

1 **A comparison of predictive performance of joint species distribution models for**
2 **presence-absence data**

3 David P. Wilkinson^{1*}, Nick Golding^{2,3,4}, Gurutzeta Guillera-Arroita^{1,5}, Reid Tingley^{6,7},
4 Michael A. McCarthy¹

5

6 1. School of Agriculture, Food and Ecosystem Sciences, University of Melbourne, Parkville,
7 3010, Victoria, Australia

8 2. Telethon Kids Institute, Perth Children's Hospital, 15 Hospital Ave, Nedlands, 6009,
9 Western Australia, Australia

10 3. Curtin University, Kent St, Bentley, 6102, Western Australia, Australia

11 4. Melbourne School of Population and Global Health, University of Melbourne, Parkville,
12 3010, Victoria, Australia

13 5. Pyrenean Institute of Ecology, Spanish National Research Council, Jaca, 22700, Huesca,
14 Spain

15 6. EnviroDNA Pty Ltd, 95 Albert Street, Brunswick, Victoria 3056 Australia

16 7. School of Biological Sciences, Monash University, Clayton, 3800, Victoria, Australia

17 *Corresponding author: dwilkinson@unimelb.edu.au

18 ORCiDs:

19 D.P.W. 0000-0002-9560-6499

20 N.G. 0000-0001-8916-5570

21 G.G.A. 0000-0002-8387-5739

22 R.T. 0000-0002-7630-7434

23 M.A.M. 0000-0003-1039-7980

24 **Acknowledgements:** We thank Peter Vesk, Brendan Wintle and Jian Yen for insightful
25 discussions. N.G. was supported by an Australian Research Council (ARC) Discovery Early
26 Career Researcher Award (DE180100635). G.G.A. is currently supported by a ‘Ramón y
27 Cajal’ grant (RYC2020-028826-I) funded by the Spanish Ministry of Science and
28 Innovation, the Agencia Estatal de Investigación (10.13039/501100011033) and “ESF
29 Investing in your future”.

30 **Data availability statement:** The code and data for this analysis can be found on GitHub and
31 is archived by Zenodo (Doi90/JSDM_Prediction; Wilkinson, 2019).

32 **Conflicts of interest:** No conflict of interest is declared.

33 **Author contribution statement:** All the authors conceived the ideas and methodology;
34 D.P.W. implemented the analysis; D.P.W. led writing the manuscript but all the authors
35 contributed significantly throughout and gave final approval before submission.

36 **ABSTRACT**

- 37 1. While there has been substantial literature on the evaluation of predictions from single
38 species distribution models, the topic of prediction has only recently begun to be
39 addressed for joint species distribution models (JSDMs). These studies have covered
40 only limited aspects of prediction: limited selection of models being compared, limited
41 number of evaluation metrics, and/or not comparing the different prediction types
42 available to JSDMs.
- 43 2. In this study, we perform a large-scale comparison of the predictive performance of
44 eight model types: two stacked species distribution models (SSDMs) and six JSDMs.
45 We fit these models to 22 real and simulated datasets, make four types of JSDM
46 predictions, and evaluate up to 32 metrics from five different classes that quantify
47 different aspects of performance of predictions about species distributions and the
48 community assemblage process.
- 49 3. We found that likelihood-based metrics indicated the JSDMs were better fit to the data
50 than the standard SSDM, but most other metric classes showed the SSDM
51 outperforming the JSDMs by generally small amounts. The spatial and non-spatial
52 implementations of the hierarchical multivariate probit regression model with latent
53 factors typically performed better than the other JSDMs, but overall still performed
54 worse than the SSDM. The SSDM predictions constrained with the spatially-explicit
55 species assemblage modelling framework (SESAM) consistently outperformed both
56 the standard SSDM and all JSDMs for both species- and community-level metrics.
- 57 4. Our results indicate that despite the additional inference they provide about the
58 community assemblage process by accounting for the residual association between
59 species, JSDMs generally yield worse predictions than stacked single species models
60 when evaluated at either the species or community level. The performance of the

61 SESAM framework suggests that exploring similar approaches to constrain JSDM
62 predictions is an interesting future avenue of research.

63 **Keywords:** biotic interactions, community assemblage, evaluation metrics, joint species
64 distribution models, prediction, species richness

65 **1. INTRODUCTION**

66 Species distributions are influenced by the response of a species to both the abiotic and biotic
67 conditions it encounters, but the development of species distribution models (SDMs) has
68 historically not accounted for the (biotic) effect of species interactions, largely focusing on
69 single-species approaches. Over the past decade, joint species distribution models (JSDMs)
70 have seen rapid development for modelling multiple species simultaneously while accounting
71 for both environmental responses and residual species associations (Kissling et al., 2012;
72 Warton et al., 2015; Wilkinson et al., 2019; Wisz et al., 2013). Research on JSDMs initially
73 focused on the development of different statistical modelling approaches (Clark et al., 2017;
74 Golding & Purse, 2016; Harris, 2015; Hui, 2016; Ovaskainen, Roy, et al., 2016; Pollock et al.,
75 2014), then extended the framework to account for additional factors such as the effect of
76 spatial scale or environmental gradients on species associations (Ovaskainen, Abrego, et al.,
77 2016; Thorson et al., 2016; Tikhonov et al., 2017). Only recently has the question of prediction
78 using JSDMs begun to be explored in any detail (Norberg et al., 2019; Wilkinson et al., 2021;
79 Zhang et al., 2018; Zurell et al., 2019).

80 A model's ability to accurately predict species distributions or communities needs to be
81 evaluated to be confident that the model performs well in practice. While there is substantial
82 literature on the evaluation of single species SDMs (Fielding & Bell, 1997; Lawson et al., 2014;
83 Liu et al., 2009), there is little research into the evaluation of the multi-species predictions of
84 JSDMs. Whilst single species model evaluation metrics can also be applied to the predictions
85 of multi-species models, the multivariate nature of JSDMs also invites the evaluation of
86 predictions using the community dissimilarity indices widely used in community ecology
87 (Legendre & De Cáceres, 2013).

88 The community assemblage process can be viewed as the result of two processes: (1) an abiotic
89 filter where species-environment relationships influence which species can occur in a given
90 environment, and (2) a biotic filter where between-species relationships influence which
91 species are more or less likely to co-occur (Cornell & Harrison, 2014; Götzenberger et al.,
92 2012). Some accounting of species interactions inside the single-species framework is possible,
93 by using additional species as predictor variables (Araújo & Luoto, 2007; Meier et al., 2010;
94 Zhang et al., 2020) or constraining predicted distributions to that of another species it depends
95 on (Schweiger et al., 2012), but this only reflects unidirectional interactions where the direction
96 of the relationship is already known (Kissling et al., 2012; Pollock et al., 2014; Wisz et al.,
97 2013).

98 Stacked species distribution models (SSDMs) represent community assemblages by stacking
99 the predictions of multiple *independent* single-species models (Guisan & Rahbek, 2011;
100 Thuiller et al., 2015). In contrast, joint species distribution models can capture both the biotic
101 and the abiotic factors impacting species and thus link distribution modelling and community
102 ecology. Wilkinson *et al.* (2021) detailed the different ways that JSDMs can make predictions:
103 environment-only marginal predictions at the species-level, joint predictions at the community-
104 level predictions that leverage the additional information on residual species occurrence, and
105 conditioning both marginal or joint on the known occurrence state of some species at a site.

106 JSDMs have been presented as a modelling approach with the potential to make better
107 predictions, particularly at the community level, than methods that do not account for the
108 residual correlations between species. However, the use of JSDMs for prediction has only
109 begun to be addressed in the literature (Gelfand & Shirota, 2021; Norberg et al., 2019;
110 Ovaskainen, Roy, et al., 2016; Poggiato et al., 2021; Wilkinson et al., 2021; Zhang et al., 2018;
111 Zurell et al., 2019). These studies have been limited in what aspects of JSDM predictive
112 performance they have considered. Only Norberg *et al.* (2019) considered prediction types that

113 account for the residual correlations between species in practice, while conditional prediction
114 types have only been theoretically discussed (Gelfand & Shirota, 2021; Poggiato et al., 2021).
115 Only Norberg *et al.* (2019) has compared multiple JSDM implementations, and they all have
116 compared predictive performance for only a select few evaluation metrics. In this study we
117 compare the predictive performance of six JSDMs and two SSDMs when fit to two real and 20
118 simulated datasets. We use 32 metrics over five metric classes to evaluate SSDM predictions
119 against the four different prediction types available to JSDMs.

