

1 **Fast model-based ordination with copulas**

Gordana C. Popovic*, Francis K. C. Hui†, David I. Warton*,

2 *School of Mathematics and Statistics and the Evolution & the Ecology Research Centre, UNSW
Sydney, NSW 2052, Australia

†Research School of Finance, Actuarial Studies & Statistics, Australia National University, Acton,
ACT 2601, Australia

3 **Header:** Fast model-based ordination

4 **Keywords:** ordination, Gaussian copula, abundance data, high-dimensional data, mul-
5 tivariate analysis, species interactions]

6 **Word count:** 6359

7 **ABSTRACT**

8 1. Visualising data is a vital part of analysis, allowing researchers to find patterns, and
9 assess and communicate the results of statistical modeling. In ecology, visualisation
10 is often challenging when there are many variables (often for different species or
11 other taxonomic groups) and they are not normally distributed (often counts or
12 presence-absence data). Ordination is a common and powerful way to overcome
13 this hurdle by reducing data from many response variables to just two or three, to
14 be easily plotted.

15 2. Ordination is traditionally done using dissimilarity-based methods, most commonly
16 non-metric multidimensional scaling (nMDS). In the last decade however, model-
17 based methods for unconstrained ordination have gained popularity. These are
18 primarily based on latent variable models, with latent variables estimating the un-
19 derlying, unobserved ecological gradients.

20 3. Despite some major benefits, a major drawback of model-based ordination meth-
21 ods is their speed, as they typically taking much longer to return a result than
22 dissimilarity-based methods, especially for large sample sizes.

23 4. We introduce copula ordination, a new, scalable model-based approach to uncon-
24 strained ordination. This method has all the desirable properties of model-based
25 ordination methods, with the added advantage that it is computationally far more
26 efficient. In particular, simulations show copula ordination is an order of magnitude
27 faster than current model-based methods, and can even be faster than nMDS for
28 large sample sizes, while being able to produce similar ordination plots and trends
29 as these methods.

30 INTRODUCTION

31 Visualisation is a vital part of working with data in ecology, both for exploration, and
32 to support conclusions from statistical analyses. Many studies in ecology however collect
33 multivariate and discrete data, which are difficult to visualise. In this article, we focus on
34 ordination (Zuur *et al.*, 2007; Legendre & Legendre, 2012), a general and very commonly
35 applied class of visualisation methods which aim to collapse or project multivariate data
36 onto a small number (often two or three) dimensions in ways that preserve as much of
37 the underlying structure as possible, such that it can be plotted to look for prevailing
38 patterns. In this paper our focus is on ordination methods for multivariate abundance
39 data e.g. species assemblages, although in principle the ideas discussed here would apply
40 more generally to any multivariate analysis.

41 There are two broad categories of ordination methods in the current literature, tradi-
42 tional dissimilarity-based methods and model-based methods. Dissimilarity-based ordi-
43 nation methods begin by calculating a dissimilarity or distance matrix between sites (or
44 more generally, observational units). Afterwards, these are collapsed to a small number
45 of dimensions for plotting using an algorithm that attempts to preserve and display infor-
46 mation about these relative distances, with the most widely used example of this being
47 non-metric multidimensional scaling (MDS; Kruskal, 1964); see also van der Maaten &
48 Hinton (2008) for a modern alternative. By contrast, model-based ordination methods
49 explicitly model the underlying distribution of the response, using an extension of gener-

50 alised linear mixed models known as latent variable models (or factor analytic models). As
51 the name suggests, these models account for the correlation between taxa via the inclusion
52 of a small number of latent variables, which are assumed to be random effects (Walker
53 & Jackson, 2011; Hui *et al.*, 2015; Warton *et al.*, 2015). As they model the correlation
54 between taxa, latent variables also act as natural ordination axes reflecting unobserved
55 covariates. Both model-based and distance-based methods often give qualitatively simi-
56 lar ordinations e.g., see *Application to data* later on, as well as Hui *et al.* (2015), among
57 others.

58 By directly modelling the data, model-based methods account for both the natural vari-
59 ation and the associated mean-variance relationships present in multivariate abundance
60 data (Warton & Hui, 2017). Furthermore, we can use standard statistical tools to check
61 the assumptions underlying the model, perform model selection, quantify the uncertainty
62 in the estimated correlations between taxa, predict to existing and/or new sites, and incor-
63 porate extensions to (for instance) account for spatio-temporal correlations or imperfect
64 detection (e.g., Warton *et al.*, 2015; Ovaskainen *et al.*, 2017; Hui *et al.*, 2018; Tobler *et al.*,
65 2019). All of this is generally more challenging to accomplish using dissimilarity-based
66 ordination procedures.

67 One major drawback of most currently implemented model-based methods for ordination
68 is that they tend to be considerably slower than dissimilarity-based methods, particularly
69 when sample size (both the number of sites and/or taxa) is small to moderate. While
70 there have been some efforts to overcome this burden (e.g., Niku *et al.*, 2019a; Tikhonov
71 *et al.*, 2020a), computation currently remains a major bottleneck when applying model-
72 based approaches for ordination to large multivariate abundance datasets.

73 The primary reason why computation presents a major challenge for model-based ap-
74 proaches to ordination is that they tend to be built using hierarchical models with latent
75 variables, formulated via a series of conditional distributions, from which it is difficult
76 to calculate the marginal likelihoods or posterior distributions needed for estimation and
77 inference (Walker & Jackson, 2011; Hui *et al.*, 2015; Warton *et al.*, 2015; Tikhonov *et al.*,
78 2020b). To overcome this challenge, we present here an alternative family of multivariate

79 models for ordination based on copulas. Copulas have only recently attracted interest
80 in community ecology (Popovic *et al.*, 2018; Anderson *et al.*, 2019; Popovic *et al.*, 2019).
81 They can be considered as an extension of generalized linear models (GLMs, McCullagh
82 & Nelder, 1989), and hence retain all the desirable statistical properties of model-based
83 methods. Importantly, in contrast to hierarchical models, copulas are based on a *marginal*
84 approach to modelling multivariate abundance data, so the relevant marginal likelihoods
85 are much easier to compute, and the model is much faster to fit (see Fieberg *et al.*, 2009;
86 Muff *et al.*, 2016, for examples of comparisons between conditional versus marginal ap-
87 proaches to modeling). In this paper, we propose a Gaussian copula model with a latent
88 variable structure as a novel, fast model-based ordination method for community ecology.
89 Simulations show our method is an order of magnitude faster than existing model-based
90 methods, while producing largely similar ordination patterns and conclusions.

91 In the remainder of the article, we introduce copulas specifically for discrete data (given
92 the prevalence of non-continuous responses in community ecology), and the proposed
93 Gaussian copula latent variable model for ordination. The proposed method is imple-
94 mented in the R package `ecoCopula`, via the `cord` function. We illustrate an example of
95 copula-based ordination on an example multivariate abundance dataset, and conduct a
96 simulation study to compare `cord` to existing model-based and distance based methods
97 of ordination in terms of both accuracy and speed.

