

Model-based ordination for species with unequal niche widths

Bert van der Veen^{1,2,3}

Francis K.C. Hui⁴

Knut A. Hovstad⁵

Erik B. Solbu¹

Robert B. O'Hara^{2,3}

¹Department of Landscape and Biodiversity, Norwegian Institute of Bioeconomy research, Trondheim, Norway

²Department of Mathematical Sciences, Norwegian University of Science and Technology, Trondheim, Norway

³Centre of Biodiversity Dynamics, Norwegian University of Science and Technology, Trondheim, Norway

⁴Research School of Finance, Actuarial Studies and Statistics, Australian National University, Canberra, Australia

⁵The Norwegian Biodiversity Information Centre, Trondheim, Norway

Summary

1. It is common practice for ecologists to examine species niches in the study of community composition. The response curve of a species in the fundamental niche is usually assumed to be quadratic. The center of a quadratic curve represents a species' optimal environmental conditions, and the width its ability to tolerate deviations from the optimum.

2. Most multivariate methods assume species respond linearly to the environment of the niche, or with a quadratic curve that is of equal width and height for all species. However, it is widely understood that some species are generalists who tolerate deviations from their optimal environment better than others. Rare species often tolerate a smaller range of environments than more common species, corresponding to a narrow niche.

3. We propose a new method, for ordination and fitting Joint Species Distribution Models, based on Generalized Linear Mixed-Effects Models, which relaxes the assumptions of equal tolerances and equal maxima.

4. By explicitly estimating species optima, tolerances, and maxima, per ecological gradient, we can better predict change in species communities, and understand how species relate to each other.

29 **keywords:** model-based ordination, unimodal response, niche model, unconstrained quadratic ordina-
30 tion, joint species distribution model.

31 Introduction

32 One of the key topics addressed by community ecology is what causes changes in community composition. In
33 order to explore species niches, species communities are surveyed at locations with different environmental
34 conditions. Species tolerances are then reflected in the resulting multivariate dataset, as differences in
35 occurrences or abundances between locations. The most favorable environmental conditions for species are
36 represented by the optimum of the niche, where species exhibit their maximum abundance or probability of
37 occurrence. Deviation from the optimum reflects increasingly unfavorable conditions.

38 Correspondence Analysis (CA) is often used to summarize community data, as it implicitly approximates
39 the fit of a quadratic model, with the additional assumptions of equally spaced optima, sites that are well
40 within the range of species optima, equal tolerances, and equal or independent maxima (ter Braak 1985).
41 The combination of assuming equally spaced optima, equal maxima, and equal tolerances, gives an early
42 niche model, called the species packing model (MacArthur & Levins 1967). The relationship of the species
43 packing model to CA has added to its popularity among applied ecologists (Wehrden *et al.* 2009).

44 Recent advances in the estimation of species niches have focussed on performing ordination with explicit
45 statistical models, such as Generalized Linear Latent Variable Models (GLLVMs; Warton *et al.* 2015). The
46 GLLVM framework is well known for its capability to fit Joint Species Distribution Models (JSDMs; Pollock
47 *et al.* 2014; Ovaskainen *et al.* 2017; Tobler *et al.* 2019; Zurell *et al.* 2020). In the context of JSDMs, GLLVMs
48 assume species abundances are correlated due to similarity in response to ecological gradients, modelled with
49 covariates or latent variables respectively. Latent variables can be understood as combinations of missing
50 covariates, so that GLLVMs allow us to parsimoniously model species distributions. They are equivalent to
51 ordination axes, representing complex ecological gradients (Halvorsen 2012). Recently, the use of GLLVMs
52 to perform model-based ordination has increased in popularity (Inoue *et al.* 2017; Björk *et al.* 2018; Lacoste
53 *et al.* 2019; Damgaard *et al.* 2020).

54 With intercepts included for row standardization, GLLVMs fit the species packing model (Jamil & ter
55 Braak 2013; Hui *et al.* 2015), though with maxima that are equal for the latent variables. Existing GLLVMs
56 assume that latent variables are linear, just as all classical ordination methods (Jamil & ter Braak 2013).
57 However, it is widely understood that species have unequal tolerances and maxima, so that the assumptions
58 of linear latent variables, and equal tolerances, are unlikely to hold in practice.

59 In this paper, our goal is to overcome the assumptions of equal tolerances, and equal maxima, by formu-

60 lating a GLLVM with quadratic latent variables. To our knowledge, there has been no attempt to implement
61 such a GLLVM until now. Although seemingly a straightforward extension, the quadratic term explicitly
62 allows species niches to be estimated without constraints on the parameters. This means that optima, tol-
63 erances, and maxima per latent variable, as well as the lengths of ecological gradients, can all be explicitly
64 estimated. Explicitly estimating the combination of these three parameters gives unique insight into reasons
65 for species low detectability, whether it is due to low abundance or probability of occurrence (maxima),
66 a high degree of habitat specialization (tolerance), or due to unsuitable observed environmental conditions
67 (optima). In combination with knowledge of the study system, these parameters can help ecologists to deter-
68 mine why certain species are rare. Additionally, due to the model-based nature of the proposed ordination
69 method, it is possible to calculate confidence intervals for each set of parameters, providing unparalleled
70 benefits for inference when using ordination. In the context of JSDMs, the quadratic GLLVM models latent
71 species distributions, without covariates in the model. When covariates are included, the quadratic GLLVM
72 partitions species distributions in observed (fixed effects) and latent or unobserved (random effects), similar
73 to the partitioning of fixed and random effects in mixed-effects models when covariates are included.

74 In contrast to classical ordination methods, GLLVMs model the latent variables as unobserved, treating
75 them as random rather than fixed (Walker & Jackson 2011), which consequently have to be integrated
76 over in the likelihood. Here, we develop a variational approximations (VA) implementation after Hui *et al.*
77 (2017) and Niku *et al.* (2019a), to perform calculations quickly and efficiently. In addition to presenting
78 the quadratic GLLVM, we perform simulations to evaluate the accuracy of the VA implementation, and the
79 capability of the quadratic GLLVM to retrieve the true species-specific parameters and latent variables. We
80 use two real world datasets to demonstrate use and interpretation of the proposed quadratic GLLVM: 1) a
81 small dataset of hunting spiders in a Dutch dune ecosystem (van der Aart & Smeek-Enserink 1974), and 2)
82 a larger dataset on Swiss alpine plant species on a strong elevation gradient (D’Amen *et al.* 2018).

83 Model formulation

84 The ecological niche is here described by a quadratic function involving three parameters; the optimum \mathbf{u}_j ,
85 the tolerance \mathbf{t}_j , and the maximum \mathbf{c}_j . The optimum \mathbf{u}_j is the location on the ecological gradient where a
86 species exhibits its highest abundance or probability of occurrence (the maximum \mathbf{c}_j). The tolerance \mathbf{t}_j is a
87 measure of the width or breadth of the niche, and indicates if a species is a generalist or specialist.

88 Consider an $n \times p$ matrix of observations, where y_{ij} denotes the response of species $j = 1 \dots p$ at site
89 $i = 1 \dots n$. Then in the quadratic GLLVM, we assume that, conditional on a vector \mathbf{z}_i of $q = 1 \dots d$ latent
90 variables where $d \ll p$, the responses y_{ij} at site i are independent observations from a distribution whose

91 mean, denoted here as $E(y_{ij}|\mathbf{z}_i)$, is modelled as:

$$\begin{aligned} g\{E(y_{ij}|\mathbf{z}_i)\} &= \sum_{q=1}^d \left\{ c_{jq} - \frac{(z_{iq} - u_{jq})^2}{2t_{jq}^2} \right\} \\ &= \sum_{q=1}^d \left(c_{jq} - \frac{u_{jq}^2}{2t_{jq}^2} + \frac{z_{iq}^2}{2t_{jq}^2} - \frac{z_{iq}u_{jq}}{t_{jq}^2} \right), \end{aligned} \quad (1)$$

92 where $g\{\cdot\}$ is a known link function (e.g. the log-link when the responses are assumed to be Poisson, negative-
93 binomial, or gamma distributed, the probit-link when the responses are assumed to be Bernoulli or ordinal
94 distributed, and the identity-link for responses that are assumed to be Gaussian distributed).

