

Generalized Linear Latent Variable Models for Multivariate Count and Biomass Data in Ecology

Jenni Niku^{*1}, David I. Warton^{2,3}, Francis K.C. Hui⁴, and Sara Taskinen¹

¹Department of Mathematics and Statistics, University of Jyväskylä, Finland

²School of Mathematics and Statistics and Evolution & Ecology Research Centre, The
University of New South Wales, Sydney, Australia

³School of Mathematics and Statistics, The University of New South Wales, Sydney,
Australia

⁴Mathematical Sciences Institute, The Australian National University, Australia

Abstract

In this paper we consider generalized linear latent variable models that can handle overdispersed counts and continuous but non-negative data. Such data are common in ecological studies when modelling multivariate abundances or biomass. By extending the standard generalized linear modelling framework to include latent variables, we can account for any covariation between species not accounted for by the predictors, notably species interactions and correlations driven by missing covariates. We show how estimation and inference for the considered models can be performed efficiently using the Laplace approximation method, and use simulations to study the finite-sample properties of the resulting estimates. In the overdispersed count data case, the Laplace approximated estimates perform similarly to the estimates based on variational approximation method, which is another method that provides a closed form approximation of the likelihood. In the biomass data case, we show that ignoring the correlation between taxa affects the regression estimates unfavourably. To illustrate how our methods can be used in unconstrained ordination and in making inference on environmental variables, we apply them to two ecological datasets: abundances of bacterial species in three arctic locations in Europe and abundances of coral reef species in Indonesia.

Keywords: biomass; Laplace approximation; ordination; overdispersed count; species interactions

^{*}Corresponding author: Jenni Niku, *email:* jenni.m.e.niku@jyu.fi

1 Introduction

In many studies in community ecology, multivariate abundance data are often collected, comprising the records of a large number of interacting species at a set of observational units or sites. Such data are characterized by two main features. First, the data are high-dimensional in that the number of species, many of which may interact, is often close to or exceeding the number of sites. Second the data almost always are not or cannot be suitably transformed to be normally distributed. Instead, the most common types of responses recorded include presence-absence records, overdispersed species counts, biomass (non-negative, continuous data often with large number of zeros, representing the total mass of a species found at a site), and heavily discretized percent cover data.

As a motivating example, we consider data on diversity of plant-associated bacteria (Nissinen et al., 2012). The data consists of counts of 1276 interacting bacteria species measured from different habitats (bulk soil) in 56 sites across three locations. The study design is explained in Section 5.1 in detail. This example, which is by no means an extreme case, exhibit both of the above characteristics, with the number of species approximately 23 times that of the number of sites, and the counts being highly overdispersed with nearly half of the species present at ten or fewer sites.

Multivariate abundance data are often collected to answer a number of key questions concerning the species community. In our motivating dataset for instance, Nissinen et al. (2012) were interested in performing an ordination to visualize whether sites are similar in terms of their species composition, which could be helpful in planning future sampling designs as well as identifying the drivers of microbial community composition such as soil physiochemical properties. They were also interested in conducting multivariate inference on the associations between climate zone, environment and soil microflora on microbial communities associated with plant or with particular plant species. Such analyses have important implications to help in interpreting drivers of biological associations (bacteria-plant) as well as abiotic factors (Männistö et al., 2007; Chu et al., 2010). A model-based analysis of such data poses some major challenges not just due to the high-dimensionality and non-normality of the data, as previously discussed, but also because of the (potentially) complex between species interactions. Analogous to longitudinal data, while the observational units (sites) are often independent by design, we cannot assume that species within a unit are independent: species responses are likely to be correlated due to a host of ecological reasons, such as biotic interactions, phylogeny and missing covariates (Araújo and Luoto, 2007; Morales-Castilla et al., 2015). Ignoring the correlation between species responses may result in inflated Type I errors and too narrow confidence intervals when assessing the significance of one or more predictors in the model, and too narrow prediction intervals when extrapolating key community quantities such as species richness into new sites and/or under various

58 climate scenarios (Warton et al., 2015, 2016).

59 Over the past few years, the above challenges have spurred a variety of work into model-based joint
60 analysis of multivariate abundance data. One promising approach, as reviewed by Warton et al. (2015), is
61 generalized linear latent variable models (GLLVMs, Moustaki and Knott, 2000). This rich class of models
62 extend the basic generalized linear model framework by including one or more latent variables, with corre-
63 sponding factor loadings, as a parsimonious method of modeling any residual correlation between species not
64 accounted by the covariates. Warton et al. (2015) showed how GLLVMs overcome the challenges discussed
65 above to offer a viable approach for analyzing multivariate abundance data. Specifically, by using a factor
66 analytic type approach based on rank reduction to model the high-dimensional between species covariance
67 matrix, GLLVMs offer a viable method of constructing model-based (residual) ordination and biplots, as well
68 as conducting multivariate inference such as hypothesis testing of environmental and/or treatment effects,
69 environment-by-trait interactions, and how species interactions vary at different spatial and temporal scales;
70 see Letten et al. (2015) and Ovaskainen et al. (2016a) for recent applications of GLLVMs to multivariate
71 abundance data.

72 While a promising approach, one of the major and outstanding challenges with using GLLVMs is compu-
73 tationally efficient estimation and inference. Since the responses are not normally distributed, the marginal
74 likelihood, which involves integrating out the unknown latent variables, does not possess a closed form. This
75 problem in general has attracted much attention in the statistical literature, and below we review several
76 of the well-known methods proposed to overcome this issue. In Moustaki (1996) and Moustaki and Knott
77 (2000), GLLVMs for mixtures of binary and normal responses were fitted using Gauss-Hermite quadrature.
78 This was expanded upon by Rabe-Hesketh et al. (2002), who proposed adaptive Gauss-Hermite quadrature
79 to fit GLLVMs, allowing for normal, binomial, gamma, and Poisson distributed responses. While quadra-
80 ture in general works well for simple latent variable models, the method scales poorly with the number of
81 latent variables, and becomes computationally impractical if the number of latent variables is moderate e.g.,
82 exceeds two. Another drawback is that the method of Rabe-Hesketh et al. (2002) is only available in the
83 proprietary software STATA. More recently, Hui et al. (2016) proposed a fast variational approximation
84 method to approximate the likelihood in the case of binary, ordinal and overdispersed count data. While
85 quick, the method is rather case specific, offering only a closed approximation for specific combinations of
86 response distributions and link functions. Furthermore, little is known about the theoretical properties of
87 variational approximations as a framework e.g., the convergence rate and asymptotic normality of Gaussian
88 variational approximation estimates has been derived in only specific cases such as Poisson mixed models
89 with a random intercept (Hall et al., 2011a,b).