98 **MATERIALS AND METHODS**

99 **Model-based methods for ordination**

100 Multivariate abundance data, also known as multi-species or community composition
101 data, consist of abundance measurements (generally presence/absence, counts, cover or
102 biomass) simultaneously collected for a large number of taxa. Due to the discrete nature
103 of such data, each taxon is typically modelled with a distribution appropriate to its
104 characteristics *via* a generalised linear model (GLM; McCullagh & Nelder, 1989) or some

105 variation thereof e.g., the negative binomial (Caraka *et al.*, 2018) or zero-inflated count
106 distribution for overdispersed species counts, the binomial distribution (Golding *et al.*,
107 2015) for presence-absences, and the Tweedie or Gamma distribution for biomass (Blakey
108 *et al.*, 2016) *etc.*

109 In order to model correlations between the taxa, as is central for ordination plots, it is
110 necessary to construct a multivariate distribution for all taxa jointly. By far the most
111 popular method of achieving this is via a hierarchical modelling framework, where latent
112 variables (assumed to be Gaussian) form the ordination axes; see Appendix for math-
113 ematical details. However, estimation of a hierarchical model involves computationally
114 burdensome steps, e.g., Markov Chain Monte Carlo sampling (in the case of Bayesian
115 methods, Hui, 2021; Tikhonov *et al.*, 2020b), or numerical integration of the marginal
116 likelihood (in the case of likelihood-based methods, Niku *et al.*, 2019b). In this paper, we
117 consider an alternative model-based framework that is much faster to fit to multivariate
118 abundance data – Gaussian copula models.

119 **Gaussian copula ordination**

120 In this section, we introduce copulas as a means of constructing a joint distribution for
121 multivariate abundance data, with a focus on ordination. Copulas are a flexible class of
122 models which allow the modelling of data using any set of marginal distributions, coupled
123 with the covariance structure of any desired multivariate distribution. The term “copula”
124 comes from this idea of coupling a *marginal model* with a *multivariate model* that can
125 handle covariance across response variables. To construct an ordination of multivariate
126 abundances, we will couple marginal GLMs with a *factor analysis* for ordination, which
127 we have implemented as the `cord` function in the `ecoCopula` package.

128 For the marginal model, we need to specify a separate model for each taxon that accounts
129 for key properties of the data. A key feature of abundance data is its mean-variance
130 relationship, we account for this using a GLM with an appropriate choice of distribution.
131 This marginal model allows us to calculate cumulative probabilities for each taxon, we

132 write as $F_j(\cdot)$ the cumulative distribution function (or *CDF*) assumed for the j th taxon.
133 In a Gaussian copula model, abundances are mapped to copula values z that have a
134 normal (or *Gaussian*) distribution, such that we can specify a multivariate model which
135 assumes multivariate normality. For count data, and letting y_{ij} denote the observed count
136 for site i and taxon j , these copula values z_{ij} satisfy

$$F_j(y_{ij} - 1) \leq \Phi(z_{ij}) < F_j(y_{ij}), \quad (1)$$

137 where $F_j(\cdot)$ is the CDF assumed under the marginal GLM for the j th taxon, and $\Phi(\cdot)$ is the
138 CDF of the standard normal. More generally we can replace $F_j(y_{ij} - 1)$ in equation (1) with
139 $F_j(y_{ij}^-)$, the left limit of F_j at y_{ij} . We note that these values are related to Dunn-Smyth
140 residuals (Dunn & Smyth, 1996) which are often used for residual analysis in community
141 ecology. Specifically, these residuals r_{ij} satisfy the equation $\Phi(r_{ij}) = F(y^-) + uf(y)$, where
142 u is a simulated value from the uniform distribution between zero and one, and $f(\cdot)$ is
143 the probability density/mass function corresponding to the CDF $F(\cdot)$.

144 For the multivariate model, and hence ordination, we assume a factor analytic formulation,
145 meaning we assume copula values z_{ij} 's are multivariate normally distributed and satisfy

$$z_{ij} = \boldsymbol{\xi}_i^\top \boldsymbol{\lambda}_j + \epsilon_{ij},$$

146 where $\boldsymbol{\xi}_i$ denote the latent scores for site i and are assumed to be independent normal,
147 $\boldsymbol{\lambda}_j$ are the corresponding factor loadings for taxon j , and ϵ_{ij} are independent Gaussian
148 errors, with variance σ_j^2 for taxon j . We construct an ordination of observations (sites)
149 by plotting the estimated factor scores, and can construct a biplot by overlaying the
150 estimated factor loadings, as is typically done in unconstrained ordination.

151 For unconstrained ordination, $F_j(\cdot)$ is a GLM with only an intercept. To construct a
152 residual ordination (Warton *et al.*, 2015), *i.e.* ordination controlling for some measured
153 variables, we can include these as predictors in the GLM.

154 We estimate the Gaussian copula ordination model by maximum likelihood, which in

155 this case has no closed form due to the inequalities in equation (1). The maximisation
156 employs an efficient iterative procedure where a factor analysis (we used the `factanal`
157 function in R) is applied to iteratively reweighted Dunn-Smyth residuals. This makes the
158 resulting procedure computationally efficient because the factor analysis likelihood has a
159 closed form, and so the process of estimating the multivariate model is fast; see Appendix
160 for details.

161 We conclude this section by highlighting out the two main assumptions required in spec-
162 ifying the Gaussian copula model: (1) the CDF F_j has been specified correctly for all
163 taxa $j = 1, \dots, D$; (2) the unobserved vector characterising the copula, ξ_i , comes from
164 a multivariate Gaussian distribution. We require assumption (1) in order to be able to
165 map from any set of variables to a set of Gaussian variables. But note that just because
166 a set of variables is Gaussian does not necessarily mean they are multivariate Gaussian
167 – hence the need for assumption (2), which essentially requires the further constraints
168 of linearity and equal variance when looking at each copula variable as a function of
169 other responses. Importantly, both these assumptions can be checked, say, by plotting
170 Dunn-Smyth residuals and applying residual diagnostic tools.

171 APPLICATION TO DATA

172 The hunting spider data is a well-known ecological dataset popularised in ter Braak (1986).
173 It consists of counts of hunting spiders caught in pitfall traps, with $D = 12$ species found
174 at $n = 28$ sites (van der Aart & Smeenk-Enserink, 1974). The data also contain six
175 environmental variables thought to be associated with spider abundance, namely: dry
176 soil mass; percent cover of bare sand; percent cover of fallen leaves or twigs; percent cover
177 of moss; percent cover of herb layer and reflection of the soil surface with a cloudless sky.
178 The primary question of interest in the original collection of the data was to identify the
179 main environmental drivers of abundance of the species studied.

180 Copula ordination

181 It is straightforward to carry out model-based ordination on the spider data using the
182 `ecoCopula` package. Figure 1 presents an example of the estimated ordinations (uncon-
183 strained and residual) for these data. There are two key coding steps when constructing
184 an ordination plot using the `ecoCopula` package: 1) fitting separate marginal models to
185 each spider species (to obtain the CDF F_j), for which we provide the function `stackedsdm`;
186 2) then applying the function `cord`, which takes the output from the marginal model and
187 returns the ordination to be plotted, assuming a Gaussian copula model.

188 In the first step, the `stackedsdm` function takes as input the matrix of multivariate abun-
189 dances, a formula (which defaults to an intercept-only model, for an unconstrained ordi-
190 nation), and optionally a data frame storing covariates. Because it was written specifi-
191 cally for abundance data, it defaults to negative binomial regression if a family argument
192 has not been entered, as does `manyglm` from the `mvabund` package (Wang *et al.*, 2020),
193 which could also be used with `cord`. The two main distinctions from `manyglm` are that
194 `stackedsdm`: can take vector family input, which is useful if different response variables
195 were collected in different ways; and it can take advantage of parallel processing to speed
196 up fitting of the marginal models for very large datasets.

