

Bayesian Estimation of Co-Occurrence *Affinity* With Dyadic Regression

Abstract

1) Estimating underlying co-occurrence relationships between pairs of species has long been a challenging task in ecology as the extent to which species co-occur is partially dependent on their prevalence. While recent work has taken large steps towards solving this problem, the next question is how to assess the factors that influence co-occurrence.

2) Here, I show that a recently proposed co-occurrence metric can be improved upon by assigning Bayesian priors to the latent co-occurrence relationships being estimated. In the context of analysing the factors that affect co-occurrence relationships, I demonstrate the need for a generalised linear model (GLM) that takes raw data (co-occurrences and species prevalence) not derived quantities (co-occurrence metrics) as its data. Next, I show the form that such a GLM should take in order to perform Bayesian inference while accounting for non-independence of dyadic matrix data (e.g. distance and co-occurrence matrices).

17 3) I then present 3 example analyses to highlight the types of scientific questions these meth-
18 ods can help answer, using existing data sets - measuring the effects of trait dissimilarity
19 among dung beetle species and relatedness between ant species on co-occurrence, and
20 constructing co-occurrence networks of bacteria found in cystic fibrosis patient sputum
21 samples.

22 4) Finally, I present the software package CooccurrenceRegression.jl, which provides a
23 straightforward interface for researchers to put these methods into practice.

24 **1 Introduction**

25 The analysis of patterns of co-occurrence between taxa is an important and active area of
26 ecological research (Gotelli and McCabe 2002; Barberán et al. 2012; Williams, Howe, and
27 Hofmockel 2014; Kraan, Thrush, and Dormann 2020; Zhu et al. 2023). Mainali et al. (2022)
28 have recently shown that of the numerous ways of measuring co-occurrence relationships be-
29 tween species pairs (reviewed in (Mainali et al. 2022)) the correct and unbiased method is
30 to make use of Fisher’s noncentral hypergeometric distribution, as suggested by Veech (2013)
31 and Griffith et al. (2016). Mainali et al. (2022) derived a co-occurrence metric called *affinity*
32 (or $\hat{\alpha}$) based on this distribution. This is a significant step forward in the analysis of co-
33 occurrence relationships. There are still several challenges remaining, however. For instance,
34 when species pairs of very low or very high prevalence are analysed with this method, they will
35 often be assigned very high or low affinity scores, but with low confidence (high p-value and

wide confidence interval). For downstream analysis then the researcher may (i) treat all data points as equal, (ii) remove data above some high p-value threshold or (iii) devise a scheme to weight data appropriately. (i) wastes information and can yield misleading results. (ii) again, wastes information and could severely bias results depending on the reasons for the differences in prevalence/p-values. If an appropriate, unbiased scheme for (iii) should be devised, then this would be a welcome development. However, adopting a Bayesian approach that builds upon the work of Mainali et al. (2022) not only yields more accurate estimates of the *affinity* between species, but also naturally propagates uncertainty through the analysis, i.e., it accounts for the differing levels of confidence we have about the co-occurrence relationships between different species pairs. In the present work I illustrate this point and provide the model description and code to use this Bayesian method in practice.

A significant advantage of this framework is that it allows co-occurrence data to be analysed as response data in a Bayesian general linear model (GLM). That is, by supplying species occurrence data as the prevalence of individual species, the number of times a given species pair co-occur and the total number of sites considered, it is now no more complicated to construct a regression model (with a Fisher's noncentral hypergeometric likelihood function) than it would be to perform binomial regression with count data. Similarly, this method makes the best use of all available information and weighs data points appropriately. Thus, just as it is not appropriate to convert count data to proportions and conduct linear regression, it is no longer best practice to summarise co-occurrence data as point estimates for linear regression.

56 2 Materials and Methods

57 2.1 Simulation data

58 In the following, maximum likelihood estimates of co-occurrence affinity were obtained using
59 the R (R Core Team 2014) package `co-occurrenceAffinity` (Mainali et al. 2022). All other
60 analyses and visualisations were carried out in Julia (Bezanson et al. 2017). All models were
61 constructed in the probabilistic programming language `Turing` (Ge, Xu, and Ghahramani
62 2018) with MLE estimates of regression coefficients fit to data using the Nelder-Mead method
63 (Nelder and Mead 1965) and MAP estimates of α fit using L-BFGS (Liu and Nocedal 1989)
64 implemented in `Optim` (Mogensen and Riseth 2018) and results visualised in `Makie` (Danisch
65 and Krumbiegel 2021).

66 **Maximum *a posteriori* vs maximum likelihood pairwise affinity estimates**

67 Mainali et al. (2022) show that the log of the odds ratio term in Fisher’s noncentral hypergeo-
68 metric distribution (a quantity they term α) can be used to appropriately describe the extent
69 to which two species will tend to co-occur more or less than would be expected based just on
70 the prevalence of the species. They go on to propose that the maximum likelihood estimate
71 of this parameter $\hat{\alpha}$ or *affinity* should be used as a pairwise co-occurrence metric. However,
72 Maximum likelihood estimates can yield values of positive or negative infinity. This causes
73 difficulties for downstream analyses (one cannot do something as simple as calculating the
74 mean of a set containing infinite values). Furthermore, we know *a priori* that an infinitely

75 large or small affinity is not sensible for most cases of interest to ecologists. Bayesian analysis
 76 uses prior knowledge to avoid the estimation of physically or biologically implausible values.
 77 Mainali et al. (2022) reassign these infinite estimates an absolute value of $\log(2N^2)$, where N
 78 is the total number of sample sites (from an argument made based on the Jeffreys' prior for
 79 the beta distribution). While this figure comes from sound argument, there are at least two
 80 problems with this approach. Firstly, not all data are treated the same way, i.e., no regulari-
 81 sation is applied to finite affinity estimates, only to these extreme values. Secondly, the value
 82 $\log(2N^2)$ is only a function of N , not the species prevalence. Thus, it is not influenced by our
 83 actual state of knowledge about the species in question.

84 To make these ideas more concrete and show the practical implications, I simulated species
 85 pairs using the affinity model. For each pair there are $N = 30$ sites they can inhabit, species
 86 A has a prevalence mA and species B has prevalence mB . The number of sites at which they
 87 co-occur k was drawn from a Fisher's noncentral hypergeometric distribution

$$k \sim \text{fnchypg}(mA, N - mA, mB, e^\alpha),$$

88 with 10 draws per combination of mA and mB for each of 41 different values of α . Then,
 89 given the values for N , k , mA and mB I estimated α using two methods. Firstly, I used
 90 the original maximum likelihood estimate (MLE) of Mainali et al. (2022). Next, I obtained
 91 maximum *a posteriori* (MAP) estimates with a Gaussian prior $N(0, 3)$ for α . Note that these
 92 are not strongly regularising priors as when exponentiated in the likelihood function a standard

deviation of $3 \approx 20$ and two standard deviations $6 \approx 403$ which is a very large odds ratio for most applications. In order to compare the two methods, for each combination of mA and mB I calculated the deviation between the inferred and ground truth α values as the root mean squared error (RMSE).