90 The most well-known approach for estimating GLLVMs is to apply an Expectation Maximization (EM)

91 algorithm or some variant of it, as in Sammel et al. (1997) and Hui et al. (2015). In the ecology literature
92 however, with the growing popularity in hierarchical approaches to community level modeling (Cressie et al.,
93 2009; Ovaskainen et al., 2016b), most of the applications of GLLVMs have instead employed Bayesian Markov
94 Chain Monte Carlo estimation based on the complete likelihood function (Blanchet, 2014; Ovaskainen et al.,
95 2016a; Hui, 2016). A major downside of both Markov Chain Monte Carlo and the EM algorithm estimation
96 though is that they are computationally very intensive: the E-step in the EM algorithm (still) does not
97 possess a closed form, and so some form of Monte-Carlo integration is still necessary.

98 Computational efficiency is a key requirement of methods of parameter estimation, given the sizes of
99 datasets now encountered in practice in ecology. While historically most multivariate abundance datasets
100 had a few hundred variables, modern lab-based sampling and classification techniques, such as metabarcoding
101 in Yu et al. (2012) commonly result in datasets exceeding a thousand response variables, as in our microbial
102 application. As such, the most feasible maximum likelihood approaches for fitting GLLVMs in the foreseeable
103 future are those that approximate the marginal likelihood as a closed form, in particular, a variational
104 approximation (where applicable), or as in this paper, a Laplace approximation.

105 In this paper, we propose estimating and performing inference with GLLVMs using the Laplace approx-
106 imation for overdispersed count and biomass data, motivated by multivariate abundance data in ecology.
107 Although the Laplace method is a special case of adaptive Gauss-Hermite quadrature with only one quadra-
108 ture point, one of the major advantages of the Laplace approximation is that it provides a general but
109 fully closed form approximation of the likelihood, which can be maximized efficiently even for very complex
110 models applied to high-dimensional data such as overdispersed species counts in our motivating example.
111 This article is not the first to propose the Laplace approximation for GLLVMs, but the key innovation is our
112 extension particularly to handle overdispersed counts and biomass data in ecology. Huber et al. (2004) pre-
113 viously provided a Laplace approximation of the likelihood function in the general exponential family case,
114 with mixtures of binomial and normal responses serving as examples. This was extended by Bianconcini
115 and Cagnone (2012), who proposed a fully exponential Laplace approximation method for fitting GLLVMs.
116 They also treated the general exponential family case, but focused on ordinal data in simulation studies.
117 This article differs from these previous works though in that we are motivated specifically by multivariate
118 abundance data in ecology, and provides the first Laplace approximated likelihood forms for response dis-
119 tributions appropriate for overdispersed count and biomass data. More precisely, we derive forms in the
120 case of negative binomial or zero-inflated Poisson distributions for overdispersed counts and the Tweedie
121 distribution for biomass data. To our knowledge, the Laplace approximation method has not been formally
122 considered for any of these distributions so far. Notice that the two other important response types in ecol-
123 ogy, that is, presence-absence records and heavily discretized percent cover data, can be handled with the

124 tools provided by Huber et al. (2004) for binary responses and Bianconcini and Cagnone (2012) for ordinal
 125 responses, respectively.

126 The paper is organized as follows. In Section 2, we formulate the generalized linear latent variable model
 127 framework and response distributions of interest for multivariate abundance data. In Section 3, Laplace
 128 approximations of the likelihood functions are derived, and estimation and inference based on these are
 129 discussed. Section 4 provides a simulation study to compare the performance of Laplace approximation
 130 estimates to variational approximation estimates in the case of overdispersed count data. In the case of
 131 biomass data, we empirically illustrate the detrimental effect of ignoring the correlation inherent in the
 132 responses on parameter estimates. Finally, Section 5 applies the proposed Laplace approximated GLLVMs
 133 to the microbial community data (Nissinen et al., 2012) and coral community data (Warwick et al., 1990),
 134 in both cases demonstrating how common aspects of inference such as ordination can be performed within
 135 a model-based framework via the Laplace approximation.

136 2 GLLVMs for Multivariate Abundance Data

137 Let \mathbf{Y} denote a $n \times m$ response matrix, where rows $i = 1, \dots, n$ are observational units (sites) and columns
 138 $j = 1, \dots, m$ consist of m -variate correlated responses (species). For each site $\mathbf{y}_i = (y_{i1}, \dots, y_{im})'$, a k -vector
 139 of environmental covariates, denoted here as \mathbf{x}_i , may also be recorded.

140 In GLLVMs, the mean response $\mu_{ij} = E(y_{ij})$ is regressed against a vector of $d \ll m$ latent variables,
 141 denoted as \mathbf{u}_i , along with the vector of k covariates if available. That is,

$$g(\mu_{ij}) = \eta_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}_i' \boldsymbol{\beta}_j + \mathbf{u}_i' \boldsymbol{\gamma}_j, \quad (1)$$

142 where $g(\cdot)$ is a known link function, and α_i are β_{0j} denote row effects and species-specific intercepts respec-
 143 tively. While optional, row and column effects may be included to account for differences in site and species
 144 total abundance. For example, a row effect is included to ensure that the latent variables quantify differences
 145 in species composition only, as opposed to species abundance (a combination of composition and site total
 146 abundance; see Hui et al., 2015, for more details). The vectors $\boldsymbol{\beta}_j$ and $\boldsymbol{\gamma}_j$ denote species-specific regression
 147 coefficients and loadings, that is, coefficients related to the covariates and latent variables, respectively.

148 In model (1), the term $\mathbf{u}_i' \boldsymbol{\gamma}_j$ captures any residual correlation across species not accounted for by the
 149 observed covariates \mathbf{x}_i . We assume that the latent variables are drawn from independent, standard normal
 150 distributions, $\mathbf{u}_i \sim N_d(\mathbf{0}, \mathbf{I}_d)$, where \mathbf{I}_d denotes a $d \times d$ identity matrix. The purpose of the zero mean and
 151 unit variance assumption is to fix the locations and scales of the latent variables (see Chapter 5, Skrondal

152 and Rabe-Hesketh, 2004). Also, to avoid rotation invariance and ensure parameter identifiability, we set all
 153 the upper triangular elements of $m \times d$ matrix $\mathbf{\Gamma} = (\gamma_1 \cdots \gamma_m)'$ to zero, and constrain its diagonal elements
 154 to be positive (Huber et al., 2004). It is important to emphasize that these constraints do not limit the
 155 flexibility of the GLLVM to model between species correlation: there are no restrictions on the form of the
 156 residual covariance matrix induced by (1), namely $\mathbf{\Sigma}_{\text{res}} = \mathbf{\Gamma}\mathbf{\Gamma}'$, aside from it being of reduced rank d .