95 To facilitate easier estimation, and for a closer comparison to the linear GLLVM, we formulate the
96 quadratic GLLVM in matrix notation:

$$g\{E(y_{ij}|\mathbf{z}_i)\} = \beta_{0j} + \mathbf{z}_i^\top \boldsymbol{\gamma}_j - \mathbf{z}_i^\top \mathbf{D}_j \mathbf{z}_i, \quad (2)$$

97 with a species-specific intercept β_{0j} that accounts for e.g. mean abundances, and a vector of coefficients per
98 species for the linear term $\boldsymbol{\gamma}_j$. We can see a third term is added here to the existing structure of the linear
99 GLLVM, which models tolerances and maxima per species and latent variable. Specifically, we introduce a
100 diagonal matrix \mathbf{D}_j of quadratic coefficients with each diagonal element being the quadratic effect for latent
101 variable q and species j . We require \mathbf{D}_j to be a positive-definite diagonal matrix, to ensure concave curves
102 to the latent variables. Thus, $2\mathbf{D}_j$ is the precision matrix of the ecological niche (likewise $(2\mathbf{D}_j)^{-1}$ is the
103 covariance matrix). Additionally, row intercepts or covariates can be included as in Hui *et al.* (2017), or
104 species traits as in Niku *et al.* (2019a), though we have chosen to omit those terms here and focus on the
105 case of unconstrained ordination.

106 With \mathbf{D}_j being a diagonal matrix with the positive elements D_{jqj} , the vector of species maxima \mathbf{c}_j
107 with elements c_{jq} , the vector of species optima \mathbf{u}_j with elements u_{jq} , and the vector of species tolerances
108 \mathbf{t}_j with elements t_{jq} , we derive the following connections between the parameters in equations (1) and (2):
109 $\beta_{0j} = \sum_{q=1}^d c_{jq} - u_{jq}^2/(2t_{jq}^2)$, $\gamma_{jq} = -u_{jq}/t_{jq}^2$, and $D_{jqj} = 1/(2t_{jq}^2)$. Similarly, for the formulation in equation
110 (2), the parameters in equation (1) can be retrieved: $c_{jq} = \beta_{0j} + u_{jq}\gamma_{jq} - u_{jq}^2 D_{jqj}$, $u_{jq} = -\gamma_{jq}/(2D_{jqj})$, and
111 $t_{jq} = 1/\sqrt{2D_{jqj}}$.

112 Four special cases of the quadratic GLLVM, as formulated in equation (2), are worth discussing: 1)
113 $\mathbf{D}_j = \mathbf{D}$, i.e. common tolerances for species, 2) $\mathbf{D}_j = D_{11}\mathbf{I}_d$ where \mathbf{I}_d is a $d \times d$ identity matrix, i.e. equal
114 tolerances for species and latent variables, 3) when $\mathbf{D}_j = 0$ for a subset of the p species, and 4) when $\mathbf{D}_j = 0$
115 for all p species. The first case assumes tolerances to be the same across species, but not latent variables,
116 and additionally places constraints on the species maxima. This species-common tolerances model might

117 prove useful in practice, as it requires fewer observations per species than the full quadratic GLLVM, but still
118 explicitly includes quadratic latent variables. The second case can be shown to be equivalent to the linear
119 GLLVM with row intercepts as presented in Hui *et al.* (2015), which assumes tolerances to be the same for
120 all species and latent variables, and the maxima to be the same for all latent variables. In the third case,
121 some species respond to the latent variable linearly, while others exhibit quadratic responses. The fourth
122 case is the most basic GLLVM with linear latent variables, currently possible to fit with e.g. `boral` (Hui
123 2016), `HMSC-R` (Tikhonov *et al.* 2020), and `gllvm` (Niku *et al.* 2020).

124 Model interpretation

125 In this section, we derive and discuss various tools that are commonly used in the application of JSDMs and
126 ordination, such as calculating residual correlations and partitioning residual variance, calculating gradient
127 length, and visualizing the ordination, and demonstrate how they can be adapted to the proposed quadratic
128 GLLVM.

129 Residual covariance matrix

130 One aspect GLLVMs are known for is modelling species residual correlations, calculated from the residual
131 covariance matrix (Zurell *et al.* 2018; Blanchet *et al.* 2020). To facilitate calculation of the residual covariance
132 matrix, we can reparameterize all GLLVMs as a multivariate mixed-effects model with a residual term:

$$g\{E(y_{ij}|\mathbf{z}_i)\} = \beta_{0j} + \epsilon_{ij}. \quad (3)$$

133 Here, ϵ_{ij} accounts for any residual information that is not accounted for by fixed-effects in the model, such
134 as covariates or intercepts (Warton *et al.* 2015). Assuming the latent variables are independent for all sites,
135 the elements of the residual covariance matrix are given by:

$$\Sigma_{jk} = \text{cov}(\epsilon_{ij}, \epsilon_{kl}), \quad \forall i, k = 1 \dots n, j, l = 1 \dots p.$$

136 For a length p vector ϵ_i , existing JSDM implementations assume $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$, i.e. the residual term follows
137 a multivariate normal distribution. For the linear GLLVM, it is straightforward to show that $\epsilon_{ij} = \mathbf{z}_i^\top \boldsymbol{\gamma}_j$, so
138 it follows that $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\gamma}_j^\top \boldsymbol{\gamma}_j)$. In essence, GLLVMs perform a low rank approximation to the covariance
139 matrix of a residual term. The rank of this residual covariance matrix is equal to the number of estimated
140 latent variables d in the model for the linear GLLVM.

141 Turning to the quadratic GLLVM, where $\epsilon_{ij} = \mathbf{z}_i^\top \boldsymbol{\gamma}_j - \mathbf{z}_i^\top \mathbf{D}_j \mathbf{z}_i$, the elements of the residual covariance
142 matrix are:

$$\Sigma_{\text{quad},jk} = \sum_{q=1}^d (\gamma_{jq} \gamma_{kq} + 2D_{jqq} D_{kqq}), \quad (4)$$

143 for which a proof is given in Appendix S1. This can be rewritten in terms of the species optima \mathbf{u}_j and
144 tolerances t_j :

$$\Sigma_{\text{quad},jk} = \sum_{q=1}^d \{(t_{jq}^2 t_{kq}^2)^{-1} (0.5 + u_{jq} u_{kq})\}, \quad (5)$$

145 from which it follows that $\epsilon_{ij} \sim \sum_{q=1}^d t_{jq}^2 \chi^2(\frac{1}{4}, \frac{1}{4} u_{jq}^2)$, in words: the residual term follows a generalized χ^2
146 distribution (Khatri 1980).

147 Equation (4) and equation (5) additionally serve to demonstrate how to partition the residual variance
148 of the quadratic GLLVM, e.g. per latent variable, for the linear and quadratic term separately, or both.
149 Variance partitioning is commonly used in the application of ordination methods, e.g. to determine fit
150 (Økland 1999), or to explore causes of residual variance (Borcard *et al.* 1992; Økland & Eilertsen 1994).
151 Covariates can be included in the model to account for the residual variance otherwise accounted for by the
152 latent variables. The residual variance can be used to identify indicator species i.e. those species that best
153 represent an ecological gradient, or to calculate a measure of R^2 (Nakagawa & Schielzeth 2013).

154 The rank of the residual covariance matrix is double that of a linear GLLVM with the same number of
155 latent variables: $2d$. The additional quadratic term thus allows us to account for more residual correlations
156 between species, with fewer latent variables. This corresponds with the ecological notion that species often
157 respond to few major complex ecological gradients (Halvorsen 2012). From this, we see that when the number
158 of latent variables in a quadratic GLLVM exceeds $\frac{1}{2}p$, there are more parameters included than in a JSJM
159 with an unstructured residual covariance matrix. However, this is not an issue here, since for ordination
160 purposes we are only interested in cases where there are much fewer latent variables d than species p .

161 Gradient length

162 The length of an ecological gradient is of great interest to ecologists in the use of ordination, because it
163 provides a measure of beta diversity (Oksanen & Tonteri 1995). Longer gradients indicate higher diversity,
164 as spacing (i.e. dissimilarity in the species community) between sites in latent space is potentially larger. In
165 the past, it has been emphasized that short gradients are better analysed using linear ordination methods,
166 and longer with unimodal methods (ter Braak & Prentice 1988). However, the quadratic GLLVM allows

167 species to exhibit both linear and unimodal responses, and so it is appropriate for both, and it is no longer
168 required to switch ordination method as a consequence of gradient length.

