

Community Confounding In Joint Species Distribution Models

Justin J. Van Ee (✉ vanee002@colostate.edu)

Colorado State University

Jacob S. Ivan

Colorado Parks and Wildlife

Mevin B. Hooten

The University of Texas at Austin

Research Article

Keywords: models, relationships, epidemic, community

Posted Date: December 6th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1111644/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on July 18th, 2022. See the published version at <https://doi.org/10.1038/s41598-022-15694-6>.

1 Community Confounding in Joint Species

2 Distribution Models

3 **Justin J. Van Ee^{1,*}, Jacob S. Ivan², and Mevin B. Hooten³**

4 ¹Colorado State University, Department of Statistics, Fort Collins, 80523, USA

5 ²Colorado Parks and Wildlife, Fort Collins, 80526, USA

6 ³University of Texas at Austin, Department of Statistics and Data Sciences, Austin, 78712, USA

7 *vanee002@colostate.edu

8 **ABSTRACT**

Joint species distribution models have become ubiquitous for studying species-habitat relationships and dependence among species. Accounting for community structure often improves predictive power, but can also alter inference on species-habitat relationships. Modulated species-habitat relationships are indicative of community confounding: The situation in which interspecies dependence and habitat effects compete to explain species distributions. We discuss community confounding in a case study of mammalian responses to the Colorado bark beetle epidemic in the subalpine forest by comparing the inference from independent single species distribution models and a joint species distribution model. We present a method for measuring community confounding and develop a restricted version of our hierarchical model that orthogonalizes the habitat and species random effects. Our results indicate that variables associated with the severity and duration of the bark beetle epidemic suffer from community confounding. This implies that mammalian responses to the bark beetle epidemic are governed by interconnected habitat and community effects. Disentangling habitat and community effects can improve our understanding of the ecological system and possible management strategies. We evaluate restricted regression as a method for alleviating community confounding and distinguish it from other inferential methods for confounded models.

10 **Introduction**

11 Ecological datasets that provide insights about collections of organisms have become prevalent over the last decade thanks to efforts like Long Term Ecological Research Network (LTER), National Ecological Observatory Network (NEON), citizen science surveys, etc.¹. In addition, technology has improved our ability to fit modern statistical models to these datasets jointly. As a consequence, joint species distribution models (JSDM)²⁻⁴ have become popular for modeling dependence among species simultaneously with environmental drivers of occurrence and/or abundance. JSDMs provide inference for species-habitat relationship conditional on the community. Species habitat preferences are modulated by interspecies dynamics⁵⁻⁷. Therefore, we should consider the relationships between community and environmental effects in JSDMs.

18 We describe a statistical model that we used to improve our understanding of montaine mammal communities in what
19 follows. Critically, we consider potential confounding among interspecies relationships and environmental relationships in our
20 model and demonstrate how to orthogonalize these mechanisms in the model and the resulting inference compared to models
21 where species are treated independently. Unlike previous approaches that have applied restricted regression techniques similar
22 to ours, we use it in the context of well-known ecological models for species occupancy and abundance. While such approaches
23 have become well-known in spatial statistics and environmental science, they have not been adopted in settings involving the
24 multivariate analysis of community data.

25 **Joint Species Distribution Model**

26 We present a JSDM based on a multispecies extension to the Royle-Nichols model⁸. The Royle-Nichols model accounts for
27 heterogeneity in detection induced by the species' latent intensity, a surrogate related to true species abundance. Abundance
28 estimation often requires an explicit spatial region that is closed to emmigration and immigration. Intensity is a measure of how
29 frequently members of the species use a region and does not require a population closure assumption. In the *Model* section, we
30 further discuss the distinctions between abundance and intensity in the Royle-Nichols model.

31 The Royle-Nichols model utilizes occupancy survey data but provides inference distinct from the basic occupancy model⁹.
32 In the Royle-Nichols model, we estimate individual detection probability for homogeneous members of the population, whereas
33 in an occupancy model, we estimate probability of observing at least one member of the population given that the site is
34 occupied. Furthermore, the Royle-Nichols model allows us to relate habitat covariates to the latent intensity associated with a
35 species at a site, while in an occupancy model, habitat covariates are associated with the species latent probability of occupancy
36 at a site. Species intensity and occupancy may be governed by different mechanisms, and inference from an intensity model
37 can be distinct from that provided by an occupancy model¹⁰⁻¹². Cingolani et al.¹¹ proposed that, in plant communities, certain
38 environmental filters preclude species from occupying a site and an additional set of filters may regulate if a species can flourish.
39 Hence, certain covariates that were unimportant in an occupancy model may improve predictive power in an intensity model.

40 **Community Confounding**

41 Species distributions are shaped by habitat as well as competition and mutualism within the community^{7,13,14}. Community
42 confounding occurs when species distributions are explained by a convolution of habitat and interspecies effects and can lead to
43 inferential discrepancies between a joint and single species distribution model. Former studies have incorporated interspecies
44 dependence into an occupancy model¹⁵⁻¹⁸, and others have addressed spatial confounding^{1,19-21}, but none of these considered
45 community confounding.

46 We address community confounding by formulating a version of our model that orthogonalizes the habitat effects and
47 species random effects. Orthogonalizing the fixed and random effects is common practice in spatial statistics and often referred
48 to as restricted spatial regression²²⁻²⁶. Restricted regression has been applied to spatial generalized linear mixed models
49 (SGLMM) for observations \mathbf{y} , which can be expressed as

$$\mathbf{y} \sim [\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\psi}], \quad (1)$$

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}, \quad (2)$$

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{H}(\boldsymbol{\phi})), \quad (3)$$

50 where $g(\cdot)$ is a link function, $\boldsymbol{\psi}$ are additional parameters for the data model, and $\boldsymbol{\phi}$ parameterizes the random effect. In the
 51 SGLMM, prior information facilitates the estimation of $\boldsymbol{\eta}$, which would not be estimable otherwise due to its shared column
 52 space with $\boldsymbol{\beta}$ ²⁴. This is analogous to applying a ridge penalty to $\boldsymbol{\eta}$, which stabilizes the likelihood. Another method for fitting
 53 the confounded SGLMM is to specify a restricted version:

$$\mathbf{y} \sim [\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\psi}], \quad (4)$$

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\theta} + (\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}, \quad (5)$$

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{H}(\boldsymbol{\phi})), \quad (6)$$

54 where \mathbf{P}_X is the projection matrix onto the column space of \mathbf{X} . In the unrestricted SGLMM, variability in the latent mean $\boldsymbol{\mu}$
 55 in the direction of \mathbf{X} is explained by both the regression coefficients $\boldsymbol{\beta}$ in (2) and random effect $\boldsymbol{\eta}$. In the restricted model,
 56 however, all variability in the direction of \mathbf{X} is explained solely by the regression coefficients $\boldsymbol{\theta}$ in (5)²⁵. We refer to $\boldsymbol{\beta}$ as the
 57 conditional effects because they depend on $\boldsymbol{\eta}$, and $\boldsymbol{\theta}$ as the unconditional effects.

58 Restricted regression, as specified (4)-(6), was proposed by Reich et al.²². Reich et al.²² described a disease-mapping
 59 example in which the inclusion of a spatial random effect rendered one covariate effect unimportant that was important in
 60 the non-spatial model. Spatial maps indicated an association between the covariate and response, making inference from the
 61 spatial model appear untenable. Reich et al.²² proposed restricted spatial regression as a method for recovering the posterior
 62 expectations of the non-spatial model and shrinking the posterior variances which tend to be inflated for the unrestricted
 63 SGLMM.

64 Several modifications of restricted spatial regression have been proposed^{24,27-30}. All restricted spatial regression methods
 65 seek to provide posterior means $E(\theta_j|\mathbf{Y})$ and marginal posterior variances $\text{Var}(\theta_j|\mathbf{Y})$, $j = 1, \dots, p$ that satisfy the following two
 66 conditions³¹:

- 67 1. $E(\boldsymbol{\theta}|\mathbf{Y}) = E(\boldsymbol{\beta}_{NS}|\mathbf{Y})$ and,
- 68 2. $\text{Var}(\boldsymbol{\beta}_{NS,j}|\mathbf{Y}) \leq \text{Var}(\theta_j|\mathbf{Y}) \leq \text{Var}(\boldsymbol{\beta}_{Spatial,j}|\mathbf{Y})$ for $j = 1, \dots, p$,

69 where $\boldsymbol{\beta}_{NS}$ and $\boldsymbol{\beta}_{Spatial}$ are the regression coefficients corresponding to the non-spatial and unrestricted spatial models,
 70 respectively.