157 We now study specific cases of GLLVMs of key relevance to multivariate abundance data in ecology,
 158 namely, overdispersed species counts and biomass (a continuous, non-negative value typically obtained as
 159 total mass of a species at a site).

160 2.1 Species Counts

161 Species counts are often overdispersed due to their clustered nature i.e., species tend to be found in large
 162 numbers or not at all. A standard approach is to assume a negative binomial distribution for the response,
 163 $y_{ij} \sim \text{NegBin}(\mu_{ij}, \phi_j)$, where ϕ_j is a species-specific dispersion parameter, and choose $g(\cdot)$ to be the log link
 164 function. The probability density function is given by

$$f(y_{ij}|\mathbf{u}_i, \mathbf{\Psi}) = \frac{\Gamma(y_{ij} + 1/\phi_j)}{y_{ij}!\Gamma(1/\phi_j)} \left(\frac{\mu_{ij}}{1/\phi_j + \mu_{ij}}\right)^{y_{ij}} \left(\frac{1}{1 + \mu_{ij}\phi_j}\right)^{1/\phi_j}, \quad (2)$$

165 such that $E(y_{ij}) = \mu_{ij}$ and the quadratic mean-variance relationship $V(\mu_{ij}) = \mu_{ij} + \mu_{ij}^2\phi_j$. When $\phi_j \rightarrow 0$,
 166 the response variable approaches the Poisson distribution.

167 The negative binomial distribution is often appropriate when the zeros (species absences) in the data can
 168 be explained via the same environmental filtering mechanism as the non-zero counts (Warton, 2005). But
 169 if the ecological process governing most species absences is believed to be independent of the mechanism
 170 driving the non-zero counts, then a more appropriate and common choice is a zero-inflated Poisson (ZIP)
 171 model (Welsh et al., 1996; Martin et al., 2005). A ZIP model assumes that responses are either structural
 172 zeros obtained with probability p or Poisson distributed count values obtained with probability $1 - p$. If
 173 $y_{ij} \sim \text{ZIP}(p_j, \mu_{ij})$, the probability distribution function is

$$f(y_{ij}|\mathbf{u}_i, \mathbf{\Psi}) = \begin{cases} p_j + (1 - p_j) \exp(-\mu_{ij}), & \text{if } y_{ij} = 0, \\ (1 - p_j) \exp(-\mu_{ij}) \mu_{ij}^{y_{ij}} / y_{ij}!, & \text{if } y_{ij} > 0. \end{cases}, \quad (3)$$

174 where μ_{ij} is modelled as in (1) with log link function. Here we assume the probability of extra zeros
 175 is modelled for each species separately and without reference to the covariates. Under the ZIP model,
 176 $E(y_{ij}) = \mu_{ij}(1 - p_j)$ and $\text{Var}(y_{ij}) = E(y_{ij})(1 + p_j\mu_{ij})$. When $p_j = 0$, the ZIP model reduces to the
 177 Poisson model. Finally, notice the negative binomial distribution could also be extended to account for

178 extra zeros (e.g., Welsh et al., 1996). Zero-inflated negative binomial models however can often fit poorly to
 179 overdispersed count data and can suffer from convergence problems (Warton, 2005; Rodrigues-Motta et al.,
 180 2013), and so we do not pursue such a model in this article.

181 2.2 Biomass Data

182 For biomass data, which take continuous but non-negative values, an often appropriate assumption is the
 183 Tweedie distribution (Jorgensen, 1997). For a comprehensive discussion on Tweedie models and their suit-
 184 ability for biomass data, see Foster and Bravington (2013). If y_{ij} follows a Tweedie distribution, then
 185 $E(y_{ij}) = \mu_{ij}$ and $Var(y_{ij}) = \phi_j \mu_{ij}^\nu$, where ϕ_j is a species-specific dispersion parameter and ν is a power
 186 parameter controlling the shape of the distribution. The mean-variance relationship is thus explicitly defined
 187 by Taylor’s power law (Taylor, 1961), which empirically arises under a range of ecological processes (Kendal,
 188 2004).

189 The Tweedie distribution does not possess an explicit analytic form, but the density function can be
 190 evaluated numerically. For a typical power parameter value, $1 < \nu < 2$, a Tweedie random variable follows
 191 a compound Poisson distribution, and the probability distribution function can be written as

$$f(y_{ij}; \mathbf{u}_i, \Psi) = \begin{cases} \exp\left(-\frac{\mu_{ij}^{2-\nu}}{\phi_j(2-\nu)}\right), & y = 0 \\ W(y_{ij}, \phi_j, \nu) \exp\left\{\left(\frac{y_{ij}\mu_{ij}^{1-\nu}}{1-\nu} - \frac{\mu_{ij}^{2-\nu}}{2-\nu}\right) / \phi_j\right\} / y_{ij}, & y > 0 \end{cases}, \quad (4)$$

192 where $W(y_{ij}, \phi_j, \nu) = \sum_{k=1}^{\infty} W_k$, and

$$W_k = \frac{y_{ij}^{-k\alpha} (\nu - 1)^{\alpha k}}{\phi_j^{k(1-\alpha)} (2 - \nu)^k k! \Gamma(-k\alpha)}$$

193 with $\alpha = (2-\nu)/(1-\nu)$. The function $W(y_{ij}, \phi_j, \nu)$ can be evaluated numerically using the method described
 194 in Dunn and Smyth (2005). Foster and Bravington (2013) and Dunstan et al. (2013) noted that a Tweedie
 195 distribution is equivalent to the distribution obtained by summing a Poisson number of gamma random
 196 variables. Such a parametrization makes it particularly suitable for example in analyzing marine data, e.g.,
 197 the total weight of a fish species at a site can be considered as the sum of the individual fish weights, where
 198 the number of fish caught is given by a Poisson random variable and the weight of each fish follows a gamma
 199 distribution.