169 To determine gradient length from the proposed quadratic GLLVM, we define the ecological gradients
170 \tilde{z}_i , as a function of the latent variables z_i , but with a diagonal covariance matrix \mathbf{G} of size $d \times d$. First, for a
171 species-common tolerances model, we note that the quadratic term in equation (2), i.e. $z_i^\top \mathbf{D} z_i$, can instead
172 be written as $\sum_{q=1}^d z_{iq}^2 D_{qq}$, so that $\tilde{z}_{iq} = z_{iq} \sqrt{D_{qq}}$, and $\tilde{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{G})$, where $\mathbf{G} = \mathbf{D}$. Then, the gradient
173 length is approximately $4\mathbf{G}^{\frac{1}{2}}$ (i.e. the approximate width of a normal distribution), and $\mathbf{D} = \mathbf{I}$

174 For the species-specific tolerances model, we note that one of the uses of gradient length in the past has
175 been to rescale the latent variables so that an ordination diagram can be understood in terms of compositional
176 turnover (Hill & Gauch 1980). This requires the mean species tolerances to be one (as is the case for the
177 species-common tolerances model above), so that the covariance matrix of the ecological gradient in the
178 species-specific tolerances model is $G_{qq} = \frac{1}{p} \sum_{j=1}^p D_{jqqq}$, so the matrix of quadratic coefficients \mathbf{D}_j is scaled
179 by the inverse of the covariance matrix of the ecological gradient, \mathbf{G}^{-1} . However, we choose to use the
180 median of the species tolerances instead, as it more accurately represents gradient length with both linear
181 and quadratic responses of species in the model. In general, the proposed quadratic model allows further
182 exploration of measures of gradient length by, for example, using the mean tolerance of species with clear
183 quadratic responses, rather than the median of all tolerances.

184 The measure of gradient length calculated here, can be interpreted in the same manner as the gradient
185 length provided by Detrended Correspondence Analysis (Hill & Gauch 1980).

186 Ordination diagram

187 Usually, a biplot (Gabriel 1971) is constructed to visually inspect results from an ordination. For the
188 quadratic GLLVM, biplots tend to create an arch when the residual variance of the linear term is smaller
189 than the residual variance of the quadratic term.

190 Instead, we propose that species optima and tolerances can be plotted directly, so that species niches
191 are visualized in a two-dimensional latent space from a top-down perspective. The widths of the niches
192 can then potentially be represented as ellipses using the estimated species tolerances (i.e. providing species
193 distributions in latent space), so that co-occurrence patterns can be inferred from the (lack of) overlap
194 between ellipses. Additionally, information on sites, such as the predicted locations and prediction regions,
195 can be added (Hui *et al.* 2017). Information for the sites can be used to infer the distance of sites to the
196 species optima (i.e. the suitability of sites for species), or to the edges of species niches (see the hunting
197 spiders example below).

198 Finally, based on the discussion in the two subsections above, there are two ways of scaling the ordination
199 diagram: 1) by the residual variance per latent variable, or 2) by the mean or median tolerance. In the first
200 scaling, the diagram is scaled to draw attention to the latent variable that explains most variance in the model.
201 However, the second scaling has a more ecological intuitive interpretation. If the tolerances are assumed to
202 be common for species, the second scaling provides an ordination diagram in units of compositional turnover
203 (Gauch 1982). When the linear and quadratic terms in the model explain an equal proportion of the total
204 residual variance per latent variable, these scalings produce similar results.

205 Model estimation

206 We propose to use variational approximations (VA; Hui *et al.* 2017) for estimation and inference for the
207 quadratic GLLVM. Broadly speaking, VA is a general technique used to provide a closed-form approximation
208 to the marginal log-likelihood of a model with random effects or latent variables, when an analytical solution is
209 not available. Computationally, VA can be orders of magnitude faster than MCMC, numerical integration,
210 or even the Laplace approximation (Niku *et al.* 2019a), and without loss of accuracy (Hui *et al.* 2017).
211 However, the calculation of the VA log-likelihoods needs to be derived on a case-by-case basis. In contrast,
212 the Laplace approximation can be applied automatically in many cases (Kristensen *et al.* 2016). Note it is
213 not possible to approximate the marginal likelihood of a quadratic GLLVM with the Laplace approximation
214 (K. Kristensen, pers. comm., March 8th 2019).

215 The marginal log-likelihood of a quadratic GLLVM is given by:

$$\mathcal{L}(\Theta) = \sum_{i=1}^n \log \left\{ \int_{-\infty}^{\infty} \prod_{j=1}^p f(y_{ij} | z_i, \Theta) h(z_i) dz_i \right\}, \quad (6)$$

216 where $f(y_{ij} | z_i, \Theta)$ is the distribution of the species responses given the latent variables. As mentioned
217 previously, and as per Hui *et al.* (2015), we assume the distribution of the latent variables $h(z_i)$ to be
218 multivariate standard normal i.e. $h(z_i) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. The vector Θ includes all parameters in the model
219 $\Theta = \{\beta_{01} \dots \beta_{0j}, \gamma_{11} \dots \gamma_{jq}, D_{111} \dots D_{jqj}\}^T$. Equation (6) can be straightforwardly modified if covariates are
220 also included in the quadratic GLLVM.

221 In VA, we construct a lower bound to equation (6), by assuming that the posterior distribution of the
222 latent variables can be approximated by a closed form distribution e.g., a multivariate normal distribution.
223 We then minimize the Kullback-Leibler divergence between this approximate closed-form distribution (also
224 known as the variational distribution) and the true posterior distribution. Hui *et al.* (2017) showed that,
225 for GLLVMs with linear latent variables, the optimal variational distribution is multivariate normal $z_i \sim$

226 $\mathcal{N}(\mathbf{a}_i, \mathbf{A}_i)$, with mean \mathbf{a}_i and covariance matrix \mathbf{A}_i , so we will adopt this choice here as well.

227 In Appendix S1 we provide information on calculating approximate confidence intervals for the parame-
228 ters. In Appendix S2 we provide derivations for the log-likelihood of common response types in community
229 ecology, such as count data (Poisson, and negative-binomial with quadratic mean-variance relationship,
230 and both assuming a log-link function), binary data and ordinal data (both with probit-link function), as
231 well as positive continuous data (gamma, with log-link function) and continuous data (Gaussian, with an
232 identity-link function).

233 Simulation study

234 To assess how well the proposed model retrieves the true latent variables \mathbf{z}_i , optima \mathbf{u}_j , tolerances \mathbf{t}_j , and
235 maxima \mathbf{c}_j , we performed simulations for six response distributions; 1) Gaussian, 2) gamma, 3) Poisson, 4)
236 negative-binomial, 5) Bernoulli, and 6) ordinal. The R-code used for the simulations is provided in Appendix
237 S3. For each of the distributions, we simulated 1000 datasets with different numbers of sites and species.
238 A consequence of the negative-only third term in the quadratic GLLVM, is that the model often simulates
239 a large number of zeros (more so than the linear GLLVM), providing a challenge in testing its accuracy,
240 especially for small datasets. First, to study the accuracy of the VA approximation, we simulated datasets
241 of $p = 20$ to 100 species in increments of 10, while keeping the number of sites constant at $n = 100$. Hui
242 *et al.* (2017) argued that the VA log-likelihood is expected to converge to the true likelihood as $p \rightarrow \infty$
243 (i.e. for a large number of species), thus this will allow us to study the finite sample properties of the VA
244 approximation for the proposed model.

245 Second, to explore the sample size required to accurately estimate the species-specific parameters e.g.,
246 species optima \mathbf{u}_j , tolerances \mathbf{t}_j , and maxima \mathbf{c}_j , we simulated datasets of $n = 20$ to 100 sites in increments of
247 10, while keeping the number of species constant at $p = 100$. For each dataset, we compared 12 combinations
248 of initial values and fitting algorithms (see Appendix S4: Fitting, for details), and picked the model with the
249 highest log-likelihood (see Appendix S5: Fig. S1 for the frequency at which different types of initial values
250 and fitting algorithm were used in the best models per distribution).

251 As a true model, we considered a quadratic GLLVM with $d = 2$ latent variables, which was constructed
252 as follows. First, the species-specific intercepts β_{0j} were simulated as $\text{Uniform}(-1, 1)$, which corresponds
253 to species with low abundance or occurrence. Next, the true coefficients corresponding to the linear terms
254 in the model $\boldsymbol{\gamma}_j$, were simulated independently as $\text{Uniform}(-5, 5)$, and the true quadratic coefficients as
255 $\text{Uniform}(-5, -0.5)$. The true latent variables were simulated as $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For the Gaussian, negative-
256 binomial, and gamma distribution, the dispersion parameter for all species was set equal to one. For the

257 ordinal distribution we assumed six classes with the true cut-offs being 0, 1, 2, 3, 4, 5, meaning that species
258 were most often absent (category 1), while they were rarely very abundant (category 6).