71 The inferential impacts of spatial confounding on the regression coefficients has been debated. Hodges and Reich²³ outlined
 72 five viewpoints on spatial confounding and restricted regression in the literature and refuted the two following views:

- 73 1. Adding the random effect $\boldsymbol{\eta}$ corrects for bias in $\boldsymbol{\beta}$ resulting from missing covariates.

74 2. Estimates of β in a SGLMM are shrunk by the random effect and hence conservative.

75 The random effect η can increase or decrease the magnitude of β , and the change may be galvanized by mechanisms not related
76 to missing covariates. Therefore, we cannot assume the regression coefficients in the SGLMM will exceed those of the restricted
77 model, nor should we regard the estimates in either model as biased due to misspecification. Confounding in the SGLMM
78 causes $\text{Var}(\beta_j|\mathbf{Y}) \geq \text{Var}(\theta_j|\mathbf{Y})$, $j = 1, \dots, p$, because of the shared column space of the fixed and random effects. Thus, we
79 refer to the conditional coefficients as conservative with regards to their credible intervals, not their posterior expectations.

80 Reich et al.²² argued that restricted spatial regression should always be applied because the spatial random effect is generally
81 added to improve predictions and/or correct the fixed effect variance estimate. While it may be inappropriate to orthogonalize a
82 set of fixed effects in an ordinary linear model, orthogonalizing the fixed and random effect in a spatial model is permissible
83 because the random effect is generally not of inferential interest. Paciorek³² provided the alternative perspective that, if
84 confounding exists, it is inappropriate to attribute all contested variability in \mathbf{y} to the fixed effects. Hanks et al.²⁵ discussed
85 factors for deciding between the unrestricted and restricted SGLMM on a continuous spatial support. The restricted SGLMM
86 leads to improved computational efficiency, but the unconditional effects are less conservative under model misspecification
87 and more prone to type-S errors: The Bayesian analogue of Type I error. Fitting the unrestricted SGLMM when the fixed
88 and random effect are truly orthogonal does not introduce bias, but it will increase the fixed effect variance. Given these
89 considerations, Hanks et al.²⁵ suggested a hybrid approach where the conditional effects, β , are extracted from the restricted
90 SGLMM. This is possible because the restricted SGLMM is a reparameterization of the unrestricted SGLMM. This hybrid
91 approach leads to improved computational efficiency but yields the more conservative parameter estimates.

92 Restricted regression has also been applied in time series applications. Dominici et al.³³ debiased estimates of fixed effects
93 confounded by time using restricted smoothing splines. Without the temporal random effect, Dominici et al.³³ asserted all
94 temporal variation in the response would be wrongly attributed to temporally correlated fixed effects. Houseman et al.³⁴ used
95 restricted regression to ensure identifiability of a nonparametric temporal effect and highlighted certain covariate effects that
96 were more evident in the restricted model (i.e., the unconditional effects' magnitude was greater). Furthermore, restricted
97 regression is implicit in restricted maximum likelihood estimation (REML). REML is often employed for debiasing the estimate
98 of the variance of \mathbf{Y} in linear regression and fitting linear mixed models that are not estimable in their unrestricted format³⁵.
99 Because REML is generally applied in the context of variance and covariance estimation, considerations regarding the effects
100 of REML on inference for the fixed effects are lacking in the literature.

101 In ecological science, multispecies models often include random effects to account for interspecies dependence. Interspecies
102 dependence in the random effects can be characterized by a covariance matrix. Unlike a spatial or temporal random effect,
103 we consider species random effects to be inferentially important, rather than a tool for improving predictions or catch-all for
104 missing covariates. A restriction approach in a multispecies model attributes contested variation between the fixed effects
105 (habitat information) and random effect (community information) to the fixed effect.

106 We demonstrate restricted regression in multispecies models by comparing inference between a restricted and unrestricted

107 model for the camera trap data. Furthermore, we discuss community confounding and its relevance to restricted regression in
 108 the context of multispecies models. We present a method for measuring confounding and highlight its inferential utility. We
 109 also discuss other inferential methods for confounded models and consider their appropriateness in the multispecies context.

110 In what follows, we motivate and formulate a multispecies model. Referencing our multispecies model, we describe
 111 community confounding and propose methods to detect and alleviate it. We then apply our model to data collected from 2013-14
 112 during a Colorado Parks and Wildlife camera trap survey. Finally, we discuss our findings and highlight their implications for
 113 inference in models that exhibit community confounding.

114 **Model**

115 **Royle-Nichols Link Function**

116 We specify a model for analyzing multispecies binary detection data, y_{ijk} , arising from a Bernoulli process with probability
 117 of success p_{ijk} , where $i = 1, \dots, n$, $j = 1, \dots, J_i$, and $k = 1, \dots, K$ correspond to sites, occasions, and species, respectively.
 118 Occupancy data of this form have traditionally been analyzed in a latent variable framework^{9,36,37}. Let $z_{ik} \sim \text{Bern}(\psi_{ik})$ be
 119 an indicator on whether species k occupies site i . Given the site is occupied, we detect species k on occasion j with some
 120 probability p_{ijk} , such that $(y_{ijk}|z_{ik} = 1) \sim \text{Bern}(p_{ijk})$, but if species k is absent from the site, we have zero probability of
 121 detecting it, $P(y_{ijk} = 0|z_{ik} = 0) = 1$.

122 Royle and Nichols⁸ introduced a method for analyzing occupancy data that explicitly modeled the probability of detecting
 123 species k at a site as a function of a surrogate related to the true species abundance. Assuming there are N_{ik} individuals of
 124 species k in sample region i and that all individuals in species k on the sample unit have identical detection probabilities and are
 125 detected independently of other individuals, the probability of detecting at least one of these individuals can be expressed as

$$126 \quad p_{ijk} = 1 - (1 - r_{jk})^{N_{ik}}, \quad (7)$$

126 where r_{jk} is a binomial sampling probability that a particular individual of species k is detected on occasion j . While the
 127 Royle-Nichols model facilitates inference on number of individuals of species k , N_{ik} , at each site when all the assumptions
 128 are met, we do not interpret them as such because sites are not necessarily closed in camera trap studies due to highly mobile
 129 species with home ranges larger than the sampling radius of the camera.

130 The nonlinear function of r_{jk} and N_{ik} in (7) involves more parameters than would be identifiable in a typical occupancy
 131 model, especially when the individual detection probability is heterogeneous across occasions (e.g., r_{jk} are heterogeneous). In
 132 the heterogeneous case, r_{jk} is connected to covariates with the logit link function:

$$133 \quad \text{logit}(r_{jk}) = f(\mathbf{w}_{jk}, \boldsymbol{\alpha}_k), \quad (8)$$

133 where $f(\mathbf{w}_{jk}, \boldsymbol{\alpha}_k)$ is a linear function of the covariates \mathbf{w}_{jk} and regression parameters $\boldsymbol{\alpha}_k$.

134 Modeling Interspecies Dependence

135 Following Royle and Nichols⁸, we assume $N_{ik} \sim \text{Pois}(\lambda_{ik})$, where λ_{ik} is mean intensity of species k at site i . We let $\boldsymbol{\lambda}$ denote
136 the vectorized intensities of the K species in the community stacked across all n sites. To model interspecies dependence, we
137 specify the conditional multivariate normal distribution:

$$\log(\boldsymbol{\lambda}) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}, \mathbf{H}), \quad (9)$$

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{spp} \otimes \mathbf{I}_n), \quad (10)$$

138 where \mathbf{X} is a block-diagonal matrix of the K species design matrices, $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_K)'$ is a stacked vector of species specific
139 regression coefficients, $\boldsymbol{\eta}$ represents the species random effects, and $\boldsymbol{\Sigma}_{spp}$ is a species covariance matrix, and \mathbf{H} is a matrix
140 that allows for additional covariance structures such as spatial dependence. For our purpose of comparing the single and
141 joint species distribution models, we marginalized the model over the species random effects and set $\mathbf{H} = \mathbf{0}$. The resulting
142 distribution for the species latent intensities was $\log(\boldsymbol{\lambda}) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}_{spp} \otimes \mathbf{I}_n)$.

143 The formulation given in (9) allows for dependence between all K species in the community at each site. Previous joint
144 occupancy models allowed for interspecies dependence in the probability of occupancy¹⁵⁻¹⁷, whereas our model allows for
145 dependence in the species latent intensities. Just as certain habitat features may not preclude species occupancy but can curb
146 intensity, some species may coexist in a region but not be able to jointly flourish³⁸. Hence, interspecies dependence that is
147 observed in a joint intensity model may not be present in a joint occupancy model.