120

121 **2. MATERIALS AND METHODS**

122 **2.1. MODELS**

123 The two SSDM models are generalised linear models with a probit link, fit individually to each
124 species but differ in their approach to stacking predictions. The first is a standard stacked
125 approach (SSDM) where individual species predictions are summed together to obtain species
126 richness predictions. The second approach is spatially explicit species assemblage modelling
127 (SESAM; Guisan & Rahbek, 2011) that selects species as present at a site up to a calculated
128 maximum limit (e.g. from a macroecological model).

129 All six JSDMs are based on the multivariate probit regression model of Chib and Greenberg
130 (1998). The first is the standard multivariate probit regression model (MPR) implemented in
131 the R package *BayesComm* (Golding et al., 2015). Second, the hierarchical multivariate probit
132 regression model (HPR) of Pollock *et al* (2014). Third, the multivariate probit regression with
133 latent factors (LPR) implemented in the R package *boral* (Hui, 2016). Fourth, the multivariate
134 generalised regression model (DPR) implemented in the R package *gjam* (Clark et al., 2017).
135 The fifth and sixth models are the spatial (HLR-S) and non-spatial (HLR-NS) implementations

136 of the hierarchical multivariate probit regression model with latent factors implemented in the
137 R package *HMSC* (Ovaskainen, Roy, et al., 2016). All of these models are implemented using
138 a Bayesian framework and fit using Markov chain Monte Carlo (MCMC) sampling, using their
139 default or suggested settings as defined in their source articles or the documentation of the
140 software implementing them. Model equations, default priors, and MCMC regimes are defined
141 in greater detail in Wilkinson *et al* (2019) and in Appendix S1.

142 **2.2. PREDICTION METHODS**

143 Single-species models generate environment-only predictions that ignore species associations.
144 For SSDMs, a prediction is obtained independently for each species using the estimated
145 regression coefficients and the corresponding measured variables. This provides a predicted
146 probability of presence for each species at each site. Binary predictions can be generated by
147 taking draws from a Bernoulli distribution, using the predicted probabilities. The SESAM only
148 provides binary predictions and constrains them to a site-specific species richness upper limit
149 using an estimated species richness, either from a macroecological model or alternatively, as
150 we have done here, using the probability rank rule (selecting species in decreasing order of
151 probability, considering as present as many as the total sum of probabilities rounded down).

152 JSDMs, in contrast, provide additional information on residual correlations between species
153 occurrence in their model outputs which can inform predictions. Marginal JSDM predictions
154 are species-level, environment-only predictions where the information about species co-
155 occurrence is ignored. Predictions are again probabilistic (a probability for each species at each
156 site), and binary predictions can be generated by drawing Bernoulli samples, as for the standard
157 SSDM. To quantify prediction uncertainty in binary predictions, Bernoulli samples can be
158 drawn for each species, site, and sample of predicted probability from the model posterior in a
159 Bayesian model fitting framework.

160 Joint predictions are obtained from the estimated joint probability distribution (a multivariate
161 normal distribution) over plausible community assemblages. This probability distribution
162 defines the probability of occurrence of an observed or hypothesised assemblage, or can be
163 used to simulate binary community assemblage predictions. Averaging the presence-absence
164 of individual species across simulated community assemblages approximates the marginal
165 distribution.

166 Conditional predictions include the additional information of *known* occurrence states for one
167 or more species in the community. The known occurrence state of a species truncates the
168 multivariate normal distribution on one axis. As the total probability must sum to one, this
169 affects the probability of the remaining possible community assemblages.

170 Conditional marginal predictions are similar to the conditional predictions in that they truncate
171 the multivariate normal distribution based on the occurrence state of one or more species. But
172 rather than representing the joint distribution of the remaining species, the distribution over the
173 remaining species is marginalised, to yield a single probability of presence for each species.
174 Thus, these are predictions conditional on one or more species but marginal to the remainder.

175 An in-depth explanation of these prediction methods can be found in Wilkinson *et al.*
176 (Wilkinson et al., 2021) or Appendix S2.

177 **2.3. DATASETS**

178 For this comparison we have used two real datasets, on frogs and eucalypts, and twenty
179 simulated presence-absence datasets. The frog dataset comprises 9 species, 104 sites, and 3
180 covariates from Melbourne, Australia (Parris, 2006). The eucalypt dataset comprises 12
181 species, 458 sites, and 7 covariates from Grampians National Park, Australia (Pollock et al.,
182 2014). The simulated datasets all have 10 species, 100 sites, and 5 covariates (3 continuous, 2

183 binary), with the species correlation matrix generated using three latent factors. Species
184 presence absence data was generated using the *HMSC::communitySimul* function (Blanchet et
185 al., 2019).

186 For all datasets, the continuous variables were standardised. For evaluation, we implemented
187 two different cross-validation approaches depending on the dataset. For the eucalypts and half
188 of the simulated datasets, we used five-fold spatial block cross-validation using the
189 *blockCV::spatialBlock* function (Valavi et al., 2019). For the frog and half of the simulated
190 datasets, we implemented five-fold random cross-validation, using the *caret::createFolds*
191 function (Kuhn et al., 2019), as the spatial scale of the frog dataset was too small for practical
192 implementation of spatial block cross-validation.

193 **2.4. EVALUATION METRICS**

194 The metrics available for evaluating JSDM predictions can be broadly classified into five
195 groups based on the aspects of performance they consider. Threshold-independent metrics
196 evaluate continuous predicted probabilities against observed presence-absence data. A
197 common example in SDMs is the Area Under the Receiver Operating Characteristic Curve
198 (AUC), but also includes root mean square error (RMSE), and the coefficient of determination
199 (R^2).

200 Threshold-dependent metrics compare binary predictions against observed presence-absence
201 data. Predicted probabilities are converted to presences if they exceed a set threshold value or
202 absences if they do not. A confusion matrix that contrasts observed and predicted species
203 occurrence states can then be built. Example metrics here include precision, sensitivity, and
204 true/false positive/negative rates. However, thresholding predictions is a contentious topic in
205 the SDM literature (Freeman & Moisen, 2008; Guillera-Aroita et al., 2015; Lawson et al.,

206 2014; Liu et al., 2005). In addition to debates about the use of thresholds in general, there are
207 also debates about how to determine the threshold value. Thresholds are commonly set to an
208 arbitrary value of 0.5 (Freeman & Moisen, 2008), but an alternative is to make the threshold
209 equivalent to the observed prevalence of the species (Hanberry & He, 2013). A logical
210 extension of this debate for JSDMs is species-specific or community-wide thresholds.
211 However, Lawson *et al* (2014) showed that we can evaluate threshold-dependent metrics
212 without thresholding predictions by calculating a probabilistic confusion matrix. We have used
213 this approach here to avoid any influence of threshold choice impacting our analysis.

214 As JSDMs are multi-species in nature we can use additional evaluation metrics from
215 community ecology in the form of community dissimilarity indices. These metrics compare
216 how dissimilar our observed and predicted species assemblages are. Common examples are
217 Bray-Curtis dissimilarity and Jaccard distance. These metrics are restricted to evaluating binary
218 predictions. To evaluate these metrics on probabilistic predictions, we simulated binary
219 community assemblages from the appropriate probability distribution. For JSDMs, a
220 community assemblage was drawn per posterior sample; for SSDMs, the same number of
221 community assemblages were simulated.

222 Species richness metrics consider a model's ability to predict a single aspect of community
223 composition: the number of species present at a site. We consider species richness difference -
224 the predicted richness minus the observed richness.

225 Likelihood-based metrics assess model fit by computing the probability of observing a given
226 community assemblage, assuming a particular model structure, and given the set of model
227 parameter estimates representing the prediction. It is common to work with the log of the
228 likelihood for numerical stability reasons. The independent log-likelihood represents the
229 typical log-likelihood metric used in SSDMs. This metric independently assesses each species

230 across all sites– computing the probability of observing a species’ presence/absence
231 observations– and combines them into a single metric, assuming the species’ distributions to
232 be independent. The joint log-likelihood simultaneously assesses all species, as an assemblage,
233 at each site and accounts for the correlation structure encoded in the core JSDM formulation –
234 the multivariate probit model.

235 More detail on the metrics including how they are calculated, which prediction types they are
236 appropriate for, and how to interpret them can be found in Appendix S3.