197 In the second step, the `cord` function applies the `factanal` function to iteratively reweighted
198 Dunn-Smyth residuals from the marginal model fit. By default two factors are fitted, but
199 other options can be specified via the `nlv` argument. The `cord` function was written so
200 it could be applied to marginal model fits constructed using `stackedsdm` or `manyglm`. In
201 principle, other functions could be used to fit the marginal models too, as long as they
202 have a `residuals` function and store fitted values as a matrix response.

203 Competing ordination methods

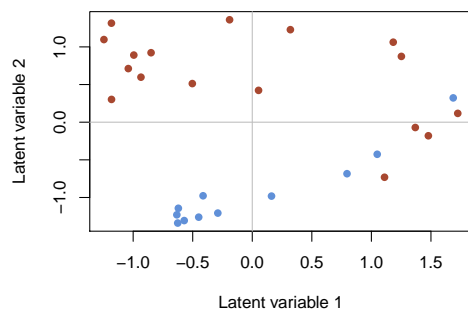
204 There are several competing methods commonly used for ordination in community ecology.
205 Hierarchical model-based ordination methods e.g., R packages `Hmsc` (Tikhonov *et al.*,
206 2020c), `boral`, and `gllvm` use latent variables to define the ordination axes. For ordination

Does community abundance change with the cover of fallen leaves and twigs?

```
library(ecoCopula)
cols=ifelse(spider$x["fallen.leaves"] > 0, "#6090d8", "#a84830")
```

Unconstrained ordination

```
spider$abund %>%
  stackedsdm(~ 1,
             data = spider$x) %>%
  cord() %>%
  plot(site.col=cols)
```



Residual ordination

```
spider$abund %>%
  stackedsdm(~ soil.dry + reflection,
             data = spider$x) %>%
  cord() %>%
  plot(site.col=cols)
```

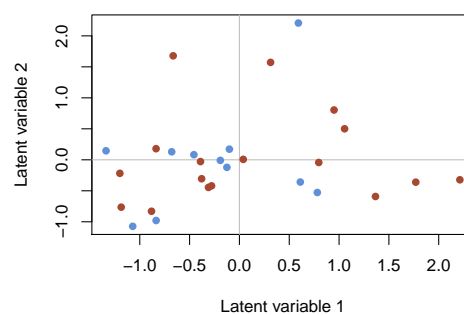


Figure 1: An illustration of estimated ordinations for the hunting spider data using the `ecoCopula` package. Sites are coloured according to presence (blue) or absence (red) of fallen leaves or twigs. In the unconstrained ordination (left), we see a clear separation of the sites based on this covariate, with sites with fallen leaves or twigs having predominantly negative values on the second ordination axis. On the right we have performed a residual ordination, controlling for moisture and reflectivity of the soil. The sites with and without fallen leaves and twigs are no longer well separated, indicating that moisture and reflectivity of the soil has explained some of the effect of fallen leaves and twigs on spider abundances.

207 purposes, these all effectively fit the same type of model but by slightly different means
208 of estimation. Not surprisingly then, speed is arguably the main distinguishing feature
209 between these approaches, with the Bayesian methods being slower than their likelihood-
210 based counterparts. From the literature and our experience, we found that `gllvm` was
211 the fastest existing package that implemented a hierarchical approach to ordination. For
212 example, carrying out a model-based unconstrained ordination on the spider data with
213 `gllvm` takes 1.06 seconds, compared to 25.82 seconds for `boral` and 79.48 seconds for
214 `Hmsc` (each with 10000 MCMC samples).

215 There are many ordination methods in common use in community ecology that do not
216 use a model-based framework, with nMDS being the most popular at the time of writing.
217 We include this in our comparison (using the `metaMDS` function in the `vegan` package,
218 applied after fourth root transformation using the Bray-Curtis dissimilarity). In addition,
219 we considered another popular approach in detrended correspondence analysis (DCA;
220 Hill & Gauch, 1980), as implemented using the `decorana` function in the `vegan` package
221 (Oksanen *et al.*, 2019).

222 Applying the range of ordination methods discussed above to the spider data, we see
223 that they give qualitatively similar results (e.g., in Fig 2), despite using quite different
224 formulations.

225 **SIMULATION STUDY**

226 To quantitatively compare the Gaussian copula ordination method to existing methods
227 in terms of speed and accuracy, we conducted a simulation study using the spider data
228 as the basis for simulating multivariate abundance data.

229 **Simulation design**

230 For assessing ordination score recovery, we simulated datasets that mimicked properties
231 of the hunting spider data, using one of three approaches: 1) the `simulate` function in the

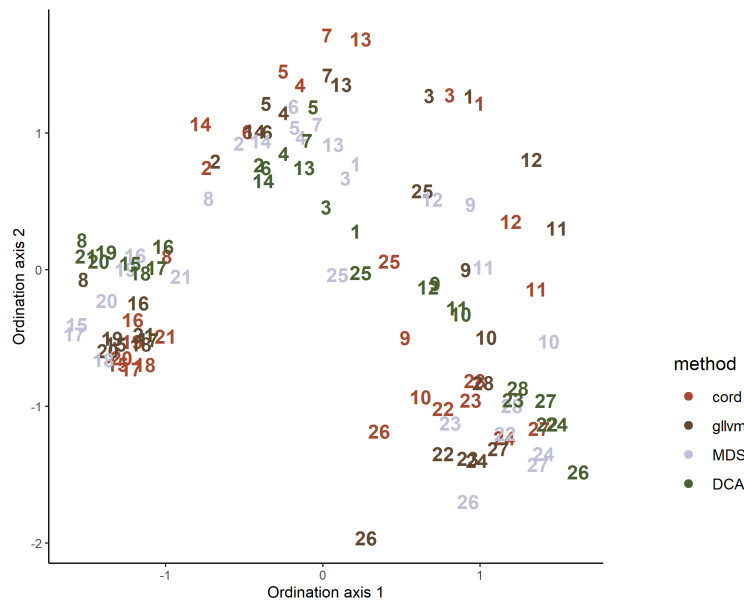


Figure 2: Unconstrained ordination scores estimated on the `cord` function in the `ecoCopula` package (which uses the proposed Gaussian copula method, red), the `gllvm` package (which is computationally the most efficient package for implementing ordination via hierarchical modeling currently available, brown), nMDS (based on applying the fourth root transformation and the Bray-Curtis dissimilarity, grey), and DCA (green). All scores have been rotated and scaled using `vegan::procrustes`. Qualitatively, it can be seen that these ordinations are fairly similar, despite being estimate using very different approaches.

232 `gllvm` package, which simulates abundances from a latent variable model i.e, a hierarchical
233 modelling approach; 2) the `simulate` function for `cord` models in the `ecoCopula` package,
234 which is based on our proposed Gaussian copula model, and; 3) `compas` (Minchin, 1987),
235 which simulates multivariate count data based on generating unimodal curves representing
236 species responses to one or more indirect gradients. To simulate from `gllvm` and `cord`, we
237 fit the corresponding model to the spider data without covariates, and then simulate from
238 the resulting model fit treating the estimated parameters and predicted latent variables
239 as the true values. To simulate using `compas`, we use the scores obtained from applying
240 nMDS to the spider data (assuming a fourth root transformation and Bray-Curtis dissim-
241 ilarity), and then apply the `compas` function from the `CommEcol` package (Melo, 2019),
242 and using default values for all arguments as appropriate. For each case, we simulated
243 100 multivariate count datasets.