148 Scheffe³⁹ stipulated that the levels of a random effect are draws from a population, and the draws are not of interest in
149 themselves but only as samples from the larger population, which is of interest. In more recent literature, the term random effect
150 is used more broadly. Hodges and Clayton⁴⁰ categorized modern definitions of a random effect into three different varieties.
151 The definition commonly used in spatial statistics is, the levels of the effect arise from a meaningful population, but they are the
152 whole population and these particular levels are of interest. We adopt this definition for the species random effects in (9). In
153 practice, some levels of the population will likely not be included in the species random effects. For example, in Ivan et al.⁴¹,
154 cameras were baited and arranged to capture all members of the mammalian community, but several species were excluded
155 from the species random effects due to a lack of detections.

156 Priors

157 We used normal priors for the regression coefficients in both the intensity, $\boldsymbol{\beta}$, and detection, $\boldsymbol{\alpha}$, process and choose the conjugate
158 Inverse-Wishart prior for the species covariance matrix $\boldsymbol{\Sigma}_{spp}$. A more general alternative to the Inverse-Wishart prior is to apply
159 a Cholesky decomposition, $\boldsymbol{\Sigma}_{spp} = \mathbf{L}\mathbf{D}^{-1}\mathbf{L}'$, where \mathbf{L} is lower diagonal with ones along the diagonal and \mathbf{D} is diagonal with
160 positive diagonal elements, and specify priors for the lower diagonal elements of \mathbf{L} and diagonal elements of \mathbf{D} ⁴². We found the
161 Inverse-Wishart prior suitable for our inferential goals, but see Chan and Jeliazkov⁴² for alternative covariance matrix priors.

The joint posterior distribution associated with our model is

$$[\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\lambda}, \mathbf{N}, \boldsymbol{\Sigma}_{spp} | \mathbf{Y}] \propto \prod_{k=1}^K \left(\prod_{i=1}^n \left(\prod_{j=1}^{J_i} \left([y_{ijk} | N_{ik}, \boldsymbol{\alpha}_k] \right) [N_{ik} | \lambda_{ik}] \right) [\boldsymbol{\alpha}_k] [\boldsymbol{\beta}_k] \right) [\boldsymbol{\lambda} | \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K, \boldsymbol{\Sigma}_{spp}] [\boldsymbol{\Sigma}_{spp}]. \quad (11)$$

162 See Appendix A for the full statement of our Royle-Nichols multispecies model.

163 Community Confounding

164 Restricted Regression Approach

165 We fit a restricted version of the Royle-Nichols multispecies model that orthogonalizes the fixed effects and random effect. We
166 express the species latent intensity process conditionally as

$$\log(\boldsymbol{\lambda}) \sim \mathcal{N}(\mathbf{X}\mathbf{B} + (\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}, \mathbf{H}), \quad (12)$$

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{spp} \otimes \mathbf{I}_n), \quad (13)$$

167 where \mathbf{P}_X is the projection matrix onto the column space of \mathbf{X} . This specification forces the species covariance matrix to
168 explain patterns in the community that are orthogonal to the fixed effects. As in (9), we marginalized over the species random
169 effects but did not set $\mathbf{H} = \mathbf{0}$ since the resulting normal distribution would be degenerate because $(\mathbf{I} - \mathbf{P}_X)(\boldsymbol{\Sigma}_{spp} \otimes \mathbf{I}_n)(\mathbf{I} - \mathbf{P}_X)$
170 is not invertible. To remedy the singularity of the covariance matrix, we applied a ridge penalty by letting $\mathbf{H} = \tau^2 \mathbf{I}$. The
171 resulting distribution was $\log(\boldsymbol{\lambda}) \sim \mathcal{N}(\mathbf{X}\mathbf{B}, \tau^2 \mathbf{I} + (\mathbf{I} - \mathbf{P}_X)(\boldsymbol{\Sigma}_{spp} \otimes \mathbf{I}_n)(\mathbf{I} - \mathbf{P}_X))$. For implementation, we kept the hierarchical
172 representation in equations (12) and (13) and obtained samples of $(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}$ in our Markov Chain Monte Carlo (MCMC)
173 through conditioning by Kriging^{25,43}.

174 In the restricted model, we set $\tau^2 = 2.25$. This choice was supported by the asymptotic equivalence between Poisson
175 and logistic regression. In a generalized occupancy model, the latent probability of occupancy is specified as $\text{logit}(\psi_i) \sim$
176 $\mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}, \tau^2)$.⁴⁴ investigated the relation between the prior on $\boldsymbol{\beta}$ and induced prior on the latent probability of success ψ_i in
177 logistic regression; their work showed that specifying an uninformative normal prior on $\boldsymbol{\beta}$, i.e. setting τ^2 large, induces a
178 U-shaped prior for ψ_i with most of the density concentrated near 0 and 1. Broms et al.¹⁶ recommended setting $\tau^2 = 2.25$ in
179 occupancy models, which results in a relatively flat prior for ψ . For rare species, λ_i in (12) is analogous to ψ_i , and specifying a
180 variance of $\tau^2 = 2.25$ is minimally informative.

181 Baddeley⁴⁵ motivated the asymptotic equivalence of Poisson and logistic regression in a spatial context where counts of
182 points from a non-homogeneous Poisson process are recorded in a lattice; they showed that, as the grid cells of the lattice
183 become infinitesimally small, the inference yielded from Poisson and logistic regression are equivalent. This result can
184 be applied more generally to any dataset where there is a high proportion of zero counts. We demonstrate the asymptotic
185 equivalence between Poisson and logistic regression in the Royle-Nichols model in Appendix D.

186 Measuring Confounding

187 Hefley et al.²⁶ showed how to assess confounding in SGLMM models by computing the Pearson correlation coefficient between
188 each pair of covariates and eigenvectors from the spectral decomposition of the spatial covariance matrix. We propose another
189 approach relevant to our method that aids in interpretation. We compute the coefficient of determination of each covariate for
190 species k regressed on the estimated latent intensities (no intercept) of the $K - 1$ other species in the community. Because the
191 latent intensities are unknown, the coefficients of determination of all covariates are derived quantities and can be computed at
192 each iteration of the MCMC algorithm:

$$R^{2(l)}(\mathbf{x}_k) = \frac{SSR^{(l)}(\mathbf{x}_k)}{SST(\mathbf{x}_k)} = \frac{\left(\mathbf{\Lambda}_{-k}^{(l)} \hat{\boldsymbol{\theta}}_{-k}^{(l)} - \bar{\mathbf{x}}_k\right)' \left(\mathbf{\Lambda}_{-k}^{(l)} \hat{\boldsymbol{\theta}}_{-k}^{(l)} - \bar{\mathbf{x}}_k\right)}{\left(\mathbf{x}_k - \bar{\mathbf{x}}_k\right)' \left(\mathbf{x}_k - \bar{\mathbf{x}}_k\right)}, \quad (14)$$

193 where $\bar{\mathbf{x}}_k = (\bar{x}_k, \dots, \bar{x}_k)$ is the mean of the covariate \mathbf{x}_k for species k repeated n times, $\mathbf{\Lambda}_{-k}^{(l)} = \left(\boldsymbol{\lambda}_1^{(l)}, \dots, \boldsymbol{\lambda}_{k-1}^{(l)}, \boldsymbol{\lambda}_{k+1}^{(l)}, \dots, \boldsymbol{\lambda}_K^{(l)}\right)$ is
194 a matrix of the $K - 1$ other species intensities sampled for MCMC iteration l , and $\hat{\boldsymbol{\theta}}_{-k}^{(l)}$ are estimated regression coefficients
195 relating the estimated species intensities at iteration l to \mathbf{x}_k . The posterior mean $E\left(R^2(\mathbf{x}_k) | \mathbf{Y}\right)$ provides a measure of community
196 confounding for the covariate \mathbf{x}_k and can help identify which fixed effects will vary between the unrestricted and restricted
197 models. We demonstrate this approach in the following case study.

198 Camera Trap Survey

199 Study Area

200 We analyzed data arising from a study area comprised of subalpine forests in the state of Colorado between 2590 and 3660
201 m elevation (Figure 1). Sites were restricted to public lands managed by the United States Forest Service, National Park
202 Service, Bureau of Land Management, and Colorado State Forest Service. Forests in our study area were primarily composed
203 of Lodgepole pine (*Pinus contorta*), Engelmann spruce (*Picea engelmannii*), and subalpine fir (*Abies lasiocarpa*). Lodgepole
204 pine was dominant at lower elevations as well as higher elevations that were drier and/or on south-facing slopes; high elevation
205 regions that had cool north-facing slopes were co-dominated by Engelmann spruce and subalpine fir. Lodgepole pine is
206 restricted to the northern two-thirds of Colorado, so all sites in the southern region of the study area were Engelmann spruce,
207 subalpine fir co-dominated. Quaking aspen (*Populus tremuloides*), Douglas-fir (*Pseudotsuga menziesii*), bristlecone pine
208 (*Pinus aristata*), limber pine (*Pinus flexilis*), and blue spruce (*Picea pungens*) were also present at some sites. Mean July and
209 January temperature across the study area were 14°C and -6.1°C respectively. All camera data were collected during summers
210 2013-2014.