3 The Laplace approximation for GLLVMs

Consider again a $n \times m$ matrix, \mathbf{Y} , of observed responses and GLLVMs as defined in equation (1). Write $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)'$, $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0m})'$, $\mathbf{B} = (\boldsymbol{\beta}_1 \dots \boldsymbol{\beta}_m)'$ and $\boldsymbol{\Gamma} = (\gamma_1 \dots \gamma_m)'$, and collect all the model parameters as a vector $\boldsymbol{\Psi} = (\boldsymbol{\alpha}, \boldsymbol{\beta}_0, \text{vec}(\mathbf{B}), \text{vec}(\boldsymbol{\Gamma}), \boldsymbol{\Phi})$, where without loss of generality $\boldsymbol{\Phi}$ is used to denote any nuisance parameters depending on the assumed distribution, i.e., ϕ_1, \dots, ϕ_m for the negative binomial and Tweedie distributions and p_1, \dots, p_m for the ZIP distribution. Here $\text{vec}(\cdot)$ is the vectorizing operator, which stacks the columns of a matrix in a column vector. Conditionally on latent variables \mathbf{u}_i , the responses y_{i1}, \dots, y_{im} at site i are assumed to be independent, such that $f(\mathbf{y}_i, \mathbf{u}_i, \boldsymbol{\Psi}) = \prod_{j=1}^m f(y_{ij} | \mathbf{u}_i, \boldsymbol{\Psi}) h(\mathbf{u}_i)$, where $h(\mathbf{u}_i) = N_d(\mathbf{0}, \mathbf{I}_d)$. The marginal distribution of \mathbf{y}_i is obtained by integrating over the distribution of \mathbf{u}_i , leading to the log-likelihood function

$$l(\boldsymbol{\Psi}) = \sum_{i=1}^n \log\{f(\mathbf{y}_i, \boldsymbol{\Psi})\} = \sum_{i=1}^n \log \left(\int \prod_{j=1}^m f(y_{ij} | \mathbf{u}_i, \boldsymbol{\Psi}) h(\mathbf{u}_i) d\mathbf{u}_i \right). \quad (5)$$

For the distributions discussed in Section 2, as well as for non-normally distributed responses in general, the marginal likelihood in (5) involves a d -dimensional integral, which cannot be solved analytically. We propose to overcome this by applying a Laplace approximation to $l(\boldsymbol{\Psi})$. The Laplace approximation for the log-likelihood in the case of the general exponential family is given in Huber et al. (2004), and is reviewed in the Appendix A. Here we focus on response types and distributions discussed in Section 2, which are frequently collected in ecology.

Consider first the negative binomial distribution which, for fixed dispersion parameters ϕ_j , is a member of the exponential family. Thus a Laplace approximation for the log-likelihood function can be derived directly from the general result of Huber et al. (2004).

Theorem 1. *The Laplace approximation \tilde{l} of the log-likelihood function in negative binomial GLLVM in (2) is given by*

$$\begin{aligned} \tilde{l}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) = \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{ \boldsymbol{\Gamma}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) \} + \sum_{j=1}^m \left\{ y_{ij} \hat{\eta}_{ij} - \left(y_{ij} + \frac{1}{\phi_j} \right) \log \{ 1 + \phi_j \exp(\hat{\eta}_{ij}) \} \right. \right. \\ \left. \left. + y_{ij} \log(\phi_j) + \log \Gamma \left(y_{ij} + \frac{1}{\phi_j} \right) - \log(y_{ij}!) - \log \Gamma \left(\frac{1}{\phi_j} \right) \right\} - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right), \end{aligned}$$

where

$$\boldsymbol{\Gamma}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) = \sum_{j=1}^m \frac{(\phi_j y_{ij} + 1) \exp(\hat{\eta}_{ij})}{\{1 + \phi_j \exp(\hat{\eta}_{ij})\}^2} \boldsymbol{\gamma}_j \boldsymbol{\gamma}_j' + \mathbf{I}_d,$$

with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_j + \hat{\mathbf{u}}_i' \boldsymbol{\gamma}_j$, and $\hat{\mathbf{u}}_i$ is the maximum of

$$Q(\boldsymbol{\Psi}, \mathbf{u}_i) = \sum_{j=1}^m \left\{ y_{ij} \eta_{ij} + y_{ij} \log(\phi_j) - \left(y_{ij} + \frac{1}{\phi_j} \right) \log \{1 + \phi_j \exp(\eta_{ij})\} + \log \Gamma \left(y_{ij} + \frac{1}{\phi_j} \right) - \log(y_{ij}!) - \log \Gamma \left(\frac{1}{\phi_j} \right) \right\} - \frac{\mathbf{u}'_i \mathbf{u}_i}{2}.$$

219 If the dispersion parameters ϕ_j are unknown as is usually the case, they can be estimated jointly with
220 the other model parameters by maximizing $\tilde{l}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i)$.

221 Next, for a ZIP model, the Laplace approximation of the log-likelihood function is given as follows. Note
222 that this is not part of the exponential family and so we cannot directly use results from Huber et al. (2004).

Theorem 2. *The Laplace approximation \tilde{l} of the log-likelihood function for the zero-inflated Poisson GLLVM in (3) is given by*

$$\begin{aligned} \tilde{l}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) = & \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{ \boldsymbol{\Gamma}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) \} + \sum_{j=1}^m \left(\log(p_j + (1-p_j) \hat{A}_{ij}) I_{(y_{ij}=0)} \right. \right. \\ & \left. \left. + \{ \log(1-p_j) - \exp(\hat{\eta}_{ij}) + y_{ij} \hat{\eta}_{ij} - \log(y_{ij}!) \} I_{(y_{ij}>0)} \right) - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right), \end{aligned}$$

where $A_{ij} = \exp\{-\exp(\eta_{ij})\}$,

$$\begin{aligned} \boldsymbol{\Gamma}(\boldsymbol{\Psi}, \hat{\mathbf{u}}_i) = & \sum_{j=1}^m \left(\exp(\hat{\eta}_{ij}) I_{(y_{ij}>0)} - \left(\frac{(1-p_j) \hat{A}_{ij} \exp(\hat{\eta}_{ij}) (\exp(\hat{\eta}_{ij}) - 1)}{p_j + (1-p_j) \hat{A}_{ij}} \right. \right. \\ & \left. \left. - \frac{(1-p_j)^2 \hat{A}_{ij}^2 \exp(2\hat{\eta}_{ij})}{(p_j + (1-p_j) \hat{A}_{ij})^2} \right) I_{(y_{ij}=0)} \right) \boldsymbol{\gamma}_j \boldsymbol{\gamma}'_j + \mathbf{I}_d, \end{aligned}$$

with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_j + \hat{\mathbf{u}}_i' \boldsymbol{\gamma}_j$ and $\hat{A}_{ij} = \exp\{-\exp(\hat{\eta}_{ij})\}$, and $\hat{\mathbf{u}}_i$ is the maximum of

$$Q(\boldsymbol{\Psi}, \mathbf{u}_i) = \sum_{j=1}^m \left(\log(p_j + (1-p_j) A_{ij}) I_{(y_{ij}=0)} + \{ \log(1-p_j) - \exp(\eta_{ij}) + y_{ij} \eta_{ij} - \log(y_{ij}!) \} I_{(y_{ij}>0)} \right) - \frac{\mathbf{u}'_i \mathbf{u}_i}{2}.$$

