

A network approach for inferring species associations from co-occurrence data

Naia Morueta-Holme, Benjamin Blonder, Brody Sandel, Brian J. McGill, Robert K. Peet, Jeffrey E. Ott, Cyrille Violle, Brian J. Enquist, Peter M. Jørgensen and Jens-Christian Svenning

N. Morueta-Holme (morueta-holme@berkeley.edu), B. Blonder, B. Sandel and J.-C. Svenning, Section for Ecoinformatics and Biodiversity, Dept of Bioscience, Aarhus Univ., Ny Munkegade 114, DK-8000 Aarhus C, Denmark. NM-H also at: Integrative Biology, Univ. of California – Berkeley, CA, USA. BB also at: Environmental Change Inst., School of Geography and the Environment, Univ. of Oxford, Oxford, UK. – B. J. McGill, School of Biology and Ecology/Sustainability Solutions Initiative, Univ. of Maine, Orono, ME, USA. – R. K. Peet and J. E. Ott, Dept of Biology, Univ of North Carolina, Chapel Hill, NC 27599-3280, USA. – C. Violle, CEFU UMR 5175, CNRS – Univ. de Montpellier – Univ. Paul-Valéry Montpellier – EPHE -1919 route de Mende, FR-34293 Montpellier, Cedex 5, France. – B. J. Enquist, Dept of Ecology and Evolutionary Biology, Univ. of Arizona, Tucson, AZ, USA. – P. M. Jørgensen, Missouri Botanical Garden, PO Box 299, St Louis, MO 63166-0299, USA.

Positive and negative associations between species are a key outcome of community assembly from regional species pools. These associations are difficult to detect and can be caused by a range of processes such as species interactions, local environmental constraints and dispersal. We integrate new ideas around species distribution modeling, covariance matrix estimation, and network analysis to provide an approach to inferring non-random species associations from local- and regional-scale occurrence data. Specifically, we provide a novel framework for identifying species associations that overcomes three challenges: 1) correcting for indirect effects from other species, 2) avoiding spurious associations driven by regional-scale distributions, and 3) describing these associations in a multi-species context. We highlight a range of research questions and analyses that this framework is able to address. We show that the approach is statistically robust using simulated data. In addition, we present an empirical analysis of > 1000 North American tree communities that gives evidence for weak positive associations among small groups of species. Finally, we discuss several possible extensions for identifying drivers of associations, predicting community assembly, and better linking biogeography and community ecology.

An unresolved question in ecology is how species assemblages come together at different spatial scales from regional species pools to form local communities (Weiher et al. 2011, HilleRisLambers et al. 2012). Some subsets of species may be consistently associated or dissociated in assemblages due to multiple processes including chance, species interactions, and the indirect effect of showing the same (or opposite) response to environmental conditions. While the effect of broad-scale environmental filtering and dispersal limitation can be assessed using niche modeling techniques (De Marco et al. 2008, Guisan and Rahbek 2011, Normand et al. 2011), the influence of species interactions and local environmental conditions on community assembly is more challenging to measure and incorporate into predictions (Kissling et al. 2012, Pottier et al. 2013, Thuiller et al. 2013, Wisz et al. 2013, Araújo and Rozenfeld 2014). If species interact with each other or share resource or scenopoetic requirements (Soberón 2007) not adequately described

by broad-scale models, then stacking independent species distribution models to predict species assemblages (sensu Guisan and Rahbek 2011, Calabrese et al. 2014) will provide misleading predictions of fine-scale community assembly. Thus, a better understanding of species associations across scales could improve predictions of the dynamics of local community composition in changing environments.

The goal of this paper is to improve the tools needed to detect interspecific associations from co-occurrence data. We first briefly describe the development of co-occurrence methods and then draw from different lines of research to present a more complete and flexible general framework for inferring species associations that overcomes multiple challenges faced by previous approaches.

From experiments to co-occurrence methods

Efforts to infer species associations and their role in structuring communities have a long history. Traditionally, associations have been derived from small-scale field observations

NM-H and BB contributed equally to this project.

(MacArthur 1958, Bullock et al. 2000) or manipulative experiments (Stewart and Aldrich 1951). Such methods may be suitable for studying associations among a few species at local scales. However, addressing whether small-scale associations occur consistently across large regions raises major practical issues. Experimental data is time-consuming to obtain even at small scales (e.g. Bullock et al. 2000, Callaway et al. 2002b), so these approaches are unfeasible for large numbers of species because the number of possible interactions grows rapidly with assemblage size. For example, $n = 100$ species have 4950 possible pairwise interactions and 161 700 possible three-way interactions.

An alternative approach for detecting associations is to gather occurrence data from field observations and analyze species co-occurrence patterns. Such approaches were widely used during discussions of the ‘checkerboard’ paradigm for forbidden combinations of species on islands, with extensive debate over the need for null models to compare observed association patterns to random expectation (Diamond 1975, Connor and Simberloff 1979). The recent increase in availability of occurrence data has brought a renewed interest in using correlations and algorithms to infer species associations and the mechanisms behind them (Bruehlheide 2000, Blick and Burns 2009, Blois et al. 2014). Since the 1970s, statistical approaches have been refined to simultaneously analyze co-occurrence patterns of multiple species pairs (Gotelli and Ulrich 2010). The main idea is to create a community matrix where rows are species, columns are sites and elements represent the observed presence/absence or abundance of each species at each site. The matrix is then compared to a set of randomized matrices in order to detect non-random co-occurrence patterns (Connor and Simberloff 1979, Gotelli and Ulrich 2010). Modern updates to null-model-based co-occurrence approaches can test the effects of environmental drivers, species interactions, or both in structuring communities; for approaches based on randomized null models see Gotelli et al. 2010, Ulrich et al. 2012; for approaches based on analytical null models see Araújo et al. 2011, Veech 2013.

Three missing ideas

One challenge that most previous co-occurrence approaches ignore is the potential effect of other species on a particular pairwise association, i.e. indirect effects (Brown et al. 2004, Harris 2015). If two competing species share a positive (or negative) relationship to a third species, their occurrences could be positively correlated, and thus an indirect effect (correlation) is inferred when the true effect (partial correlation) is actually negative (Brown et al. 2004, Schäfer and Strimmer 2005). A similar idea is implemented in ecology for joint species distribution models (JSDMs; Pollock et al. 2014), where species associations are inferred after accounting for the environment. However, JSDMs usually do not resolve indirect effects from other species (but see the inversion approach of Harris 2015). The indirect association problem is well known in other fields, such as association detection in large genomic and cell-signaling datasets. One solution is to use Gaussian graphical models (Schäfer and Strimmer 2005) or modifications of them, which estimate

the partial correlations, e.g. between each gene pair after taking into account the remaining genes (Dobra et al. 2004, Friedman et al. 2008). The approach ensures estimation of conditional instead of joint associations and has recently been extended to ecological association networks (Harris 2015).

A second challenge is the need for more robust null models. Species are not distributed randomly across sites, but rather have regional geographic distributions that are constrained by climate and dispersal limitations. In the past, most null models have been defined by simply resampling the observed species \times site community matrix \mathbf{O} (Gotelli and Ulrich 2010, Borthagaray et al. 2014). However, this approach fails to account for the additional broad-scale constraints, introducing unrealistic null expectations of spatial independence across randomized sites that disregard spatial autocorrelation in species’ distributions (Legendre 1993, Lennon 2000). In the framework we present here, we suggest that alternative null models can be defined to simulate a specific process. By incorporating regional structure that is contingent on e.g. climate-based species distribution models, we can thus gain more confidence that associations are the outcome of species interactions rather than shared environmental requirements.