211 Sampling Design

212 The primary goal of Ivan et al.⁴¹ was to assess mammalian responses to bark beetle outbreaks, thus sites were randomly
213 selected to facilitate inference on the beetle outbreak covariates. Beetle outbreak covariates included the number of years since
214 the initial outbreak (YSO) and the severity of the outbreak measured by mean overstory mortality (severity). The sample of

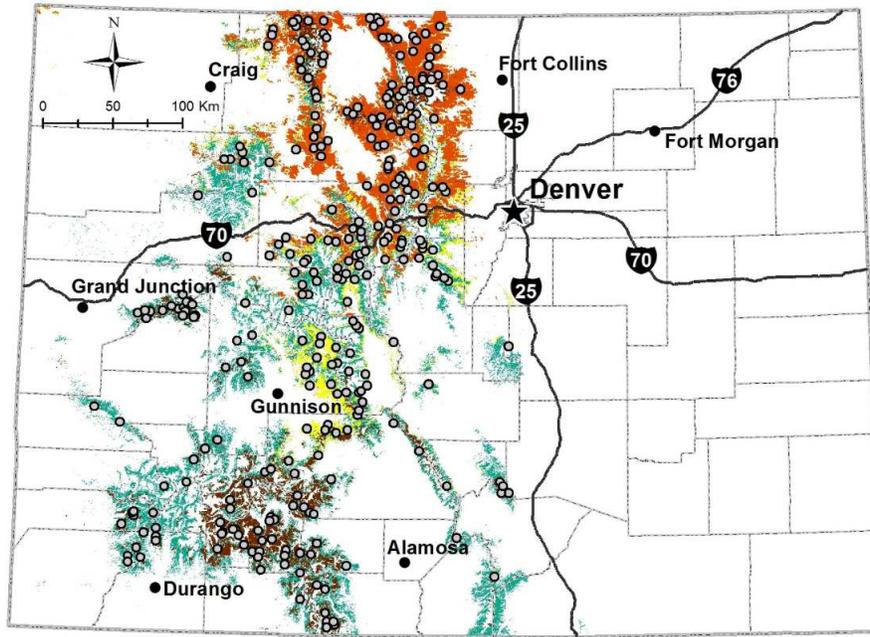


Figure 1. Randomly selected sampling sites (gray circles) where passive infrared game cameras were deployed in spruce-fir (green) and lodgepole pine (yellow) forests in Colorado, USA, 2013–2014. Brown and orange are the approximate extents of spruce beetle and mountain pine beetle impacts in spruce-fir and lodgepole pine forests, respectively, as of 2014. Reprinted from "Mammalian responses to changed forest conditions resulting from bark beetle outbreaks in the southern Rocky Mountains," by J. S. Ivan, 2018, *Ecosphere*, 9(8), e02369. Copyright [Year] by the Ecological Society of America. Reprinted with permission.

215 $n = 300$, 1 km² sites was evenly split across the two dominant forest types, spruce-fir and lodgepole pine. Additional habitat
 216 covariates were collected at each site, and a description of these is included in Appendix B.

217 Passive infrared camera traps (Reconyx PC800, Holmen, Wisconsin, USA) were deployed near the center of each site.
 218 Cameras were approximately 0.5 m above the ground and pointed toward a lure tree 4–5 m away⁴⁶. The lure tree was baited
 219 with minimal amounts of peanut butter and commercial rabbit lure. This setup was designed to maximize detections of both
 220 large and small-bodied mammals in the local community while minimizing attraction of individuals from outside the sampling
 221 region of the site. The sampling regions were likely not closed to immigration/emigration; thus, we interpret elevated detections
 222 at a site as more individuals using, as opposed to occupying, that site⁴⁷. For additional details regarding the sampling design
 223 and study area see Ivan et al.⁴¹.

224 Model Fitting

225 We fit the Royle-Nichols joint species model to the multispecies camera trap data binned into 20 two-day occasions because
 226 simulations showed this was the number of replications needed to identify a quadratic effect of occasion on individual detection
 227 probability. Not all cameras were operational for the entire 40 day sampling period, and thus the number of occasions varied
 228 from 7-20. We discarded four sites at which the camera was operational for less than one occasion. We also discarded another
 229 12 sites that had been infested by bark beetles for more than 10 years. Ivan et al.⁴¹ truncated the bark beetle infestation covariate

230 at 10 years because estimates of response curves beyond 10 years would be unreliable with so few sites. The final sample size
231 was $n = 284$ sites.

232 For consistency, our joint occupancy model included the 13 species for which Ivan et al.⁴¹ performed a single species
233 analysis; several rare species were excluded from analysis due to insufficient detections. We note, however, that these rare
234 species parameters may be identifiable in the joint model as has been the case in previous studies^{2,48–52}. Ivan et al.⁴¹ used
235 a sequential procedure similar to that described in Lebreton et al.⁵³ to select the covariates in the occupancy and detection
236 processes for each species. We adopted their detection model and used the same covariates but a different set of basis functions
237 for YSO. Ivan et al.⁴¹ treated YSO as a grouping variable and considered probability of use response curves that allowed for
238 cubic associations and delayed responses to bark beetle infestation. Multiple response curves were model averaged to produce
239 predictive YSO response curves for each species. We used orthogonal polynomial basis functions for the YSO variable in the
240 species intensity models. The basis functions included a linear (YSO1) and quadratic (YSO2) effect. Appendix E provides a
241 full description of the intensity and detection models.

242 We fit the model using a MCMC algorithm. To improve mixing and predictive ability, we regularized the coefficients β and α
243 with slightly informative priors: $\beta \sim \mathcal{N}(\mathbf{0}, 10\mathbf{I})$ and $\alpha \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ⁵⁴. We specified a vague prior of $\Sigma^{-1} \sim \text{Wishart}(15, (15\mathbf{I})^{-1})$
244 for the species variance-covariance matrix⁵⁵. We used Gibbs sampling based on conjugate priors for parameters Σ_{spp} and β
245 and Metropolis-Hastings updates for \mathbf{N} , λ , and α . Derivations of the conjugate full-conditional distributions are provided in
246 Appendix C with details about the Metropolis-Hastings updates. We tuned the Metropolis-Hastings updates so that acceptance
247 rates varied between 20-40% for α , \mathbf{N} , and λ . All continuous covariates were scaled to have mean 0 and variance 1. We ran the
248 MCMC algorithm for $L = 20000$ iterations, and the first 5000 iterations were discarded as burn-in. We fit both the unrestricted
249 and restricted models expressed in (9) and (12), respectively. The *Results* section presents inference for regression coefficients
250 from the unrestricted joint species distribution model and contrasts it with that from Ivan et al.⁴¹ single species distribution
251 models to highlight inferential discrepancies caused by community confounding. Inference from the restricted joint species
252 distribution model is omitted because it was similar in sign and significance to Ivan et al.⁴¹ single species distribution models
253 but differed in effect size because habitat covariates model intensity rather than occupancy.

254 **Code Availability**

255 All algorithms and code for fitting the unrestricted and restricted joint species distribution models are available in the
256 Supplementary Information files. All MCMC algorithms and analyses were coded in R 4.0.3.

257 **Results**

258 Ivan et al.⁴¹ fit single species distribution models to infer changes in mammalian use of stands impacted by the bark beetle
259 epidemic. The impact of bark beetle damage was measured by years since initial infestation (YSO) and severity of outbreak
260 quantified by mean overstory mortality (DeafConif). Moose, elk, and mule deer exhibited positive associations with bark beetle
261 activity; red squirrels, golden-mantled ground squirrels, chipmunks, and coyotes exhibited negative associations; American

262 martens, black bears, snowshoe hares, and porcupines showed no associations; and red foxes and yellow-bellied marmots
 263 showed mixed associations. Inference from the joint model, Figures 2-4, was largely consistent with that of the single species
 264 models. All three ungulates species had a positive association with severity but the YSO effects Ivan et al.⁴¹ noted for mule deer
 265 and moose were not observed (Figure 2). Both red and golden-mantled ground squirrels had negative associations with severity.
 266 Furthermore, use of subalpine stands decreased as a function YSO for red squirrels and coyotes. We did not, however, observe
 267 any significant associations between YSO and stand use for golden-mantled ground squirrels and chipmunks, contrary to the
 268 inference from the single species models. Consistent with inference from the single species models, red foxes had a negative
 269 association with severity but positive association with YSO1, whereas yellow-bellied marmots had a positive association with
 270 severity but negative association with YSO1.