259 We measured performance of the quadratic GLLVM by the prediction of the latent variables z_i and the
260 species optima u_j . The species optima are a function of both the linear and quadratic coefficients and should
261 provide a good overall measure of performance for retrieving the true species-specific parameters, in addition
262 to being of specific interest to ecologists. Though it is common to measure the performance of ordination
263 methods using the Procrustes error (Peres-Neto & Jackson 2001), we chose to use the Median Absolute Error
264 (MAE) instead, as we often observed a highly skewed error distribution for the species optima. Additionally,
265 interpreting the MAE is more intuitive, as it measures the deviation from the truth in the same units as the
266 coefficients of interest. We excluded the first optimum of the second latent variable as this was fixed to zero
267 for reasons of parameter identifiability (Hui *et al.* 2015), and excluded optima that could not be estimated.
268 Since the quadratic GLLVM allows species to exhibit linear responses, which have infinite optima, we chose
269 to remove all optima larger than 10 and smaller than -10, i.e. for those species that lacked a sufficiently
270 strong quadratic signal in the simulated datasets. Including these optima would result in a biased view of
271 the accuracy of the optima that can be estimated by the model. This process resulted in a vector of optima,
272 which we then used to calculate the MAE. For clarity and transparency, we additionally present the number
273 of optima removed for each of the datasets, to further provide an impression of the data requirements of the
274 proposed quadratic GLLVM.

275 For all of the models fitted to Gaussian and gamma response datasets, typically none or only a few optima
276 were excluded, meaning that the median number excluded was zero. In general, and not surprisingly, more
277 optima were excluded for models fitted to datasets where n/p was small and for discrete distributions. For
278 example, when $n = 20$ sites and $p = 100$ species, the median number of optima excluded for datasets with
279 Poisson responses was 4 (2 - 8, first and third quartiles), for datasets with negative-binomial responses this
280 was 7 (5 - 10), for datasets with Bernoulli responses this was 31 (24 - 38), and for datasets with ordinal
281 responses this was 17 (13 - 23). In contrast, for datasets where n/p was large, considerably less optima were
282 excluded across all response types. For example, when $n = 100$ and $p = 100$, for Poisson and negative-
283 binomial response datasets the median number of excluded optima was zero, while for Bernoulli response
284 datasets the median number of optima excluded was 4 (2 - 5), and for ordinal response datasets this was 2
285 (1 - 3).

286 The MAE per distribution and for the different sized datasets is presented in Figure 1 (see Appendix S5:
287 Fig. S2 for the same figure with all species optima). As expected, the quadratic GLLVM was more accurate
288 for datasets with larger p and larger n . For all distributions, the latent variables were often better retrieved
289 than the species optima. This is not surprising, as the species optima are a function of two parameters,

290 particularly the inverse of the quadratic coefficients, so that a small change in the quadratic coefficients
 291 can result in a large change in the species optima. When fitted to Gaussian or gamma response datasets,
 292 regardless of the dimensions of the data, the model performed best. The accuracy of the estimated species
 293 optima was only slightly lower for the Poisson distributed datasets with 70 or more sites, while the latent
 294 variables were accurately estimated even with small p . These results are consistent with the results above
 295 regarding the number of excluded optima. Although the accuracy of species optima for negative-binomial
 296 response datasets seems similar to that of Poisson response datasets, this is inconsistent with the number
 297 of excluded optima reported above, as that was considerably larger for the negative-binomial. The model
 298 was not accurate for Bernoulli or ordinal response datasets with small n and p . However, when the number
 299 of sites and species increased above 40, the performance of the quadratic GLLVM in these cases improved
 300 considerably, consistent with the number of excluded species optima.

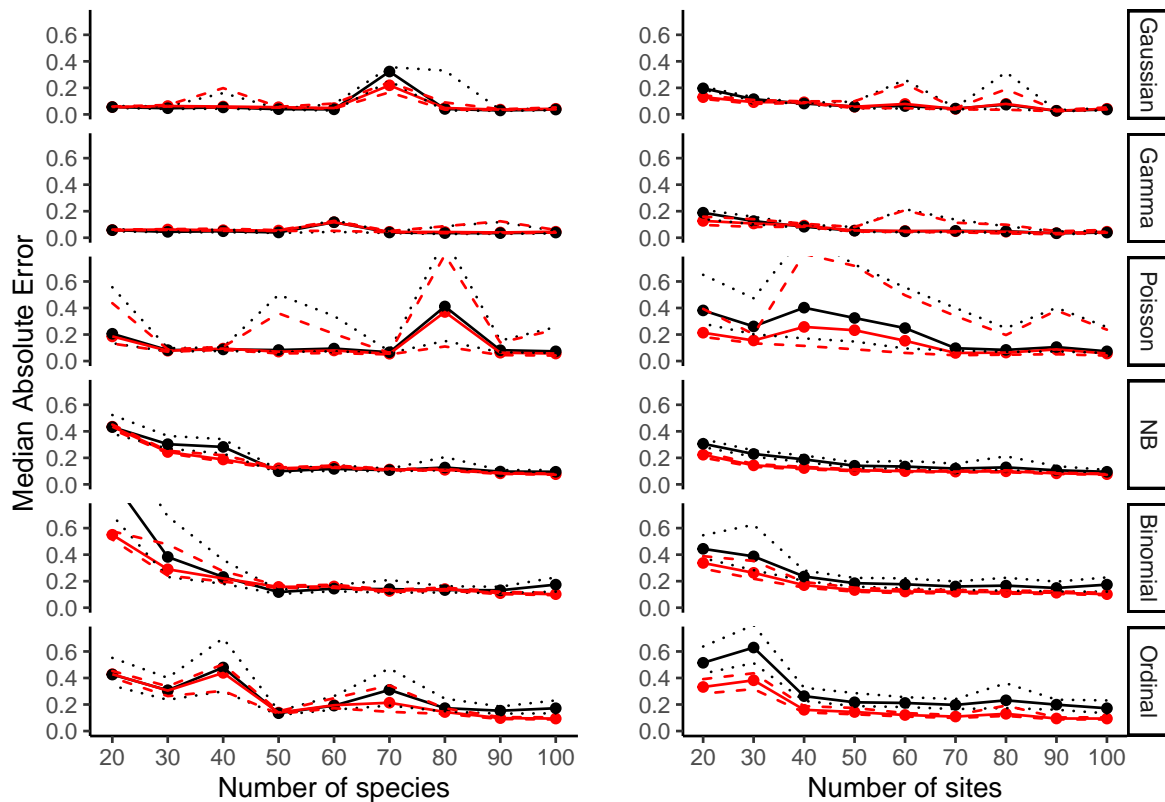


Figure 1: Simulation results for the 1000 best fitting quadratic GLLVMs across six different response distributions, with the MAE calculated based on optima that could be estimated (optima outside the range $(-10,10)$ were excluded). The left column shows simulations where the number of sites was kept constant at $n = 100$, and analogous for the right column with $p = 100$. The figure includes the median MAE for species optima (black) and latent variables (red), with the first and third quartiles represented as dotted (optima) and dashed (latent variables) lines.

301 Applications to real data

302 We applied the proposed quadratic GLLVM to two different datasets: 1) the classical hunting spiders dataset
303 collected by van der Aart & Smeek-Enserink (1974) in Dutch dunes, available in the `mvabund` R package
304 (Wang *et al.* 2012), and 2) a dataset of plants in the Swiss Alps (available in the dryad database; D’Amen
305 *et al.* 2017).

306 Hunting spiders

307 For the hunting spiders dataset, van der Aart & Smeek-Enserink (1974) used pitfall traps to collect spiders
308 over a 60 week period, resulting in a dataset of counts for each of the $n = 28$ sites and $p = 12$ species. It
309 has been used in the testing of ordination methods before (e.g. ter Braak 1985, 1986; Yee 2004; Hui *et al.*
310 2015), providing some reference results for comparison here. To find the model that best fitted the hunting
311 spiders dataset, and to limit the number of required model fits, we first performed model selection using
312 Akaike’s Information Criterion (AIC; Burnham & Anderson 2002) on linear GLLVMs with $d = 1$ to 5 latent
313 variables, and with Poisson distributions. Though the hunting spiders dataset exhibits overdispersion in the
314 linear GLLVM (Hui *et al.* 2015), the quadratic GLLVM models overdispersion with the latent variables (see
315 Appendix S2: Negative-Binomial: overdispersed counted responses). Second, we fitted a linear GLLVM with
316 random row intercepts (i.e. equal tolerances), a quadratic GLLVM with common tolerances, and a quadratic
317 GLLVM with unequal tolerances, to determine which model structure was most suited for the data. Third,
318 with the best model structure from step two, we again tested for the optimal number of latent variables,
319 after which we explored different sets of initial values and fitting algorithms to find the model that maximizes
320 the VA log-likelihood (see Appendix S5 for the results).