237 **2.5. MODEL PREDICTION COMPARISONS**

238 For the standard SSDM, we obtained binary and probabilistic predictions of community
239 assemblages, and for SESAM only binary predictions. For the JSDMs, we evaluated nine
240 prediction types: binary and probabilistic marginal predictions, binary joint predictions, binary
241 conditional predictions for low-, middle-, and high-prevalence known species scenarios, and
242 probabilistic conditional marginal predictions for low-, middle-, and high-prevalence known
243 species scenarios. For real datasets, the low-prevalence species were randomly selected from
244 those within the bottom 20% of prevalence, medium-prevalence species the middle 30%, and
245 high-prevalence the top 20%. For simulated datasets, this was the species with lowest, median,
246 and highest prevalence. We evaluated a suite of 32 evaluation metrics in total, but some applied
247 only to binary or probabilistic prediction types.

248 **2.6. ANALYSIS OF RESULTS**

249 We used linear mixed effects models (MEMs) to analyse the predictive performance of the
250 eight models for the nine prediction methods. We fit an MEM to each combination of
251 evaluation metric and prediction type, to assess the relationship between the response variable
252 (the evaluation statistic) and three explanatory variables (model, dataset, and cross-validation

253 fold), with a random effect on the intercept of either species or site (depending on the test
254 statistic’s calculation method). The performance results for any given model are relative to the
255 standard SSDM approach, which is set as the reference class. A partial interaction between
256 model and dataset was included for the HPR and DPR models to account for observed patterns
257 in the residuals. The MEMs explicitly considered different residual variances for the eight
258 model types to account for evident inhomogeneity of variance. We assessed whether the model
259 residuals met the model assumptions of being normally-distributed with homogenous variance
260 (after accounting for inhomogeneity between model types) with a Kolmogorov-Smirnov test
261 (Massey Jr, 1951), hereafter referred to as a KS-test. The KS-test determines if the distribution
262 of the residuals is significantly different from a normal distribution with the same mean and
263 standard deviation. As we performed a large number of comparisons, we used a Bonferroni-
264 corrected p-value, $\frac{0.05}{405} = 1.2 * 10^{-4}$ to consider whether normality assumptions were violated;
265 to reduce the sensitivity of these tests (Dunn, 1961).

266 **2.7. SOFTWARE**

267 All models were fit using R v3.5.2 (R Core Team, 2018). R packages for model fitting include
268 *BayesComm* v0.1-2 (Golding & Harris, 2015), *boral* v1.7 (Hui, 2018), *gjam* v2.2.5 (Clark &
269 Taylor-Rodríguez, 2018), and *HMSC* v2.2-0 (Blanchet et al., 2019). R packages required for
270 prediction and prediction evaluation are *mvtnorm* v1.0-10 (Genz et al., 2019), *tmvtnorm* v1.4-
271 10 (Wilhelm & B G, 2015), *TruncatedNormal* v1.0 (Zdravko, 2015), *Metrics* v0.1.4 (Hamner
272 & Frasco, 2018), *caret* 6.0-84 (Kuhn et al., 2019), *vegan* v2.5-5 (Oksanan et al., 2018), and
273 *psych* v1.8.12 (Revelle, 2018). Analyses were run on The University of Melbourne’s Spartan
274 HPC infrastructure (Meade et al., 2017).

275 **3. RESULTS**

276 No model outperformed the others across all prediction types, but there were consistent trends
277 within each class of validation metrics. The relative performance of the models compared to
278 the standard SSDM ($(\text{metric}_{\text{model}} - \text{metric}_{\text{SSDM}}) / |\text{metric}_{\text{SSDM}}| * 100$) are summarised in Figure
279 1 for marginal predictions and Figure 3 for joint predictions. Likelihood-based metrics showed
280 that the JSDMs were better fit to the data than the SSDM. Both threshold-dependent and
281 threshold-independent metrics indicated better performance by the SSDM, although the
282 difference was small for the majority of metrics. The SSDM generally outperformed JSDMs
283 for community dissimilarity metrics, with a greater difference for joint prediction types than
284 marginal ones. The JSDMs almost always overpredicted species richness compared to the
285 SSDM for marginal predictions, but HLR-S and HLR-NS had more accurate estimates than the
286 SSDM for most joint prediction types. The SESAM model outperformed both the SSDM and
287 JSDMs for threshold-dependent and community dissimilarity metrics for binary prediction
288 types (Figure 3).

289 Only a selection of figures is presented in the main article. Forest plots for the absolute value
290 of model performance for each evaluation metric are presented in Appendix S4. Individual
291 heatmaps for the different prediction types are presented in Appendix 5.

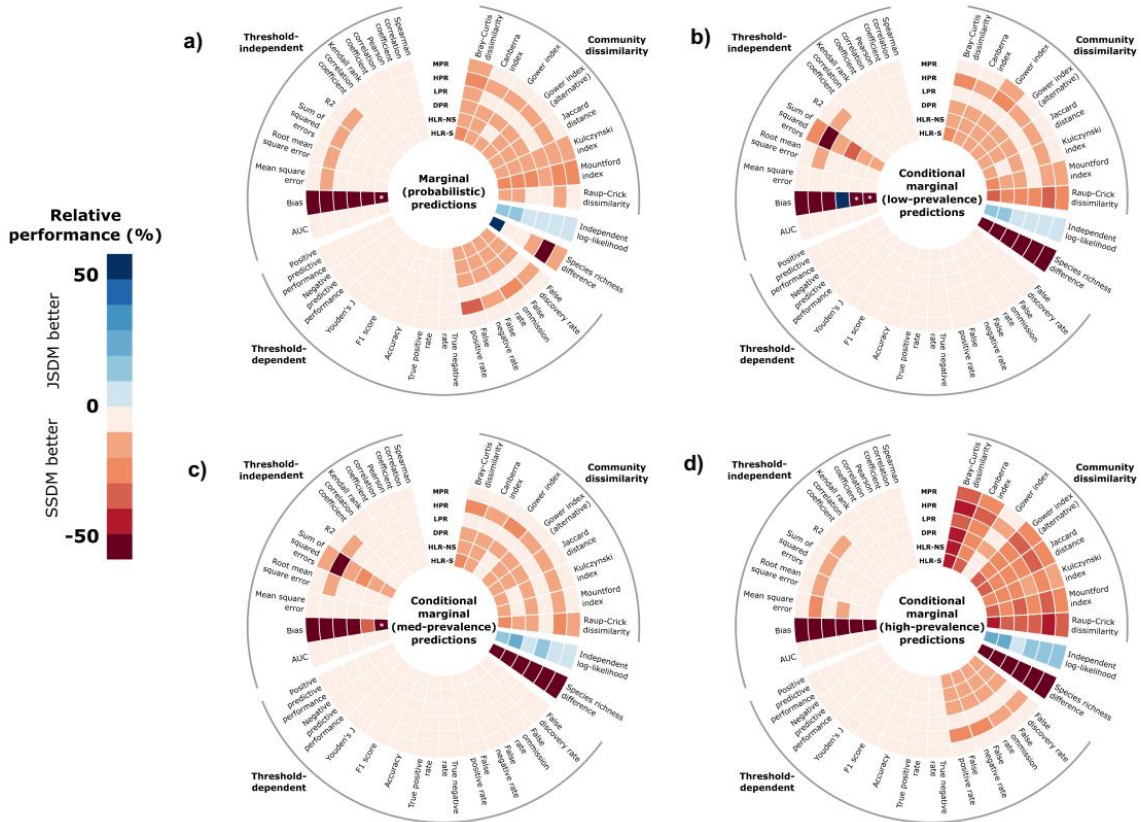
292 As the JSDM underperformance results are unexpected we performed checks to ensure the
293 MEMs were not returning erroneous results. We found no evidence to suggest our results are
294 an artefact of the model fitting process. We present the result of these checks in Appendix 6.