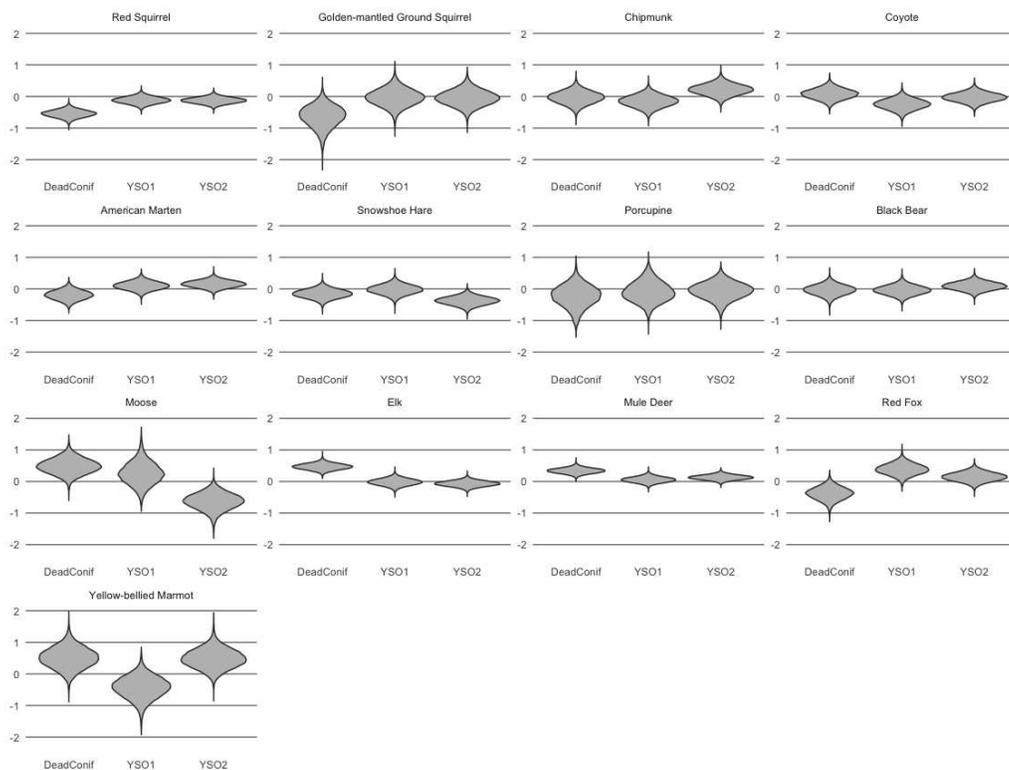


Figure 2. Violin plot for marginal posterior histograms of infestation regression parameters. Estimates are from the unrestricted joint species distribution model. DeadConif is the overstory mortality percentage, a proxy for severity of bark beetle infestation. YSO1 is the linear effect of the number of years since a site was infested with bark beetles. YSO2 is the quadratic effect.

271 The inferential discrepancies between the joint and single species models are indicative of community confounding.
 272 Covariates that were important in the single species model no longer indicate an effect when accounting for the community. For
 273 example, the posterior distribution of the effect pertaining to the indicator of whether the site was in a federally designated
 274 wilderness (WILD) overlapped zero in the joint species distribution model for Yellow-bellied Marmots (Figure 3). Ivan et al.⁴¹
 275 single species distribution model, however, identified the covariate as helpful quantity for predicting the distribution of the
 276 Yellow-bellied Marmot. The covariate could have been important in the single species model as a surrogate measure of other
 277 species that the American marten avoids or pursues. After we accounted for the interspecies dependence in the joint model
 278 (Figure 4), the WILD variable was no longer helpful as a predictor, although the contribution to that signal on species intensity
 279 was still accounted for.

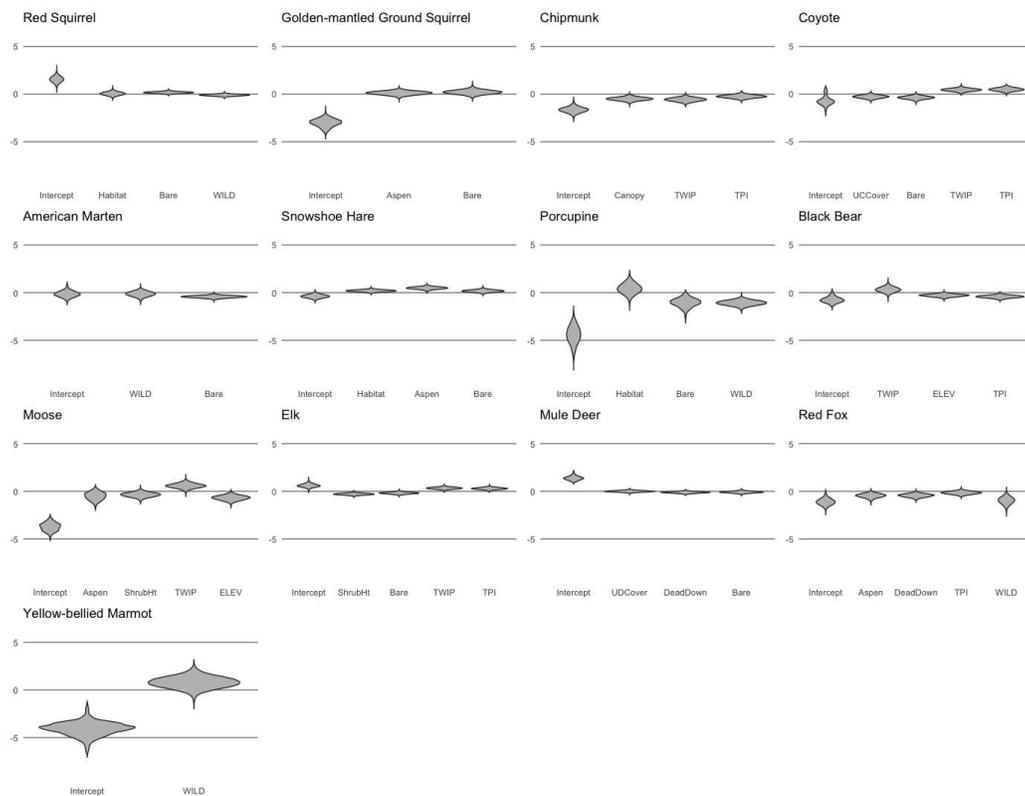


Figure 3. Violin plot for marginal posterior histograms of base habitat regression parameters. Estimates are from the unrestricted joint species distribution model. Appendix B provides a description of habitat covariates.

280 The restricted regression approach we described in the *Community Confounding* section estimates the unconditional habitat
 281 effects, which would be similar to those obtained from a single species model. Irrespective of whether inference on conditional
 282 or unconditional effects is desired, it may be helpful to compute measures of community confounding to understand which fixed
 283 effects will vary between an unrestricted and restricted model. For example, severity (DeadConif) was the most confounded
 284 covariate for yellow-bellied marmots (Table 1). Inference on the effect of this covariate varied across the two approaches;
 285 the posterior distribution of the regression coefficient did not overlap zero in the unrestricted model 0.53 (0.03, 1.03) but did
 286 contain zero for the restricted model 0.33 (-0.14, 0.78) (90% credible intervals).

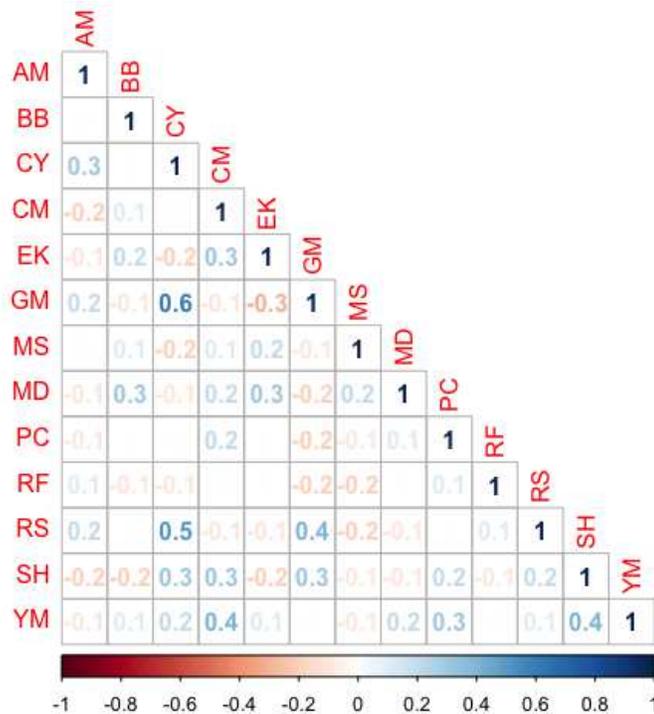


Figure 4. Posterior mean of species correlation matrix. Estimates are from the unrestricted joint species distribution model. AM = American Marten, BB = Black Bear, CY = Coyote, CM = Chipmunk spp., Ek = Elk, GM = Golden-mantled Ground Squirrel, MS = Moose, MD = Mule Deer, PC = Porcupine, RF = Red Fox, RS = Red Squirrel, SH = Snowshoe Hare, YM = Yellow-bellied Marmot.