223 Finally for the Tweedie distribution, we have the following result.

224 **Theorem 3.** *A Laplace approximation \tilde{l} of the log-likelihood function in Tweedie GLLVM in (4),*

is given by

$$\begin{aligned} \tilde{l}(\Psi, \hat{\mathbf{u}}_i) &= \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{ \Gamma(\Psi, \hat{\mathbf{u}}_i) \} + \sum_{j=1}^m \left[\left\{ \log \hat{W}(y_{ij}, \phi_j, \nu) - \log(y_{ij}) \right\} I_{(y_{ij}=0)} \right. \right. \\ &\quad \left. \left. + \frac{1}{\phi_j} \left(\frac{y_{ij} \exp\{(1-\nu)\hat{\eta}_{ij}\}}{1-\nu} - \frac{\exp\{(2-\nu)\hat{\eta}_{ij}\}}{2-\nu} \right) \right] - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right), \end{aligned}$$

where

$$\Gamma(\Psi, \hat{\mathbf{u}}_i) = \sum_{j=1}^m \frac{1}{\phi_j} [(2-\nu) \exp\{(2-\nu)\hat{\eta}_{ij}\} - y_{ij}(1-\nu) \exp\{(1-\nu)\hat{\eta}_{ij}\}] \boldsymbol{\gamma}_j \boldsymbol{\gamma}_j' + \mathbf{I}_d$$

with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_j + \hat{\mathbf{u}}_i' \boldsymbol{\gamma}_j$, and $\hat{\mathbf{u}}_i$ is the maximum of

$$\begin{aligned} Q(\Psi, \mathbf{u}_i) &= \sum_{j=1}^m \left[\left\{ \log \hat{W}(y_{ij}, \phi_j, \nu) - \log(y_{ij}) \right\} I_{(y_{ij}=0)} + \frac{1}{\phi_j} \left(\frac{y_{ij} \exp\{(1-\nu)\eta_{ij}\}}{1-\nu} - \frac{\exp\{(2-\nu)\eta_{ij}\}}{2-\nu} \right) \right] \\ &\quad - \frac{\mathbf{u}'_i \mathbf{u}_i}{2}. \end{aligned}$$

225 Note a common power parameter ν is used for all species. This is done mainly for reasons of stability, as
 226 there is typically very little information within each species to estimate the power parameter, and previous
 227 studies have shown that most species tend to have very similar values of ν (Dunstan et al., 2013).

228 3.1 Estimation and Inference

229 In all of the cases above, the Laplace approximated likelihood has a fully closed form, and therefore parameter
 230 estimates, $\hat{\Psi}$, and predictions of the latent variables, $\hat{\mathbf{u}}_i$ for the GLLVM are easily obtained by using standard
 231 quasi-Newton optimization routines available in R and alternately maximizing $\tilde{l}(\Psi, \hat{\mathbf{u}}_i)$ and $Q(\Psi, \mathbf{u}_i)$ until
 232 convergence. For this, we have developed an R package `gllvm`, which is now available on GitHub and
 233 implements the framework proposed in this paper among other functionalities.

234 For Laplace's method, the asymptotic error is of order $O(m^{-1})$, where m is the number of species. The
 235 method is therefore well suited and provides a good approximation for high-dimensional abundance data
 236 where m/n is often close to or exceeds one. As discussed in Huber et al. (2004) the Laplace approximated
 237 estimates solve the M -estimation equations, thus their consistency and asymptotic normality follow under
 238 general assumptions (Chapters 6.2-6.3, Huber and Ronchetti, 2009). Furthermore, the asymptotic standard
 239 errors for $\hat{\Psi}$ are easy to compute as the observed information matrix (negative Hessian) is obtained as part
 240 of the estimation process. This allows us to construct confidence intervals as well as conduct Wald tests
 241 for the model parameters. Likelihood ratio tests are also readily available, although with the small sample

242 sizes as well as the fact that removing a covariate from the model actually removes m coefficients, their use
243 requires careful consideration. In our examples, we use instead the corrected Akaike information criterion
244 for variable selection, although this is by no means the only information criterion one could employ.

245 Regarding ordination, similar to Hui et al. (2015) we can construct an ordination plot using predicted
246 latent variables from the fitted GLLVM. The asymptotic standard errors for $\hat{\mathbf{u}}_i$ are easily obtained in a similar
247 fashion as those for $\hat{\Psi}$, and can be used for example in constructing prediction regions around ordination
248 points. In particular if $d = 2$, then $\hat{\mathbf{u}}_i$ is a pair of coordinates representing the position of the site i in a latent
249 two-dimensional indirect gradient space. Furthermore, the coefficients γ_j quantify how each species response
250 relates to the latent variables. Therefore, we can construct a model-based biplot, where the site ordinations
251 give an indication of how species composition differs across sites, while plotting the species loadings identify
252 the indicator species characterizing the sites.

253 In Section 5, we illustrate how the model-based inference discussed above using GLLVMs can be applied,
254 using two ecological datasets.

255 4 Simulation studies

256 To evaluate the finite-sample properties of estimates obtained using the Laplace approximation method, we
257 performed two simulation studies on overdispersed count and biomass data. Details on the simulation setups
258 as well as example R code are given in Appendix C.

259 4.1 Overdispersed counts

260 In the overdispersed count data case, we compared the Laplace approximation estimates to those given by
261 variational approximation method (Hui et al., 2016). To our knowledge, this is the only other maximum
262 likelihood based method currently available which can handle negative binomial GLLVMs in a computation-
263 ally feasible manner. In Hui et al. (2015, 2016), MCMC based methods and the EM algorithm were used in
264 estimation and inference respectively, but we found these methods to be computationally so intensive that
265 they could not be included for comparison. For instance, in our initial testing with the simulation setup (d)
266 below, MCMC based method took approximately 12 hours to fit the negative binomial GLLVM.

267 The simulation setup was as follows. We simulated $K = 1500$ datasets according to the negative binomial
268 model using four different sample sizes and dimensions: (a) $n = 100$ and $m = 50$, (b) $n = 50$ and $m = 100$,
269 (c) $n = 50$ and $m = 500$ and (d) $n = 50$ and $m = 1000$. Note that especially response matrices with
270 $m \gg n$ typically arise with multivariate abundance data in ecology. As a mean model, we used $\log(\mu_{ij}) =$
271 $\alpha_i + \beta_{0j} + \mathbf{u}'_i \boldsymbol{\gamma}_j$, meaning no covariates were included in the model. The true latent variables, \mathbf{u}_i , were

272 generated from the mixture of bivariate normal distributions all having covariance matrices $0.5I_2$, means
 273 $(-1, 1)$, $(2, 1.5)$ and $(0.5, -1.5)$, and proportions 0.4, 0.3 and 0.3, respectively. The sites thus exhibit a
 274 clustering on a latent variable space. The population parameters γ_j were generated so that all the elements
 275 in both columns were generated independently from a uniform distribution $U(-2, 2)$. The population row
 276 parameters α_i and species-specific parameters β_{0j} were generated from a uniform distribution $U(-1, 1)$, and
 277 the dispersion parameters were set to $\phi_j = 1$ for all species j .