Regression analysis

Often, we do not simply wish to report co-occurrence relationships, but to measure how they change with some other variables of interest. For many types of data there are well understood and regularly used probability distributions which can be used in Bayesian and frequentist GLMs. For co-occurrence data this has not been the case. Given the issues with deriving point estimates highlighted above it is unlikely that simply fitting a linear model to such point estimates of co-occurrence affinity will yield reliable results. Thus, rather than supplying a second inference model with uninformative point estimates from a previous model, we can provide our regression model with all the data on species prevalence and co-occurrences instead. To demonstrate the impact of this I simulated 41 sets of predictor data \vec{x} , each consisting of 30 draws from $N(0,1)$. Affinity values were generated by multiplying the predictor data by a regression coefficient β . For each affinity value - species prevalence mA and mB values were chosen randomly between 1 and 29 and a k value was drawn from the Fisher's noncentral hypergeometric distribution as above. For each generated data set pairwise affinity values were estimated by the MAP and MLE methods. Then linear regression analysis was conducted on these point estimates $\alpha \sim \beta\vec{x}$. Additionally, I obtained a maximum likelihood estimate of a

113 GLM of the form

$$\begin{aligned}\vec{k}_i &\sim \text{fnchypg}(\vec{mA}_i, N - \vec{mA}_i, \vec{mB}_i, e^\alpha) \\ \alpha &= \beta_0 + \beta \vec{x}_i,\end{aligned}\tag{1}$$

114 where \vec{k} , \vec{mA} and \vec{mB} are vectors containing the values of k , mA and mB respectively and β_0
115 is the intercept.

116 **Dyadic regression**

117 In many if not most cases when working with co-occurrence data it will be in the form of a
118 square co-occurrence matrix similar to the distance and dissimilarity matrices used to record
119 e.g. phylogenetic distances between species or community dissimilarities between sample sites.
120 As with these other types of matrices, if we wish to perform regression analysis treating each
121 entry in the matrix as data point, we must account for non-independence of data coming from
122 the same row or column, e.g. same site, species etc.. To deal with this we can include a random
123 effect λ for each species (Clarke, Rothery, and Raybould 2002; Gompert et al. 2014). Thus,
124 whereas the GLM above contains the term

$$\alpha_i = \beta_0 + \beta \vec{x}_i$$

125 we would now have

$$\alpha_i = \beta_0 + \lambda_i + \lambda_j + \beta X_{ij}.$$

Where X is a (possibly dissimilarity or distance) matrix in which element X_{ij} is a quantity of interest relating species i to species j . Given a similarly arranged matrix K which holds the k values for all species pairs and a vector \vec{m} containing species prevalence, the dyadic GLM becomes

$$\begin{aligned} K_{ij} &\sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha) \\ \alpha &= \beta_0 + \lambda_i + \lambda_j + \beta X_{ij} \end{aligned} \tag{2}$$

for each index ij in either the lower or upper triangle of K . While Maximum likelihood estimates of the λ , β and β_0 parameters can be obtained, we still need a way to properly account for our uncertainty in our estimates. Bayesian inference provides an intuitive framework for this, grounded in probability theory. Thus, we can construct a Bayesian GLM by assigning priors to the unknown parameters. Assuming Gaussian priors for all parameters gives

$$\begin{aligned} K_{ij} &\sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha) \\ \alpha &= \beta_0 + \lambda_i + \lambda_j + \beta X_{ij} \\ \beta &\sim \text{N}(0, \beta_\sigma) \\ \beta_0 &\sim \text{N}(0, \beta_{0\sigma}) \\ \lambda_i &\sim \text{N}(0, \lambda_\sigma) \end{aligned} \tag{3}$$

135 where β_σ , $\beta_{0\sigma}$ and λ_σ are to be chosen according to the specifics of the system being analysed.
136 This forms the basic model structure used for analyses in the examples with real data.

137 **2.2 Examples with real data**

138 In the following examples all explanatory variables were standardised to have mean 0 and
139 standard deviation 1, in order to make results comparable. Inference was performed via Markov
140 chain Monte Carlo (MCMC), using the No U-turns (NUTS) sampler (Homan and Gelman 2014)
141 implemented the Turing probabilistic programming language (Ge, Xu, and Ghahramani 2018),
142 with 8 chains of 1000 iterations in each case. Point estimates are reported with 95% credible
143 intervals (CI) calculated from quantiles of MCMC samples and probability of direction (PD),
144 the posterior probability that the reported effect is in the estimated direction.

145 **Change in co-occurrence affinity owing to trait dissimilarity amongst dung beetles**

146 For the first example I used previously published occurrence data for dung beetles along an
147 altitudinal gradient at Serra do Cipó, State of Minas Gerais, Brazil (Nunes et al. 2016). The
148 data set includes occurrence data for 56 beetle species, across 7 sites, ranging in elevation
149 from 800 to 1400 meters above sea level, as well as mean biomass (dried weight) in grammes
150 and the functional guild for each species. In order to assess how functional trait dissimilarity
151 affected co-occurrence affinity we employed 3 different regression models 1) combining the two
152 functional traits into a single Gower's distance (Gower 1971); 2) including both traits and their

153 interaction as explanatory variables; 3) focussing on the effect of biomass *within* functional
 154 guild. Thus, model 1 was identical in form to equation Equation 3

$$K_{ij} \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha)$$

$$\alpha = \beta_0 + \lambda_i + \lambda_j + \beta X_{ij}$$

$$\beta \sim \text{N}(0, 1)$$

$$\beta_0 \sim \text{N}(0, 4)$$

$$\lambda_i \sim \text{N}(0, 1).$$

155 Model 2 was

$$K_{ij} \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha)$$

$$\alpha = \beta_0 + \lambda_i + \lambda_j + \beta_1 X1_{ij} + \beta_2 X2_{ij} + \beta_3 X3_{ij}$$

$$\beta_k \sim \text{N}(0, 1)$$

$$\beta_0 \sim \text{N}(0, 4)$$

$$\lambda_i \sim \text{N}(0, 1),$$

156 where $X3 =$ the element-wise product $X1 \odot X2$.