295 **3.1. PROBABILISTIC PREDICTIONS**

296 **3.1.1. SINGLE SPECIES METRICS**

297 The SSDM outperformed all JSDMs in marginal predictions for almost every threshold-
298 independent metric (Figure 1). The only exception to this was the bias metric for DPR in low-

299 prevalence conditional marginal predictions. The relative performance of the JSDBMs for the
300 bias metric suggests very poor performance of the JSDBMs, but on an absolute scale these
301 differences are actually quite small. As the optimum bias value is 0, good performing models
302 may show high relative differences in this metric when they actually have little absolute
303 difference (Figure 2). Across all probabilistic predictions, the SSDM had a mean bias of -
304 $2.7 \cdot 10^{-4}$, while the JSDBMs average $-1.6 \cdot 10^{-3}$, i.e. all models were largely unbiased. For all
305 other metrics, the average relative JSDBM performance across all marginal prediction types was
306 -6.3% [-20.7, -0.3]. HPR was the worst performer for the error-based metrics with a mean
307 relative difference of -18.3% [-61.8, -2.2]. For low- and medium-prevalence scenarios of
308 conditional marginal predictions, there was a large relative difference between SSDM and the
309 JSDBMs in sum of squared errors (-26.9% [-62.0, -13.0]). The mean difference over all marginal
310 prediction types for AUC between the SSDM and the worst performing JSDBM is only 0.01.



311

312

313 Figure 1: Heat maps showing the relative performance of the different JSDMs compared to the SSDM
 314 for all metrics applicable to *probabilistic* marginal predictions (see Appendix S4 for binary marginal
 315 predictions). Each ring represents a particular model, while each ray represents a different evaluation
 316 metric, which are clustered by class. Values in blue indicate better performance relative to the SSDM,
 317 while red values indicate worse performance relative to the SSDM. The four prediction types shown
 318 are a) probabilistic marginal, b) low-prevalence conditional marginal, c) medium-prevalence
 319 conditional marginal, and d) high-prevalence conditional marginal prediction. These heat maps were
 320 generated using Circos (Krzywinski et al., 2009).

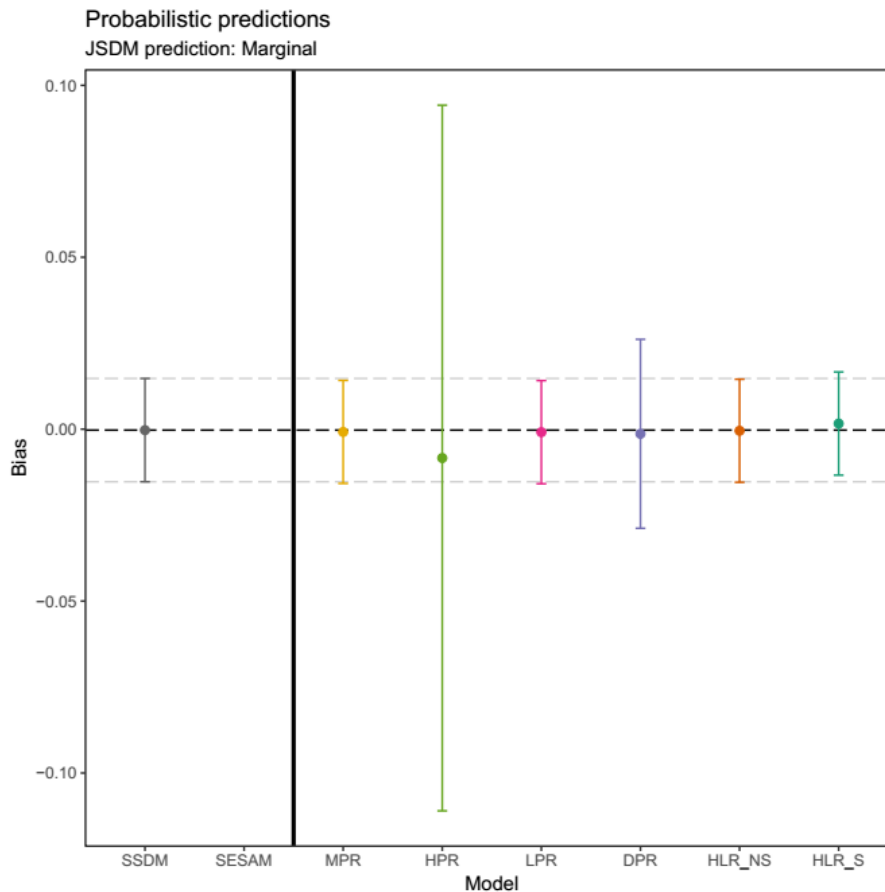
321

322 JSDMs performed worse than the SSDM for threshold-dependent metrics but differences were
 323 small. Across all probabilistic predictions, the relative JSDM performance was -3.9% [-18.2, -
 324 0.2]. The only notable relative differences were in HPR, DPR, HLR-S, and HLR-NS for the

325 false rate metrics for probabilistic marginal predictions (-17.8% [-28.7, -13.6]) and high-
326 prevalence conditional marginal predictions (-16.5% [-25.8, -13.5]).

327 **3.1.2. COMMUNITY DISSIMILARITY METRICS**

328 The SSDM outperformed the JSDMs for all community dissimilarity metrics on marginal
329 prediction types (Figure 1). The mean relative JSDM performance was -14.5% [-25.5, -4.9] for
330 probabilistic marginal predictions, -14.7% [-34.4, -5.7] for low-prevalence, -12.4% [-22.4, -
331 6.1] for medium-prevalence, and -27.8% [-44.1, -6.1] for high prevalence conditional marginal
332 predictions. For high-prevalence predictions the largest differences were for Bray-Curtis
333 dissimilarity (-40.4% [-46.9, -34.6]) and Raup-Crick dissimilarity (-36.2% [-42.5, -30.9]),
334 while the smallest was for Gower index (-9.3% [-17.8, -5.15]). The HPR model was the worst
335 performer overall with a mean relative difference of -23.1% [-43.7, -14.7] across all marginal
336 prediction types. LPR and MPR were the best JSDMs with a mean relative performance for all
337 marginal prediction types of -12.7% [-32.5, -4.8], compared to -19.7% [-41.2, -7.8] for the
338 other JSDMs.



339

340 Figure 2: Model performance for the bias (mean error) evaluation metric for probabilistic marginal
 341 predictions. Performance estimates are shown as the mean and 95% confidence intervals, over species
 342 in each dataset, after accounting for dataset and fold using an MEM. The dashed dark grey line
 343 corresponds to the SSDM mean bias, and the dashed light grey line corresponds to 95% confidence
 344 intervals from SSDM predictions. SSDMs are shown to the left of the black vertical line, and JSDMs
 345 to the right. This figure is illustrative of the absolute metric plots provided for all metrics and prediction
 346 types as supplementary information in Appendix S4. SESAM is not plotted here as it does not make
 347 predictions for individual species.

348 3.1.3. SPECIES RICHNESS

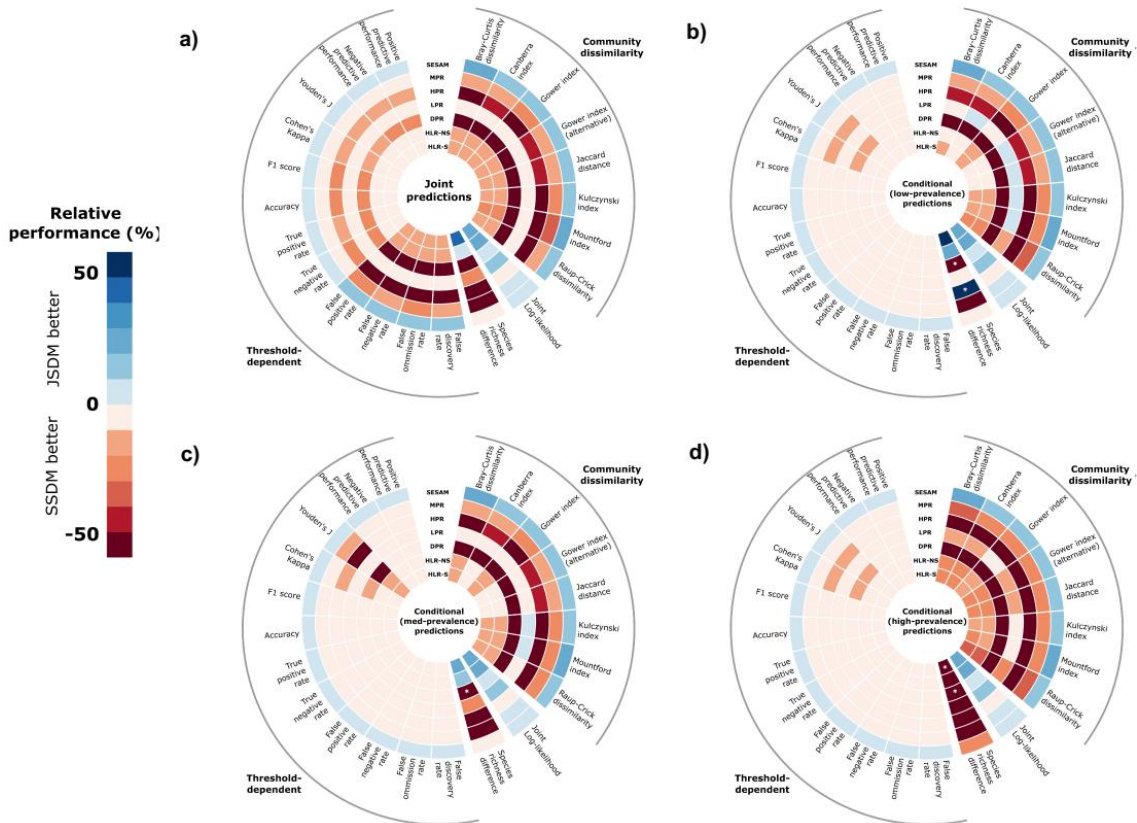
349 The SSDM generally outperformed JSDMs in species richness difference estimates for
 350 marginal predictions. However, the difference in performance is not as dramatic as suggested
 351 by the relative performance results in Figure 1 as relative differences can be exaggerated where
 352 the metric's optimum value is 0, so we discuss the absolute differences here. The SSDM had a
 353 mean species richness difference of 0.06 [0.02, 0.08] across the MEMs fitted for marginal