A third missing idea is the incorporation of network theory. Network approaches have become common in the study of protein interactions, social structure, etc. (Newman 2010), but have not been applied widely to non-bipartite species association networks. The central idea is that any individual association between species can be better understood in the context of the network of associations between all other species. For example, species with many associations may be preferentially linked to those with few associations, suggesting a scenario where some species act as hubs or keystones. Thus, the network of species associations is a useful way to visualize the community as well as quantify changes in species and multi-species patterns via node- and network-level statistics to answer more sophisticated questions about associations. Network statistics at node-level provide insight into species’ roles. For example, the unweighted degree of each species characterizes the number of its association partners. The ratio between the number of positive and negative associations per species can provide a measure of a species’ role in the network (e.g. as an attractive ‘aggregator’, if a species has more positive than negative associations, or a repulsive ‘segregator’ of other species, if otherwise). Network-level statistics also provide insight into the overall structure of species assemblages. For instance, modularity gives a measure of the overall structure of the network, indicating the amount of division of the network into clusters of nodes that are densely connected to each other, but sparsely connected to nodes in other clusters (Newman 2010, Borthagaray et al. 2014). Higher modularity indicates that groups of species are more likely to be mutually associated. Additionally, overall non-random numbers of links assigned to each node in the network can be tested by comparing the degree distribution to a binomial distribution, which in the limit of many species is equivalent to a chi-square test for deviation from a Poisson distribution.

The approach we suggest thus builds on recent attempts to infer species associations from occurrence data and network theory (Borthagaray et al. 2014). Here we pro-

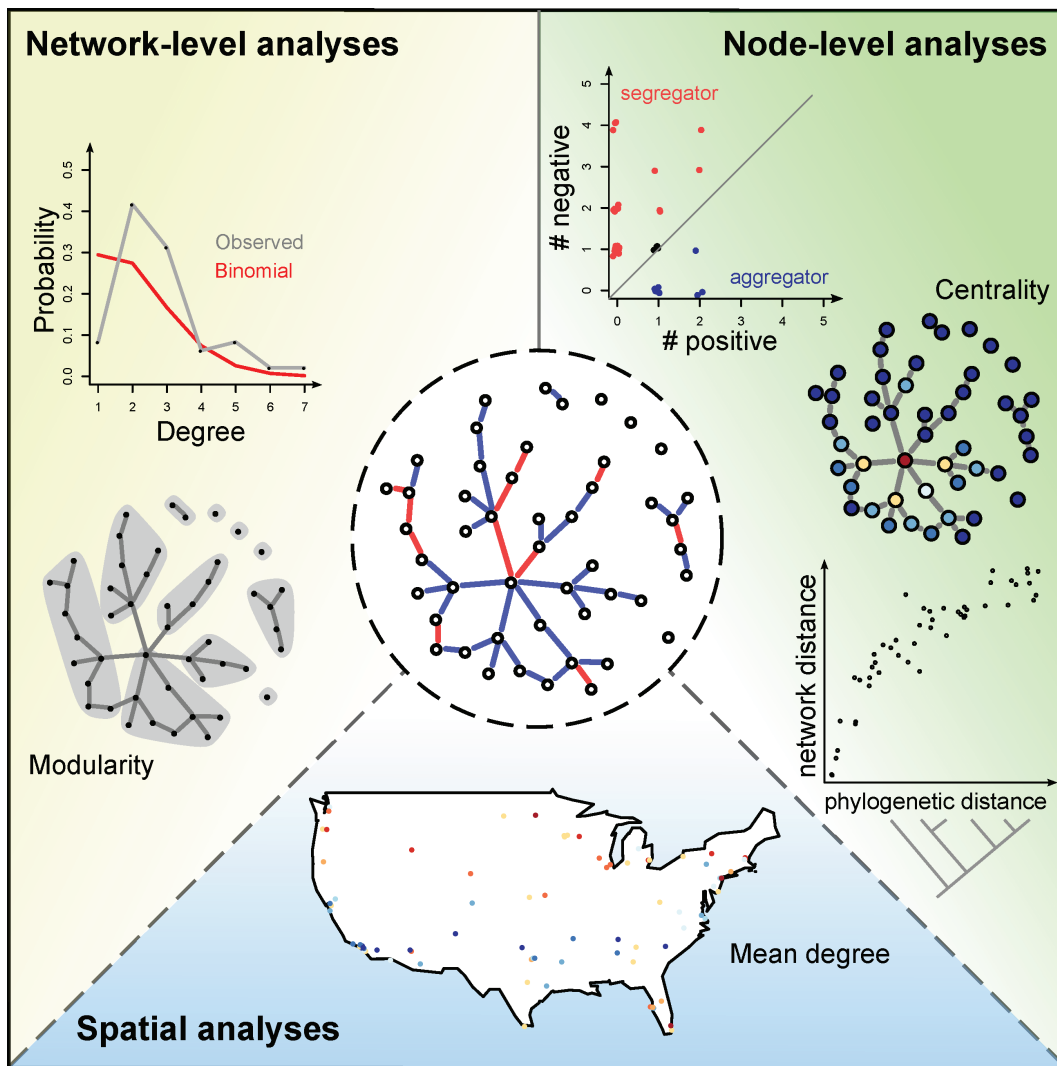
vide a more flexible framework that includes both positive and negative associations, and has the ability to test for deviations from a range of regional-scale null models. The framework is useful for systems in which associations are possible across any pair of the study species (i.e. not bipartite networks such as pollinator–plant interactions, Bascompte 2010). We currently implement the framework for symmetric associations (i.e. (+/+), (0,0) or (-,-)). For asymmetric associations like predation (+/-) or commensalism (+/0), see Discussion.

Interpreting associations

Our approach can identify multi-species modules (i.e. non-random groups) of positively or negatively associated species and assess the importance of particular species in shaping the modules. If we focus on plants, positive associations, or

aggregations of species, may be caused by biotic interactions such as facilitation between nurse trees and seedlings, nitrogen fixation by certain species improving soil fertility for other species, or through shared pollinators, seed dispersers, or facultative mutualisms with endophytic fungi (Afkhami et al. 2014). Alternatively, positive associations may be indicative of shared local environmental requirements or effects of stabilizing niche differences (Chesson 2000, Lasky et al. 2014), or reflect historical dispersal dynamics such as the expansion from glacial refugia (Svenning and Skov 2007). Negative associations, or segregation of species, may also be driven by biotic interactions through competition, or alternatively reflect different local requirements (e.g. variation in microclimate or edaphic conditions).

Once we have inferred the patterns of species associations, we can proceed to determine which hypothesized underlying drivers are most important (Box 1), with the aim of



Box 1. Examples of possible analyses to explore species associations based on the network. Network-level analyses (left) give information on the overall structure of the network. Examples include 1) testing for overall-nonrandom number of links comparing the degree distribution to e.g. a binomial distribution, and 2) quantifying modularity, i.e. groups of species more likely to be associated. Node-level analyses (right) serve to 1) identify the roles of individual species (e.g. as aggregators or segregators by looking at the proportion of negative and positive links per species), 2) quantify the centrality of each species, or 3) test for correlations of node metrics such as network distance and hypothesized drivers such as phylogenetic or functional trait distance. Finally, a hybrid of network and node-level measures can be used for spatial analyses (bottom) to test for trends in e.g. mean degree of communities across space and their correlation to climate gradients.

understanding the principles that cause associations to occur. Functional trait predictors may help identify mechanisms. For instance, if positive associations reflect complementary niches where each species occupies a different niche in relation to resource use (Bolnick et al. 2011), then phylogenetic (Webb et al. 2002) or functional (Weiher et al. 2011) distances should be largest between positively associated species (Violle et al. 2011). Likewise, by inferring spatial patterns of community-weighted mean associations (Box 1) we can test whether associations are influenced by climate gradients, and are e.g. more common in warmer areas (Brown et al. 1996, Schleuning et al. 2012), or whether positive associations become more prevalent at higher elevations (Callaway et al. 2002a).

A framework for detecting associations

Drawing on the multiple lines of research described above, we here propose a general framework to uncover species associations from co-occurrence patterns. The associations we identify are those not explained by broad-scale climate gradients. The main idea of our approach is to pair broad- and local-scale co-occurrence information with the Gaussian graphical model approach, spatially explicit regional null models (Guisan and Zimmermann 2000, Gotelli and Ulrich 2010, Borthagaray et al. 2014) and network theory (Newman 2010). Our framework is implemented as the ‘netassoc’ R package.