Table 1. Posterior means of coefficients of determination for each species. Only maximum of posterior means shown for each species.

Species	Covariate	$E(R^2 \mathbf{Y})$
Black Bear	TWIP	0.62
American Marten	DeadConif	0.54
Coyote	DeadConif	0.53
Snowshoe Hare	DeadConif	0.53
Red Fox	DeadConif	0.52
Elk	DeadConif	0.51
Moose	DeadConif	0.51
Porcupine	DeadConif	0.51
Yellow-bellied Marmot	DeadConif	0.51
Golden-mantled Ground Squirrel	DeadConif	0.49
Red Squirrel	Bare	0.47
Chipmunk	DeadConif	0.46
Mule Deer	YSO1	0.44

Discussion

Species distributions are shaped by habitat as well as competition and mutualism within the community, and these effects are likely confounded. Because habitat and interspecies effects operate simultaneously and can be equally influential, restricted regression that gives priority to the habitat effects may not be appropriate always. Alternative approaches for adjusting the estimates of confounded effects have their own caveats. Consider a simpler case where the latent intensities, $\boldsymbol{\lambda}$, of the K species in our community were known. We could construct K regression models for predicting each species intensity as follows:

$$\boldsymbol{\lambda}_k = \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\Lambda}_{-k} \boldsymbol{\eta}_k + \boldsymbol{\varepsilon}, \quad (15)$$

where $\boldsymbol{\Lambda}_{-k} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k-1}, \boldsymbol{\lambda}_{k+1}, \dots, \boldsymbol{\lambda}_K)$ is a matrix of the $K - 1$ other species intensities. If \mathbf{X}_k and $\boldsymbol{\Lambda}_{-k}$ were highly collinear, principal component analysis (PCA) might be applied. PCA decomposes the variation explained by \mathbf{X}_k and $\boldsymbol{\Lambda}_{-k}$ into $p = p_1 + p_2$ principal components, $\boldsymbol{\Gamma}_k = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p)$, where p_1 and p_2 are the number of columns of \mathbf{X}_k and $\boldsymbol{\Lambda}_{-k}$ respectively.

The regression model

$$\boldsymbol{\lambda}_k = \mathbf{W}_k \boldsymbol{\theta} + \boldsymbol{\varepsilon}, \quad (16)$$

$$\mathbf{W}_k = (\mathbf{X}_k, \boldsymbol{\Lambda}_{-k}) \boldsymbol{\Gamma}_k, \quad (17)$$

improves stability and can recover the posterior means and variances of $\boldsymbol{\beta}_k$ and $\boldsymbol{\eta}_k$ in (15). In practice, inference on $\boldsymbol{\beta}_k$ and $\boldsymbol{\eta}_k$ is often adjusted by truncating off the last $p - r$, for $r < p$, eigenvectors of $\boldsymbol{\Gamma}_k$ and employing the new design matrix

$$\mathbf{W}_k^* = (\mathbf{X}_k, \boldsymbol{\Lambda}_{-k}) \boldsymbol{\Gamma}_k^*, \quad (18)$$

$$\boldsymbol{\Gamma}_k^* = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_r). \quad (19)$$

299 By retaining only the first r principal components, the smallest sources of variation are ignored in the estimation of β_k and η_k .
 300 Jeffers⁵⁶ implemented this approach truncating off the last 7 of 13 principal components to adjust the estimates of regression
 301 coefficients relating various tree characteristics to maximum compressive strength. Other studies have selected a subset of
 302 principal components based on their strength of association with the response variable^{57–60}. In some cases, the coefficient
 303 estimates from these reduced rank approaches appeared more tenable than those from the full rank specifications based on
 304 known physical relationships between the predictors and response. Thus, PCA regression can offer both computational and
 305 inferential improvements. PCA regression, however, does not explicitly address confounding, and interpretation of the adjusted
 306 regression coefficients is unclear.

307 Another dubious approach for adjusting parameter estimates in a confounded model is model averaging^{61,62}. Consider
 308 an ensemble of S models with posteriors $[\theta|\mathbf{y}, \mathcal{M}_1], \dots, [\theta|\mathbf{y}, \mathcal{M}_S]$ and a particular parameter θ_l present in all S models. The
 309 marginal posterior distribution of θ_l across all models is given by

$$[\theta_l|\mathbf{y}] = \sum_{s=1}^S [\theta_l|\mathbf{y}, \mathcal{M}_s] P(\mathcal{M}_s|\mathbf{y}), \quad (20)$$

310 where $P(\mathcal{M}_s|\mathbf{y})$ is the posterior model probability of model s . One pitfall of model averaging is that the contribution of the
 311 posterior distribution of θ_l in model \mathcal{M}_s is weighted by the posterior model probability $P(\mathcal{M}_s|\mathbf{y})$, which is not necessarily
 312 indicative of the parameter's importance in the model⁶³. Furthermore, the interpretation of the parameter θ_l often depends
 313 on which of the other θ_{-l} parameters are included in the model⁶⁴. Some models in the ensemble may include parameters
 314 confounded with θ_l while others do not. Averaging θ_l across the ensemble makes no explicit adjustments for confounding and
 315 complicates inference^{63–65}. Restricted regression is appealing in that, unlike PCA regression and model averaging, confounding
 316 is addressed by attributing all contested sources of variation to the fixed effects.

317 Recently, concerns regarding the coverage properties of the fixed effects estimator under restricted regression have been
 318 expressed^{31,66}. For example, Zimmerman and Ver Hoef⁶⁶ showed that applying any restricted regression method to a SGLMM
 319 leads to frequentest coverage of the fixed effects that is lower than the corresponding non-spatial model. Similarly, Khan and
 320 Calder³¹ found that when fitting a restricted version of the SGLMM with an intrinsic conditional autoregressive prior, credible
 321 intervals of the fixed effects from the restricted model were generally nested inside those yielded by the non-spatial model.
 322 Given these results, both Zimmerman and Ver Hoef⁶⁶ and Khan and Calder³¹ recommended reverting to inference from the
 323 non-spatial model, rather than that of the restricted SGLMM, when inference from the unrestricted SGLMM appears untenable.
 324 This practice may not be appealing for multispecies models in which removing the species random effects would prevent
 325 inference on interspecies dependence. Because the random effect η is rarely of interest in spatial applications, there has been
 326 little investigation on the inferential impacts of restricted regression on η . Such investigation, however, may be helpful in
 327 determining the appropriateness of restricted regression for multispecies models.

328 In summary, we specified a joint species distribution model that accounts for interspecies dependence at the intensity level.
 329 We examined how inference from the joint model differed from that of the single species models in Ivan et al.⁴¹. Confounding in

330 multispecies models is unique from spatial and time series applications in that the random effect is almost always of inferential
331 interest, and hence, adjustments to the regression coefficients, β , and random effects, η , should both be considered. Using
332 unconditional effects may not be appropriate in all settings, but alternative methods for adjusting parameter estimates, like PCA
333 regression and model averaging, are also not universally applicable.

334 There may not exist a general remedial method for handling confounding in multispecies models, but confounding measures
335 should be investigated because they can provide new insights into ecological systems. For example, we discovered that the
336 severity of bark beetle outbreak was heavily confounded with the species random effects (Table 1). This suggests that effects
337 of bark beetle outbreaks are complex with components potentially related to both the mammalian community and habitat.
338 Consequently, restoration of the subalpine forest may not coincide with recovery of the mammalian community. Changes in
339 mammalian use of subalpine stands impacted by bark beetles may be temporary reactions spurred by changes to the intensities
340 of competitors, predators, and mutualists. Thus, management practices aimed at restoring previous habitat conditions might not
341 immediately trigger recovery of species negatively impacted by the outbreak. Understanding the implications of any particular
342 management strategy rely on our ability to disentangle habitat and community effects.