321 The best model of step one included $d = 2$ latent variables, the best model from step two included unequal
322 tolerances, and the best model from step three included $d = 3$ latent variables. The results for the first two
323 latent variables of the final model fit, which explained most residual variation, are presented in Figure 2.

324 We used the residual variance to determine which latent variables explained most variation i.e. were most
325 important to consider for inference. For the quadratic GLLVM, the first and second latent variables explained
326 most variation in the model; 40% and 57% respectively. Overall, the quadratic GLLVM explained four and
327 a half times more residual variation than a linear GLLVM with the same number of latent variables. The
328 lengths of the first two ecological gradients were 5.46 (4.28-6.64, 95% confidence interval), and 3.35 (3.14-
329 3.55). The confidence interval of the gradient length for the third ecological gradient included zero, so we
330 do not present results of that here.

331 Ter Braak (1985) and Yee (2004) both visualized quadratic curves of the first latent variable using

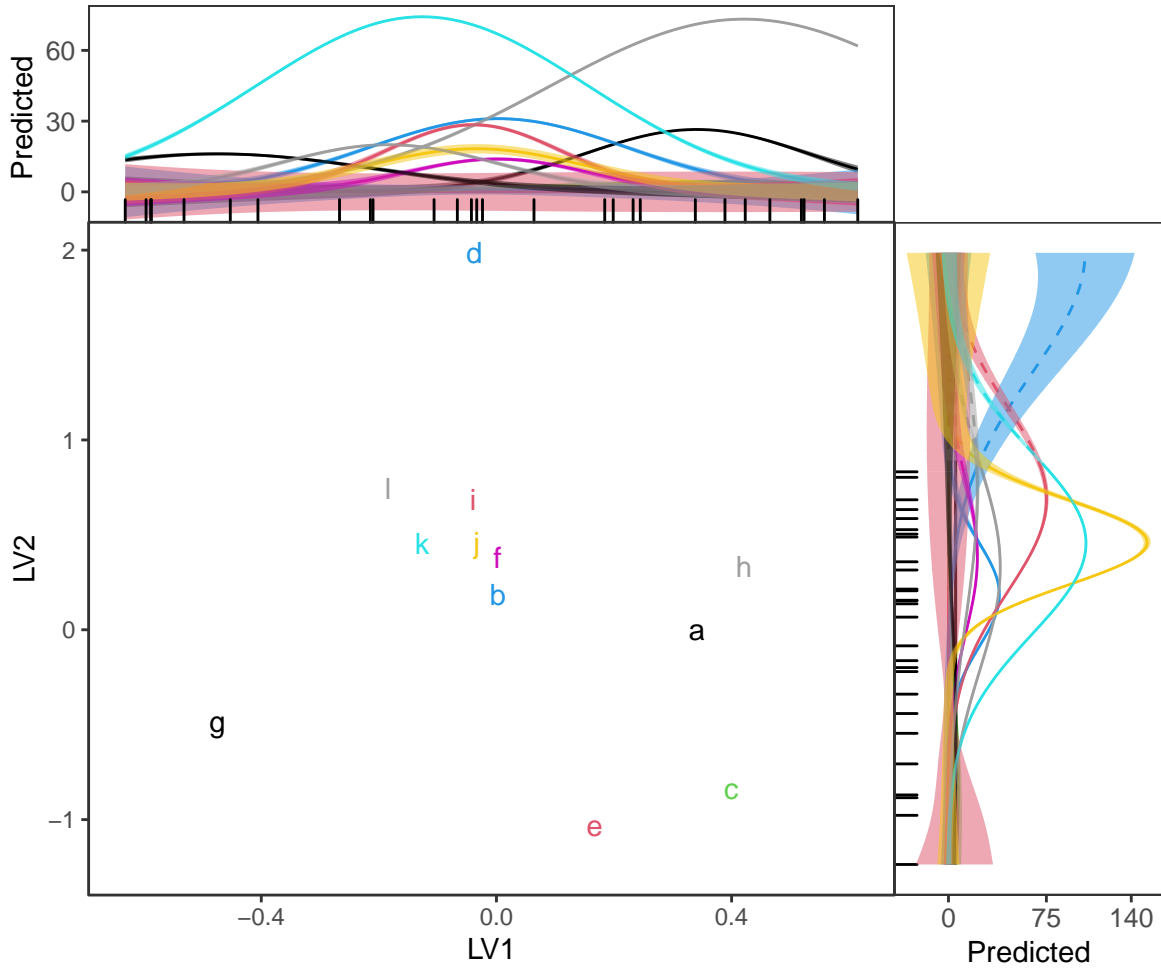


Figure 2: Ordination plot for the first two latent variables of the final quadratic GLLVM fit to the hunting spiders dataset, scaled by the residual variances. Species optima are shown as letters, indicating the following species: a = *Alopecosa accentuata*, b = *Alopecosa cuneata*, c = *Alopecosa fabrilis*, d = *Arctosa lutetiana*, e = *Arctosa perita*, f = *Alonia albimana*, g = *Pardosa lugubris*, h = *Pardosa monticola*, i = *Pardosa nigriceps*, j = *Pardosa pullata*, k = *Trochosa terricola*, l = *Zora spinimana*. Species quadratic curves are included as side panels, with dashed lines indicating unobserved parts of species niches, and with bands representing 95% confidence intervals. Site locations and prediction regions have not been included, in favor of readability.

332 variations of Poisson regression and generalized additive models, respectively. There are clear similarities
333 between the species response curves for the first latent variable in Figure 2, and the corresponding response
334 curves described by ter Braak (1985) and Yee (2004). Though ter Braak (1985) concluded all species exhibited
335 unimodal curves on the first latent variable (without formal testing), the species *Alopecosa fabrilis*, *Arctosa*
336 *perita* and *Pardosa lugubris* had confidence intervals, for the quadratic coefficients, that include zero in the
337 quadratic GLLVM. For the following conclusions, species with confidence intervals of quadratic coefficients
338 that crossed zero were excluded.

339 Turning to the species niches, on the first latent variable all optima were observed (i.e. within the
340 range of the latent variable), and on the second latent variable only the optimum of *Arctosa lutetiana* was
341 unobserved. On the first latent variable, *Aulonia albimana* had the lowest maximum and *Trochosa terricola*
342 the highest. On the second latent variable, *Alopecosa accentuata* had the lowest maximum, and *Pardosa*
343 *pullata* the highest. On the first latent variable it was possible to distinguish that *Trochosa terricola* and
344 *Pardosa monticola* had wider niches than *Pardosa nigriceps*, *Aulonia albimana*, and *Arctosa lutetiana*, and
345 the confidence intervals of these two groups did not overlap. On the second latent variable, *Pardosa pullata*
346 had the most narrow niche, and *Pardosa monticola* the widest, and the confidence intervals of the tolerances
347 for these species did not overlap. *Alopecosa cuneata* was more tolerant to changes in the environment than
348 *Pardosa pullata* but less than *Pardosa monticola*. Additionally, *Alopecosa fabrilis*, *Trochosa terricola*, and
349 *Pardosa nigriceps* were more tolerant to changes in the environment than *Pardosa pullata*, though it was not
350 possible to say if this was more or less than *Pardosa monticola*. Overall, *Arctosa lutetiana* had the smallest
351 tolerance across all three latent variables.

352 Overall, due to a combination of low maxima and low tolerances, *Arctosa lutetiana* is predicted to be
353 most prone to changes in the environment of the first latent variable, and for the second latent variable
354 *Arctosa perita*.

355 **Swiss alpine plants**

356 In the second application, $n = 912$ plots of $4 m^2$ each were used to record binary data on $p = 175$ plant
357 species. Plots were located on a strong elevation gradient ranging from 375 meters to 3210 meters above
358 sea level (D'Amen *et al.* 2018). We excluded 72 plots without any presences, and 103 plots with less than
359 six presences, though it is possible to run the model including these plots, so that the final dataset included
360 $n = 737$ plots. Species with less than 20 presences were excluded by the original study (D'Amen *et al.*
361 2018), though it would not have presented a problem for the quadratic GLLVM had we included those here.
362 Instead of selecting the optimal number of latent variables, we directly fitted the model to the data, using

363 the Bernoulli distribution and with $d = 2$ latent variables, for the purpose of creating an ordination diagram.
364 We tested different sets of initial values and retained the model that had the highest log-likelihood.