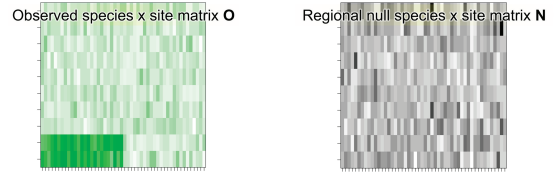
We first show how the co-occurrence framework can infer a network of species associations by pairing regional occurrence data with local-scale assemblage data. Second, we perform a simulation analysis demonstrating that the framework has acceptable error rates in realistic use cases. Third, we demonstrate a range of network-based analyses that can describe associations and test hypothesized mechanisms. Lastly, we illustrate the framework with an application to the trees of eastern United States using a large dataset of local co-occurrences and regional occurrences.

In Box 2 we describe the main data inputs and methodological choices for the approach. Briefly, the first step is to identify non-random species associations for an observed and a null dataset. To do this, we compare the partial correlation coefficients inferred for observed local-scale data to those inferred for regional-scale null expectations calculated from independent data. We then interpret these effect sizes within a network framework.

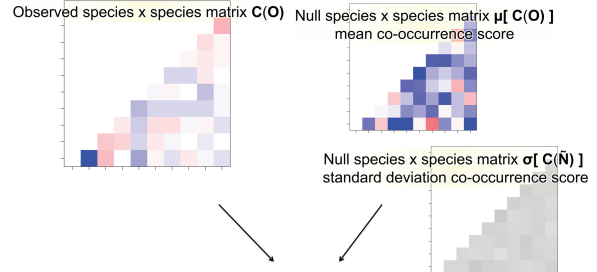
Network construction

To compute the species association network we need 1) the observed co-occurrences, i.e. a set of observed presence–absence or abundance data for n species found in a random sample of m sites, and 2) the expected co-occurrences for the same species and sites based on a null model. The null expectation can be for presence or abundance of the species. From 1) we generate \mathbf{O} , the observed species \times site community matrix. From 2) we generate \mathbf{N} , the expected abundance (or presence) patterns at each local site as predicted by a chosen regional species distribution model

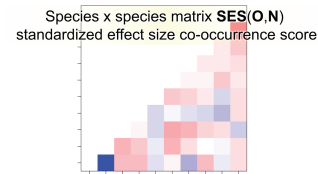
Step 1. Obtain observed and expected community data



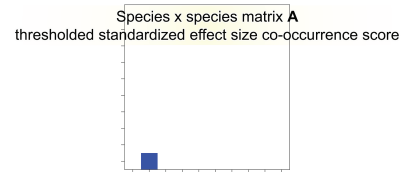
Step 2. Calculate observed and null co-occurrence scores



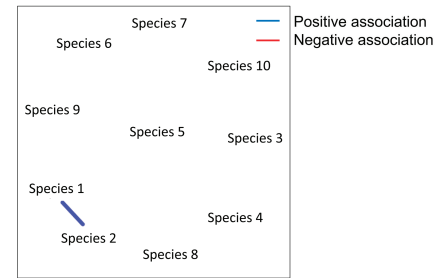
Step 3. Calculate strength and direction of each association



Step 4. Threshold values to detect significant associations



Step 5. Represent association matrix as weighted network



Box 2. Concepts and data flow underlying the framework. Each step illustrates how the network of non-random species association is derived from the observed co-occurrence matrix and the expectation based on a chosen null model. See detailed description in main text.

3). Both the \mathbf{O} and \mathbf{N} matrix will have n rows (species) and m columns (sites).

We first infer the association strength between species i and j as entries A_{ij} in the $n \times n$ matrix \mathbf{A} . We calculate an $n \times n$ covariance matrix Σ for each of \mathbf{O} and \mathbf{N} . From the inverse of this covariance matrix we obtain standard partial correlations between species, i.e. as

$$C_{ij}(\mathbf{M}) = \frac{-\Sigma_{ij}^{-1}(\mathbf{M})}{\sqrt{\Sigma_{ii}^{-1}(\mathbf{M}) \cdot \Sigma_{jj}^{-1}(\mathbf{M})}} \quad (1)$$

where \mathbf{M} is the $n \times m$ input community matrix. C_{ij} represents the effect of species i on species j after correcting for the effects of all other species. It is zero if species i is conditionally independent of species j . Equation 1 represents the fundamental mathematical approach taken when constructing Gaussian graphical models (Schäfer and Strimmer 2005) for inferring linear associations between random variables.

We calculate $C(\mathbf{O})$ as well as for $K \gg 1$ resamples of \mathbf{N} , $C(\tilde{\mathbf{N}})$. This distribution simply simulates a weighted lottery model of community assembly where species enter a community based only on their overall abundance in the regional pool (i.e. the probabilities in \mathbf{N}). To do so, the resamples preserve the total number of individuals within each site, weighting the sampling by the expected abundance of each species based on the original \mathbf{N} matrix.

We then determine if the observed association between each species pair is positive or negative by comparing the observed co-occurrence score to the distribution of expected co-occurrences across $\tilde{\mathbf{N}}$. We calculate a standard effect size $SES_{ij}(\mathbf{O}, \tilde{\mathbf{N}})$, i.e. by rescaling by the mean (μ) and standard deviation (σ) of the null distributions:

$$SES_{ij}(\mathbf{O}, \tilde{\mathbf{N}}) = \frac{C_{ij}(\mathbf{O}) - \mu[C_{ij}(\tilde{\mathbf{N}})]}{\sigma[C_{ij}(\tilde{\mathbf{N}})]} \quad (2)$$

Finally, we distinguish between significant and non-significant associations. We first calculate a two-tailed p-value for species i and j as the number of replicates in which the absolute observed association strength is smaller than the absolute null association strength divided by the total number of replicates. We then correct for multiple comparisons by specifying a false discovery rate, α , and performing a Benjamini–Hochberg correction (Benjamini and Hochberg 1995) on each p-value, producing a new set of p-values p_{ij}^* . The false discovery rate conceptualizes the type I error rate by controlling the expected proportion of false positives, i.e. the proportion of incorrect rejections of the null hypotheses across multiple comparisons.

Together, this process yields a species-by-species adjacency matrix \mathbf{A} with n rows and n columns (number of species):

$$A_{ij} = \text{If} \begin{cases} p_{ij}^* < \alpha & SES_{ij}(\mathbf{O}, \tilde{\mathbf{N}}) \\ p_{ij}^* \geq \alpha & 0 \end{cases} \quad (3)$$

This matrix \mathbf{A} is treated as the adjacency matrix (showing which species are connected to each other) used to define the species association network, such that a significantly positive or negative association between species i and j is established if A_{ij} is nonzero. The network of associations is used for all subsequent analyses.

A few important decision points in the framework

While our statistical framework for describing species associations is general, the user must choose between multiple definitions and parameters specific to the system and taxa

being studied. For example, one important choice in the analysis is the type of null model. The \mathbf{N} matrix shows the expected abundance or presence/absence at each local site for each species. Multiple methods can be used to define \mathbf{N} . For example, a leave-one-out LOESS model on occurrence data can be used to calculate the expected abundance at one site from a distance interpolation of the observed abundances at all other sites. Such a model indicates the expected community produced by a dispersal–environment model. Alternatively, MaxEnt or other species distribution models can be used to calculate the expected abundance based on only broad-scale climate. Other approaches are possible, like calculating the expected occurrence from stacking species’ regional geographic ranges sourced from expert-drawn range maps or from mechanistic regional models that predict abundance patterns across space (e.g. demographic or trait-based dispersal models; Jongejans et al. 2008).

A second set of choices that the user must define in the analysis is how to estimate the inverse covariance matrix Σ^{-1} , which can be difficult in practice. Two situations can arise: first, the number of species can be much larger than the number of sites; second, most sites can contain very few species. Both cases can lead to Σ becoming singular (i.e. non-invertible) because of very large covariances between some species pairs. A range of shrinkage estimators for Σ^{-1} have been developed that provide a robust approach to resolving this general problem (Hoerl and Kennard 1970, Schäfer and Strimmer 2005, Friedman et al. 2008). All the shrinkage estimators increase estimator bias in exchange for reduced mean squared error by introducing additional offset parameters that ‘shrink’ coefficient estimates and so force the existence of a matrix inverse. These offset parameters can be estimated by cross-validation approaches. A full survey of these methods is beyond the scope of this article, but some popular options include the James–Stein type shrinkage estimator (Schäfer and Strimmer 2005), the graphical lasso (L_1 -regularization; Friedman et al. 2008), or ridge regression (L_2 -regularization; Hoerl and Kennard 1970). We caution against using the graphical lasso because it produces sparse inverse covariance matrices that can produce singular null partial correlation distributions $C_{ij}(\tilde{\mathbf{N}})$, and instead recommend using the James–Stein estimator because of its good performance and low computational cost (Schäfer and Strimmer 2005).