343 Experimental methods and modeling techniques for alleviating confounding have been proposed in ecology. Hefley
344 et al.⁶⁷ showed that replicate populations can help disentangle confounded fixed and random effects. In the context of
345 multispecies models, replication involves analyzing several communities simultaneously, which is often infeasible. Hefley
346 et al.⁶⁷ also recommended explicit population models rather than phenomenological regression-based models for analysis of
347 temporally confounded count data. Similarly, Fieberg et al.⁶⁵ advocated for mechanistic models guided by causal diagrams
348 for analyzing temporally confounded animal movement data. An avenue of future research for multispecies modeling is to
349 compare inference from phenomenological regression-based models, such as the one proposed here, with that of models that
350 explicitly include ecological mechanisms such as competitive exclusion, mutualism, and predation. Because community and
351 temporal confounding have the same mathematical framework, mechanistic models are a promising solution for confounded
352 multispecies data.

353 References

- 354 1. Altwegg, R. & Nichols, J. D. Occupancy models for citizen-science data. *Methods Ecol. Evol.* **10**, 8–21 (2019).
- 355 2. Hui, F., Warton, D., Foster, S. & Dunstan, P. To mix or not to mix: Comparing the predictive performance of mixture
356 models vs. separate species distribution models. *Ecology* **94**, 1913–9, DOI: [10.1890/12-1322.1](https://doi.org/10.1890/12-1322.1) (2013).
- 357 3. Warton, D. et al. So many variables: Joint modeling in community ecology. *Trends Ecol. & Evol.* **30**, 766–779, DOI:
358 [10.1016/j.tree.2015.09.007](https://doi.org/10.1016/j.tree.2015.09.007) (2015).
- 359 4. Tobler, M. W. et al. Joint species distribution models with species correlations and imperfect detection. *Ecology* **100**, e02754,
360 DOI: <https://doi.org/10.1002/ecy.2754> (2019). <https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1002/ecy.2754>.

- 361 **5.** Estevo, C. A., Nagy-Reis, M. B. & Nichols, J. D. When habitat matters: Habitat preferences can modulate co-occurrence
362 patterns of similar sympatric species. *PLoS one* **12**, e0179489 (2017).
- 363 **6.** Steen, D. A. *et al.* Snake co-occurrence patterns are best explained by habitat and hypothesized effects of interspecific
364 interactions. *J. Animal Ecol.* **83**, 286–295 (2014).
- 365 **7.** Wisz, M. S. *et al.* The role of biotic interactions in shaping distributions and realised assemblages of species: Implications
366 for species distribution modelling. *Biol. Rev.* **88**, 15–30 (2013).
- 367 **8.** Royle, J. & Nichols, J. Estimating abundance from repeated presence-absence data or point counts. *Ecology* **84**, 777–790
368 (2003).
- 369 **9.** MacKenzie, D. I. *et al.* Estimating site occupancy rates when detection probabilities are less than one. *Ecology* **83**,
370 2248–2255 (2002).
- 371 **10.** Orrock, J. L., Pagels, J. F., McShea, W. J. & Harper, E. K. Predicting presence and abundance of a small mammal species:
372 The effect of scale and resolution. *Ecol. Appl.* **10**, 1356–1366 (2000).
- 373 **11.** Cingolani, A. M., Cabido, M., Gurvich, D. E., Renison, D. & Díaz, S. Filtering processes in the assembly of plant
374 communities: Are species presence and abundance driven by the same traits? *J. Veg. Sci.* **18**, 911–920 (2007).
- 375 **12.** Dibner, R. R., Doak, D. F. & Murphy, M. Discrepancies in occupancy and abundance approaches to identifying and
376 protecting habitat for an at-risk species. *Ecol. Evol.* **7**, 5692–5702, DOI: <https://doi.org/10.1002/ece3.3131> (2017).
377 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ece3.3131>.
- 378 **13.** Bascompte, J. Mutualistic networks. *Front. Ecol. Environ.* **7**, 429–436 (2009).
- 379 **14.** Van Dam, N. How plants cope with biotic interactions. *Plant Biol.* **11**, 1–5 (2009).
- 380 **15.** Tobler, M. W., Zúñiga Hartley, A., Carrillo-Percastegui, S. E. & Powell, G. V. N. Spatiotemporal hierarchical modelling
381 of species richness and occupancy using camera trap data. *J. Appl. Ecol.* **52**, 413–421, DOI: [10.1111/1365-2664.12399](https://doi.org/10.1111/1365-2664.12399)
382 (2015). <https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/1365-2664.12399>.
- 383 **16.** Bross, K. M., Hooten, M. B. & Fitzpatrick, R. M. Model selection and assessment for multi-species occupancy models.
384 *Ecology* **97**, 1759–1770, DOI: [10.1890/15-1471.1](https://doi.org/10.1890/15-1471.1) (2016). [https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/](https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/15-1471.1)
385 [15-1471.1](https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/15-1471.1).
- 386 **17.** Rota, C. T. *et al.* A multispecies occupancy model for two or more interacting species. *Methods Ecol. Evol.* **7**, 1164–1173,
387 DOI: [10.1111/2041-210X.12587](https://doi.org/10.1111/2041-210X.12587) (2016). <https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.12587>.
- 388 **18.** Maphisa, D. H., Smit-Robinson, H. & Altwegg, R. Dynamic multi-species occupancy models reveal individualistic habitat
389 preferences in a high-altitude grassland bird community. *PeerJ* **7**, e6276 (2019).
- 390 **19.** Clark, A. E. & Altwegg, R. Efficient Bayesian analysis of occupancy models with logit link functions. *Ecol. Evol.* **9**,
391 756–768 (2019).

- 392 **20.** Broms, K. M., Johnson, D. S., Altwegg, R. & Conquest, L. L. Spatial occupancy models applied to atlas data show
393 Southern Ground Hornbills strongly depend on protected areas. *Ecol. Appl.* **24**, 363–374 (2014).
- 394 **21.** Johnson, D. S., Conn, P. B., Hooten, M. B., Ray, J. C. & Pond, B. A. Spatial occupancy models for large data sets. *Ecology*
395 **94**, 801–808 (2013).
- 396 **22.** Reich, B. J., Hodges, J. S. & Zadnik, V. Effects of residual smoothing on the posterior of the fixed effects in disease-mapping
397 models. *Biometrics* **62**, 1197–1206 (2006).
- 398 **23.** Hodges, J. S. & Reich, B. J. Adding spatially-correlated errors can mess up the fixed effect you love. *The Am. Stat.* **64**,
399 325–334, DOI: [10.1198/tast.2010.10052](https://doi.org/10.1198/tast.2010.10052) (2010). <https://doi.org/10.1198/tast.2010.10052>.
- 400 **24.** Hughes, J. & Haran, M. Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *J.*
401 *Royal Stat. Soc. Ser. B (Statistical Methodol.* **75**, 139–159, DOI: <https://doi.org/10.1111/j.1467-9868.2012.01041.x> (2013).
- 402 **25.** Hanks, E. M., Schliep, E. M., Hooten, M. B. & Hoeting, J. A. Restricted spatial regression in practice: geostatistical
403 models, confounding, and robustness under model misspecification. *Environmetrics* **26**, 243–254, DOI: [10.1002/env.2331](https://doi.org/10.1002/env.2331)
404 (2015). <https://onlinelibrary.wiley.com/doi/pdf/10.1002/env.2331>.
- 405 **26.** Hefley, T. J., Hooten, M. B., Hanks, E. M., Russell, R. E. & Walsh, D. P. The bayesian group lasso for confounded spatial
406 data. *J. Agric. Biol. Environ. Stat.* **22**, 42–59 (2017).
- 407 **27.** Bradley, J. R., Holan, S. H., Wikle, C. K. *et al.* Multivariate spatio-temporal models for high-dimensional areal data with
408 application to longitudinal employer-household dynamics. *Annals Appl. Stat.* **9**, 1761–1791 (2015).
- 409 **28.** Murakami, D. & Griffith, D. A. Random effects specifications in eigenvector spatial filtering: A simulation study. *J. Geogr.*
410 *Syst.* **17**, 311–331 (2015).
- 411 **29.** Thaden, H. & Kneib, T. Structural equation models for dealing with spatial confounding. *The Am. Stat.* **72**, 239–252
412 (2018).
- 413 **30.** Prates, M. O., Assunção, R. M., Rodrigues, E. C. *et al.* Alleviating spatial confounding for areal data problems by
414 displacing the geographical centroids. *Bayesian Analysis* **14**, 623–647 (2019).
- 415 **31.** Khan, K. & Calder, C. A. Restricted spatial regression methods: Implications for inference. *J. Am. Stat. Assoc.* 1–13
416 (2020).
- 417 **32.** Paciorek, C. The importance of scale for spatial-confounding bias and precision of spatial regression estimators. *Stat. Sci.*
418 *A review journal Inst. Math. Stat.* **25**, 107–125, DOI: [10.1214/10-STS326](https://doi.org/10.1214/10-STS326) (2010).
- 419 **33.** Dominici, F., McDermott, A. & Hastie, T. J. Improved semiparametric time series models of air pollution and mortality. *J.*
420 *Am. Stat. Assoc.* **99**, 938–948 (2004).
- 421 **34.** Houseman, E. A., Coull, B. A. & Shine, J. P. A nonstationary negative binomial time series with time-dependent covariates:
422 Enterococcus counts in Boston Harbor. *J. Am. Stat. Assoc.* **101**, 1365–1376 (2006).