Table 1: Average biases, root mean squared errors (RMSEs), coverage probabilities of 95% confidence intervals and mean CI widths for GLLVM estimates based on Laplace approximation and variational approximation methods. The true models were negative binomial GLLVMs with (a) $n = 100$ and $m = 50$, (b) $n = 50$ and $m = 100$, (c) $n = 50$ and $m = 500$ and (d) $n = 50$ and $m = 1000$.

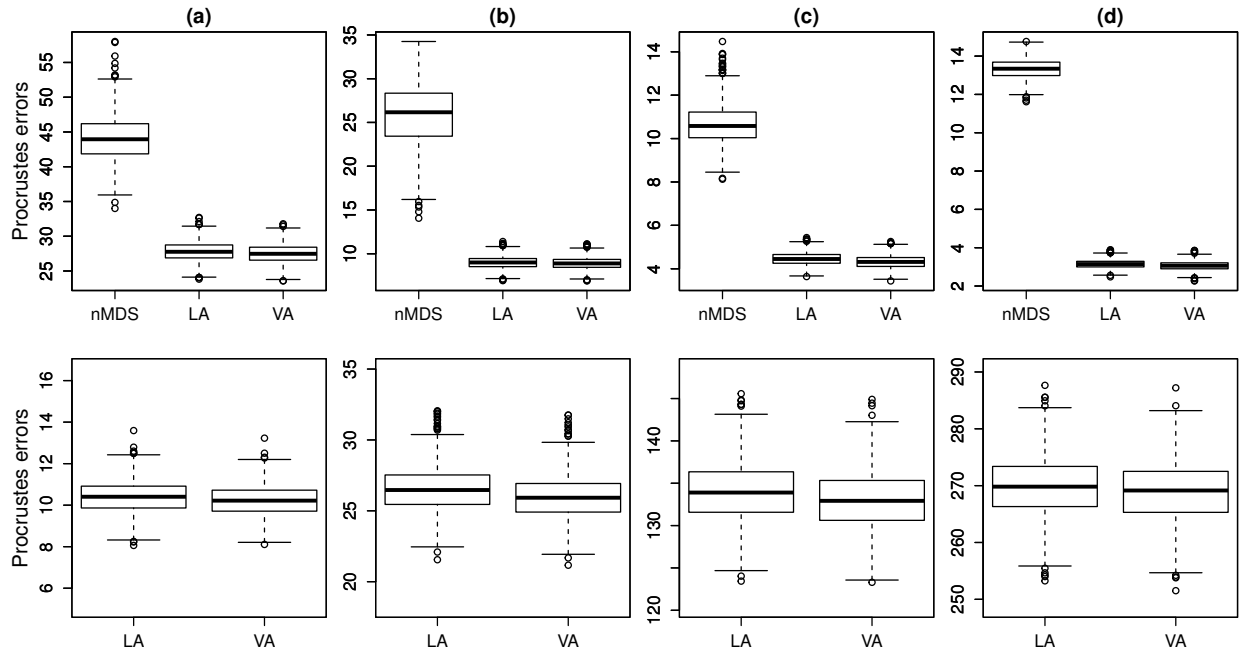
		Laplace				Variational			
		Bias	RMSE	Coverage	CI width	Bias	RMSE	Coverage	CI width
(a)	β_0	0.07	0.23	0.86	0.68	0.07	0.20	0.96	0.92
	α	-0.11	0.51	0.72	1.14	-0.11	0.39	0.85	1.14
	ϕ	-0.08	0.32	0.96	1.26	-0.03	0.30	0.96	1.16
(b)	β_0	0.10	0.30	0.88	0.97	0.10	0.30	0.95	1.27
	α	-0.19	0.30	0.95	1.21	-0.19	0.29	0.95	1.08
	ϕ	-0.12	0.42	0.97	2.30	-0.09	0.40	0.99	2.33
(c)	β_0	0.13	0.32	0.87	1.02	0.13	0.32	0.93	1.29
	α	-0.22	0.28	0.95	1.12	-0.22	0.27	0.96	1.12
	ϕ	-0.10	0.41	0.98	2.30	-0.10	0.40	0.98	2.31
(d)	β_0	0.15	0.32	0.84	0.99	0.15	0.33	0.86	1.15
	α	-0.24	0.45	0.74	1.11	-0.25	0.41	0.60	1.11
	ϕ	-0.10	0.41	0.98	2.30	-0.10	0.40	0.98	2.31

278 Table 1 lists the average biases, root mean squared errors, coverage probabilities of 95% confidence inter-
 279 vals and mean confidence interval widths for estimates of α_i , β_{0j} and ϕ_j , when the Laplace and variational
 280 approximation methods were used to fit the models assuming negative binomial distributed responses. Re-
 281 sults indicate that both methods performed similarly, with slight but noticeable biases especially for the row
 282 parameter α_i when $n \ll m$. In some cases the coverage probabilities were a lot smaller or higher than the
 283 designated level 0.95. Notice that instead of using here large-sample theory, more accurate intervals could
 284 have be obtained using, for instance, resampling based methods. This approach was however not considered
 285 due to large computational burden, and we reserve this for avenue for future empirical research.

286 To evaluate the performance of estimated γ_j and predicted latent variables, \mathbf{u}_i , the mean Procrustes
 287 errors between the estimated and true parameter values were computed (Bartholomew et al., 2011, Chapter
 288 8.4). The Procrustes error can be thought of as the mean squared error of two matrices after accounting for
 289 differences in rotation and scale. The boxplots of Procrustes errors based on Laplace approximation method
 290 and variational approximation method are given in Figure 1. To compare the performances of model based

291 ordination methods to a classical algorithmic based ordination method, non-metric multidimensional scaling
 292 (nMDS), the mean Procrustes errors between the true latent variables and nMDS ordination points were
 293 also computed. As seen in Figure 1, both model based ordination methods strongly outperform nMDS. The
 294 results based on Laplace approximation method and variational approximation method are almost equally
 295 good.

Figure 1: Comparative boxplots of Procrustes errors between true and estimated ordination points (first row) and true and estimated parameters $\hat{\gamma}_j$ (second row). Ordination points (and parameters $\hat{\gamma}_j$ when applicable) are obtained from non-metric multidimensional scaling (nMDS) and negative binomial GLLVM fitted using Laplace approximation method (LA) and variational approximation method (VA). The true model in each plot was negative binomial GLLVM with (a) $n = 100$ and $m = 50$, (b) $n = 50$ and $m = 100$, (c) $n = 50$ and $m = 500$ and (d) $n = 50$ and $m = 1000$.