354 prediction results. For probabilistic marginal predictions, the JSDBMs had a similar mean
355 species richness difference of 0.09 [0.02, 0.15], with HLR-S outperforming the SSDM with a
356 value of just 0.007. For low-, medium-, and high-prevalence conditional marginal predictions,
357 the JSDBMs had a difference of 0.86 [0.79, 0.93], 0.63 [0.57, 0.70], and 0.34 [0.28, 0.39] species
358 respectively. HLR-S exhibited the smallest differences of all of the JSDBMs for each marginal
359 prediction type.

360 **3.2. BINARY PREDICTIONS**

361 **3.2.1. SINGLE SPECIES METRICS**

362 For threshold-dependent metrics evaluated on joint predictions, the JSDBMs performed worse
363 than the SSDM (Figure 3). For binary marginal predictions (Appendix S4), the JSDBMs had a
364 mean relative difference of -3.4% [-6.9, -0.7] from the SSDM except for the false rate metrics
365 for HPR, DPR, HLR-S, and HLR-NS with a mean of -17.9% [-29.0, -13.5]. The mean relative
366 difference between JSDBMs and SSDMs for conditional predictions was -4.1% [-15.0, -0.2]
367 across all metrics, with a larger relative difference for HPR and DPR on the Cohen's Kappa
368 and Youden's J metrics of -21.8 [-66.7, -11.9]. For all JSDBM joint predictions, there was a large
369 relative difference from SSDMs of -36.3 [-97.7, -3.5] for false rate metrics. This was strongest
370 for HPR and DPR with means of -70.4% and -88.7% respectively. LPR had a comparatively
371 small relative difference of -4.0%. Across all threshold-dependent metrics for joint prediction,
372 HPR and DPR had a large relative difference of -35.7% [-82.2, -12.0] and -45.2% [-100.9, -
373 15.0] respectively.



374

375 Figure 3: Heat maps showing the relative performance of the different JSDMs and SESAM compared
 376 to the SSDM for all metrics applicable to joint (binary) predictions. Each ring represents a particular
 377 model, while each ray represents a different evaluation metric, which are clustered by class. Values in
 378 blue indicate better performance relative to the SSDM, while red values indicate worse performance
 379 relative to the SSDM. The four prediction types shown are a) joint, b) low-prevalence conditional, c)
 380 medium-prevalence conditional, and d) high-prevalence conditional prediction.

381 3.2.2.COMMUNITY DISSIMILARITY METRICS

382 For community dissimilarity metrics on binary marginal predictions, the SSDM outperformed
 383 the JSDMs. The relative performance difference was largest for HPR, DPR, HLR-S, and HLR-
 384 NS (-14.2% [-24.1, -10.0]), and smallest for MPR and LPR (-2.8% [-4.1, -2.0]). For joint
 385 prediction types, HPR and DPR performed worst with a mean relative performance difference
 386 of -75.3% [-132.7, -44.7]. LPR performed best with a mean difference of -4.3% [-20.5, 6.3],

387 and performed slightly better than the SSDM for several metrics in low- or medium-prevalence
388 scenarios.

389 **3.2.3. SPECIES RICHNESS**

390 The results for species richness were mixed. Across all joint prediction types, the mean
391 difference between predicted and observed species richness was 0.06 [0.01, 0.09] for the
392 SSDM. HLR-S had a small mean species richness difference of 0.02 [-0.01, 0.04] while HLR-
393 NS performed similarly to the SSDM with a mean difference of 0.06 [0.03, 0.09]. LPR, HPR,
394 and MPR exhibited slightly larger mean species richness differences than the SSDM with
395 values of 0.07 [0.02, 0.1], 0.09 [-0.03, 0.17], and 0.14 [0.1, 0.17] respectively. DPR was the
396 only JSDM to underpredict species richness with a mean difference value of -0.12 [-0.38, 0.14].

397 **3.3. LIKELIHOOD-BASED METRICS**

398 Both likelihood metrics (Figures 1 and 3) indicate that the JSDMs were better fit to the out-of-
399 sample data than the SSDMs, under the assumptions of the univariate and multivariate probit
400 models. For the probabilistic marginal predictions, the JSDMs on average performed 12.2%
401 [7.9, 18.4] better than the SSDM, and for the low-, medium-, and high-prevalence conditional
402 marginal predictions the JSDMs on average performed 12.5% [8.0, 19.0], 13.6% [8.5, 20.4],
403 and 15.2% [10.4, 22.1] better than the SSDM. HLR-S and HLR-NS outperformed the other
404 JSDMs across all probabilistic prediction types with a mean relative independent log likelihood
405 of 19.3% [17.0, 22.0] compared to 9.4% [7.9, 12.0] for the remaining models. Results were
406 similar for the joint log-likelihood except that HPR performed worse relative to the SSDM by
407 7.8%. Relative to the SSDM, HLR-S and HLR-NS performed better than the other JSDMs by
408 an additional ~10%.

409 **3.4. SESAM PERFORMANCE**

410 SESAM outperformed the SSDM and all JSJM joint predictions for the threshold-dependent
411 and community dissimilarity metrics. SESAM had a mean relative performance difference of
412 3.5% [0.7, 15.6] across threshold-dependent metrics for all joint prediction types. This
413 performance is strongest in the false rate metric for binary marginal and joint predictions with
414 a mean relative difference of 13.2% [10.0, 16.7]. For the community dissimilarity metrics,
415 SESAM had a mean relative performance across all joint prediction types of 17.8% [14.4,
416 22.5]. SESAM and SSDM richness predictions were functionally equivalent, as expected, as
417 SESAM uses SSDM predictions to set its species richness limit.

418 **4. DISCUSSION**

419 **4.1. SINGLE SPECIES METRICS**

420 The SSDM routinely outperformed the JSJMs in marginal prediction for all threshold-
421 independent metrics, but the magnitude of the difference was generally small on the absolute
422 scale of the metric. This differs from Zurell *et al* (2019) which found the JSJM performed
423 much worse than the stacked model. However, that study used a stacked ensemble SDM while
424 we have just used a generalised linear model based SDM. Ensemble models have the potential
425 to outperform single models (Dormann *et al.*, 2018; Hao *et al.*, 2019), which may have
426 increased the observed difference in performance between the JSJM and SSDM. We also
427 observed small differences on the absolute scale for both common error-based metrics (e.g.
428 bias, MSE, RMSE) and the correlation metrics. The only exception to this was the SSE metric
429 for low- and medium-prevalence conditional marginal metrics, but all models, including the
430 SSDM, had significantly larger errors here compared to the other marginal prediction types.
431 This suggests that, for species-level predictions, the JSJMs and SSDM are performing
432 similarly in terms of threshold-independent metrics.