Finally, we also recommend log-transforming abundance data as $f(x, a): x \mapsto \log(x + a) - \log(a)$ for some small number a , e.g. 10^{-6} . This transformation can improve normality of the distribution of abundance data, which can otherwise take either zero or very large values.

Testing the framework with a simulation analysis

To measure the expected performance of the network framework, we simulated co-occurrence matrices with known associations and determined how well network-detected associations matched these.

Consider a scenario involving n species distributed across m sites, of which a fraction h are unsuitable. We first generated an $n \times m$ expected species-by-site matrix \mathbf{N} , all of whose abundance entries were independently and identically distributed according to a hurdle model, such that N_{ij} was zero with probability h and Poisson-distributed (mean λ) with probability $1 - h$. If any of the marginal sums of \mathbf{N} were zero (i.e. a site with no species or a species with no sites; problematic only for small n and m), we re-generated \mathbf{N} until all marginal sums were non-zero. We then independently generated an $n \times m$ observed matrix \mathbf{O} via the same hurdle process.

As a next step, we assumed that there were Z associations in the ‘true’ association network. We chose the Z associations by generating a random graph with n vertices and Z edges (Erdős and Rényi 1959), with weight w_z ($z \in \{1, 2, \dots, Z\}$) set to either 1 or -1 with equal probability. To model this, we iterated over all associations z ; for each pair of species i_z and j_z for which a true association exists, we chose a random fraction f of sites $\{m_z\}$; at each of these sites we either increased the abundance for both species (when $w_z > 0$) or increased for one species and decreased for the other (when $w_z < 0$) by a factor s proportional to the mean abundance of both species at these sites:

$$\begin{aligned} \mathbf{O}_{i_z, m_{\{m_z\}}} &:= \mathbf{O}_{i_z, \{m_z\}} + w_z^2 \cdot s \cdot \left(\mathbf{O}_{i_z, \{m_z\}} + \mathbf{O}_{j_z, \{m_z\}} \right) / 2 \\ \mathbf{O}_{j_z, m_{\{m_z\}}} &:= \mathbf{O}_{j_z, \{m_z\}} + w_z \cdot s \cdot \left(\mathbf{O}_{i_z, \{m_z\}} + \mathbf{O}_{j_z, \{m_z\}} \right) / 2 \end{aligned} \quad (4)$$

This process effectively increased the covariance between species when the species were positively associated and decreased it when the species were negatively associated, with the parameters h and s controlling the strength of the association. We used $h = 0.2$ and $s = 0.2$ in this analysis.

Next, we applied our network framework using the matrices for all possible parameter combinations of $n = 10, 100$; $m = 10, 100, 1000$, $f = 0, 0.5$, and $Z = 10, 50, 100$ (note that some combinations were not possible, e.g. when $n = 10$, we cannot simulate values of $Z > 45$ because they would exceed the number of links in a fully connected network). We used a James–Stein type shrinkage estimator (Schäfer and Strimmer 2005) with significant associations inferred at the $\alpha = 0.05$ level. We set the number of null replicates to 1000 and repeated the entire analysis for each parameter combination 10 times.

In order to calculate error rates for our method, we compared the inferred networks’ structure to the true network’s structure. We counted a true positive association if it was detected for the correct pair of species and had the correct sign; as a true negative if it was not detected for a pair of species for which an association did not exist. A false positive association was counted if it was detected but was either the incorrect sign or the pair of species did not have a true association. Similarly, a false negative was counted if it was not detected but the pair of species did have a true association. These counts allowed us to calculate the positive predictive value (PPV; true positives divided by true positives plus false positives) and the negative predictive value (NPV; true negatives divided by true negatives plus false negatives) as summary statistics.

To determine the sensitivity of the method to different parameters, we constructed a random forest regression

model for NPV and PPV. The method generates an ensemble of regression trees and therefore allows for multi-way interactions between variables. We calculated the importance of each variable as the residual sum of squares caused by splitting on the variable of interest, averaged over all trees. Random forest models were built using the randomForest R package (Liaw and Wiener 2002).

Across all parameter combinations, PPV took a mean value of 12 ± 26 SD %, while NPV took a mean value of 87 ± 26 SD % (Supplementary material Appendix 1, Fig. A1). That is, the method was better at detecting the absence of associations than the presence of associations. A random forest model predicting PPV as a function of m , n , f , and Z explained 56% of the variation in the data and showed that m had the largest impact (14.6) on PPV, followed by Z (8.1), with smaller contributions of n (1.7) and f (1.3). Partial dependence plots indicated that larger values of each parameter led to higher values of PPV. A similar random forest model for NPV explained 53% of the variation in the data, and showed that m (11.1) and Z (8.7) had the largest impacts, with smaller contributions of n (1.5) and f (2.6). Partial dependence plots indicated that smaller values of each parameter led to higher values of NPV. Additionally, datasets with large fractions of zero-abundance records did not challenge the model’s ability to infer associations.

Overall, this analysis indicates that the method trades off between successfully detecting true associations (high PPV) and successfully detecting the absence of false associations (high NPV). Better PPV occurs for large datasets, while better NPV occurs for small datasets. Increasing the value of the false discovery rate α can further control this tradeoff. We repeated analyses for $\alpha = 0.5$ (results not shown), which approximately doubled PPV and halved NPV in all cases. Thus, the method can yield acceptable performance in a wide range of realistic use cases.

Empirical test of the framework

After establishing the robustness of the framework through the simulations, we can apply it to real datasets and analyze the inferred associations. Here we present an illustration of the approach using communities from temperate forests in North America. In particular, we use the framework to test whether 1) there are positive and/or negative associations among tree species that cannot be explained either from a regional dispersal–environmental model, or from broad-scale climate gradients; and whether 2) such associations can be explained by phylogenetic relatedness, functional trait similarity, or environmental gradients. We predict that across null models, association networks will have non-random structure and high modularity, and that positive (negative) associations will be more common among less (more) closely related species or in warmer (colder) environments.

Species data

We chose local community data to represent eastern North America. We extracted community-level tree species abundance data from the Forest Inventory and Analysis (FIA)

database (Gray et al. 2012) (see Supplementary material Appendix A1 for query details). We selected 5138 0.07 ha plots from the eastern USA (east of the 100°W meridian), and subsampled to a maximum of three plots for each 10 000 km² to reduce spatial sampling biases. Each plot consists of four 7.3 m radius subplots located 36.6 m from each other. All plots included were surveyed in the field between 2004 and 2008, used standard sampling protocols, and were marked as natural stands without evidence of artificial regeneration or human disturbance. If a plot was surveyed more than once in the time period chosen, we only included the newest survey. We excluded FIA taxa that were not identified to the species level. A total of $m = 1009$ plots out of the 5138 plots and $n = 137$ tree species were included in the analyses.

Point occurrence data for the null models came from BIEN, the Botanical Information and Ecology Network (Enquist et al. 2009, < <http://bien.nceas.ucsb.edu/bien/> >) for each of the 137 tree species.

Alternative species distribution models

We computed results based on three definitions of the null model. First, we used the commonly used random-swap algorithm, where the local community matrix \mathbf{O} is randomized to create the null matrix \mathbf{N} (Connor and Simberloff 1979, Gotelli and Ulrich 2010), keeping row and column sums fixed (i.e. total species and site abundances). Second, we used a leave-one-out LOESS regression of plot abundances with a span parameter of 0.2 (mirroring ter Steege et al. 2013) to compute the expected abundance at each plot from its spatial position and those of the remaining 5137 plots in the full FIA dataset. Third, we also used species distribution models created with the algorithm MaxEnt (Phillips and Dudík 2008) and based on the BIEN point occurrence data to estimate the climatic potential range of each species. As a test case, we used 19 bioclimatic layers representing ‘current’ climate (average 1950–2000 conditions) as model predictors, extracted from WorldClim 1.4 at 30-arc second resolution (Hijmans et al. 2005). Each model was fit using default parameters to both North and South America to capture the climatic range of the full New World distribution of each species. We converted the model suitability scores to expected abundance values by standardizing them so that the summed suitability scores for each species equaled the total number of individuals across all plots. This procedure assumes a linear relationship between suitability and abundance.