244 For each simulated dataset, we then compared the performance of the estimated ordination
245 scores using four competing methods: `gllvm`, `cord`, nMDS, and DCA. Performance was

246 assessed using the Procrustes distance between the true and estimated ordination scores.
247 Note we did not consider other packages such as `Hmsc` or `boral`, as we found they tended
248 to produce similar answers to `gllvm`, but, as is the case above in *Competing ordination*
249 *methods* took considerably longer to produce the ordination.

250 Next, as an assessment of computation time between the methods, we simulated new
251 multivariate count datasets with varying numbers of sites and/or taxa. We then fitted the
252 same four ordination methods as above to obtain ordination scores, tracking computation
253 time in seconds required for each method. For each combination of number of sites n and
254 number of taxa D , we simulated 20 datasets, and averaged computation time across the
255 20 fits for each method. Note that because computation speed rather than ordination
256 accuracy was of interest at this point, multivariate abundance data were simulated from
257 the `gllvm` simulation model only.

258 **Results**

259 **Ordination score recovery**

260 As expected, the precise performance of each method, while qualitatively similar (Fig 2),
261 depended strongly on how the data were generated. When data were simulated using
262 `gllvm` i.e., a latent variable model, then `gllvm` performed best at recovering ordination
263 scores, followed by the proposed Gaussian copula method using `cord`. Similarly, when the
264 data were simulated from a Gaussian copula model, `cord` performed best at recovering
265 the scores followed by the hierarchical modelling approach using `gllvm`. When data were
266 simulated using `compas`, none of the methods stood out in terms of ordination recovery.
267 The two model-based approaches to ordination clearly outperformed nMDS and DCA
268 when the data generation process was model-based (even if it was the incorrect model),
269 but when the data generation process was more “neutral” i.e., `compas`, there was no clear
270 difference between the four approaches.

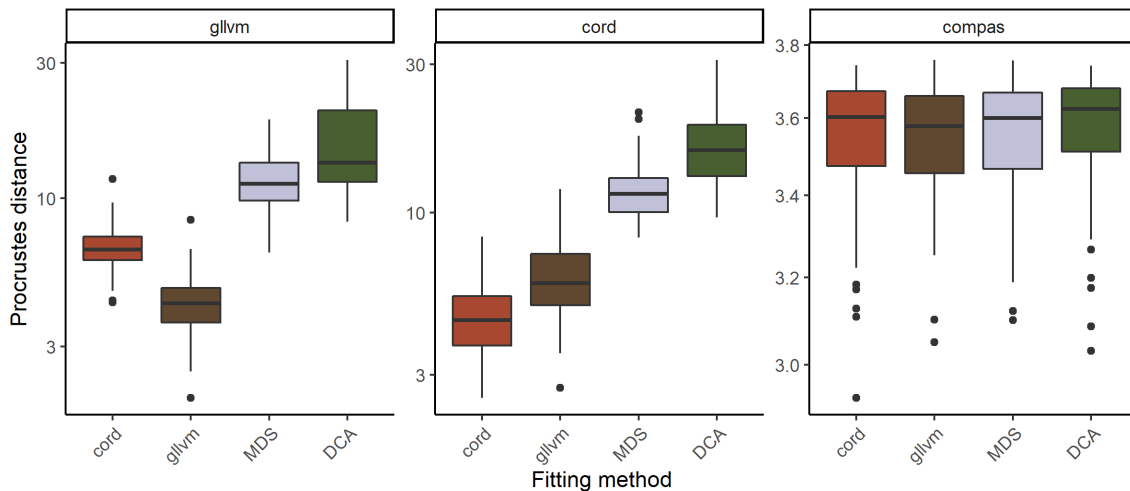


Figure 3: Comparative boxplots assessing ordination score recovery of the four ordination methods fitted to data generated from a hierarchical modelling approach (`gllvm`), the proposed Gaussian copula model (`cord`), and a more “neutral” approach (`compas`). Results showed that `gllvm` and `cord` performed best when data were generated from their respective models, while for `compas` all four methods performed similarly.

271 Computation speed

272 In all simulations, the proposed Gaussian copula method i.e., `cord`, was noticeably faster
273 than `gllvm` (Fig 4), and indeed we anticipate that as an alternative approach to or-
274 dination it will likely be faster than any other currently available software that uses a
275 hierarchical/latent variable approach for ordination. For a small number of sites n , `cord`
276 was slower than the dissimilarity-based methods we considered. Interestingly though,
277 `cord` scaled better with the number of sites n compared to other methods, and for large
278 enough n e.g., approximately $n \geq 200$ it was faster than nMDS (under its default settings
279 on `vegan`, Fig 4). Dissimilarity-based methods were consistently the fastest methods when
280 the number of taxa D became large.

281 DISCUSSION

282 In this article, we have proposed a computationally efficient approach to model-based
283 ordination which performs well at recovering ordination scores in community ecology, and
284 is faster and more scalable than all existing model-based ordination methods in the litera-
285 ture. Interestingly, it was even faster than the popular dissimilarity-based method nMDS

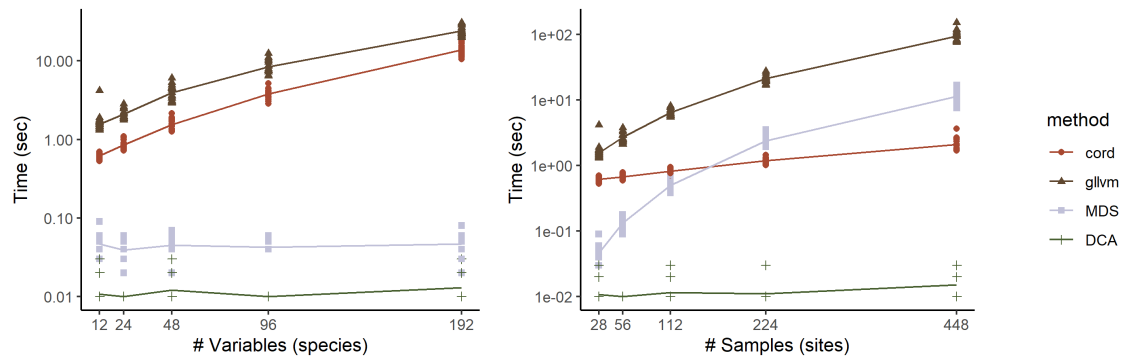


Figure 4: Comparison of computation times when constructing an ordination using different approaches. The proposed Gaussian copula methods `cord` was noticeably faster than the fastest model-based competitor `gllvm`, and also faster than nMDS for large values of n . Dissimilarity-based methods were faster when the number of taxa D was large and/or the number of sites n was small.