433 JSDBMs performed slightly worse than, but generally similar to, the SSDM for all threshold-
434 dependent metrics on both marginal and joint prediction types but, like for threshold-
435 independent metrics, differences were minor. The relative performance of HPR, DPR, HLR-S,
436 and HLR-NS suggests a large difference between them and the SSDM for the false rate metrics
437 in probabilistic marginal and high-prevalence conditional marginal predictions, but on the
438 absolute scale the performance difference is $\leq 4\%$ for all metrics. The same trend can be seen
439 in the joint prediction type exhibit larger relative differences for the false rate metrics. On the
440 absolute scale these differences are $\leq 4\%$ for most metrics, and usually $\leq 1\%$ for MPR and LPR.
441 The exception to this is that both HPR and DPR perform poorly in the joint predictions, with
442 difference on the absolute scale of 15.5% and 19.6% respectively.

443 **4.2.COMMUNITY DISSIMILARITY METRICS**

444 One of the purported benefits of JSDBMs over single species models is that accounting for
445 residual species co-occurrences during the model fitting process they will better predict
446 community composition by accounting for species associations. These species associations,
447 which could include species interactions or shared responses to unmeasured environmental
448 variables, provide information on how likely species are to co-occur beyond their response to
449 the measured variables. However, our community dissimilarity metrics results show the
450 JSDBMs predicting worse than the SSDM, although the difference was again minor on the
451 absolute scale. Only LPR outperformed the SSDM for some conditional predictions, but just
452 by a mean of 0.004, so the performance can in practice be considered identical.

453 **4.3.LIKELIHOOD**

454 In almost all cases JSDBMs outperformed the SSDM for both of the likelihood-based metrics,
455 indicating they were better fit to the in-sample data. This is consistent with Norberg *et al*

456 (2019), who also found better likelihood metrics for JSDBMs when they contrasted model pairs
457 that were identical in all aspects except accounting, or not, for residual correlations. The JSDBMs
458 exhibiting better likelihood performance is expected as they include additional useful model
459 parameters (i.e. a residual correlation structure). The estimation of regression coefficients is
460 fairly robust to the covariance matrix, so JSDBMs and SSDMs would estimate similar species
461 niches(Chib & Greenberg, 1998; Poggiato et al., 2021), therefore a possible explanation for
462 JSDBMs performing better on likelihood-based metrics and worse on the other classes is that the
463 JSDBMs have overfit the covariance matrix and are explaining some noise in the data (Poggiato
464 et al., 2021). This could also explain why we see greater differences in performance between
465 the JSDBMs and the SSDM for the community-level joint predictions that leverage this
466 information. Conditioning these predictions with known information still improves the
467 predictions, which suggests the estimated correlations are too strong rather than sign-switched.
468 A potential solution would be penalising the residual covariance matrix to prevent overfitting
469 (Pichler & Hartig, 2021).

470 **4.4. SPECIES RICHNESS**

471 The species richness difference metric presented mixed results where the SSDM did not
472 consistently perform better or worse than the JSDBMs. All models generally over-predicted
473 species richness, with the exception of DPR for conditional predictions, but the mean difference
474 was minimal at $\sim \leq 0.1$ species per site. HLR-S and HLR-NS were the best performing JSDBMs
475 overall. For all prediction types, except the conditional marginal, HLR-S outperformed the
476 SSDM and HLR-NS performed equivalently to it. As HLR-S was the only JSDBM to regularly
477 outperform the SSDM, it suggests that the effect of spatial scale and/or spatially-driven
478 unmeasured variables could be a potential driver of this result. Species co-occurrence can be
479 driven by several factors, including species interactions or shared responses to environmental

480 conditions, which operate at different scales. Two species can have the same broad
481 environmental condition preferences but tend to rarely co-occur at a finer scale. If the estimated
482 correlations between species are driven by shared responses to environmental variables at
483 larger spatial scales this can lead to higher species richness estimates at the site level
484 (Ovaskainen, Roy, et al., 2016). All JSDMs had much larger species richness difference
485 estimates for the probabilistic conditional marginal predictions compared to the binary
486 conditional marginal predictions, but it is unclear why. This result contrasts previous works
487 that indicate binary predictions are more prone to overprediction than probabilistic ones
488 (Calabrese et al., 2014; Thuiller et al., 2015; Zurell et al., 2019).

489 **4.5. SESAM PERFORMANCE**

490 The SESAM predictions outperformed the SSDM and JSDMs for all evaluation metrics in the
491 threshold-dependent and community dissimilarity metric classes. This suggests that
492 constraining the number of species predicted at a site can improve predictive performance at
493 both the species- and community-level, possibly acting as a carrying capacity proxy (Guisan
494 & Rahbek, 2011). These results are consistent with Zurell *et al* (2019) who found that
495 predictions constrained with SESAM's probability rank rule performed better than SSDMs for
496 both species- and community-level metrics. Zurell *et al* (2019) found a small benefit to species
497 richness metrics for SESAM predictions compared to an SSDM, but they used a
498 macroecological model to set the species richness limit compared to using the summed SSDM
499 predictions as we have in this study (thus being unsurprising that we identified no difference
500 between the two in this aspect of predictive performance).

501 An interesting avenue of research not considered in our study is exploring whether constraining
502 JSDM predictions can improve their performance. The superior performance of SESAM in our
503 results is in concordance with that of Zurell *et al* (2019) who suggested that choosing how

504 species predictions are combined into community-level predictions is potentially more
505 important than choosing the underlying model used to generate them. JSDMs could be
506 predicting likely community assemblages in the absence of limiting factors like site carrying
507 capacities or dispersal limits, and thus could potentially benefit from the application of a
508 constraining framework. Research into how to incorporate suitable constraints into the
509 prediction process itself rather than applying them post-hoc is suggested.

510 **5. CONCLUSION**

511 While there were consistent trends within evaluation metric classes, we did not find evidence
512 to suggest that any one model outperformed all of the others across all prediction types. The
513 likelihood metrics indicated that the JSDMs were better fit to the data, but SSDMs generally
514 outperformed all of the JSDMs in the rest of evaluation metrics. On the absolute scale, the
515 difference in performance between models was generally small. HLR-S and HLR-NS were the
516 best performing JSDMs and were able to outperform the SSDM for most species richness
517 difference estimates in joint prediction types and generally had the smallest difference in
518 performance from the SSDM when underperforming. The SESAM model consistently
519 outperformed both the JSDMs and the SSDM for both binary species- and community-level
520 metrics which suggests that the application of frameworks to constrain JSDM prediction types
521 should be evaluated in the future.

522

523 **REFERENCE LIST**

- 524 Araújo, M. B., & Luoto, M. (2007). The importance of biotic interactions for modelling
525 species distributions under climate change. *Global Ecology and Biogeography*, 16(6),
526 743–753. <https://doi.org/10.1111/j.1466-8238.2007.00359.x>
- 527 Blanchet, F. G., Tikhonov, G., & Norberg, A. (2019). *HMSC: Hierarchical Modelling of*
528 *Species Community* (2.2-0). github.com/guibranchet/HMSC
- 529 Calabrese, J. M., Certain, G., Kraan, C., & Dormann, C. F. (2014). Stacking species
530 distribution models and adjusting bias by linking them to macroecological models.
531 *Global Ecology and Biogeography*, 23(1), 99–112. <https://doi.org/10.1111/geb.12102>
- 532 Chib, S., & Greenberg, E. (1998). Analysis of multivariate probit models. *Biometrika*, 85(2),
533 347–361. <https://doi.org/10.1093/biomet/85.2.347>
- 534 Clark, J. S., Nemergut, D., Seyednasrollah, B., Turner, P. J., & Zhang, S. (2017). Generalized
535 joint attribute modeling for biodiversity analysis: median-zero, multivariate, multifarious
536 data. *Ecological Monographs*, 87(1), 34–56. <https://doi.org/10.1002/ecm.1241>
- 537 Clark, J. S., & Taylor-Rodríguez, D. (2018). *gjam: Generalised Joint Attribute Modelling*
538 (2.2.5). <https://cran.r-project.org/package=gjam>
- 539 Cornell, H. V., & Harrison, S. P. (2014). What Are Species Pools and When Are They
540 Important? *Annual Review of Ecology, Evolution, and Systematics*, 45(1), 45–67.
541 <https://doi.org/10.1146/annurev-ecolsys-120213-091759>
- 542 Dormann, C. F., Calabrese, J. M., Guillerá-Arroita, G., Matechou, E., Bahn, V., Bartoń, K.,
543 Beale, C. M., Ciuti, S., Elith, J., Gerstner, K., Guelat, J., Keil, P., Lahoz-Monfort, J. J.,
544 Pollock, L. J., Reineking, B., Roberts, D. R., Schröder, B., Thuiller, W., Warton, D. I.,
545 ... Hartig, F. (2018). Model averaging in ecology: a review of Bayesian, information-
546 theoretic, and tactical approaches for predictive inference. *Ecological Monographs*,
547 88(4), 485–504. <https://doi.org/10.1002/ecm.1309>
- 548 Dunn, O. J. (1961). Multiple Comparisons Among Means. *Journal of the American*
549 *Statistical Association*, 56(293), 52–64.
- 550 Fielding, A. H., & Bell, J. F. (1997). A review of methods for the assessment of prediction
551 errors in conservation presence/absence models. *Environmental Conservation*, 24(1),
552 38–49. <https://doi.org/10.1017/S0376892997000088>
- 553 Freeman, E. A., & Moisen, G. G. (2008). A comparison of the performance of threshold
554 criteria for binary classification in terms of predicted prevalence and kappa. *Ecological*
555 *Modelling*, 217(1–2), 48–58. <https://doi.org/10.1016/j.ecolmodel.2008.05.015>
- 556 Gelfand, A. E., & Shirota, S. (2021). The role of odds ratios in joint species distribution
557 modeling. *Environmental and Ecological Statistics*. [https://doi.org/10.1007/s10651-021-](https://doi.org/10.1007/s10651-021-00486-4)
558 00486-4
- 559 Genz, A., Bretz, F., Miwa, T., Mi, X., Leisch, F., Scheipl, F., & Hothorn, T. (2019).
560 *mvtnorm: Multivariate Normal and t Distributions* (1.0-10). [http://cran.r-](http://cran.r-project.org/package=mvtnorm)
561 [project.org/package=mvtnorm](http://cran.r-project.org/package=mvtnorm)