We used 1000 resamples $\hat{\mathbf{N}}$ of the expected species \times site matrix \mathbf{N} . As in the simulation analysis, we used a James–Stein shrinkage estimator for the inverse covariance matrix, and specified an overall false discovery rate of $\alpha = 0.05$ to exclude non-significant associations. Modules in the network were inferred using a standard fast-greedy algorithm (Clauset et al. 2004).

Potential predictors of species’ associations

We obtained data on phylogenetic relatedness between species and functional trait values to illustrate how the network statistics can be used to test hypotheses of drivers of species associations.

We calculated phylogenetic distances from a phylogeny for all the trees of eastern USA using Phylocom’s ‘phyloomatic’ tool (Webb et al. 2008). We used the R20120829 backbone tree, with branch lengths adjusted by fossil constraints (Gastauer and Meira-Neto 2013). We then computed the distance between each pair of species.

We also obtained measurements of four traits thought to underlie major axes of ecological strategy variation (Westoby et al. 2002): maximum height (m), specific leaf area (SLA; cm² g⁻¹), seed mass (g) and wood density (g cm⁻³). Functional trait data were extracted from the BIEN database. There was good coverage for trait data: 99% for height, 72% for SLA, 93% for seed mass, and 81% for wood density. We log₁₀-transformed each trait value to reduce skewness, then rescaled all values by subtracting means and dividing by standard deviations. We then computed trait differences between all pairs of species.

We used linear regression to determine whether the links between each species pair (link strength if adjacent; 0 if not adjacent) was predicted by pairwise phylogenetic or trait distance between the species pair. We did not correct for non-independence of predictor distances (e.g. Mantel-type test). This approach should lead to an increased rate of falsely rejecting the null hypothesis, meaning that failing to reject the null hypothesis is more likely to reflect a true absence of relationship.

Network structure of trees of eastern USA

The species association network based on the random-swap algorithm showed a non-random structure (chi-square test for Poisson degree distribution, $p < 10^{-45}$), identifying many positive and negative associations (Table 1). On the other hand, the overall structure of networks based on the

Table 1. Descriptive statistics for association networks constructed from the same dataset using different null models. We report an overall p-value for deviations from a Poisson distribution of edges, and for the subsets of the network for positive or negative associations, the mean \pm standard deviation of degree, as well as the number of non-singleton modules (i.e. with more than one member) detected in the network.

Regional model	p-value	Degree		Positive modules		Negative modules	
		Positive	Negative	Number	Mean size	Number	Mean size
Swap	2.44×10^{-46}	20.31 ± 10.81	56.26 ± 15.42	12	11.4	5	27.4
LOESS	0.230	0.09 ± 0.31	0 ± 0	5	2.2	0	NA
MaxEnt	0.565	0.82 ± 0.89	0 ± 0	24	3.2	0	NA

two regional null models – the LOESS dispersal–environment model, and the MaxEnt climate model – was not different from random. This result indicates that the overall distribution of local associations between trees in North America can mainly be explained by broad-scale drivers. Looking at the individual associations, we only found a few positive interactions in both of these networks. Interestingly, all five modules of 2–3 associations identified using the LOESS regional model also appeared when applying the MaxEnt model, although sometimes with an additional species in the module (Supplementary material Appendix 1, Fig. A2). We found additional associations deviating from the MaxEnt regional model, and modules were larger in general (Table 1, Fig. 1). In the random-swap network most species were inferred to be ‘segregators’, while in the MaxEnt and LOESS network most

species were inferred to be ‘aggregators’ (Supplementary material Appendix 1, Fig. A3).

What predicts species’ associations?

The positive and negative associations we found were not predicted by phylogenetic or trait distances (Supplementary material Appendix 1, Fig. A4) regardless of null model. The network structure explained by phylogeny in any network was no more than 0.05%, and the variation explained by all four traits together was no more than 0.12% in any network.

We did find weak spatial gradients in abundance-weighted mean values of degree for the species comprising each community (Supplementary material Appendix

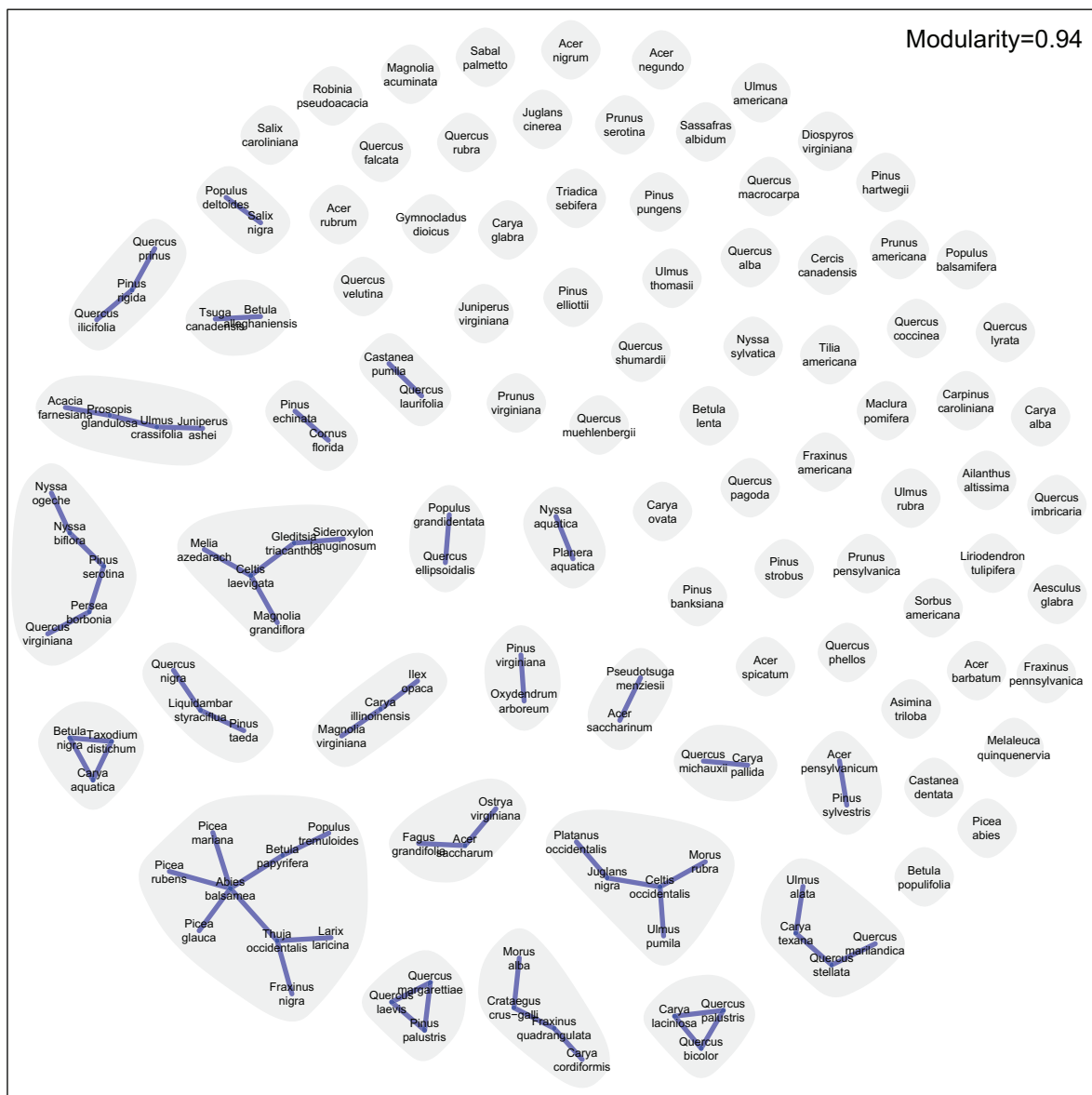


Figure 1. Empirical association network for North American trees. The network was constructed using the MaxEnt regional null model and a shrinkage inverse covariance estimator using 1000 null replicates and a false discovery rate of $\alpha = 0.05$. Gray envelopes indicate distinct modules. Positive associations are shown as blue lines; negative would be shown as red lines but none were found.