157 Model 3 was

$$K_{ij} \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha) \quad \text{for } (i, j) \in \mathcal{S}$$

$$\alpha = \beta_0 + \lambda_i + \lambda_j + \beta X_{ij}$$

$$\beta \sim \text{N}(0, 1)$$

$$\beta_0 \sim \text{N}(0, 4)$$

$$\lambda_i \sim \text{N}(0, 1),$$

158 where $\mathcal{S} = \{(i, j) : g_i = g_j\}$ is the set of indices denoting species pairs from the same functional
 159 guild.

160 **Change in co-occurrence affinity owing to trait dissimilarity amongst ants**

161 For the second example I used two data sets, one containing global ant species occurrence
 162 data and phylogenetic tree (Economo et al. 2018) and containing ant occurrence data and
 163 trait measurements from two nature reserves near Hong Kong (Wong et al. 2021). The
 164 trait data consisted of species mean values for body size, leg length, head width, mandible
 165 length, pronotum width and scape length. In order to analyse the relationship between trait
 166 dissimilarity and co-occurrence affinity I first subset the data to work with only those species
 167 that occurred in both data sets, so that I had trait measurements and occurrence data at both
 168 scales for all of them. Separate models were then used for the two geographic scales, following
 169 the same basic model structure as above for the global data set

$$K_{ij} \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha)$$

$$\alpha = \beta_0 + \lambda_i + \lambda_j + \sum_{k=1}^n \beta_k X_{k_{ij}}$$

$$\beta_k \sim \text{N}(0, 1)$$

$$\beta_0 \sim \text{N}(0, 4)$$

$$\lambda_i \sim \text{N}(0, 1),$$

170 where here $n = 6$ traits. However, the the local scale data set used dat from two separate
 171 nature reserves and from two types of sites: those where invasive species *Solenopsis invicta*
 172 was present and those where it was absent. I pooled the data from these four site types, with
 173 random intercepts for each, giving the model

$$K_{ijl} \sim \text{fnchypg}(\vec{m}_{il}, N_l - \vec{m}_{il}, \vec{m}_{jl}, e^\alpha)$$

$$\alpha = \beta_{0_l} + \lambda_{il} + \lambda_{jl} + \sum_{k=1}^n \beta_k X_{k_{ijl}}$$

$$\beta_k \sim \text{N}(0, 1)$$

$$\beta_{0_l} \sim \text{N}(0, 4)$$

$$\lambda_{il} \sim \text{N}(0, 1),$$

174 In order to determine the relationship between phylogenetic distance and co-occurrence affinity
 175 I used 3 models: 1) measuring the relationship at genus level; 2) at species level with a
 176 single regression coefficient; 3) at species level with a hierarchical model. The phylogenetic
 177 tree was read and traversed using phylo.jl (Reeve, Borregaard, and Harris 2024). Genus

178 level phylogenetic distances were calculated as the tree distance between most recent common
 179 ancestor node (MRCA) of one species and the MRCA of another. Prior to conducting the
 180 species level analysis, I removed all genera with less than 25% unique phylogenetic distances.
 181 This was to reduce any bias introduced by species having artificially low phylogenetic distance
 182 due to lack of resolution in the tree. 61 genera remained for analysis. These analyses were
 183 conducted only across the global geographic scale, as there was insufficient phylogenetic data
 184 to calculate distances between many of the species at the local scale.

185 Analysing the effect of phylogenetic distance on co-occurrence at the genus level took an
 186 identical form to Equation 3 with same priors as model 1 of the dung beetle analysis. Similarly,
 187 the pooled species level model was identical to model 3 of the dung beetle analysis. The
 188 hierarchical model was given by

$$K_{ij} \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha) \quad \text{for } (i, j) \in \mathcal{S}$$

$$\alpha = \beta_0 + \lambda_i + \lambda_j + \beta_{g_i} X_{ij}$$

$$\beta_k \sim \text{N}(\mu, \sigma)$$

$$\beta_0 \sim \text{N}(0, 4)$$

$$\lambda_i \sim \text{N}(0, 1)$$

$$\sigma \sim \text{gamma}(2, 1/10)$$

$$\mu \sim \text{N}(0, 1)$$

189 where $\mathcal{S} = \{(i, j) : g_i = g_j\}$ is the set of indices denoting species pairs from the same genus

190 and g is a vector of indices mapping each species to it's genus.

191 **Bacterial co-occurrence networks derived from cystic fibrosis patient sputum samples**

192 To demonstrate another use for Bayesian estimation of co-occurrence affinity I used microbiome
193 data from cystic fibrosis patient sputum samples (Quinn et al. 2019) to construct co-occurrence
194 networks. The data originally consisted of 4148 unique sequences of 100 nucleotides of the
195 v4 region of the bacterial 16S rRNA (primers: 515f GTGCCAGCMGCCGCGGTAA; 806r
196 GGACTACHVGGGTWTCTAAT). All sequences were identified to genus level using the RDP
197 classifier (Wang et al. 2007) implemented in AssignTaxonomy.jl. Prior to analysis, I removed
198 all sputum samples which were outliers in terms of read depth ($8000 < \text{read depth} < 16000$) to
199 reduce the impact of differences in read depth on affinity estimates e.g. rare taxa co-occurring
200 in bigger samples. After combining occurrence data relating to sequences from the same genera
201 and limiting the data to a single sample from any one patient, the final data set consisted of
202 a presence/absence matrix of 287 genera across 62 patients.

203 The affinity of each genus pair was inferred separately:

$$k \sim \text{fnchypg}(\vec{m}_i, N - \vec{m}_i, \vec{m}_j, e^\alpha)$$

$$\alpha \sim N(0, 4)$$

204 The same computational methods were used as for all the other examples, with the exception
205 that I used only a single MCMC chain for each genus pair. Results were visualised both as an
206 adjacency matrix and as networks.

3 Results

3.1 Simulation data

Maximum *a posteriori* vs maximum likelihood pairwise affinity estimates

Figure 1 shows that only when both mA and mB were equal to 15 was the RMSE approximately equivalent between MLE and MAP methods. Whenever one or both species had a high or low prevalence, and particularly as the absolute value of α became larger, the MLE method produced very poor estimates, and the extreme estimates were always the same $\log(2N^2) = 7.496$. By contrast, for the MAP values, the prior provides regularisation which can be overcome by increasing confidence in the data, which is a function of mA and mB . Thus, the models *best guess* when $mA = 15$, $mB = 5$ and $k = 5$ is higher than the equivalent situation when $mB = 1$ and $k = 1$. Neither of these methods is perfect, however. When asked for one, a model will give you its best guess point estimate, but we can make better use of the data we have collected if we can utilise not only the point estimates but also our uncertainty around them.