286 for datasets with a moderate number of taxa and a moderate to large number of sites.
287 Two other key features of the proposed approach include the straightforward capability
288 to consider a large range of data distributions for the responses, to have different distri-
289 butions for different taxa i.e., multivariate abundance data with mixed responses, and the
290 possibility of residual ordination, *i.e.* controlling for environmental factors in the ordina-
291 tion. The proposed Gaussian copula method is implemented in the R package `ecoCopula`,
292 which is now available on CRAN. The package also contains a `simulate` function, which
293 allows users to simulate data with similar characteristics to the observed data assum-
294 ing a Gaussian copula model; this is useful for sample size calculations for study design,
295 diagnostics, and visual inference, among other potential applications.

296 The speed of the copula approach to ordination, relative to hierarchical models, is driven
297 by the use of marginal models, and the choice of a Gaussian distribution for the multi-
298 variate model. This meant that the latent variable model was estimated *via* a Gaussian
299 likelihood function, which is generally well-behaved and quick to compute (largely be-
300 cause for Gaussian likelihoods, we can work directly with the sample covariance matrix
301 rather than the full dataset). While we only explored this here in the context of or-
302 dination, we expect this to be generally true across other covariance models (such as
303 graphical modelling; Popovic *et al.*, 2019), and anticipate a wealth of potential extensions
304 and applications of Gaussian copula models in ecology.

305 For example, further ordination functionality could include an extension of redundancy
306 analysis to handle data from any marginal distribution, since this can be understood as
307 a form of reduced rank regression on multivariate Gaussian data (Bach & Jordan, 2005;
308 Stoica & Viberg, 1996), and therefore can be applied to the estimated copula model using
309 a similar estimation algorithm to that used in `ecoCopula` (Popovic *et al.*, 2018). Fast
310 model-based ordination can also be combined with high-dimensional data visualisation
311 and visual inference tools such as those implemented in the `tourr` package (Wickham
312 *et al.*, 2011) to visualise three or more latent variables, and we are currently undertaking
313 research on this front. In addition, marginal models can be extended in a variety of
314 ways, from modelling non-linear relationships using generalised additive models (Wood,
315 2017), to accounting for detection (Tobler *et al.*, 2019), other response types (e.g., plant
316 cover data, Damgaard *et al.*, 2020), and explaining changes in abundance with traits
317 via a fourth-corner model (Brown *et al.*, 2014). Spatially and/or temporally indexed
318 multivariate abundance data can be modelled using more structured latent scores in the
319 Gaussian copula model (*e.g.* Thorson *et al.*, 2015, 2016); this can be interpreted as a set
320 of unobserved and spatially smoothed environmental predictions or ordination axes.

321 It should be noted that, along of the lines of what has been discussed in the article, han-
322 dling such spatial and/or temporal multivariate abundance data though a copula model is
323 expected to be computationally far more efficient and scalable than hierarchical or latent
324 variable approaches (Hui *et al.*, In press), since it can be straightforwardly done through
325 the covariance of z_{ij} 's while maintaining the multivariate normality assumption. Finally,
326 copula models can be used as a basis for performing likelihood-based inference about
327 multivariate abundance data, for example to test for effects of a one or more treatment
328 or environmental covariates. In summary, the copula-based framework has considerable
329 potential in community ecology, due to both its flexibility and its advantages in computa-
330 tional efficiency and scalability, and we look forward to seeing how this literature develops
331 in the future.

332 **ACKNOWLEDGEMENTS**

333 GCP was supported by the Australia Postgraduate Award and ARC Discovery Project
334 project (DP180103543) awarded to DIW. FKCH was supported an ARC Discovery Early
335 Career Research Fellowship (DE200100435).

336 **AUTHOR CONTRIBUTIONS**

337 GCP led the writing and analysis. All authors conceived and developed ideas, contributed
338 critically to the manuscript editing, and gave final approval for publication.

339 **DATA ACCESSIBILITY**

340 The hunting spider data are available as part of the `ecoCopula` package.

341 APPENDIX

342 Ordination with hierarchical models

343 Let \mathbf{y} be the response matrix with elements y_{ij} for site $i = 1, 2, \dots, N$ and taxon $j =$
344 $1, 2, \dots, D$. Then a typical hierarchical modelling approach to ordination models the
345 conditional mean of each element, denoted as μ_{ij} , via an extension of GLMs, where a
346 small number $d \ll D$ of latent variables $\boldsymbol{\xi}_i$ and an (optional) error ϵ_{ij} are included on the
347 linear predictor scale,

$$g(\nu_{ij}) = \eta_{ij} = \mathbf{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{\xi}_i^\top \boldsymbol{\lambda}_j + \epsilon_{ij},$$

348 where \mathbf{x}_i is a vector of observed predictors for site i , $\boldsymbol{\beta}_j$ are the corresponding taxon-
349 specific coefficients, $\boldsymbol{\xi}_i$ denote the latent variables, $\boldsymbol{\lambda}_j$ are the corresponding factor load-
350 ings. The quantity $g(\cdot)$ denotes some known link function relating the mean to the linear
351 predictor. If included, the error terms are assumed to be independent across responses,
352 such that the variance-covariance matrix $\boldsymbol{\Psi}$ of the vector $\boldsymbol{\epsilon}_i = (\epsilon_{i1}, \dots, \epsilon_{iD})$ is diagonal.
353 On the linear predictor scale, the covariance and hence correlation between the taxa is
354 driven by the factor loadings. In particular, if we define $\boldsymbol{\eta}_i = (\eta_{i1}, \dots, \eta_{iD})$, then the
355 variance-covariance at site i is then given by $\text{Cov}(\boldsymbol{\eta}_i) = \boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi}$, where $\boldsymbol{\Lambda}$ is the
356 loadings matrix formed by stacking the $\boldsymbol{\lambda}_j$'s as row vectors. The hierarchical nature of
357 the latent variable model derives from the fact that it can be formulated conditionally.
358 That is, we first assume a distribution for the latent variables e.g., the elements of \mathbf{x}_i are
359 assumed to follow a standard Gaussian distribution. Then, conditionally on the latent
360 variables, we assume a distribution for the responses as given by the GLM framework,
361 where the mean of the response is related to the measured covariates, latent variables,
362 and errors (if appropriate) as defined above.

363 With regards to ordination, from the above specification we see that the latent variables
364 can be understood as modelling unmeasured environmental gradients, while the factor
365 loadings quantify the response of each taxon to each of these underlying gradients (similar
366 to regression coefficients for a measured covariate). If the number of latent variables d is

367 small (typically $d = 2$ or 3), then the estimated loadings and latent scores can be plotted
368 separately, or can be combined into a *biplot* to visualise relationships between different
369 sites, and between different taxa and between taxa and sites.

370 Depending on what additional information is available, there are a wealth of extensions
371 that could be made to the hierarchical latent variable framework to further the ordi-
372 nation analysis e.g., the latent variables could be spatially and/or temporally varying
373 (Thorson *et al.*, 2015), and/or included at different scales (Björk *et al.*, 2018). However
374 for simplicity, we will focus on the specification as given above, noting that the proposed
375 approach to model-based ordination via copulas detailed below may also be extended in
376 such directions as avenues of future research.

377 **Copulas for multivariate abundance data**

378 **Copulas**

379 Copula modelling is rooted in Sklar’s Theorem (Sklar, 1959), which shows that any multi-
380 variate distribution can be written as a copula model, as long as the marginal and copula
381 distributions are appropriately chosen. From a modelling perspective, the main appeal of
382 copulas is that the covariance structure can be specified independently of the marginal
383 distributions, making them very flexible. In addition, their marginal specification makes
384 them computationally efficient, fast to fit, and readily parallelizable.