1, Fig. A5). For each of the random-swap and MaxEnt networks, mean degree was negatively correlated with mean annual temperature but not with mean annual precipitation (multiple regression; both $p < 0.002$, both $R^2 < 0.06$). The LOESS network was not correlated with either mean annual temperature or precipitation ($p = 0.26$).

Can we infer non-random species associations?

Our simulation analysis showed that the co-occurrence framework is indeed able to identify known associations within the parameter regimes we explored. However, there is a trade-off, where large datasets lead to better success in detecting true associations, and smaller datasets are better at identifying the absence of false associations. When including explicit regional null models, the framework is able to pinpoint which species are associated after broad-scale drivers have been accounted for. This not only allows testing for multiple regional models. Another advantage is that spatial autocorrelation in species distributions can be taken into consideration. Indeed, when looking at the results from our case study, it is apparent that the random-swap algorithm identifies spurious positive and negative associations that stem from ignoring the spatial dependence of species occurrences. This important implication shows that the null model from a hypothesized process such as dispersal or broad-scale climate patterns gives better control of the type of associations identified.

Another challenge addressed by the framework is that of indirect effects from multiple species on pairwise associations. The result is that when applying the framework to tree communities of eastern North America, we only found a small number of positive associations deviating from either the LOESS or the MaxEnt based regional models after correcting for indirect effects (Fig. 1). Ecologically, this may reflect that tree species distributions are largely controlled by environment and dispersal, with little importance of interspecific interactions, even if the latter matter for local abundances. The associations we do identify are those deviating from the broad-scale expectations. The LOESS model can be interpreted as a dispersal and environment model, since it simply interpolates abundance as a function of distance and implicitly includes environmental conditions that are spatially autocorrelated. Alternatively, the MaxEnt model computes the expected co-occurrence as a function of broad-scale climate variables, ignoring dispersal constraints, and instead representing the potential climate range of each species. It is thus not entirely surprising that the associations identified with the dispersal–environment model are all a subset of those identified with the climate-only model. Contrary to our predictions, functional traits, phylogenetic relatedness and environmental gradients did not correlate with the associations found across the networks. However, natural history, successional dynamics and missed environmental drivers could explain at least some of the associations identified. Indeed, local habitat requirements not fully captured by the broad scale environmental gradients tested here seem to explain several

associations. For instance, *Carya aquatica* and *Taxodium distichum* are both species of Coastal Plain, strongly associated with large river backswamps that are periodically flooded. A similar habitat is preferred by *Nyssa aquatica* and *Planera aquatica*, although this species pair is found more up-stream than the first ones, in areas where flooding is less prolonged. *Quercus palustris*, *Q. bicolor* and *Carya laciniosa* also prefer swamp habitats, though more from interior flatlands on more calcareous soils. Habitat preference does not seem to explain the association between *Carya pallida* and *Quercus michauxii*, which prefer dry/sandy and swampy soils, respectively. These two species rarely occur together, but both are largely confined to the southeastern Coastal Plain and lower Piedmont regions. This geographic signal could instead be driven by dispersal limitation. A few boreal modules such as the one around *Abies balsamea* are identified with the MaxEnt regional model but disappear in the LOESS model. These modules may represent local mosaics of boreal and temperate stands at the boreal-temperate transition zone, reflecting local environmental conditions and/or priority effects (Pastor and Mladenoff 1992).

Given that the MaxEnt model only covers environmental constraints, species might be simulated to co-occur that have similar climatic requirements but are allopatric. In such a case we would expect more negatively associated species in the network derived from the MaxEnt null model. Instead, the modularity of networks based on the MaxEnt models was much higher than that of LOESS-based networks, and in neither case did we find negative associations. The fact that we only identify positively associated groups of species unexplained by broad-scale climate conditions is consistent with the results of the simulation study of Araújo and Rozenfeld (2014), who found that the effect of positive (but not negative) dependencies between species scaled up to biogeographical scales and should be accounted for in range models under climate change.

One implication of our results – in particular the larger amount of positive associations deviating from the climate-only model – is that we cannot rely solely on stacked species distribution models (SDMs) (sensu Guisan and Rahbek 2011, Calabrese et al. 2014) to predict the composition of local communities, but need to take species associations into account (Araújo and Rozenfeld 2014) – whether driven by biotic interactions, dispersal, or local environmental filtering. In cases where e.g. local environmental or dispersal data are available, these can be integrated to filter broad-scale predictions and improve the performance of SDMs for communities (Boulangeat et al. 2012).

We note that even low error rates can translate into high absolute numbers of incorrect associations for datasets with large species numbers. Given that 100 species have 4950 potential pairwise interactions, our simulation error rates mean that, on average, approximately 25 links will be inferred that do not actually exist. These type 2 error rates are very high for small datasets, but settle to near zero for datasets with at least 100 sites, with performance further increasing when species co-occurrence patterns have high association weights. Thus, the algorithm does not often miss real associations, but even then, these type 2 error rates translate into approximately 500 real links that are missed. These rates and

uncertainties are similar to other association-inference approaches (Morales-Castilla et al. 2015) and suggest that these frameworks are most useful for reducing the space of possible associations to significantly more manageable numbers (cf. Fig. 2 in Morales-Castilla et al. 2015).

Inferring drivers of associations — a complex challenge

In our example with North American tree communities, the associations we identified can best be explained by ecological (and possibly geographic) groupings. It is thus not surprising that they could not be predicted by markers like functional trait similarity and phylogenetic distance. Similarly, we only found a weak negative correlation between community mean degree and mean annual temperature for the associations identified with the MaxEnt null model, possibly reflecting the geographic signal of boreal communities. The application still shows how network metrics can be used to test for drivers of associations, although in this particular case, more direct measures such as species' flooding tolerance from Ellenberg values would have been more useful.

Still, we can imagine several issues that could pose obstacles to the prediction of associations. One issue is that multiple competing processes can be at play. Competition and facilitation may cancel each other out (Callaway and Pennings 2000). Species associations may also vary across environmental gradients (Callaway et al. 2002a, Pottier et al. 2013) or across temporal scales (Blois et al. 2014, Martorell and Freckleton 2014), such that the scale of input data used would be critical. Even in networks of biotic interactions between plants and pollinators, pairwise associations do not always correlate with hypothesized drivers, even when correlations are found with metrics of overall network structure (Olito and Fox 2015). We therefore expect that disentangling the processes driving patterns of co-occurrence will remain an ongoing challenge.

Extensions of the network framework

The increased availability of regional species occurrence and local-scale co-occurrence data across large extents has been a strong driver behind recent attempts to disentangle the processes promoting species' associations through biotic interactions, dispersal limitations and environmental filtering (Blick and Burns 2009, Ulrich et al. 2012, Blois et al. 2014). We have presented a flexible framework to infer species positive and negative associations that deviate from expected broad-scale processes, and illustrated how network statistics can be used to test hypothesized drivers of associations. The approach can be readily applied to other datasets and systems for any other type of potential species associations across trophic levels.

There is an ongoing debate on the directionality of associations, with some studies finding more positive than negative associations (Blick and Burns 2009 and references therein, Blois et al. 2014), while others (including a meta-analysis) have found the opposite (Azaele et al. 2010, Gotelli and Ulrich 2010). Inconsistent associations appear to be the norm rather than the exception in the literature, changing when using species-based rather than individual-based

models (Blick and Burns 2009), or when looking for patterns across time in the paleo-record (Blois et al. 2014). Some of the disparate results seen across studies might be the unintended outcome of methodological differences. We suggest that sensitivity analyses and consensus approaches should be used to achieve robust results. Importantly, partial correlations should be widely implemented, as well as taking into account spatial autocorrelation as part of the null model to avoid identifying non-existent interactions.