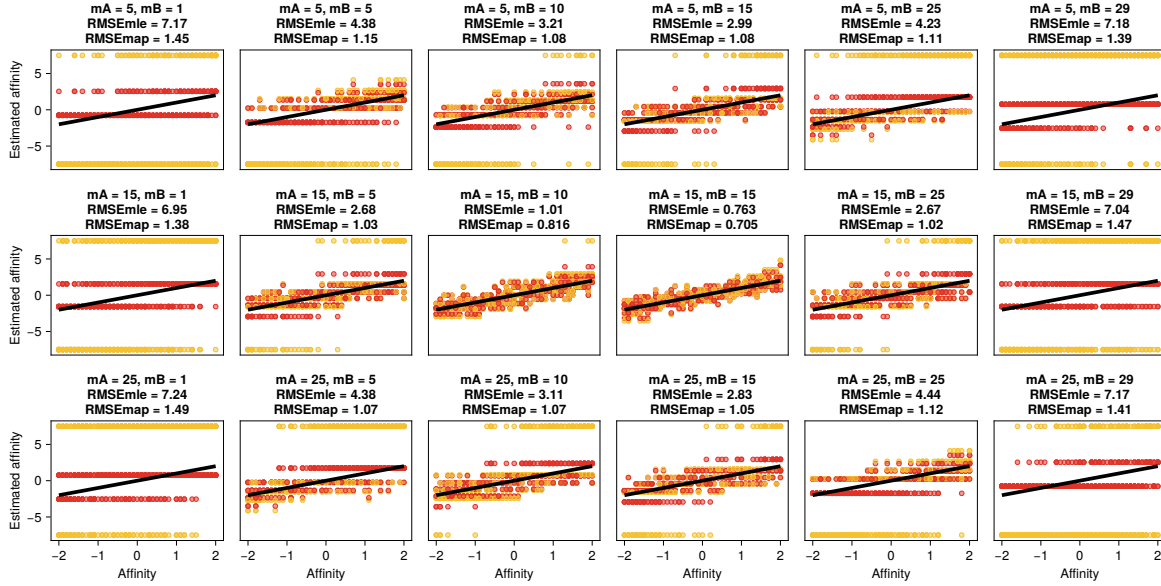


Figure 1: Actual and estimated affinity values for a range of species prevalences. Yellow points are estimates using the original maximum likelihood method and red points are maximum *a posteriori* estimates. Black lines indicate the actual affinity values used to generate the data. In each panel is shown the root mean squared error (RMSE) for both types of estimate.

220 Regression analysis

221 The results in Figure 2 show how poorly fitting a linear model to $\hat{\alpha}$ point estimates does,
 222 typically overestimating the absolute value of β by a large margin. Using MAP estimates of
 223 α does better here, exhibiting the opposite behaviour of slightly underestimating the absolute
 224 value of β . However, by cutting out the step of generating point estimates for each pair the
 225 GLM retains all pertinent information and accurately recaptures the parameters of the data
 226 generating model. It is of course expected that the GLM should be able to discover the correct
 227 parameter values, since they were generated by an identical model. What is important is the

way the other two models fail by comparison, and of course the fact that we now have the correct likelihood function for such a co-occurrence GLM.

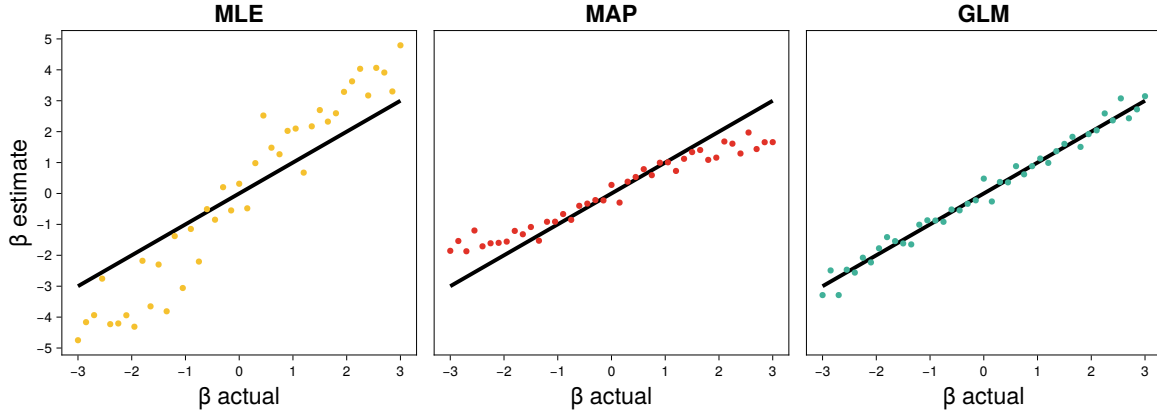


Figure 2: Estimated regression coefficients β according to three different methods: fit linear model to $\hat{\alpha}$ values, fit linear model to MAP estimates of α , fit GLM to raw data. Black lines indicate the true β values.

3.2 Examples with real data

Change in co-occurrence affinity owing to trait dissimilarity amongst dung beetles

The first example used occurrence data for dung beetles, along an altitudinal gradient at Serra do Cipó, State of Minas Gerais, Brazil (Nunes et al. 2016). Given the strong abiotic selection imposed by the elevation gradient, I would expect to find greater co-occurrence affinity between more similar beetles. However, there was very little evidence of Gower's distance (based on both functional guild and mean biomass) having any effect on co-occurrence affinity 0.05 (CI[-0.138, 0.229], PD = 0.703) (Figure 3, top row). By separating trait dissimilarity into functional guild membership and mean biomass, I found that dissimilarity in mean biomass had

239 a significant negative relationship with co-occurrence affinity, with a regression coefficient of
 240 -0.342 (CI[-0.621, -0.064], PD = 0.993) (Figure 3, middle row). There was also a significant in-
 241 teraction between the two traits - interaction strength = 0.562 (CI[0.273, 0.879], PD > 0.999),
 242 implying that the negative relationship between biomass dissimilarity and co-occurrence affin-
 243 ity is only apparent *within* functional guilds, since the combined effect of biomass and the
 244 interaction term (i.e., $\beta_{\text{biomass}} + \beta_{\text{biomass_in_different_functional_guild}}$) was 0.223 (CI[-0.097, 0.556],
 245 PD = 0.912). In fact, the effect was so completely confined to within functional guild that
 246 the weak evidence of increased affinity with functional guild dissimilarity: $\beta = 0.176$ (CI[-
 247 0.2, 0.541], PD = 0.827) combined with the fact that most species pair were not from the
 248 same functional guild led to the slightly positive estimate of the effect of Gower's distance on
 249 affinity. Thus, I present an alternative analysis where I consider only the impact of biomass
 250 dissimilarity between species of the same functional guild on co-occurrence affinity (Figure 3,
 251 bottom row) - resulting in a slightly increased absolute effect: -0.42 (CI[-0.729, -0.11], PD =
 252 0.997). These results fit the expected pattern of similar species being found together across
 253 an environmental gradient but also highlight the possibility of such effects being nuanced in
 254 such a way that a combination of expert knowledge and flexible inference method such as the
 255 one proposed here may be needed in order to detect them.