385 In more detail, copulas can be viewed as providing a means of starting from a common
386 non-Gaussian data distributions (Poisson, negative binomial, Tweedie, *etc*), and then ex-
387 tending them to multiple, *correlated* responses in a flexible manner. The idea is to map
388 from any data distribution to a convenient copula family that facilitates multivariate mod-
389 elling, most commonly, the multivariate Gaussian distribution as is done in our main text
390 with Gaussian copula method. This mapping occurs using the probability integral trans-
391 form (PIT), a basic property of continuous probability distributions which allows observed
392 data values y to be converted to uniformly distributed values between zero and one via its
393 corresponding cumulative distribution function (CDF), $F(y)$. The PIT hence provides a

394 way to transform between any two continuous distributions, in our case, between the de-
395 sired data distribution and the copula distribution, by going through the standard uniform
396 distribution itself. The latter copula distribution is where the multivariate properties of
397 the data, including covariances and correlations, are straightforwardly formulated and
398 estimated. We highlight that while there are a wide variety of possible choices of copula
399 distributions (see Nelsen, 2007, for an introduction to copulas), in our case we will focus
400 on the using the multivariate Gaussian distribution as the copula distribution, given its
401 simplicity and relative easy of interpretation in the context of community ecology.

402 **Discrete Gaussian copulas**

403 To construct a copula model, we require a set of marginal distributions, one for each
404 taxon, along with a multivariate (copula) distribution to model the correlations between
405 taxa. As ecologists often have discrete data, the marginal distributions $F_j(\cdot)$ are, like
406 in the hierarchical case, chosen to account for key properties of the data e.g., using a
407 negative binomial or zero-inflated count GLM for overdispersed counts. Note however
408 that when the marginal distributions are discrete, the PIT properties do not strictly hold,
409 in the sense that the transformation is no longer a unique mapping. Nevertheless, we
410 can still use this concept to derive a discrete Gaussian copula with marginal distributions
411 $F_j; j = 1, \dots, D$ as a function an unobserved multivariate Gaussian variable.

412 To understand how this discrete Gaussian copula can be formulated, it is useful to first
413 recap the idea of a randomized quantile or Dunn-Smyth residual (Dunn & Smyth, 1996)
414 as discussed in the main text. The Dunn-Smyth residual is illustrated in the left panel of
415 Fig 5. Specifically, instead of calculating $F(y)$ as in the PIT, we first simulate a random
416 value between $F(y^-)$ and $F(y)$ where y^- generically denotes the previous observable value
417 of a response y e.g., for counts these must occur at integer values, and then transform
418 this value using the inverse CDF of a standard normal random variable. Mathematically,
419 this can be written as $r_{ij} = \Phi^{-1}[F(y_{ij}^-) + uf(y_{ij})]$, where u is a simulated value from the
420 uniform distribution between zero and one and $f(\cdot)$ is the corresponding density function.
421 For discrete data y , it can be subsequently shown that the Dunn-Smyth residual r_{ij} follows

422 a standard univariate Gaussian distribution, provided the CDF $F(y)$ has been specified
 423 correctly. Dunn-Smyth residuals have been growing in popularity in ecology (Warton
 424 *et al.*, 2017) and are used in popular packages such as `mvabund` (Wang *et al.*, 2020), `boral`
 425 (Hui, 2021), `gllvm` (Niku *et al.*, 2019b) and `DHARMa` (Hartig, 2020), largely for the purposes
 426 of residual analysis. Here however, we use the the idea of a Dunn-Smyth residual to map
 427 from a discrete distribution (Fig 5, left panel), through the uniform distribution (centre
 428 panel), to the copula distribution of our choice (right panel). We note here that for each
 429 value of y , there are many possible values of r , however for each r , there is only one value
 430 of the response y that could have produced it. This is a consequence of the fact that, as
 431 remarked previously, the PIT is no longer a unique mapping for discrete data.

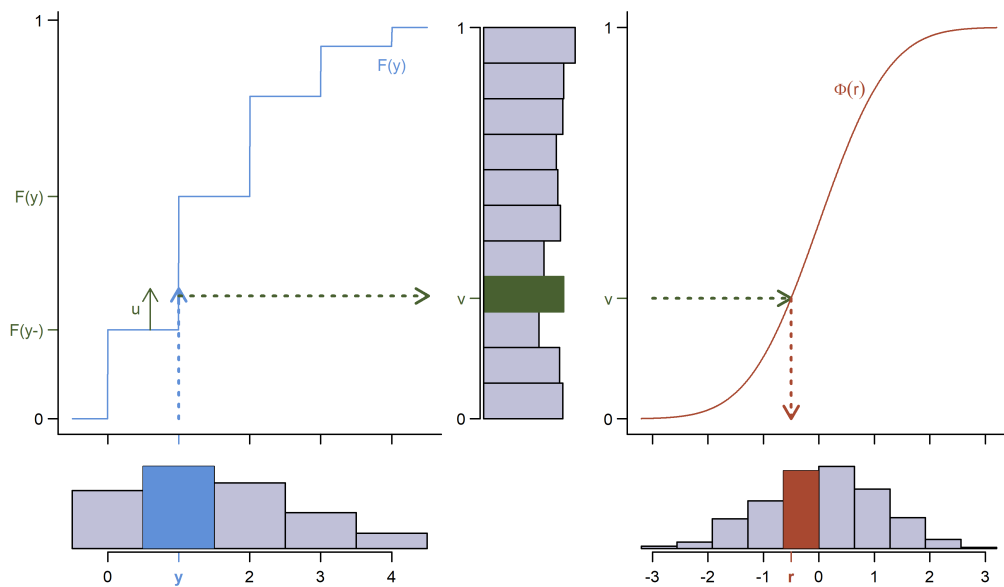


Figure 5: To generate Dunn-Smyth residuals r for a discrete observation y with CDF $F(y)$ (left panel), we first simulate a uniform value u , (green) between $F(y^-)$ and $F(y)$. The resulting $v = F(y^-) + u$ is uniform between zero and one (centre panel). We then transform v through the inverse of the Gaussian CDF, $r = \Phi^{-1}(v)$.