The framework as currently implemented only accounts for symmetric associations, i.e. those representing (+/+), (0,0) or (-,-) interactions. This limitation is imposed because the matrix **A** is derived from **C**, which is derived from Σ^{-1} , which is symmetric. The coefficients in **C** could instead be calculated using other approaches that allow for non-symmetric outcomes, and so also capture effects like predation (+/-) and commensalism (+/0). The network framework allows for such directed linkages, and the R package does allow arbitrary user-specified functions to be used for calculating **C**. However, most association metrics are symmetric (cf. Janson and Vegelius 1981) and the few that are not (e.g. Somers' D (Somers 1962), Goodman and Kruskal's lambda (Goodman and Kruskal 1972)) cannot be used with abundance data, do not take both positive and negative values, and do not account for indirect associations. Other methods such as the excess co-occurrence approach of Araújo and Rozenfeld (2014), the Markov network approach suggested by Harris (2015), and the directed partial correlation coefficients proposed for microarray studies (Yuan et al. 2011) are still experimental. Developing association metrics that capture all possible ecological integration types should be a priority.

For cases in which many strong associations are identified, the network framework could potentially be extended to provide predictions of local community composition, which in turn can be tested in new communities. For instance, with information on the identity of only one species in a community, a 'network-crawling' approach could be used to predict the identity of the remaining species by allowing species that are close to the one known species in the network more likely to be predicted in the local community (or less likely for species negatively associated). Such an extension would provide transparent and powerful tests of the predictive ability of association networks derived either solely from occurrence data, or combined with experimental or field-based interaction data.

More sophisticated analyses could be applied to our framework to provide a more thorough test of the drivers of associations observed. Indeed, a recent study found modularity analyses useful to identify biological attributes driving the connection of modules of co-occurring species (Borthagaray et al. 2014). Such additional analyses can be readily applied with our framework, potentially combined with other methods to identify network modules (Leger et al. 2015) to e.g. investigate the relationship of traits and phylogenetic relationships for species within each module, and thus explore the definitions of the scale of study.

Looking forward

Mechanistic understandings of the drivers of species distributions across scales are needed to better predict the

consequences of rapid global change. The general framework we propose here provides a novel tool to infer local-scale associations that can affect species' distributions. We have presented a flexible framework linking co-occurrence and null-models to network theory for inferring species associations that is easily applicable to other systems and datasets. This approach can resolve challenges related to assumptions of spatial independence in species' distributions and indirect effects of multiple species on associations. Variations in the implementation of our approach can be used to test for association patterns that do not follow the expectation from broad-scale climate, dispersal or other hypothesized processes. With network metrics and analyses we can test for drivers of the patterns, and the approach can be potentially extended to predictions of local community assembly. Combined with field experiments and other methods, our framework can be a powerful tool to move beyond extant concepts in the analysis of co-occurrence data to improve assembly predictions across scales and help merge community ecology and biogeography.

Acknowledgements – This study was conducted as a part of the BIEN Working Group (Principal Investigators: Brian J. Enquist, Brad Boyle, Richard Condit, Steven Dolins, Robert K. Peet, and Barbara M. Thiers) supported by the National Centre for Ecological Analysis and Synthesis, a center funded by the National Science Foundation (NSF Grant EF-0553768), the Univ. of California, Santa Barbara, and the State of California. The BIEN Working Group was also supported by The iPlant Collaborative (NSF Grant DBI-0735191). We thank all the contributors for the invaluable data provided to the BIEN (<<http://bien.nceas.ucsb.edu/bien/people/data-contributors/>>). We are grateful to Irena Šimová for preparing the trait data and to Brad Boyle for guidance on SQL queries. We thank D. Harris for constructive comments that improved this manuscript, in particular for pointing us towards the use of partial correlation methods. We further acknowledge support to NMH by an EliteForsk Award, the Aarhus Univ. Research Foundation, and the Villum Foundation. BB and J-CS acknowledge financial support from Centre for Informatics Research on Complexity in Ecology (CIRCE), funded by the Aarhus Univ. Research Foundation under the AU Ideas, the European Research Council (ERC Starting Grant #310886 'HISTFUNC', and the Danish Council for Independent Research – Natural Sciences (grant #12-125079). BB was also supported by a United States National Science Foundation graduate research fellowship. CV was supported by a Marie Curie International Outgoing Fellowship within the 7th European Community Framework Program (DiversiTraits project, no. 221060) and by the ERC Starting Grand Project 'Ecophysiological and biophysical constraints on domestication in crop plants' (Grant SRC-StG-2014-639706-CONSTRAINTS).

References

- Afkhami, M. E. et al. 2014. Mutualist-mediated effects on species' range limits across large geographic scales. – *Ecol. Lett.* 17: 1265–1273.
- Araújo, M. B. and Rozenfeld, A. 2014. The geographic scaling of biotic interactions. – *Ecography* 37: 1–10.
- Araújo, M. B. et al. 2011. Using species co-occurrence networks to assess the impacts of climate change. – *Ecography* 34: 897–908.
- Azaele, S. et al. 2010. Inferring plant ecosystem organization from species occurrences. – *J. Theor. Biol.* 262: 323–329.
- Bascompte, J. 2010. Structure and dynamics of ecological networks. – *Science* 329: 765–766.
- Benjamini, Y. and Hochberg, Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. – *J. R. Stat. Soc. Ser. B* 57: 289–300.
- Blick, R. and Burns, K. C. 2009. Network properties of arboreal plants: are epiphytes, mistletoes and lianas structured similarly? – *Perspect. Plant Ecol. Evol. Syst.* 11: 41–52.
- Blois, J. L. et al. 2014. A framework for evaluating the influence of climate, dispersal limitation, and biotic interactions using fossil pollen associations across the late Quaternary. – *Ecography* 37: 1095–1108.
- Bolnick, D. I. et al. 2011. Why intraspecific trait variation matters in community ecology. – *Trends Ecol. Evol.* 26: 183–192.
- Borthagaray, A. I. et al. 2014. Inferring species roles in metacommunity structure from species co-occurrence networks. – *Proc. Biol. Sci.* 281: 20141425.
- Boulangier, I. et al. 2012. Accounting for dispersal and biotic interactions to disentangle the drivers of species distributions and their abundances. – *Ecol. Lett.* 15: 584–593.
- Brown, J. H. et al. 1996. The geographic range: size, shape, boundaries, and internal structure. – *Annu. Rev. Ecol. Syst.* 27: 597–623.
- Brown, J. H. et al. 2004. Constraints on negative relationships. – In: Taper, M. L. and Lele, S. R. (eds), *The nature of scientific evidence: statistical, philosophical, and empirical considerations*. Univ. of Chicago Press, pp. 298–324.
- Bruelheide, H. 2000. A new measure of fidelity and its application to defining species groups. – *J. Veg. Sci.* 11: 167–178.
- Bullock, J. M. et al. 2000. Geographical separation of two *Ulex* species at three spatial scales: does competition limit species' ranges? – *Ecography* 23: 257–271.
- Calabrese, J. M. et al. 2014. Stacking species distribution models and adjusting bias by linking them to macroecological models. – *Global Ecol. Biogeogr.* 23: 99–112.
- Callaway, R. M. and Pennings, S. C. 2000. Facilitation may buffer competitive effects: indirect and diffuse interactions among salt marsh plants. – *Am. Nat.* 156: 416–424.
- Callaway, R. M. et al. 2002a. Positive interactions among alpine plants increase with stress. – *Nature* 417: 844–848.
- Callaway, R. M. et al. 2002b. Epiphyte host preferences and host traits: mechanisms for species-specific interactions. – *Oecologia* 132: 221–230.
- Chesson, P. 2000. Mechanisms of maintenance of species diversity. – *Annu. Rev. Ecol. Syst.* 31: 343–366.
- Clauset, A. et al. 2004. Finding community structure in very large networks. – *Phys. Rev. E* 70: 066111.
- Connor, E. F. and Simberloff, D. 1979. The assembly of species communities: chance or competition? – *Ecology* 60: 1132–1140.
- De Marco, P. et al. 2008. Spatial analysis improves species distribution modelling during range expansion. – *Biol. Lett.* 4: 577–580.
- Diamond, J. M. 1975. Assembly of species communities. – In: Cody, M. L. and Diamond, J. M. (eds), *Ecology and evolution of communities*. Harvard Univ. Press, pp. 342–444.
- Dobra, A. et al. 2004. Sparse graphical models for exploring gene expression data. – *J. Multivar. Anal.* 90: 196–212.
- Enquist, B. J. et al. 2009. The Botanical Information and Ecology Network (BIEN): cyberinfrastructure for an integrated botanical information network to investigate the ecological impacts of global climate change on plant biodiversity. – iPlant Collaborative, <www.iplantcollaborative.org>.
- Erdős, P. and Rényi, A. 1959. On random graphs I. – *Publ. Math.* 6: 290–297.
- Friedman, J. et al. 2008. Sparse inverse covariance estimation with the graphical lasso. – *Biostatistics* 9: 432–441.
- Gastauer, M. and Meira-Neto, J. A. A. 2013. Avoiding inaccuracies in tree calibration and phylogenetic community analysis using Phylocom 4.2. – *Ecol. Inform.* 15: 85–90.