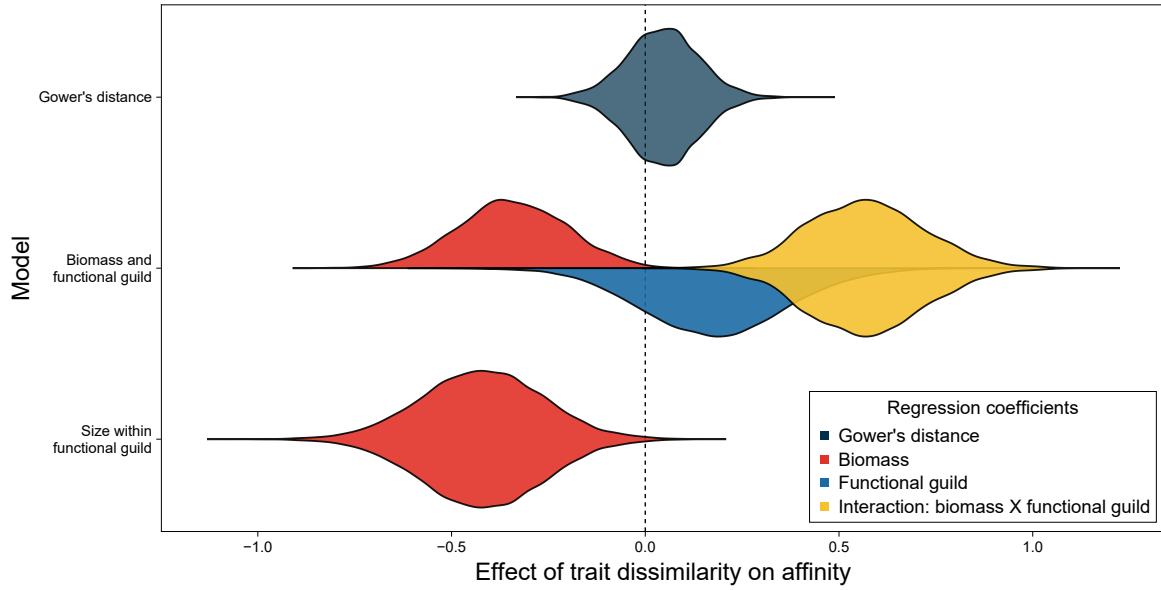


Figure 3: Posterior estimates of change in co-occurrence affinity owing to trait dissimilarity amongst dung beetles, along an altitudinal gradient at Serra do Cipó, State of Minas Gerais, Brazil.

Change in co-occurrence affinity owing to trait dissimilarity and relatedness amongst ants

The extent to which functional trait similarity and phylogenetic relatedness lead to taxa co-occurring or being overdispersed are expected to change depending on the geographic scale and (in the case of relatedness) phylogenetic scale under investigation (Webb et al. 2002). Thus, for the next example I made use of ant species occurrence data from both the global scale (Economio et al. 2018) and a smaller scale data set from two nature reserves near Hong Kong (Wong et al. 2021). First, I used mean values for a range of traits (body size, leg length, head width, mandible length, pronotum width and scape length) measured at the Hong Kong sites to see how trait dissimilarity affected affinity at both geographic scales. In general, there

was a poor match between the results at different scales. Body size was the only trait that appeared to have the same impact at both scales, with weak evidence of more dissimilar species tending to positively co-occur - $\beta = 0.039$ (CI[-0.1, 0.176], PD = 0.706) and 0.097 (CI[-0.045, 0.242], PD = 0.908) for local and global scales respectively. The only strong evidence of an effect on affinity in either direction were a negative effect of pronotum width dissimilarity at the global scale: -0.221 (CI[-0.429, -0.013], PD = 0.981), and of head width at the local scale -0.259 (CI[-0.472, -0.042], PD = 0.992) (Figure 4 a).

For the next analysis I used the phylogenetic tree from the global data set (Economo et al. 2018) to ask how phylogenetic distance affects affinity at both species and genus level. The species level analysis investigated the effect of phylogenetic distance on affinity *within* genus (as with the above example of biomass within functional guild). However, with 61 genera to analyse I was also able to employ a hierarchical model with random slopes for each genus. (This modification is simple within the proposed modelling framework and is an option available in the software package `CooccurrenceRegression.jl`). Here, again different scales yielded different results. There was strong evidence for a negative effect of phylogenetic distance on affinity at the genus level: -0.162 (CI[-0.171, -0.155], PD > 0.999) and of a positive relationship at the species level in the non-hierarchical (pooled) model: 0.071 (CI[0.03, 0.115], PD > 0.999). However, there was less certainty in the estimate from the hierarchical model: 0.059 (CI[-0.01, 0.12], PD = 0.955) (Figure 4 b). Here, by treating co-occurrence analysis as just another Bayesian GLM it was straightforward to employ the specific model structure desired to answer the exact question I was interested in - namely a hierarchical structure to measure the effect of

relatedness on affinity *within* genus. In this case, while it is reassuring that the pooled model and hierarchical model yield qualitatively similar results, the hierarchical model provides the more reliable results as it naturally accounts for clusters (here genera) in the data (McElreath 2018).

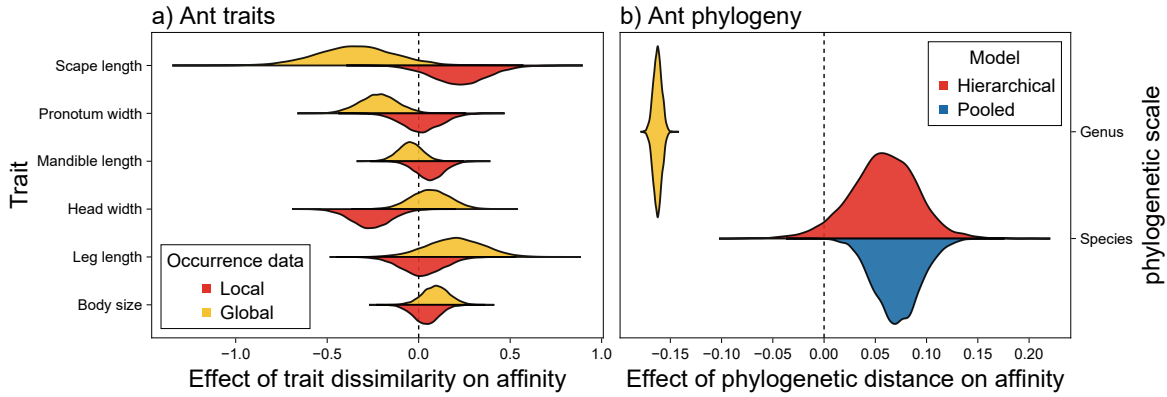


Figure 4: Posterior estimates of change in co-occurrence affinity between ant species owing to a) trait dissimilarity, b) phylogenetic distance.