432 Based on the idea of Dunn-Smyth residuals, we now discuss how to fit a Gaussian copula
 433 model to discrete data. For site $i = 1, \dots, n$, let $\mathbf{z}_i = (z_{i1}, z_{i2}, \dots, z_{iD})$ denote a D -
 434 dimensional unobserved multivariate vector, arising from the copula distribution $f(\mathbf{z}_i)$.
 435 Then the conditional distribution of y_{ij} given \mathbf{z}_i is equal to one when $\Phi^{-1}(F_j(y_{ij}^-)) \leq$

436 $z_{ij} < \Phi^{-1}(F_j(y_{ij}))$, and zero otherwise. In other words, if we know z_{ij} then we also know
 437 y_{ij} , as it must be equal to whichever the value of the response that could have produced
 438 z_{ij} , analogous to the idea of Dunn-Smyth residuals above. More formally,

$$f(y_{ij}|\mathbf{z}_i) = \mathbb{1}(\Phi^{-1}(F_j(y_{ij}^-)) \leq z_{ij} < \Phi^{-1}(F_j(y_{ij}))),$$

439 where $\mathbb{1}(\cdot)$ denotes the indicator function. To find the marginal distribution of \mathbf{y}_i , we first
 440 write the joint distribution of \mathbf{y}_i and \mathbf{z}_i ,

$$f(\mathbf{y}_i, \mathbf{z}_i) = f(\mathbf{y}_i|\mathbf{z}_i)f(\mathbf{z}_i) = \prod_{j=1}^D \mathbb{1}(\Phi^{-1}(F_j(y_{ij}^-)) \leq z_{ij} < \Phi^{-1}(F_j(y_{ij})))f(\mathbf{z}_i),$$

441 and then integrate out the unobserved vector \mathbf{z}_i , as follows

$$\begin{aligned} f(\mathbf{y}_i) &= \int f(\mathbf{y}_i, \mathbf{z}_i)d\mathbf{z}_i = \int \prod_{j=1}^D \mathbb{1}(\Phi^{-1}(F_j(y_{ij}^-)) \leq \xi_{ij} < \Phi^{-1}(F_j(y_{ij})))f(\boldsymbol{\xi}_i)d\mathbf{z}_i \\ &= \int_A f(\mathbf{z}_i)d\mathbf{z}_i, \end{aligned}$$

442 where $A_i = \cap_j [\Phi^{-1}(F_j(y_{ij}^-)), \Phi^{-1}(F_j(y_{ij}))]$. That is, the marginal likelihood of the data
 443 at site i is constructed by integrating the copula distribution $f(\mathbf{z})$ over the region A , a
 444 hypercube with boundaries defined by the PIT.

445 By choosing $f(\mathbf{z}_i)$ to follow a multivariate Gaussian distribution, the above marginal
 446 likelihood then formally defines the Gaussian copula model. In practice, to fit such a
 447 model we want to first develop an approach by which to estimate this likelihood function
 448 (given data collected at a set of n sites), and then proceed find the values of model
 449 parameters (e.g., coefficients $\boldsymbol{\beta}_j$, loadings $\boldsymbol{\lambda}_j$, $\boldsymbol{\Psi}$ and so on) that maximise $f(\mathbf{y}_i)$. The main
 450 challenge lies in the former i.e., to approximate the integral. Perhaps not surprisingly,
 451 it turns out that we can use the idea of Dunn-Smyth residuals to estimate $f(\mathbf{y}_i)$ via
 452 importance sampling Popovic *et al.* (see 2018, for details). Essentially, we construct
 453 many sets of Dunn-Smyth residuals, and then take a weighted average of the likelihood
 454 evaluated at these residual values, weighted by how well each residual fits the assumed
 455 multivariate Gaussian copula distribution. Critically, the resulting approximate likelihood

456 then has the form of a weighted sum of Gaussian likelihood terms, so maximum likelihood
457 estimation of (all the parameters in) the copula model becomes relatively straightforward
458 by adopting standard algorithms for estimating Gaussian distributed data ; we refer the
459 reader to Popovic *et al.* (2018) for the mathematical details underlying this, noting that
460 we have implemented this in the R package `ecoCopula`. As mentioned previously, one
461 major advantage of this estimation procedure, and of the copula model as a whole, is
462 its computationally efficiency and scalability relative to estimating hierarchical models.
463 This largely results from the fact that, by effectively transforming the problem from
464 working with data from an complex high-dimensional multivariate discrete distribution
465 to working with Dunn-Smyth residuals from a multivariate Gaussian distribution, we can
466 use the many well-known properties and tools of the latter to perform likelihood-based
467 estimation and inference. For example, we can use the sample covariance matrix (noting
468 sampling from a multivariate Gaussian distribution is relatively straightforward) as a
469 sufficient statistic for the correlation matrix of the copula distribution. In turn, we avoid
470 the high-dimensional integrals inherent in estimating hierarchical or conditional models,
471 instead replacing it with the need to sample from a tractable copula distribution, which
472 in the case of a Gaussian copula model can be very efficiently.

473 References

- 474 Anderson, M. J., de Valpine, P., Punnett, A. & Miller, A. E. (2019). A pathway for
475 multivariate analysis of ecological communities using copulas. *Ecology and Evolution*,
476 **9**, 3276–3294.
- 477 Bach, F. R. & Jordan, M. I. (2005). *A Probabilistic Interpretation of Canonical Correlation*
478 *Analysis*. Tech. Rep. 688 , Department of Statistics, University of California, Berkeley.
- 479 Björk, J. R., Hui, F. K. C., O’Hara, R. B. & Montoya, J. M. (2018). Uncovering the
480 drivers of host-associated microbiota with joint species distribution modelling. *Molec-*
481 *ular ecology*, **27**, 2714–2724.
- 482 Blakey, R. V., Law, B. S., Kingsford, R. T., Stoklosa, J., Tap, P. & Williamson, K. (2016).
483 Bat communities respond positively to large-scale thinning of forest regrowth. *Journal*
484 *of Applied Ecology*, **53**, 1694–1703.
- 485 Brown, A. M., Warton, D. I., Andrew, N. R., Binns, M., Cassis, G. & Gibb, H. (2014).
486 The fourth-corner solution—using predictive models to understand how species traits
487 interact with the environment. *Methods in Ecology and Evolution*, **5**, 344–352.
- 488 Caraka, R. E., Shohaimi, S., Kurniawan, I. D., Herliansyah, R., Budiarto, A., Sari, S. P. &
489 Pardamean, B. (2018). Ecological show cave and wild cave: negative binomial gllvm’s
490 arthropod community modelling. *Procedia Computer Science*, **135**, 377–384.
- 491 Damgaard, C., Hansen, R. R. & Hui, F. K. (2020). Model-based ordination of pin-point
492 cover data: Effect of management on dry heathland. *Ecological Informatics*, **60**, 101155.
- 493 Dunn, P. K. & Smyth, G. K. (1996). Randomized quantile residuals. *Journal of Compu-*
494 *tational and Graphical Statistics*, **5**, 236–244.
- 495 Fieberg, J., Rieger, R. H., Zicus, M. C. & Schildcrout, J. S. (2009). Regression modelling of
496 correlated data in ecology: subject-specific and population averaged response patterns.
497 *Journal of Applied Ecology*, **46**, 1018–1025.