- 562 Golding, N., & Harris, D. J. (2015). *BayesComm: Bayesian Community Ecology Analysis*
563 (0.1-2). <https://cran.r-project.org/package=BayesComm>
- 564 Golding, N., Nunn, M. A., & Purse, B. V. (2015). Identifying biotic interactions which drive
565 the spatial distribution of a mosquito community. *Parasites & Vectors*, 8(1), 367.
566 <https://doi.org/10.1186/s13071-015-0915-1>
- 567 Golding, N., & Purse, B. V. (2016). Fast and flexible Bayesian species distribution modelling
568 using Gaussian processes. *Methods in Ecology and Evolution*, 7(5), 598–608.
569 <https://doi.org/10.1111/2041-210X.12523>
- 570 Götzenberger, L., de Bello, F., Bräthen, K. A., Davison, J., Dubuis, A., Guisan, A., Lepš, J.,
571 Lindborg, R., Moora, M., Pärtel, M., Pellissier, L., Pottier, J., Vittoz, P., Zobel, K., &
572 Zobel, M. (2012). Ecological assembly rules in plant communities—approaches, patterns
573 and prospects. *Biological Reviews*, 87(1), 111–127. <https://doi.org/10.1111/j.1469-185X.2011.00187.x>
- 575 Guillera-Arroita, G., Lahoz-Monfort, J. J., Elith, J., Gordon, A., Kujala, H., Lentini, P. E.,
576 Mccarthy, M. A., Tingley, R., & Wintle, B. A. (2015). Is my species distribution model
577 fit for purpose? Matching data and models to applications. *Global Ecology and*
578 *Biogeography*, 24(3), 276–292. <https://doi.org/10.1111/geb.12268>
- 579 Guisan, A., & Rahbek, C. (2011). SESAM - a new framework integrating macroecological
580 and species distribution models for predicting spatio-temporal patterns of species
581 assemblages. *Journal of Biogeography*, 38(8), 1433–1444.
582 <https://doi.org/10.1111/j.1365-2699.2011.02550.x>
- 583 Hamner, B., & Frasco, M. (2018). *Metrics: Evaluation Metrics for Machine Learning* (0.1.4).
584 <https://cran.r-project.org/package=Metrics>
- 585 Hanberry, B. B., & He, H. S. (2013). Prevalence, statistical thresholds, and accuracy
586 assessment for species distribution models. *Web Ecol*, 13, 13–19.
587 <https://doi.org/10.5194/we-13-13-2013>
- 588 Hao, T., Elith, J., Guillera-Arroita, G., & Lahoz-Monfort, J. J. (2019). A review of evidence
589 about use and performance of species distribution modelling ensembles like BIOMOD.
590 *Diversity and Distributions*, 25(5), 839–852. <https://doi.org/10.1111/ddi.12892>
- 591 Harris, D. J. (2015). Generating realistic assemblages with a joint species distribution model.
592 *Methods in Ecology and Evolution*, 6(4), 465–473. <https://doi.org/10.1111/2041-210X.12332>
- 594 Hui, F. K. C. (2016). boral - Bayesian Ordination and Regression Analysis of Multivariate
595 Abundance Data in r. *Methods in Ecology and Evolution*, 7(6), 744–750.
596 <https://doi.org/10.1111/2041-210X.12514>
- 597 Hui, F. K. C. (2018). *boral: Bayesian Ordination and Regression AnaLysis* (1.7).
598 <https://cran.r-project.org/package=boral>
- 599 Kissling, W. D., Dormann, C. F., Groeneveld, J., Hickler, T., Kühn, I., McNerny, G. J.,
600 Montoya, J. M., Römermann, C., Schiffers, K., Schurr, F. M., Singer, A., Svenning, J.-
601 C., Zimmermann, N. E., & O'Hara, R. B. (2012). Towards novel approaches to
602 modelling biotic interactions in multispecies assemblages at large spatial extents.

- 603 *Journal of Biogeography*, 39(12), 2163–2178. <https://doi.org/10.1111/j.1365->
604 2699.2011.02663.x
- 605 Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S. J., &
606 Marra, M. A. (2009). Circos: An information aesthetic for comparative genomics.
607 *Genome Research*, 19(9), 1639–1645. <https://doi.org/10.1101/gr.092759.109>
- 608 Kuhn, M., Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper, T., Mayer,
609 Z., Kenkel, B., R Core Team, Benesty, M., Lescarbeau, R., Ziem, A., Scrucca, L., Tang,
610 L., Candan, C., & Hunt, T. (2019). *caret: Classification and Regression Training* (6.0-
611 84). <https://cran.r-project.org/package=caret>
- 612 Lawson, C. R., Hodgson, J. A., Wilson, R. J., & Richards, S. A. (2014). Prevalence,
613 thresholds and the performance of presence-absence models. *Methods in Ecology and*
614 *Evolution*, 5(1), 54–64. <https://doi.org/10.1111/2041-210X.12123>
- 615 Legendre, P., & De Cáceres, M. (2013). Beta diversity as the variance of community data:
616 Dissimilarity coefficients and partitioning. *Ecology Letters*, 16(8), 951–963.
617 <https://doi.org/10.1111/ele.12141>
- 618 Liu, C., Berry, P. M., Dawson, T. P., & Pearson, R. G. (2005). Selecting thresholds of
619 occurrence in the prediction of species distributions. *Ecography*, 28(3), 385–393.
620 <https://doi.org/10.1111/j.0906-7590.2005.03957.x>
- 621 Liu, C., White, M., & Newell, G. (2009). Measuring the accuracy of species distribution
622 models: A review. *18th World IMACS Congress and MODSIM09 International*
623 *Congress on Modelling and Simulation: Interfacing Modelling and Simulation with*
624 *Mathematical and Computational Sciences, Proceedings, July*, 4241–4247.
- 625 Massey Jr, F. J. (1951). The Kolmogorov-Smirnov Test for Goodness of Fit. *Journal of the*
626 *American Statistical Association*, 46(253), 68–78.
627 <https://doi.org/10.1080/01621459.1951.10500769>
- 628 Meade, B., Lafayette, L., Sauter, G., & Tosello, D. (2017). *Spartan HPC-Cloud Hybrid:*
629 *Delivering Performance and Flexibility*. <https://doi.org/10.4225/49/58ead90dceaaa>
- 630 Meier, E. S., Kienast, F., Pearman, P. B., Svenning, J.-C., Thuiller, W., Araújo, M. B.,
631 Guisan, A., & Zimmermann, N. E. (2010). Biotic and abiotic variables show little
632 redundancy in explaining tree species distributions. *Ecography*, 33(6), 1038–1048.
633 <https://doi.org/10.1111/j.1600-0587.2010.06229.x>
- 634 Norberg, A., Abrego, N., Blanchet, F. G., Adler, F. R., Anderson, B. J., Anttila, J., Araújo,
635 M. B., Dallas, T., Dunson, D., Elith, J., Foster, S. D., Fox, R., Franklin, J., Godsoe, W.,
636 Guisan, A., O’Hara, B., Hill, N. A., Holt, R. D., Hui, F. K. C., ... Ovaskainen, O.
637 (2019). A comprehensive evaluation of predictive performance of 33 species distribution
638 models at species and community levels. *Ecological Monographs*, 89(3), 1–24.
639 <https://doi.org/10.1002/ecm.1370>
- 640 Oksanen, J., Guillaume Blanchet, F., Friendly, M., Kindt, R., Legendre, P., McGlinn, D.,
641 Minchin, P., O’Hara, R. B., Simpson, G. L., Solymos, P., Stevens, M. H. H., Szoecs, E.,
642 & Wagner, H. (2018). *vegan: Community Ecology Package* (2.5-5). [https://cran.r-](https://cran.r-project.org/package=vegan)
643 [project.org/package=vegan](https://cran.r-project.org/package=vegan)