290 Bacterial co-occurrence networks derived from cystic fibrosis patient sputum samples

The final example returns to the issue of simply analysing pairwise co-occurrence relationships, without any explanatory variables. Here, I constructed microbial co-occurrence networks from microbiome data derived from cystic fibrosis patients' sputum samples (Quinn et al. 2019) using the pairwise affinity estimates (median of posterior MCMC samples) but also retaining and reporting a measure of confidence in those estimates - probability of direction (PD). Figure 5 shows the inferred co-occurrence relationships between the 50 most prevalent genera in the data set.

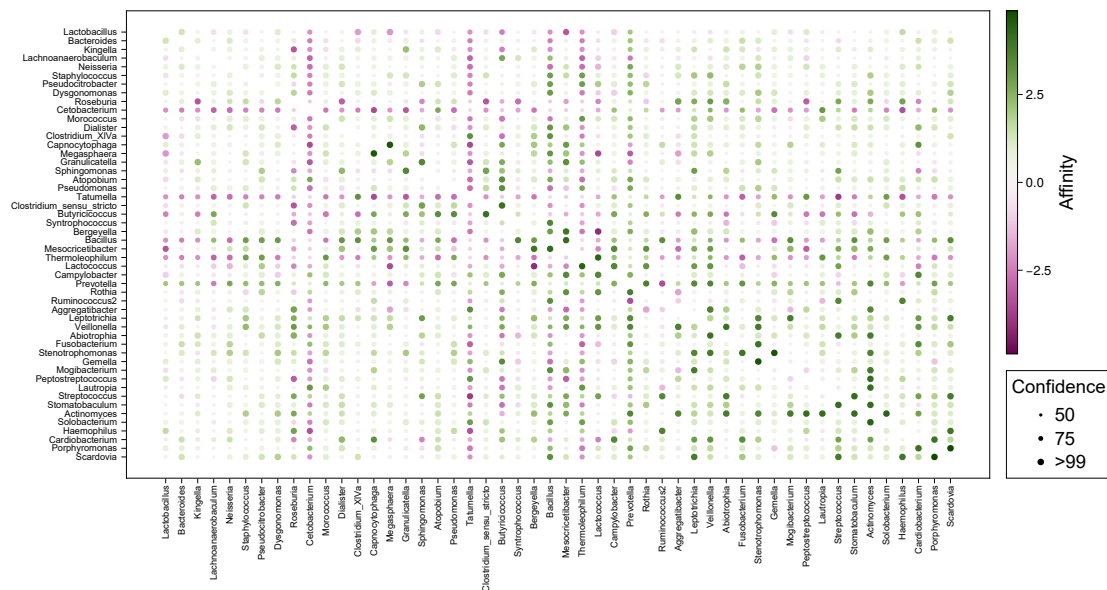


Figure 5: Affinity estimates between the 50 most prevalent genera across a set of cystic fibrosis patient sputum samples. Confidence estimates convey confidence that the affinity between the two genera is at least in the same direction as the point estimate.

Figure 6 Shows an alternative representation - networks constructed from only those genus pairs for which there is over 97.5% confidence in the direction of their affinity. In both representations it can be seen that genus pairs with higher absolute values for co-occurrence affinity do not always have higher confidence in their affinity estimates, though confidence and (absolute) affinity are clearly correlated.

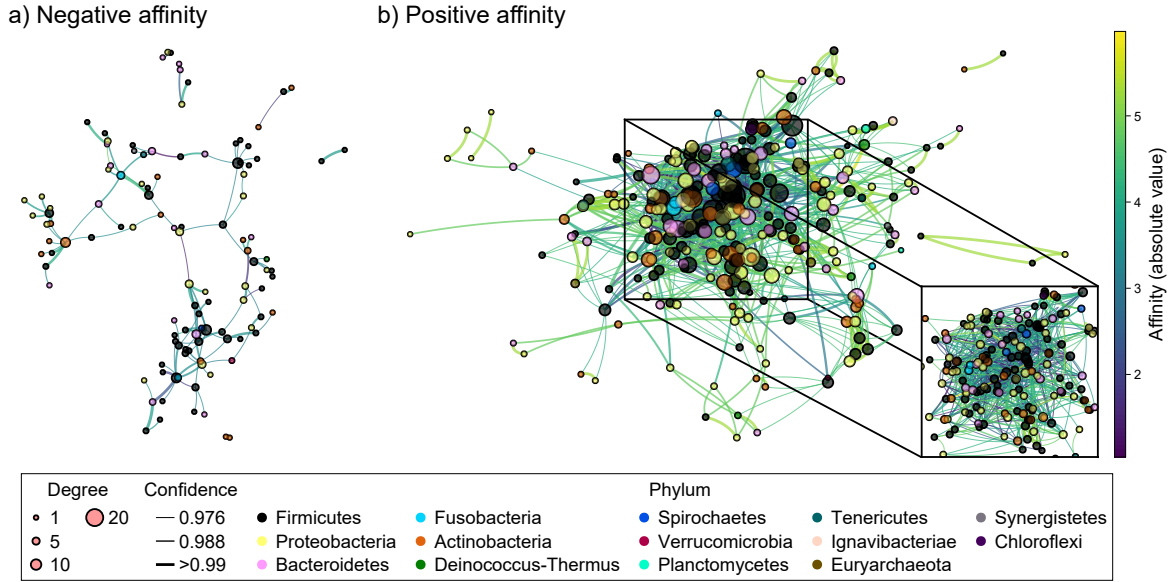


Figure 6: Microbial co-occurrence networks across a set of cystic fibrosis patient sputum samples. a) Negative affinity estimates. b) positive affinity estimates. Each vertex represents a single genus, with its size representing its *degree* (the number of genera it is linked to) and its colour representing its phylum. Confidence estimates convey confidence that the affinity between the two genera is at least in the same direction as the point estimate. Only genus pairs with greater than 97.5% confidence in the direction (positive/negative) of their affinity are included. Panel b includes an inset plot of the largest cluster of genera, showing only those links between genera contained in the cluster, with both the edges and vertices of the network reduce in scale for clarity. Here, we can see that these genera are linked together by many relatively small affinity estimates.