365 The first latent variable explained 83% of the overall residual variation in the model, of which 61%
366 was accounted for by the linear term. The length of the first ecological gradient was 3.52 (2.85-4.18, 95%
367 confidence interval). Since the first latent variable explained considerably more residual variation than the
368 second, we here focus our inference on that alone for illustration purposes. The species response curves for
369 the first latent variable are visualized in Figure 3a-c. To improve readability, species are numbered by their
370 location in the dataset, for which the corresponding names are included in Figure 4. In Figure 4 species
371 tolerances for the first latent variable are visualized, with approximate 95% confidence intervals.

372 Environmental tolerances from species of which the confidence interval for the quadratic coefficients on
373 the first latent variable did not include zero, ranged from 0.45 (*Veratrum album*) to 1.58 (*Silene vulgaris*)
374 with a median tolerance of 0.73 and a standard deviation of 0.21.

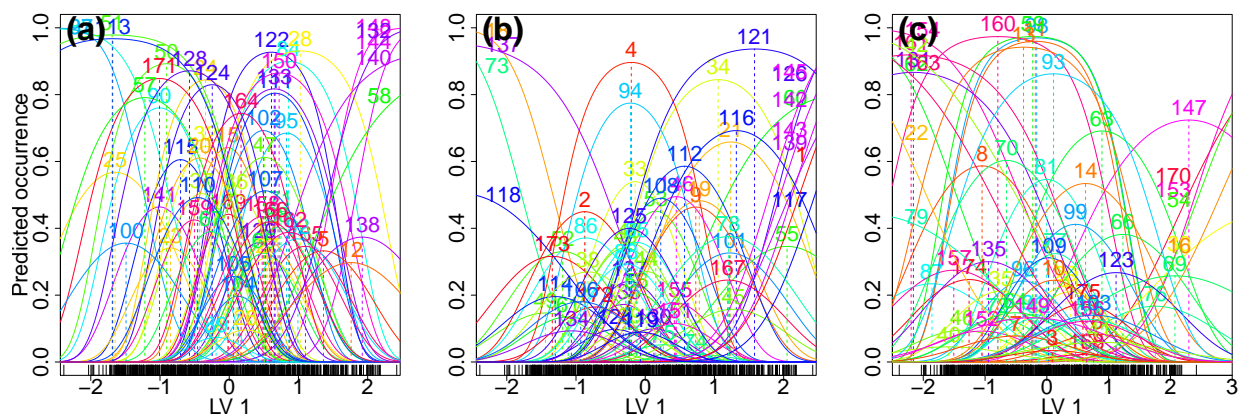


Figure 3: One-dimensional figures for the quadratic GLLVM fit to the Swiss alpine plants dataset. Each plot includes approximately one third of the species in the dataset, which have been sorted based on their variation explained, so that the first plot includes species explaining most of the variation. Plot a) represents 56% of the residual variation, plot b) represents 28% of the residual variation, and plot 2) represents 16% of the residual variation. Dashed coloured lines indicate the position of species optima, and the rug plot at the bottom indicates predicted locations of the plots. The numbers correspond with the species names in Figure 4.

375 The original dataset additionally included multiple covariates, measuring the growing degree-days above
376 zero, a moisture index, total solar radiation over the year, slope, topography, and elevation. In an attempt
377 to identify the ecological gradient represented by the first latent variable, we *post-hoc* calculated correlation
378 coefficients between the covariates and the first latent variable. From all covariates, elevation was most
379 correlated with the first latent variable (a correlation coefficient of 0.93), though this was collinear with
380 growing degree-days above zero and the moisture index. We additionally fitted two unconstrained linear

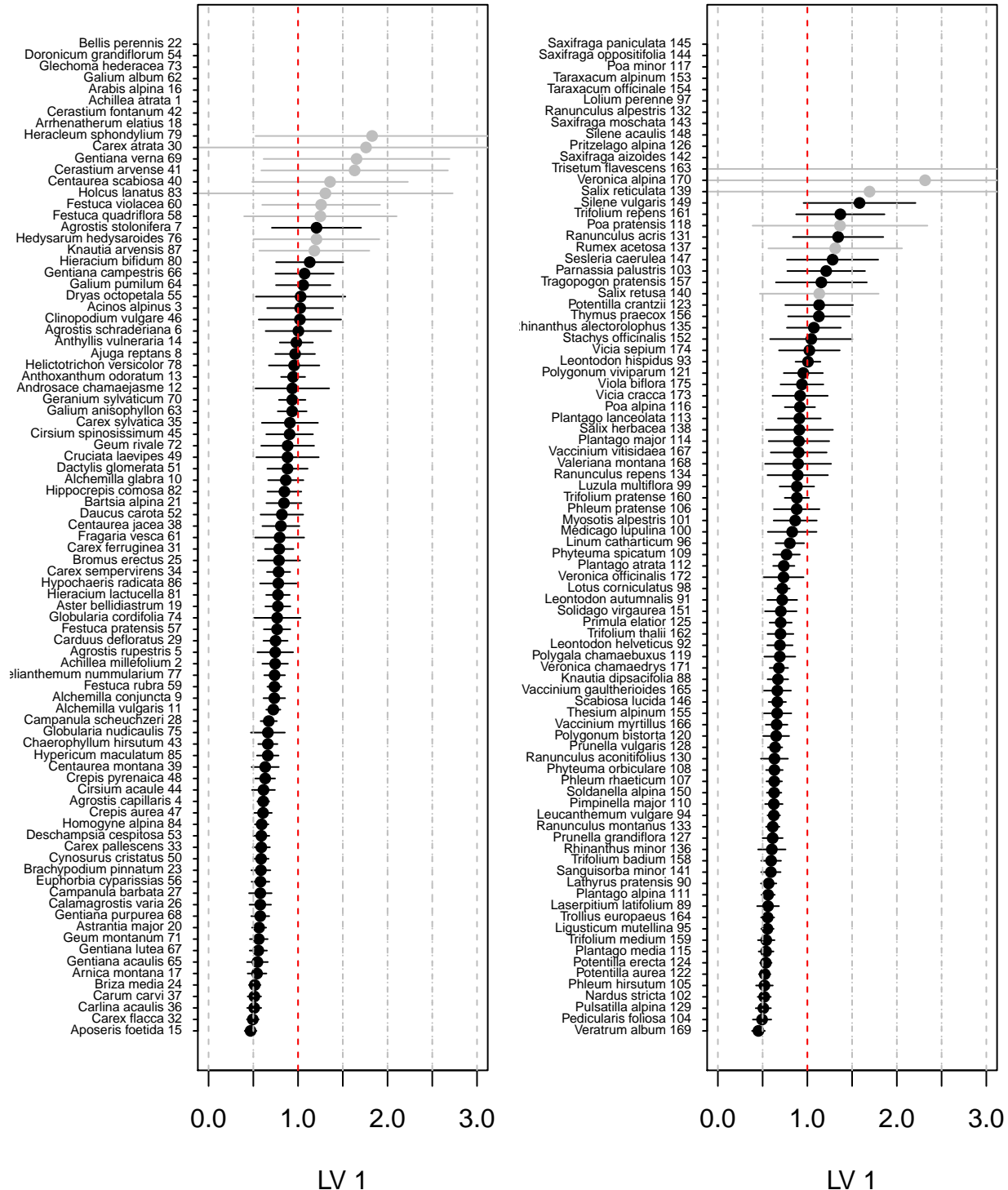


Figure 4: Species tolerances and approximate 95% confidence intervals derived using the Delta method, for the Swiss plants dataset. When optima are outside the range of the latent variable, or when tolerances cross one (indicated with a red dashed line), species have partially unobserved niches. The panels show the first and second half of species in the dataset respectively, ordered by the size of their tolerances. Species for which the confidence interval for the quadratic coefficients crosses zero are shown in grey. Species at the top of the plot, seemingly without tolerances, exhibit near linear responses, so that their tolerances are very large. Grey dashed lines are added at increments of 0.5 as visual aid.