- Goodman, L. A. and Kruskal, W. H. 1972. Measures of association for cross classifications, IV: simplification of asymptotic variances. – *J. Am. Stat. Assoc.* 67: 415–421.
- Gotelli, N. J. and Ulrich, W. 2010. The empirical Bayes approach as a tool to identify non-random species associations. – *Oecologia* 162: 463–477.
- Gotelli, N. J. et al. 2010. Macroecological signals of species interactions in the Danish avifauna. – *Proc. Natl Acad. Sci. USA* 107: 5030–5035.
- Gray, A. et al. 2012. Forest Inventory and Analysis Database of the United States of America (FIA). – *Biodivers. Ecol.* 4: 225–231.
- Guisan, A. and Zimmermann, N. E. 2000. Predictive habitat distribution models in ecology. – *Ecol. Model.* 135: 147–186.
- Guisan, A. and Rahbek, C. 2011. SESAM – a new framework integrating macroecological and species distribution models for predicting spatio-temporal patterns of species assemblages. – *J. Biogeogr.* 38: 1433–1444.
- Harris, D. J. 2015. Estimating species interactions from observational data with Markov networks. – *bioRxiv*, dx.doi.org/10.1101/018861
- Hijmans, R. J. et al. 2005. Very high resolution interpolated climate surfaces for global land areas. – *Int. J. Climatol.* 25: 1965–1978.
- HilleRisLambers, J. et al. 2012. Rethinking community assembly through the lens of coexistence theory. – *Annu. Rev. Ecol. Evol. Syst.* 43: 227–248.
- Hoerl, A. E. and Kennard, R. W. 1970. Ridge regression: biased estimation for nonorthogonal problems. – *Technometrics* 12: 55–67.
- Janson, S. and Vegelius, J. 1981. Measures of ecological association. – *Oecologia* 49: 371–376.
- Jongejans, E. et al. 2008. Dispersal, demography and spatial population models for conservation and control management. – *Perspect. Plant Ecol. Evol. Syst.* 9: 153–170.
- Kissling, W. D. et al. 2012. Towards novel approaches to modelling biotic interactions in multispecies assemblages at large spatial extents. – *J. Biogeogr.* 39: 2163–2178.
- Lasky, J. R. et al. 2014. Trait-mediated assembly processes predict successional changes in community diversity of tropical forests. – *Proc. Natl Acad. Sci. USA* 111: 5616–5621.
- Legendre, P. 1993. Spatial autocorrelation: trouble or new paradigm? – *Ecology* 74: 1659–1673.
- Leger, J.-B. et al. 2015. Clustering methods differ in their ability to detect patterns in ecological networks. – *Methods Ecol. Evol.* 6: 474–481.
- Lennon, J. J. 2000. Red-shifts and red herrings in geographical ecology. – *Ecography* 23: 101–113.
- Liaw, A. and Wiener, M. 2002. Classification and regression by randomForest. – *R News* 2: 18–22.
- MacArthur, R. H. 1958. Population ecology of some warblers of northeastern coniferous forests. – *Ecology* 39: 599.
- Martorell, C. and Freckleton, R. P. 2014. Testing the roles of competition, facilitation and stochasticity on community structure in a species-rich assemblage. – *J. Ecol.* 102: 74–85.
- Morales-Castilla, I. et al. 2015. Inferring biotic interactions from proxies. – *Trends Ecol. Evol.* 30: 347–356.
- Newman, M. 2010. Networks: an introduction. – Oxford Univ. Press.
- Normand, S. et al. 2011. Postglacial migration supplements climate in determining plant species ranges in Europe. – *Proc. R. Soc. B* 278: 3644–3653.
- Olito, C. and Fox, J. W. 2015. Species traits and abundances predict metrics of plant–pollinator network structure, but not pairwise interactions. – *Oikos* 124: 428–436.
- Pastor, J. and Mladenoff, D. J. 1992. The southern boreal-northern hardwood forest border. – In: Shugart, H. H. et al. (eds), *A systems analysis of the global boreal forest*. Cambridge Univ. Press, pp. 216–240.
- Phillips, S. J. and Dudík, M. 2008. Modeling of species distributions with Maxent: new extensions and a comprehensive evaluation. – *Ecography* 31: 161–175.
- Pollock, L. J. et al. 2014. Understanding co-occurrence by modelling species simultaneously with a joint species distribution model (JSDM). – *Methods Ecol. Evol.* 5: 397–406.
- Pottier, J. et al. 2013. The accuracy of plant assemblage prediction from species distribution models varies along environmental gradients. – *Global Ecol. Biogeogr.* 22: 52–63.
- Schäfer, J. and Strimmer, K. 2005. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. – *Stat. Appl. Genet. Mol. Biol.* 4: Article 32.
- Schleuning, M. et al. 2012. Specialization of mutualistic interaction networks decreases toward tropical latitudes. – *Curr. Biol.* 22: 1925–1931.
- Soberón, J. 2007. Grinnellian and Eltonian niches and geographic distributions of species. – *Ecol. Lett.* 10: 1115–1123.
- Somers, R. H. 1962. A new asymmetric measure of association for ordinal variables. – *Am. Sociol. Rev.* 27: 799–811.
- Stewart, R. E. and Aldrich, J. W. 1951. Removal and repopulation of breeding birds in a spruce-fir forest community. – *Auk* 68: 471–482.
- Svenning, J.-C. and Skov, F. 2007. Could the tree diversity pattern in Europe be generated by postglacial dispersal limitation? – *Ecol. Lett.* 10: 453–460.
- ter Steege, H. et al. 2013. Hyperdominance in the Amazonian tree flora. – *Science* 342: 1243092.
- Thuiller, W. et al. 2013. A road map for integrating eco-evolutionary processes into biodiversity models. – *Ecol. Lett.* 16 (Suppl. 1): 94–105.
- Ulrich, W. et al. 2012. Null model tests for niche conservatism, phylogenetic assortment and habitat filtering. – *Methods Ecol. Evol.* 3: 930–939.
- Veech, J. A. 2013. A probabilistic model for analysing species co-occurrence. – *Global Ecol. Biogeogr.* 22: 252–260.
- Violle, C. et al. 2011. Phylogenetic limiting similarity and competitive exclusion. – *Ecol. Lett.* 14: 782–787.
- Webb, C. O. et al. 2002. Phylogenies and community ecology. – *Annu. Rev. Ecol. Syst.* 33: 475–505.
- Webb, C. O. et al. 2008. Phylocom: software for the analysis of phylogenetic community structure and trait evolution. – *Bioinformatics* 24: 2098–2100.
- Weihner, E. et al. 2011. Advances, challenges and a developing synthesis of ecological community assembly theory. – *Phil. Trans. R. Soc. B* 366: 2403–2413.
- Westoby, M. et al. 2002. Plant ecological strategies: some leading dimensions of variation between species. – *Annu. Rev. Ecol. Syst.* 33: 125–159.
- Wisz, M. S. et al. 2013. The role of biotic interactions in shaping distributions and realised assemblages of species: implications for species distribution modelling. – *Biol. Rev. Camb. Phil. Soc.* 88: 15–30.
- Yuan, Y. et al. 2011. Directed partial correlation: inferring large-scale gene regulatory network through induced topology disruptions. – *PLoS One* 6: e16835.

Supplementary material (Appendix ECOG-01892 at <www.ecography.org/appendix/ecog-01892>). Appendix 1.