3.3 CooccurrenceRegression.jl

The main reason behind this work was to communicate a general methodological framework for the benefit of ecologists. In order to simplify the use of this framework I also present the Julia (Bezanson et al. 2017) package `CooccurrenceRegression.jl`, which is built on top the probabilistic programming language (PPL) `Turing.jl` (Ge, Xu, and Ghahramani 2018). With

308 this package one can recreate the single explanatory variable (Gower's distance) model from
309 the dung beetles example for N species across M sites as follows:

```
cooccurrence_regression(X,Y)
```

310 where X is an $M \times M$ array of dyadic explanatory variables e.g. a distance matrix and Y is
311 an $N \times M$ presence absence matrix with rows of sites and columns of species. Additionally,
312 one could replicate the second model from the dung beetle example by supplying a vector of
313 explanatory matrices:

```
X3 = X1 .* X2  
Xs = [X1, X2, X3]  
cooccurrence_regression(Xs,Y)
```

314 Lastly, hierarchical models can be run as follows

```
cooccurrence_regression(X,Y,g)
```

315 where g is a vector of length M assigning each species to a particular group.

316 More detail on changing priors and inference parameters can be found in the package repository

317 <https://github.com/EvoArt/CooccurrenceRegression.jl>.

3.4 Discussion

Here, I have shown how to construct a Bayesian dyadic GLM for the analysis of co-occurrence data. This builds on the work of Mainali et al. (2022) as well as Veech (2013) and Griffith et al. (2016). The identification of Fisher’s noncentral hypergeometric distribution (or mathematically equivalent formulations) as the correct distribution for modelling co-occurrence led first to null model approaches to co-occurrence analysis (Griffith, Veech, and Marsh 2016; Veech 2013), then to a useful co-occurrence metric (Mainali et al. 2022) and now to a general model capable of analysing raw co-occurrence data as a response variable even when data points are not independent (which will generally be the case). It should be noted that the co-occurrence relationships discussed here and in the works cited above are probabilistic in nature, i.e., I do not assume that either a high or low affinity between a specific pair of species implies significant ecological interaction. While co-occurrence analysis does not give the researcher direct access to ecological interaction data, it does have some bearing (Cazelles 2024) and may be used as an initial screening to identify likely interactions.

Part of the motivation for this work was the failure to recapture known regression coefficients when fitting linear models to pairwise affinity estimates (Figure 2). There may be other ways of combating this failure. For example, the removal of data points for which we have low confidence may be an option. For instance, if $N = 30$, $mA = 29$ and $mB = 1$, then a k of 1 tells us very little, since our null expectation is that k will very likely = 1. This will lead to a high affinity estimate but with a wide confidence interval and a high p-value. Using a cut-off threshold e.g. only using data for which $p < 0.05$ may lead to better results. However, these

339 data points are now *missing* from further analyses. The potential pitfalls involved in dealing
340 with missing data are numerous (Kang 2013), but it is unnecessary to risk the possible bias
341 associated with systematically removing data points when the Bayesian analysis framework
342 naturally accounts for differing levels of confidence between data (McElreath 2018).

343 The method proposed here is as flexible as any Bayesian GLM and can thus be adjusted to
344 fit the specific questions and modelling assumptions of the researcher. Here, I have shown
345 the applicability of this approach to analysing the relationship between trait dissimilarity or
346 phylogenetic distance and co-occurrence, as well as constructing co-occurrence networks. In
347 doing so, I have needed to use hierarchical models and interacting explanatory variables to
348 get accurate results. Many more model structures can be employed, as well as many more
349 applications inside and outside of ecology e.g. social networks in the humanities and gene
350 co-occurrence relationships in genetics. For simple models with a single presence/absence
351 matrix as response and one or more matrices of explanatory variables I have developed the
352 Julia (Bezanson et al. 2017) package CooccurrenceRegression.jl. However, it will serve many
353 researchers to consult accessible texts on probabilistic programming (McElreath 2018), read
354 the source code of the package, and develop their own models to suit their specific research
355 question.

356 While the focus here was on regression, retaining information on both confidence/uncertainty
357 and estimated effect size will also be important for downstream analyses of co-occurrence
358 networks e.g. generating random networks and analysing network metrics for each. However,
359 with the approach used here, the minimum confidence in any interaction is 50% (unless an

effect size threshold is used e.g. the probability of an absolute affinity value >1). I used probability of direction as a measure of confidence in the visualisations, because researchers are often interested in the binary classification of positive vs negative. However, other measures (e.g. $1/(\text{width of 95\% credible interval})$) may also be used. Future work should investigate the use of sparsity inducing spike and slab priors and their approximations (Castillo, Schmidt-Hieber, and Van der Vaart 2015), to model the assumption that many species will not interact in any meaningful way. Although, assumption may or may not be valid for the affinity metric, baring in mind that affinity measures a statistical likelihood for species to co-occur, and does not directly measure biotic interactions.

The method proposed here provides ecologist with an important new tool for the analysis of co-occurrence, and in particular discovering the relationships between co-occurrence and other variables e.g. phylogenetic distance, which is an active area of research (Goberna et al. 2019) and has been the subject of much research effort over the past few decades (Webb et al. 2002) but now has an analysis framework based on the simple application of probability theory (McElreath 2018) with correct modelling of co-occurrence probabilities (Griffith, Veech, and Marsh 2016; Veech 2013; Mainali et al. 2022) while accounting for non-independence of data in matrices of pairwise species measurements (Clarke, Rothery, and Raybould 2002; Gompert et al. 2014).

References

- Barberán, Albert, Scott T Bates, Emilio O Casamayor, and Noah Fierer. 2012. “Using Network Analysis to Explore Co-Occurrence Patterns in Soil Microbial Communities.” *The ISME Journal* 6 (2): 343–51.
- Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. “Julia: A Fresh Approach to Numerical Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- Castillo, Ismaël, Johannes Schmidt-Hieber, and Aad Van der Vaart. 2015. “Bayesian Linear Regression with Sparse Priors.” *The Annals of Statistics*, 1986–2018.
- Cazelles, Kevin. 2024. “Isolating Interactions from Co-Occurrences.” *Nature Ecology & Evolution* 8 (2): 184–85.
- Clarke, Ralph T., Peter Rothery, and Alan F. Raybould. 2002. “Confidence Limits for Regression Relationships Between Distance Matrices: Estimating Gene Flow with Distance.” *Journal of Agricultural, Biological, and Environmental Statistics* 7 (3): 361–72. <http://www.jstor.org/stable/1400681>.
- Danisch, Simon, and Julius Krumbiegel. 2021. “Makie.jl: Flexible High-Performance Data Visualization for Julia.” *Journal of Open Source Software* 6 (65): 3349. <https://doi.org/10.21105/joss.03349>.
- Economo, Evan P, Nitish Narula, Nicholas R Friedman, Michael D Weiser, and Benoit Guénard. 2018. “Macroecology and Macroevolution of the Latitudinal Diversity Gradient in Ants.” *Nature Communications* 9 (1): 1778.