381 GLLVMs with two latent variables, one of which included a random row intercept, and again calculated a
382 correlation coefficient between the latent variables and elevation. The linear GLLVM without row intercept
383 estimated the ecological gradient less successfully (highest correlation coefficient of -0.68), than when a row
384 intercept was included (highest correlation coefficient of -0.92). To test more explicitly for the effect of
385 elevation, we additionally fitted a quadratic GLLVM with elevation included as a covariate (both the linear
386 and quadratic term), and with two latent variables. Including the covariate reduced the residual variance
387 to 36% of that in the unconstrained model. The results presented here are from the unconstrained model,
388 though the covariate effect of the second model is presented in appendix S5, Figure S3.

389 We examined groups of plants at the extremes of the gradient, i.e. plants that had optima of minus two or
390 smaller, and plants with optima of two or larger, to further investigate whether the estimated latent variable
391 from the quadratic GLLVM represented an elevation gradient. This approach allowed us to distinguish two
392 groups of plants, the first indicative of lowlands (see Fig. 3). In contrast, plant species included on the
393 opposite side of the latent variable were clearly indicative of alpine conditions. Here, we focus our inference
394 on the alpine plants, as those are likely to be most affected by climate change (Walther *et al.* 2005). Of
395 the alpine species, only two had confidence intervals for the quadratic coefficients that did not include zero:
396 *Dryas octopetala* and *Sesleria caerulea*. *Dryas octopetala* had a maximum probability of occurrence of 0.35,
397 a tolerance of 1.03, and an optimum at 2.05, and is as such predicted to be most prone to future changes
398 in the environment. However, Figure 4 clearly shows some species that have more narrow tolerances, thus
399 more specialised species are likely present in the dataset, though it was not possible to conclude this due to
400 the confidence intervals of species quadratic coefficients crossing zero.

401 Discussion

402 In this article, we extended the GLLVM approach of Hui *et al.* (2015), to estimate the niches of species
403 with quadratic response curves, for unobserved ecological gradients. We fitted and performed inference
404 for the quadratic GLLVM by extending the VA approach from Hui *et al.* (2017). The relation between
405 latent variable models (i.e. unobserved ecological gradients) and ecological niches has been well described
406 for classical ordination methods (ter Braak & Prentice 1988; Jongman *et al.* 1995), yet a method (either
407 classical or model-based) to perform unconstrained (residual) ordination without limiting assumptions for
408 species tolerances and maxima has not been available until now.

409 The similarity in responses of species to unobserved environments can be assessed by examining tolerances,
410 by examining an ordination diagrams for overlap in species distributions, or by using the residual correlation
411 matrix. Determining if species exhibit fully quadratic curves in response to ecological gradients, whether

412 tolerances are common for all species per ecological gradient, or if the equal tolerances assumption is suited
413 for a dataset, comes down to a problem of model selection for the quadratic GLLVM. To that end, future
414 research can further investigate approaches such as regularization (e.g., possibly extending the approach of
415 Hui *et al.* 2018), hypothesis testing, or the use of confidence intervals of the quadratic coefficients. Similar to
416 DCA, the quadratic GLLVM provides estimates of gradient length. But in contrast to DCA, where gradient
417 length is a result of a heuristic rescaling of the ecological gradient, here it is calculated from the quadratic
418 coefficients, which are estimated with approximate maximum likelihood.

419 For datasets with 50 species and 50 sites or more, the quadratic GLLVM accurately retrieved ecological
420 gradients and species-specific parameters, though for continuous responses or counts it is possible to accu-
421 rately estimate parameters with fewer species or sites. In general, when fitting the quadratic GLLVM to
422 binary or ordinal responses, more information is required than for other data types (similarly as reported in
423 Yee 2004). However, this is conditional on the information content in a dataset, and the number of required
424 sites and species here should only be considered as a rough rule of thumb.

425 We studied the response of species to ecological gradients for hunting spiders in a Dutch dune ecosystem
426 (van der Aart & Smeek-Enserink 1974), and for Swiss alpine plants (D'Amen *et al.* 2017), with use of
427 the quadratic GLLVM. Various generalist species can be identified for both datasets, but as specialists are
428 more likely to be affected by future changes in the environment, their identification is of critical importance
429 to community ecology, to better focus recommendations for conservation efforts. We suggest that, for
430 the hunting spiders dataset, *Arctosa perita*, and *Arctosa lutetiana* are most vulnerable to changes in the
431 environment, and for the Swiss alpine plants dataset *Dryas octopetala* is most vulnerable to changes in the
432 environment.

433 Modelling rare species is often difficult in community ecology as few ordination methods have the ca-
434 pability to explicitly do so. The quadratic GLLVM has great potential for community ecology, as it can
435 simultaneously accommodate common (large tolerance and maxima i.e. a wide and high niche) and rare
436 species (small tolerance and maxima i.e. a narrow and low niche) with the quadratic term. The quadratic
437 GLLVM predicts species with unobserved optima, narrow niches, and small maxima will have the fewest
438 observations. Since the quadratic GLLVM includes two species-specific parameters per latent variable, and
439 thus requires more information in the data for accurate estimation of parameters than the linear GLLVM, it
440 potentially requires a large dataset to include sufficient information on rare species. However, the example
441 in this paper using the dataset of counts for hunting spiders (van der Aart & Smeek-Enserink 1974) shows
442 that the quadratic GLLVM can be feasible to fit even to small datasets. An advantage of GLLVMs is their
443 ability to use information from common species to improve estimation of parameters to describe the niches
444 of rare species. However, without penalization or borrowing information for estimation from more abundant

445 species, the parameters for species with few observations are not necessarily expected to be accurate.

446 The implementation of the quadratic GLLVM here is constrained to produce concave shapes only, though
447 it could instead be used to estimate species minima rather than maxima. However, we did not do that here,
448 as clear ecological foundations for such a model are lacking. An easy to use implementation based on the
449 the `gllvm` R package (Niku *et al.* 2019b) is available on github ([https://github.com/BertvanderVeen/gllvm-](https://github.com/BertvanderVeen/gllvm-1/tree/goGLLVM)
450 `1/tree/goGLLVM`), which will be included in the `gllvm` R package after publication.

451 Acknowledgements

452 Manuela D’Amen kindly provided the elevation covariate for the Swiss Alpine plants dataset. B.V. was
453 supported by a scholarship from the Research Council of Norway (grant number 272408/F40). F.K.C.H.
454 was supported by two Australian Research Council Discovery grants.

455 Authors contributions

456 B.V., K.A.H. and R.B.O. conceived the ideas. B.V., F.K.C.H. and R.B.O. designed the methodology. All
457 authors contributed to the writing, reviewing and editing of the draft and gave final approval for publication.

458 References

459 Björk, J.R., Hui, F.K.C., O’Hara, R.B. & Montoya, J.M. (2018). Uncovering the drivers of host-associated
460 microbiota with joint species distribution modelling. *Molecular Ecology*, **27**, 2714–2724.

461 Blanchet, F.G., Cazelles, K. & Gravel, D. (2020). Co-occurrence is not evidence of ecological interactions.
462 *Ecology Letters*, **23**, 1050–1063.

463 Borcard, D., Legendre, P. & Drapeau, P. (1992). Partialling out the Spatial Component of Ecological
464 Variation. *Ecology*, **73**, 1045–1055.

465 Burnham, K.P. & Anderson, D.R. (2002). *Model Selection and Multimodel Inference: A Practical*
466 *Information-Theoretic Approach*, Secondn. Springer-Verlag, New York.

467 D’Amen, M., Mod, H.K., Gotelli, N.J. & Guisan, A. (2018). Disentangling biotic interactions, environ-
468 mental filters, and dispersal limitation as drivers of species co-occurrence. *Ecography*, **41**, 1233–1244.

469 D’Amen, M., Mod, H.K., Gotelli, N.J. & Guisan, A. (2017). Disentangling biotic interactions, environ-
470 mental filters, and dispersal limitation as drivers of species co-occurrence. *Dryad*.