- 423 **35.** Corbeil, R. R. & Searle, S. R. Restricted maximum likelihood (reml) estimation of variance components in the mixed
424 model. *Technometrics* **18**, 31–38 (1976).
- 425 **36.** Hoeting, J. A., Leecaster, M. & Bowden, D. An improved model for spatially correlated binary responses. *J. Agric. Biol.*
426 *Environ. Stat.* **5**, 102–114 (2000).
- 427 **37.** Tyre, A. J. *et al.* Improving precision and reducing bias in biological surveys: Estimating false-negative error rates. *Ecol.*
428 *Appl.* **13**, 1790–1801, DOI: <https://doi.org/10.1890/02-5078> (2003). [https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.](https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/02-5078)
429 [1890/02-5078](https://doi.org/10.1890/02-5078).
- 430 **38.** Clark, J. S. *et al.* High-dimensional coexistence based on individual variation: a synthesis of evidence. *Ecol. Monogr.* **80**,
431 569–608 (2010).
- 432 **39.** Scheffe, H. *The Analysis of Variance*, vol. 72 (John Wiley & Sons, 1959).
- 433 **40.** Hodges, J. S. & Clayton, M. K. Random effects old and new. *Stat. Sci.* (2011).
- 434 **41.** Ivan, J., Seglund, A., Truex, R. & Newkirk, E. Mammalian responses to changed forest conditions resulting from bark
435 beetle outbreaks in the southern Rocky Mountains. *Ecosphere* **9**, DOI: [10.1002/ecs2.2369](https://doi.org/10.1002/ecs2.2369) (2018).
- 436 **42.** Chan, J. C.-C. & Jeliazkov, I. MCMC estimation of restricted covariance matrices. *J. Comput. Graph. Stat.* **18**, 457–480
437 (2009).
- 438 **43.** Rue, H. & Held, L. *Gaussian Markov Random Fields: Theory and Applications* (CRC press, 2005).
- 439 **44.** Hanson, T. E., Branscum, A. J., Johnson, W. O. *et al.* Informative *g*-priors for logistic regression. *Bayesian Analysis* **9**,
440 597–612 (2014).
- 441 **45.** Baddeley, A. *et al.* Spatial logistic regression and change-of-support in poisson point processes. *Electron. J. Stat.* **4**,
442 1151–1201, DOI: [10.1214/10-EJS581](https://doi.org/10.1214/10-EJS581) (2010).
- 443 **46.** Blecha, K. A. Risk-reward tradeoffs in the foraging strategy of cougar (puma concolor): prey distribution, anthropogenic
444 development, and patch selection. *Thesis. Colo. State Univ. Fort Collins, Color. USA* (2015).
- 445 **47.** MacKenzie, D. I. *et al.* *Occupancy Estimation and Modeling: Inferring Patterns and Dynamics of Species Occurrence*
446 (Academic Press, 2006).
- 447 **48.** Guisan, A., Weiss, S. & Weiss, A. GLM versus CCA spatial modeling of plant species distribution. *Plant Ecol.* **143**,
448 107–122, DOI: [10.1023/A:1009841519580](https://doi.org/10.1023/A:1009841519580) (1999).
- 449 **49.** Ovaskainen, O. & Soininen, J. Making more out of sparse data: hierarchical modeling of species communities. *Ecology*
450 **92**, 289–295 (2011).
- 451 **50.** Madon, B., Warton, D. I. & Araújo, M. B. Community-level vs species-specific approaches to model selection. *Ecography*
452 **36**, 1291–1298, DOI: [10.1111/j.1600-0587.2013.00127.x](https://doi.org/10.1111/j.1600-0587.2013.00127.x) (2013).

- 453 **51.** Ovaskainen, O., Abrego, N., Halme, P. & Dunson, D. Using latent variable models to identify large networks of species-
454 to-species associations at different spatial scales. *Methods Ecol. Evol.* **7**, 549–555, DOI: [10.1111/2041-210X.12501](https://doi.org/10.1111/2041-210X.12501)
455 (2015).
- 456 **52.** Tikhonov, G., Abrego, N., Dunson, D. & Ovaskainen, O. Using joint species distribution models for evaluating how
457 species-to-species associations depend on the environmental context. *Methods Ecol. Evol.* **8**, 443–452, DOI: [10.1111/
458 2041-210X.12723](https://doi.org/10.1111/2041-210X.12723) (2017).
- 459 **53.** Lebreton, J.-D., Burnham, K. P., Clobert, J. & Anderson, D. R. Modeling survival and testing biological hypotheses using
460 marked animals: A unified approach with case studies. *Ecol. Monogr.* **62**, 67–118, DOI: [10.2307/2937171](https://doi.org/10.2307/2937171) (1992).
- 461 **54.** Hooten, M. B. & Hobbs, N. T. A guide to Bayesian model selection for ecologists. *Ecol. Monogr.* **85**, 3–28, DOI:
462 <https://doi.org/10.1890/14-0661.1> (2015). <https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/14-0661.1>.
- 463 **55.** Chung, Y., Gelman, A., Rabe-Hesketh, S., Liu, J. & Dorie, V. Weakly informative prior for point estimation of covariance
464 matrices in hierarchical models. *J. Educ. Behav. Stat.* **40**, 136–157, DOI: [10.3102/1076998615570945](https://doi.org/10.3102/1076998615570945) (2015). [https:
465 //doi.org/10.3102/1076998615570945](https://doi.org/10.3102/1076998615570945).
- 466 **56.** Jeffers, J. N. Two case studies in the application of principal component analysis. *J. Royal Stat. Soc. Ser. C (Applied Stat.*
467 **16**, 225–236 (1967).
- 468 **57.** Hill, C. R., Fomby, T. B. & Johnson, S. R. Component selection norms for principal components regression. *Commun.*
469 *Stat. Methods* **6**, 309–334 (1977).
- 470 **58.** Kung, E. C. & Sharif, T. A. Regression forecasting of the onset of the indian summer monsoon with antecedent upper air
471 conditions. *J. Appl. Meteorol. Climatol.* **19**, 370–380 (1980).
- 472 **59.** Smith, G. & Campbell, F. A critique of some ridge regression methods. *J. Am. Stat. Assoc.* **75**, 74–81 (1980).
- 473 **60.** Jolliffe, I. T. A note on the use of principal components in regression. *J. Royal Stat. Soc. Ser. C (Applied Stat.* **31**, 300–303
474 (1982).
- 475 **61.** Roberts, H. V. Probabilistic prediction. *J. Am. Stat. Assoc.* **60**, 50–62, DOI: [10.1080/01621459.1965.10480774](https://doi.org/10.1080/01621459.1965.10480774) (1965).
- 476 **62.** Leamer, E. E. *Specification Searches: Ad Hoc Inference with Nonexperimental Data*, vol. 53 (Wiley New York, 1978).
- 477 **63.** Cade, B. S. Model averaging and muddled multimodel inferences. *Ecology* **96**, 2370–2382 (2015).
- 478 **64.** Banner, K. M. & Higgs, M. D. Considerations for assessing model averaging of regression coefficients. *Ecol. Appl.* **27**,
479 78–93 (2017).
- 480 **65.** Fieberg, J., Ditmer, M. & Freckleton, R. Understanding the causes and consequences of animal movement: A cautionary
481 note on fitting and interpreting regression models with time-dependent covariates. *Methods Ecol. Evol.* **3**, DOI: [10.1111/j.
482 2041-210X.2012.00239.x](https://doi.org/10.1111/j.2041-210X.2012.00239.x) (2012).
- 483 **66.** Zimmerman, D. L. & Ver Hoef, J. M. On deconfounding spatial confounding in linear models. *The Am. Stat.* 1–9 (2021).

484 **67.** Hefley, T. J., Hooten, M. B., Drake, J. M., Russell, R. E. & Walsh, D. P. When can the cause of a population decline be
485 determined? *Ecol. Lett.* **19**, 1353–1362 (2016).

486 **Acknowledgements**

487 This research was funded by Colorado Parks and Wildlife.

488 **Appendix**

489 All appendices referenced in the manuscript are in the Supplementary Information files.

490 **Author Contributions**

491 J.V. and M.H. wrote the manuscript and designed the model. J.I. acquired the data. All authors reviewed the manuscript.

492 **Data Availability**

493 The data are available in the Supplementary Information files.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Appendix.pdf](#)