H. 1958. Population ecology of some warblers of northeastern coniferous
328 forests. *Ecology* 39:599–619.
- 329 Murphy, K. P. 2012. *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- 330 Pollock, L. J., R. Tingley, W. K. Morris, N. Golding, R. B. O’Hara, K. M. Parris, P. A. Vesk,
331 and M. A. McCarthy. 2014. Understanding co-occurrence by modelling species
332 simultaneously with a Joint Species Distribution Model (JSDM). *Methods in Ecology and*
333 *Evolution*:n/a–n/a.
- 334 R Core Team. 2015. *R: A Language and Environment for Statistical Computing*. R
335 Foundation for Statistical Computing, Vienna, Austria.
- 336 Roughgarden, J. 1983. *Competition and Theory in Community Ecology*. The American
337 Naturalist 122:583–601.
- 338 Schäfer, J., R. Opgen-Rhein, V. Zuber, M. Ahdesmäki, A. P. D. Silva, and K. Strimmer.
339 2014. Corpcor: Efficient Estimation of Covariance and (Partial) Correlation.
- 340 Schmidt, M., and K. Murphy. 2012. Modeling Discrete Interventional Data using Directed
341 Cyclic Graphical Models. arXiv preprint arXiv:1205.2617.
- 342 Strong, D. R., D. Simberloff, L. G. Abele, and A. B. Thistle. 1984. *Ecological communities:*
343 *Conceptual issues and the evidence*. Princeton University Press.
- 344 Tjelmeland, H., and J. Besag. 1998. Markov Random Fields with Higher-order Interactions.
345 *Scandinavian Journal of Statistics* 25:415–433.
- 346 Veech, J. A. 2013. A probabilistic model for analysing species co-occurrence. *Global Ecology*
347 *and Biogeography* 22:252–260.
- 348 Whittam, T. S., and D. Siegel-Causey. 1981. Species Interactions and Community Structure
349 in Alaskan Seabird Colonies. *Ecology* 62:1515–1524.

Inferring species interactions

- 350 Xiao, X., D. J. McGlinn, and E. P. White. 2015. A strong test of the Maximum Entropy
351 Theory of Ecology. *The American Naturalist* 185:E70–E80.