- 471 Damgaard, C., Hansen, R.R. & Hui, F.K.C. (2020). Model-based ordination of pin-point cover data:
472 Effect of management on dry heathland. *bioRxiv*, 2020.03.05.980060.
- 473 Gabriel, K.R. (1971). The biplot graphic display of matrices with application to principal component
474 analysis. *Biometrika*, **58**, 453–467.
- 475 Gauch, H.G. (1982). *Multivariate Analysis in Community Ecology*. Cambridge University Press, Cam-
476 bridge.
- 477 Halvorsen, R. (2012). A gradient analytic perspective on distribution modelling. *Sommerfeltia*, **35**,
478 1–165.
- 479 Hill, M.O. & Gauch, H.G. (1980). Detrended Correspondence Analysis: An Improved Ordination Tech-
480 nique. *Classification and Ordination: Symposium on advances in vegetation science, Nijmegen, The Nether-*
481 *lands, May 1979* (ed E. van der Maarel), pp. 47–58. Advances in vegetation science. Springer Netherlands,
482 Dordrecht.
- 483 Hui, F.K.C. (2016). Boral Bayesian Ordination and Regression Analysis of Multivariate Abundance Data
484 in r. *Methods in Ecology and Evolution*, **7**, 744–750.
- 485 Hui, F.K.C., Tanaka, E. & Warton, D.I. (2018). Order selection and sparsity in latent variable models
486 via the ordered factor LASSO. *Biometrics*, **74**, 1311–1319.
- 487 Hui, F.K.C., Taskinen, S., Pledger, S., Foster, S.D. & Warton, D.I. (2015). Model-based approaches to
488 unconstrained ordination. *Methods in Ecology and Evolution*, **6**, 399–411.
- 489 Hui, F.K.C., Warton, D.I., Ormerod, J.T., Haapaniemi, V. & Taskinen, S. (2017). Variational Approxi-
490 mations for Generalized Linear Latent Variable Models. *Journal of Computational and Graphical Statistics*,
491 **26**, 35–43.
- 492 Inoue, K., Stoeckl, K. & Geist, J. (2017). Joint species models reveal the effects of environment on
493 community assemblage of freshwater mussels and fishes in European rivers. *Diversity and Distributions*, **23**,
494 284–296.
- 495 Jamil, T. & ter Braak, C.J.F. (2013). Generalized linear mixed models can detect unimodal species-
496 environment relationships. *PeerJ*, **1**, e95.
- 497 Jongman, R., ter Braak, C. & van Tongeren, O. (Eds.). (1995). *Data analysis in community and*
498 *landscape ecology*. Cambridge university press, Cambridge.
- 499 Khatri, C.G. (1980). 14 Quadratic forms in normal variables. *Handbook of Statistics*, pp. 443–469.
500 Analysis of Variance. Elsevier.
- 501 Kristensen, K., Nielsen, A., Berg, C.W., Skaug, H. & Bell, B. (2016). TMB: Automatic Differentiation
502 and Laplace Approximation. *Journal of Statistical Software*, **70**. Retrieved from [http://arxiv.org/abs/1509.](http://arxiv.org/abs/1509.00660)
503 00660

- 504 Lacoste, É., Weise, A.M., Lavoie, M.-F., Archambault, P. & McKindsey, C.W. (2019). Changes in
505 infaunal assemblage structure influence nutrient fluxes in sediment enriched by mussel biodeposition. *Science*
506 *of The Total Environment*, **692**, 39–48.
- 507 MacArthur, R. & Levins, R. (1967). The limiting similarity, convergence, and divergence of coexisting
508 species. *The American Naturalist*, **101**, 377–385.
- 509 Nakagawa, S. & Schielzeth, H. (2013). A general and simple method for obtaining R² from generalized
510 linear mixed-effects models. *Methods in Ecology and Evolution*, **4**, 133–142.
- 511 Niku, J., Brooks, W., Herliansyah, R., Hui, F.K.C., Taskinen, S. & Warton, D.I. (2019a). Efficient
512 estimation of generalized linear latent variable models. *PLOS ONE*, **14**, e0216129.
- 513 Niku, J., Brooks, W., Herliansyah, R., Hui, F.K.C., Taskinen, S. & Warton, D.I. (2020). *Gllvm: Gener-*
514 *alized linear latent variable models*.
- 515 Niku, J., Hui, F.K.C., Taskinen, S. & Warton, D.I. (2019b). Gllvm: Fast analysis of multivariate
516 abundance data with generalized linear latent variable models in r. *Methods in Ecology and Evolution*, **10**,
517 2173–2182.
- 518 Oksanen, J. & Tonteri, T. (1995). Rate of compositional turnover along gradients and total gradient
519 length. *Journal of Vegetation Science*, **6**, 815–824.
- 520 Ovaskainen, O., Tikhonov, G., Norberg, A., Blanchet, F.G., Duan, L., Dunson, D., Roslin, T. & Abrego,
521 N. (2017). How to make more out of community data? A conceptual framework and its implementation as
522 models and software. *Ecology Letters*, **20**, 561–576.
- 523 Peres-Neto, P.R. & Jackson, D.A. (2001). How well do multivariate data sets match? The advantages of
524 a Procrustean superimposition approach over the Mantel test. *Oecologia*, **129**, 169–178.
- 525 Pollock, L.J., Tingley, R., Morris, W.K., Golding, N., O’Hara, R.B., Parris, K.M., Veski, P.A. & Mc-
526 Carthy, M.A. (2014). Understanding co-occurrence by modelling species simultaneously with a Joint Species
527 Distribution Model (JSDM). *Methods in Ecology and Evolution*, **5**, 397–406.
- 528 ter Braak, C.J. (1986). Canonical Correspondence Analysis: A New Eigenvector Technique for Multi-
529 variate Direct Gradient Analysis. *Ecology*, **67**, 1167–1179.
- 530 ter Braak, C.J.F. (1985). Correspondence Analysis of Incidence and Abundance Data: Properties in
531 Terms of a Unimodal Response Model. *Biometrics*, **41**, 859–873.
- 532 ter Braak, C.J.F. & Prentice, I.C. (1988). A Theory of Gradient Analysis. *Advances in Ecological*
533 *Research* (eds M. Begon, A.H. Fitter, E.D. Ford & A. Macfadyen), pp. 271–317. Academic Press.
- 534 Tikhonov, G., Ovaskainen, O., Oksanen, J., de Jonge, M., Opedal, O. & Dallas, T. (2020). *Hmsc:*
535 *Hierarchical model of species communities*.
- 536 Tobler, M.W., Kéry, M., Hui, F.K.C., Guillera-Aroita, G., Knaus, P. & Sattler, T. (2019). Joint species

- 537 distribution models with species correlations and imperfect detection. *Ecology*, **100**, e02754.
- 538 van der Aart, P. & Smeek-Enserink, N. (1974). Correlations between distributions of hunting spiders
539 (Lycosidae, Ctenidae) and environmental characteristics in a dune area. *Netherlands Journal of Zoology*, **25**,
540 1–45.
- 541 Walker, S.C. & Jackson, D.A. (2011). Random-effects ordination: Describing and predicting multivariate
542 correlations and co-occurrences. *Ecological Monographs*, **81**, 635–663.
- 543 Walther, G.-R., Beißner, S. & Burga, C.A. (2005). Trends in the upward shift of alpine plants. *Journal*
544 *of Vegetation Science*, **16**, 541–548.
- 545 Wang, Y., Naumann, U., Wright, S.T. & Warton, D.I. (2012). Mvabund an R package for model-based
546 analysis of multivariate abundance data. *Methods in Ecology and Evolution*, **3**, 471–474.
- 547 Warton, D.I., Blanchet, F.G., O’Hara, R.B., Ovaskainen, O., Taskinen, S., Walker, S.C. & Hui, F.K.C.
548 (2015). So Many Variables: Joint Modeling in Community Ecology. *Trends in Ecology & Evolution*, **30**,
549 766–779.
- 550 Wehrden, H.V., Hanspach, J., Bruelheide, H. & Wesche, K. (2009). Pluralism and diversity: Trends in
551 the use and application of ordination methods 1990–2007. *Journal of Vegetation Science*, **20**, 695–705.
- 552 Yee, T.W. (2004). A New Technique for Maximum-Likelihood Canonical Gaussian Ordination. *Ecological*
553 *Monographs*, **74**, 685–701.
- 554 Zurell, D., Pollock, L.J. & Thuiller, W. (2018). Do joint species distribution models reliably detect
555 interspecific interactions from co-occurrence data in homogenous environments? *Ecography*, **41**, 1812–1819.
- 556 Zurell, D., Zimmermann, N.E., Gross, H., Baltensweiler, A., Sattler, T. & Wüest, R.O. (2020). Testing
557 species assemblage predictions from stacked and joint species distribution models. *Journal of Biogeography*,
558 **47**, 101–113.
- 559 Økland, R.H. (1999). On the variation explained by ordination and constrained ordination axes. *Journal*
560 *of Vegetation Science*, **10**, 131–136.
- 561 Økland, R.H. & Eilertsen, O. (1994). Canonical Correspondence Analysis with variation partitioning:
562 Some comments and an application. *Journal of Vegetation Science*, **5**, 117–126.