296 Finally, regarding computation time, the proposed Laplace approximation method averaged 13.2, 12.1,
 297 159.4 and 609.3 seconds respectively to estimate the parameters and their standard errors using models in
 298 simulation settings (a) to (d) above. This was a substantial gain on the corresponding mean computation
 299 times for variational approximation method, which averaged 56.4, 54.9, 233.4 and 650.9 seconds, respectively.
 300 The main reason for differences in computation times is that for these setups, the variational approximation
 301 needs to estimate $5n$ variational parameters (corresponding to the mean and covariance parameters in the
 302 variational distribution) on top of the model parameters.

303 **4.2 Biomass data**

304 In the case of biomass data, we used simulations to study the effect of ignoring the correlation between
 305 taxa on regression estimates. We used only Laplace approximation method to fit the models, as there are
 306 currently no alternative maximum likelihood based methods available for fitting GLLVMs to biomass data.

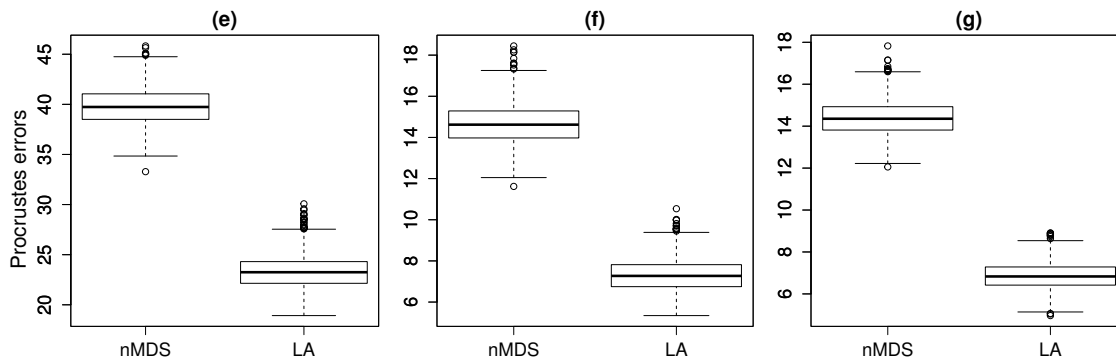
307 The simulation setup differed slightly from the one used previously for overdispersed counts. Specifically,
 308 we simulated $K = 1500$ datasets according to the Tweedie model with fixed power parameter $\nu = 1.6$ using
 309 three different sample sizes with dimensions: (e) $n = 100$ and $m = 50$, (f) $n = 50$ and $m = 100$ and (g)
 310 $n = 50$ and $m = 200$. As a mean model, we used $\log(\mu_{ij}) = \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_j + \mathbf{u}'_i \boldsymbol{\gamma}_j$, with two covariates included
 311 in the model. The true latent variables for the GLLVM, \mathbf{u}_i , were generated from a three component mixture
 312 of bivariate normal distributions all having covariance matrices $0.5I_2$, with differing means $(-1, 1)$, $(1.5, 1.5)$
 313 and $(0.5, -1.5)$, and proportions 0.4, 0.3 and 0.3 respectively. The first covariate x_{i1} was generated from the
 314 standard normal distribution and the second covariate x_{i2} from the exponential distribution with rate $\lambda = 1$.
 315 Finally, as per the overdispersed count simulation, we constructed $\boldsymbol{\gamma}_j$ such that all elements in both columns
 316 were obtained from the uniform distribution $U(-2, 2)$, while the species-specific covariate coefficients $\boldsymbol{\beta}_j$ and
 317 intercept parameters β_{0j} were chosen from the uniform distribution $U(-1, 1)$. The dispersion parameters
 were set to $\phi_j = 1$ for all species j .

Table 2: Average biases and root mean squared errors (MSEs) of Tweedie GLLVM and Tweedie GLM estimates based on Laplace approximation method. The true models were Tweedie GLLVMs with (e) $n = 100$ and $m = 50$, (f) $n = 50$ and $m = 100$ and (g) $n = 50$ and $m = 200$.

		GLLVM		GLM	
		Bias	RMSE	Bias	RMSE
(e)	β_0	0.06	0.31	1.15	1.37
	β_1	0.03	0.16	-0.09	0.18
	β_2	-0.08	0.32	0.02	0.21
	ϕ	-0.03	0.12	2.06	2.71
(f)	β_0	-0.02	0.25	0.97	1.17
	β_1	0.00	0.17	-0.20	0.32
	β_2	-0.03	0.23	0.06	0.34
	ϕ	-0.07	0.18	1.79	2.44
(g)	β_0	-0.02	0.27	0.94	1.12
	β_1	-0.00	0.17	-0.18	0.32
	β_2	-0.03	0.25	0.06	0.34
	ϕ	-0.07	0.18	1.63	2.10

318
 319 Table 2 lists the average biases and mean squared errors for regression estimates based on a Tweedie
 320 GLLVM compared to a Tweedie generalized linear model (GLM). The latter does not include any latent
 321 variables to account for residual correlation between species i.e., it assumes the species are independent after

Figure 2: Comparative boxplots of Procrustes errors between true and estimated ordination points. Ordination points are obtained from non-metric multidimensional scaling (nMDS) and Tweedie GLLVM fitted using Laplace approximation method (LA). The true model in each plot was Tweedie GLLVM with (e) $n = 100$ and $m = 50$, (f) $n = 50$ and $m = 100$ and (g) $n = 50$ and $m = 200$.



322 accounting for correlations due to the observed predictors \mathbf{x}_i . In all of the considered setups ignoring the
 323 correlation yields biased estimates with high variability, particularly for the species specific intercepts and
 324 overdispersion parameters. Additionally, Figure 2 displays the boxplots of Procrustes errors between true
 325 and predicted latent variables, as well as those between the true latent variables and ordination points given
 326 by nMDS. Again, the model based approach of GLLVM yields substantially better ordination results.

327 5 Examples

328 5.1 Microbial Community Data

329 We applied Laplace approximated GLLVMs on the bacterial species data discussed in Nissinen et al. (2012).
 330 Altogether eight different sampling sites were selected from three locations. Three of the sites were in
 331 Kilpisjärvi, Finland, three in Ny-Ålesund, Svalbard, Norway, and two in Mayrhofen, Austria. From each
 332 sampling site, several soil samples were taken and their bacterial species were recorded. The data consist of
 333 $m = 1276$ bacterial species counts measured from $n = 56$ sites. The sites can be considered as independent
 334 from each other since bacterial communities are known to be very location specific. As many of the species
 335 were observed only in few sites, we decided to exclude such rare species and considered only species present
 336 at five or more sites. This reduced the number of species to $m = 985$. In addition to bacteria counts, three
 337 continuous environmental variables (pH, available phosphorous and soil organic matter) were measured from
 338 each soil sample.