- 498 Golding, N., Nunn, M. A. & Purse, B. V. (2015). Identifying biotic interactions which
499 drive the spatial distribution of a mosquito community. *Parasites & vectors*, **8**, 367.
- 500 Hartig, F. (2020). *DHARMA: Residual Diagnostics for Hierarchical (Multi-Level / Mixed)*
501 *Regression Models*.
- 502 Hill, M. O. & Gauch, H. G. (1980). Detrended correspondence analysis: an improved
503 ordination technique. *Classification and ordination*, pp. 47–58. Springer.
- 504 Hui, F. K. C. (2021). *boral: Bayesian Ordination and Regression AnaLysis*. R package
505 version 2.0.
- 506 Hui, F. K. C., Hill, N. A. & Welsh, A. H. (In press). Assuming independence in spatial
507 latent variable models: Consequences and implications of misspecification. *Biometrics*.
- 508 Hui, F. K. C., Sara, T., Pledger, S., Foster, S. D. & Warton, D. I. (2015). Model-based
509 approaches to unconstrained ordination. *Methods in Ecology and Evolution*, **6**, 399–411.
- 510 Hui, F. K. C., Tanaka, E. & Warton, D. I. (2018). Order selection and sparsity in latent
511 variable models via the ordered factor LASSO. *Biometrics*, **74**, 1311–1319.
- 512 Kruskal, J. B. (1964). Multidimensional scaling by optimizing goodness of fit to a non-
513 metric hypothesis. *Psychometrika*, **29**, 1–27.
- 514 Legendre, P. & Legendre, L. (2012). *Numerical ecology*. Elsevier.
- 515 McCullagh, P. & Nelder, J. A. (1989). *Generalized Linear Models, Volume 37*. CRC
516 Press.
- 517 Melo, A. S. (2019). *CommEcol: Community Ecology Analyses*.
- 518 Minchin, P. R. (1987). Simulation of multidimensional community patterns: towards a
519 comprehensive model. *Vegetatio*, **71**, 145–156.
- 520 Muff, S., Held, L. & Keller, L. F. (2016). Marginal or conditional regression models for
521 correlated non-normal data? *Methods in Ecology and Evolution*, **7**, 1514–1524.
- 522 Nelsen, R. B. (2007). *An introduction to copulas*. Springer Science & Business Media.

- 523 Niku, J., Brooks, W., Herliansyah, R., Hui, F. K. C., Sara, T. & Warton, D. I. (2019a).
524 Efficient estimation of generalized linear latent variable models. *PloS one*, **14**, e0216129.
- 525 Niku, J., Hui, F. K. C., Taskinen, S. & Warton, D. I. (2019b). gllvm: Fast analysis
526 of multivariate abundance data with generalized linear latent variable models in R.
527 *Methods in Ecology and Evolution*, **10**, 2173–2182.
- 528 Oksanen, J., Blanchet, F. G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., Minchin,
529 P. R., O’Hara, R. B., Simpson, G. L., Solymos, P., Stevens, M. H. H., Szoecs, E. &
530 Wagner, H. (2019). *vegan: Community Ecology Package*.
- 531 Ovaskainen, O., Tikhonov, G., Norberg, A., Guillaume Blanchet, F., Duan, L., Dunson,
532 D., Roslin, T. & Abrego, N. (2017). How to make more out of community data? A
533 conceptual framework and its implementation as models and software. *Ecology letters*,
534 **20**, 561–576.
- 535 Popovic, G. C., Hui, F. K. C. & Warton, D. I. (2018). A general algorithm for covariance
536 modeling of discrete data. *Journal of Multivariate Analysis*, **165**, 86 – 100.
- 537 Popovic, G. C., Warton, D. I., Thomson, F. J., Hui, F. K. C. & Moles, A. T. (2019). Un-
538 tangling direct species associations from indirect mediator species effects with graphical
539 models. *Methods in Ecology and Evolution*, **10**, 1571–1583.
- 540 Sklar, M. (1959). *Fonctions de répartition à n dimensions et leurs marges*. Université
541 Paris.
- 542 Stoica, P. & Viberg, M. (1996). Maximum likelihood parameter and rank estimation in
543 reduced-rank multivariate linear regressions. *IEEE Transactions on Signal Processing*,
544 **44**, 3069–3078.
- 545 ter Braak, C. J. F. (1986). Canonical Correspondence Analysis: A New Eigenvector
546 Technique for Multivariate Direct Gradient Analysis. *Ecology*, **67**, 1167–1179.
- 547 Thorson, J. T., Ianelli, J. N., Larsen, E. A., Ries, L., Scheuerell, M. D., Szuwalski, C. &
548 Zipkin, E. F. (2016). Joint dynamic species distribution models: a tool for community

- 549 ordination and spatio-temporal monitoring. *Global Ecology and Biogeography*, **25**, 1144–
550 1158.
- 551 Thorson, J. T., Scheuerell, M. D., Shelton, A. O., See, K. E., Skaug, H. J. & Kristensen,
552 K. (2015). Spatial factor analysis: a new tool for estimating joint species distributions
553 and correlations in species range. *Methods in Ecology and Evolution*, **6**, 627–637.
- 554 Tikhonov, G., Duan, L., Abrego, N., Newell, G., White, M., Dunson, D. & Ovaskainen,
555 O. (2020a). Computationally efficient joint species distribution modeling of big spatial
556 data. *Ecology*, **101**, e02929.
- 557 Tikhonov, G., Opedal, Ø. H., Abrego, N., Lehikoinen, A., de Jonge, M. M. J., Oksanen,
558 J. & Ovaskainen, O. (2020b). Joint species distribution modelling with the R-package
559 Hmsc. *Methods in Ecology and Evolution*, **11**, 442–447.
- 560 Tikhonov, G., Ovaskainen, O., Oksanen, J., de Jonge, M., Opedal, O. & Dallas, T.
561 (2020c). *Hmsc: Hierarchical Model of Species Communities*.
- 562 Tobler, M. W., Kéry, M., Hui, F. K. C., Guillera-Arroita, G., Knaus, P. & Sattler,
563 T. (2019). Joint species distribution models with species correlations and imperfect
564 detection. *Ecology*, **100**, e02754.
- 565 van der Aart, P. & Smeenk-Enserink, N. (1974). Correlations Between Distributions of
566 Hunting Spiders (Lycosidae, Ctenidae) and Environmental Characteristics in a Dune
567 Area. *Netherlands Journal of Zoology*, **25**, 1–45.
- 568 van der Maaten, L. & Hinton, G. (2008). Visualizing Data using t-SNE. *Journal of*
569 *Machine Learning Research*, **9**, 2579–2605.
- 570 Walker, S. C. & Jackson, D. A. (2011). Random-effects ordination: describing and predict-
571 ing multivariate correlations and co-occurrences. *Ecological Monographs*, **81**, 635–663.
- 572 Wang, Y., Naumann, U., Eddelbuettel, D., Wilshire, J. & Warton, D. (2020). *mvabund:*
573 *Statistical Methods for Analysing Multivariate Abundance Data*.

- 574 Warton, D. I., Blanchet, F. G., O'Hara, R. B., Ovaskainen, O., Sara, T., Walker, S. C.
575 & Hui, F. K. C. (2015). So Many Variables: Joint Modeling in Community Ecology.
576 *Trends in Ecology & Evolution*, **30**, 766 – 779.
- 577 Warton, D. I. & Hui, F. K. C. (2017). The central role of mean-variance relationships in
578 the analysis of multivariate abundance data: a response to Roberts (2017). *Methods in*
579 *Ecology and Evolution*, **8**, 1408–1414.
- 580 Warton, D. I., Thibaut, L. & Wang, Y. A. (2017). The PIT-trap — A “model-free”
581 bootstrap procedure for inference about regression models with discrete, multivariate
582 responses. *PLOS ONE*, **12**, 1–18.
- 583 Wickham, H., Cook, D., Hofmann, H. & Buja, A. (2011). tourr: An R Package for
584 Exploring Multivariate Data with Projections. *Journal of Statistical Software*, **40**, 1–
585 18.
- 586 Wood, S. N. (2017). *Generalized additive models: an introduction with R*. CRC press.
- 587 Zuur, A. F., Ieno, E. N. & Smith, G. M. (2007). Principal component analysis and
588 redundancy analysis. *Analysing ecological data*, pp. 193–224.