- 399 Ge, Hong, Kai Xu, and Zoubin Ghahramani. 2018. “Turing: A Language for Flexible Prob-
400 abilistic Inference.” In *International Conference on Artificial Intelligence and Statistics*,
401 *AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*, 1682–
402 90. <http://proceedings.mlr.press/v84/ge18b.html>.
- 403 Goberna, Marta, Alicia Montesinos-Navarro, Alfonso Valiente-Banuet, Yannick Colin, Alicia
404 Gómez-Fernández, Santiago Donat, Jose A Navarro-Cano, and Miguel Verdú. 2019. “Incor-
405 porating Phylogenetic Metrics to Microbial Co-Occurrence Networks Based on Amplicon
406 Sequences to Discern Community Assembly Processes.” *Molecular Ecology Resources* 19
407 (6): 1552–64.
- 408 Gompert, Zachariah, Lauren K. Lucas, C. Alex Buerkle, Matthew L. Forister, James A.
409 Fordyce, and Chris C. Nice. 2014. “Admixture and the Organization of Genetic Diversity
410 in a Butterfly Species Complex Revealed Through Common and Rare Genetic Variants.”
411 *Molecular Ecology* 23 (18): 4555–73. <https://doi.org/https://doi.org/10.1111/mec.12811>.
- 412 Gotelli, Nicholas J, and Declan J McCabe. 2002. “Species Co-Occurrence: A Meta-Analysis
413 of JM Diamond’s Assembly Rules Model.” *Ecology* 83 (8): 2091–96.
- 414 Gower, John C. 1971. “A General Coefficient of Similarity and Some of Its Properties.” *Bio-*
415 *metrics*, 857–71.
- 416 Griffith, Daniel M, Joseph A Veech, and Charles J Marsh. 2016. “Cooccur: Probabilistic
417 Species Co-Occurrence Analysis in r.” *Journal of Statistical Software* 69: 1–17.
- 418 Homan, Matthew D., and Andrew Gelman. 2014. “The No-u-Turn Sampler: Adaptively
419 Setting Path Lengths in Hamiltonian Monte Carlo.” *J. Mach. Learn. Res.* 15 (1): 1593–
420 623.

- 421 Kang, Hyun. 2013. “The Prevention and Handling of the Missing Data.” *Korean Journal of*
422 *Anesthesiology* 64 (5): 402–6.
- 423 Kraan, Casper, Simon F Thrush, and Carsten F Dormann. 2020. “Co-Occurrence Patterns
424 and the Large-Scale Spatial Structure of Benthic Communities in Seagrass Meadows and
425 Bare Sand.” *BMC Ecology* 20 (1): 1–8.
- 426 Liu, Dong C, and Jorge Nocedal. 1989. “On the Limited Memory BFGS Method for Large
427 Scale Optimization.” *Mathematical Programming* 45 (1-3): 503–28.
- 428 Mainali, Kumar P., Eric Slud, Michael C. Singer, and William F. Fagan. 2022. “A Better
429 Index for Analysis of Co-Occurrence and Similarity.” *Science Advances* 8 (4): eabj9204.
430 <https://doi.org/10.1126/sciadv.abj9204>.
- 431 McElreath, Richard. 2018. *Statistical Rethinking: A Bayesian Course with Examples in r and*
432 *Stan*. Chapman; Hall/CRC.
- 433 Mogensen, Patrick Kofod, and Asbjørn Nilsen Riseth. 2018. “Optim: A Mathematical
434 Optimization Package for Julia.” *Journal of Open Source Software* 3 (24): 615. <https://doi.org/10.21105/joss.00615>.
435
- 436 Nelder, John A, and Roger Mead. 1965. “A Simplex Method for Function Minimization.” *The*
437 *Computer Journal* 7 (4): 308–13.
- 438 Nunes, Cássio Alencar, Rodrigo Fagundes Braga, José Eugênio Cortes Figueira, Frederico de
439 Siqueira Neves, and G Wilson Fernandes. 2016. “Dung Beetles Along a Tropical Altitudinal
440 Gradient: Environmental Filtering on Taxonomic and Functional Diversity.” *PLoS One* 11
441 (6): e0157442.
- 442 Quinn, Robert A, Sandeep Adem, Robert H Mills, William Comstock, Lindsay DeRight Golda-

sich, Gregory Humphrey, Alexander A Aksenov, et al. 2019. “Neutrophilic Proteolysis in the Cystic Fibrosis Lung Correlates with a Pathogenic Microbiome.” *Microbiome* 7: 1–13.

R Core Team. 2014. “R Core Team (2014). R: A Language and Environment for Statistical Computing.” *R Foundation for Statistical Computing, Vienna, Austria*. URL <Http://Www.R-Project.org/>.

Reeve, Richard, Michael Borregaard, and Claire Harris. 2024. “Phylo.jl.” Zenodo. <https://doi.org/10.5281/zenodo.12789192>.

Veech, Joseph A. 2013. “A Probabilistic Model for Analysing Species Co-Occurrence.” *Global Ecology and Biogeography* 22 (2): 252–60.

Wang, Qiong, George M Garrity, James M Tiedje, and James R Cole. 2007. “Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy.” *Applied and Environmental Microbiology* 73 (16): 5261–67.

Webb, Campbell O, David D Ackerly, Mark A McPeck, and Michael J Donoghue. 2002. “Phylogenies and Community Ecology.” *Annual Review of Ecology and Systematics* 33 (1): 475–505.

Williams, Ryan J, Adina Howe, and Kirsten S Hofmockel. 2014. “Demonstrating Microbial Co-Occurrence Pattern Analyses Within and Between Ecosystems.” *Frontiers in Microbiology* 5: 358.

Wong, Mark KL, Toby PN Tsang, Owen T Lewis, and Benoit Guénard. 2021. “Trait-Similarity and Trait-Hierarchy Jointly Determine Fine-Scale Spatial Associations of Resident and Invasive Ant Species.” *Ecography* 44 (4): 589–601.

Zhu, Wentao, Ming Zhu, Xiangbo Liu, Jingquan Xia, Hongyang Yin, and Xiubao Li.

465 2023. “Different Responses of Bacteria and Microeukaryote to Assembly Processes and
466 Co-Occurrence Pattern in the Coastal Upwelling.” *Microbial Ecology* 86 (1): 174–86.