- 644 Ovaskainen, O., Abrego, N., Halme, P., & Dunson, D. (2016). Using latent variable models
645 to identify large networks of species-to-species associations at different spatial scales.
646 *Methods in Ecology and Evolution*, 7(5), 549–555. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210X.12501)
647 210X.12501
- 648 Ovaskainen, O., Roy, D. B., Fox, R., & Anderson, B. J. (2016). Uncovering hidden spatial
649 structure in species communities with spatially explicit joint species distribution models.
650 *Methods in Ecology and Evolution*, 7(4), 428–436. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210X.12502)
651 210X.12502
- 652 Parris, K. M. (2006). Urban amphibian assemblages as metacommunities. *Journal of Animal*
653 *Ecology*, 75(3), 757–764. <https://doi.org/10.1111/j.1365-2656.2006.01096.x>
- 654 Pichler, M., & Hartig, F. (2021). A new joint species distribution model for faster and more
655 accurate inference of species associations from big community data. *Methods in Ecology*
656 *and Evolution*, 12(11), 2159–2173. <https://doi.org/10.1111/2041-210X.13687>
- 657 Poggiato, G., Münkemüller, T., Bystrova, D., Arbel, J., Clark, J. S., & Thuiller, W. (2021).
658 On the Interpretations of Joint Modeling in Community Ecology. *Trends in Ecology and*
659 *Evolution*, 36(5), 391–401. <https://doi.org/10.1016/j.tree.2021.01.002>
- 660 Pollock, L. J., Tingley, R., Morris, W. K., Golding, N., O’Hara, R. B., Parris, K. M., Vesk, P.
661 A., & McCarthy, M. A. (2014). Understanding co-occurrence by modelling species
662 simultaneously with a Joint Species Distribution Model (JSDM). *Methods in Ecology*
663 *and Evolution*, 5(5), 397–406. <https://doi.org/10.1111/2041-210X.12180>
- 664 R Core Team. (2018). *R: A Language and Environment for Statistical Computing* (3.5.2).
- 665 Revelle, W. (2018). *psych: Procedures for Personality and Psychological Research* (1.8-12).
- 666 Schweiger, O., Heikkinen, R. K., Harpke, A., Hickler, T., Klotz, S., Kudrna, O., Kühn, I.,
667 Pöyry, J., & Settele, J. (2012). Increasing range mismatching of interacting species
668 under global change is related to their ecological characteristics. *Global Ecology and*
669 *Biogeography*, 21(1), 88–99. <https://doi.org/10.1111/j.1466-8238.2010.00607.x>
- 670 Thorson, J. T., Ianelli, J. N., Larsen, E. A., Ries, L., Scheuerell, M. D., Szuwalski, C., &
671 Zipkin, E. F. (2016). Joint dynamic species distribution models: a tool for community
672 ordination and spatio-temporal monitoring. *Global Ecology and Biogeography*, 25(9),
673 1144–1158. <https://doi.org/10.1111/geb.12464>
- 674 Thuiller, W., Pollock, L. J., Gueguen, M., & Münkemüller, T. (2015). From species
675 distributions to meta-communities. *Ecology Letters*, 18(12), 1321–1328.
676 <https://doi.org/10.1111/ele.12526>
- 677 Tikhonov, G., Abrego, N., Dunson, D., & Ovaskainen, O. (2017). Using joint species
678 distribution models for evaluating how species-to-species associations depend on the
679 environmental context. *Methods in Ecology and Evolution*, 8(4), 443–452.
680 <https://doi.org/10.1111/2041-210X.12723>
- 681 Valavi, R., Elith, J., Lahoz-Monfort, J. J., & Guillera-Arroita, G. (2019). blockCV: An r
682 package for generating spatially or environmentally separated folds for k-fold cross-
683 validation of species distribution models. *Methods in Ecology and Evolution*, 10(2),
684 225–232. <https://doi.org/10.1111/2041-210X.13107>

685 Warton, D. I., Foster, S. D., De, G., Stoklosa, J., & Dunstan, P. K. (2015). Model-based
686 thinking for community ecology. *Plant Ecology*, 669–682.
687 <https://doi.org/10.1007/s11258-014-0366-3>

688 Wilhelm, S., & B G, M. (2015). *tmvtnorm: Truncated Multivariate Normal and Student t*
689 *Distribution* (1.4-10). <http://cran.r-project.org/package=tmvtnorm>

690 Wilkinson, D. P. (2019). *JSDM_Prediction v0.1.0*. Zenodo.
691 <https://doi.org/10.5281/zenodo.3514766>

692 Wilkinson, D. P., Golding, N., Guillera-Arroita, G., Tingley, R., & McCarthy, M. A. (2019).
693 A comparison of joint species distribution models for presence–absence data. *Methods*
694 *in Ecology and Evolution*, 10(2), 198–211. <https://doi.org/10.1111/2041-210X.13106>

695 Wilkinson, D. P., Golding, N., Guillera-Arroita, G., Tingley, R., & McCarthy, M. A. (2021).
696 Defining and evaluating predictions of joint species distribution models. *Methods in*
697 *Ecology and Evolution*, 12(3), 394–404. <https://doi.org/10.1111/2041-210X.13518>

698 Wisz, M. S., Pottier, J., Kissling, W. D., Pellissier, L., Lenoir, J., Damgaard, C. F., Dormann,
699 C. F., Forchhammer, M. C., Grytnes, J. A., Guisan, A., Heikkinen, R. K., Høye, T. T.,
700 Kühn, I., Luoto, M., Maiorano, L., Nilsson, M. C., Normand, S., Öckinger, E., Schmidt,
701 N. M., ... Svenning, J. C. (2013). The role of biotic interactions in shaping distributions
702 and realised assemblages of species: Implications for species distribution modelling.
703 *Biological Reviews*, 88(1), 15–30. <https://doi.org/10.1111/j.1469-185X.2012.00235.x>

704 Zdravko, B. I. (2015). *TruncatedNormal: Truncated Multivariate Normal* (1.0). [https://cran.r-](https://cran.r-project.org/package=TruncatedNormal)
705 [project.org/package=TruncatedNormal](https://cran.r-project.org/package=TruncatedNormal)

706 Zhang, C., Chen, Y., Xu, B., Xue, Y., & Ren, Y. (2018). Comparing the prediction of joint
707 species distribution models with respect to characteristics of sampling data. *Ecography*,
708 41(11), 1876–1887. <https://doi.org/10.1111/ecog.03571>

709 Zhang, C., Chen, Y., Xu, B., Xue, Y., & Ren, Y. (2020). Improving prediction of rare
710 species' distribution from community data. *Scientific Reports*, 10(1), 1–9.
711 <https://doi.org/10.1038/s41598-020-69157-x>

712 Zurell, D., Zimmermann, N. E., Gross, H., Baltensweiler, A., Sattler, T., & Wüest, R. O.
713 (2019). Testing species assemblage predictions from stacked and joint species
714 distribution models. *Journal of Biogeography*, 00, 1–13.
715 <https://doi.org/10.1111/jbi.13608>

716