339 In order to study whether the effect of environmental variables is seen in an unconstrained ordination
 340 plot, we first considered a generalized linear latent variable model with two latent variables and no predictors,

341 and constructed an ordination plot based on the predicted latent variables. Due to small sample size, the
 342 corrected Akaike information criterion, AIC_c , was used for selecting which count distribution was most
 343 appropriate for the data (Burnham and Anderson, 2002). The values for AIC_c (scaled by n and subtracted
 344 by 1942) based on the Poisson, negative binomial and ZIP models are given in the first column of Table 3,
 345 with results indicating that the negative binomial model fitted the data best. The ZIP model outperformed
 346 the model assuming Poisson counts.

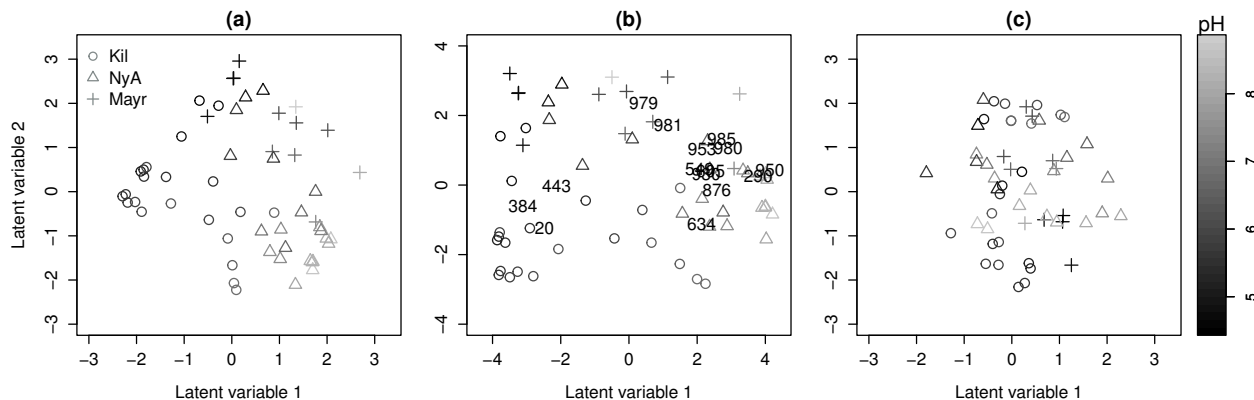
Table 3: Values of AIC_c (scaled by n and subtracted by 1942) for Poisson, negative binomial (NB) and ZIP GLLVMs (1) without covariate, (2) with pH as a covariate, (3) with pH, soil organic matter and phosphorous as covariates, (4) with pH included along with a site effect and (5) with all three soil covariates included along with a site effect.

	(1)	(2)	(3)	(4)	(5)
Poisson	771	674	463	395	244
NB	178	150	86	59	0
ZIP	630	547	377	311	189

347 The ordination of sites based on negative binomial GLLVM is plotted in Figure 3(a). The sites are
 348 coloured according to their pH values. A very clear gradient in the pH values of sites is observed, while
 349 there was less evidence of such a pattern with the two other soil variables (see Figure B1 in Appendix B). In
 350 addition, the ordination points are (also) labeled according to the sampling location (Kilpisjärvi, Ny-Ålesund
 351 and Innsbruck), and it is clear that the sites differed in terms of species composition. In Figure 3(b), a biplot
 352 based on generalized linear latent variable model is given. Here indices of the 15 species with largest factor
 353 loadings are added in the (rotated) ordination plot in Figure 3(a). The biplot suggests a small set of indicator
 354 species which prefer sites with low pH values and a larger set of indicator species for high pH sites.

355 In order to study whether the environmental variables alone are capable of explaining the variation in
 356 species composition across sites, we included them as explanatory variables in the GLLVM. Points estimates
 357 with 95% confidence intervals are plotted in Figure B2 in Appendix B, and indicate that pH value was the
 358 main covariate affecting the species composition. The corresponding ordination plots are given in Figure B3
 359 in Appendix B, and they indicate that even though the effect of environmental variables on ordination
 360 vanishes, the ordination still exhibits a sampling location effect. Several Kilpisjärvi sites in particular seem
 361 to be different from the others. To account for this, we further added the sampling location as a categorical
 362 covariate into the model. The resulting ordination plot in Figure 3(c) shows that there is no visible pattern
 363 in sampling location anymore. As the figure uses the same scale as plots in Figure 3(a), it is clear that a lot
 364 of covariation in ordination is explained by the covariates included in the model. When comparing nested
 365 models, in particular, the model with environmental covariates to the null model, and the model with all
 366 covariates to the model with environmental covariates, the deviances are 5144.6 and 4830.1, respectively,

Figure 3: (a) The ordination plot of $n = 56$ sites based on generalized linear latent variable model without any covariates assuming negative binomial distributed responses. (b) The biplot, where 15 species with the largest factor loadings (in terms of distance from the origin) are printed on top of the (rotated) site ordination to illustrate indicator species for sites with low and high pH values. (c) The ordination plot based on generalized linear latent variable model with environmental variables and sampling location as covariates. The plot (c) uses the same scale as Figure (a) to emphasize the reduction in variation. The sites in ordination plots are coloured according to their pH values and labeled according to the sampling site.



367 suggesting that about 6% of the total covariation is due to environmental covariates based on the marginal log-
 368 likelihood. Notice that changes in log-likelihood are not the only approach to quantifying variance explained,
 369 and other methods like extensions of pseudo R^2 are possible (see for instance recent work by Nakagawa and
 370 Schielzeth, 2013, for the case of generalized linear mixed models). Notice also that the corrected AIC_c picks
 371 the model with these covariates i.e., the negative binomial GLLVM with all three covariates and sampling
 372 location, as the best model (Table 3).

373 Finally, as a diagnostic tool, we plotted Dunn-Smyth residuals (Dunn and Smyth, 1996) against linear
 374 predictors for Poisson, zero-inflated Poisson and negative binomial GLLVM models with pH, soil organic
 375 matter, phosphorous and site as covariates. The plots in Figure B4 in Appendix B show residuals for 100
 376 randomly selected species to make any patterns in the plots more apparent. Specifically, the plot for the
 377 Poisson model displays a fan-shaped pattern, which means that the model is not capable of capturing the
 378 overdispersion in the data, while the plot for the ZIP model displays skew with a lowess curve showing a
 379 positive trend in residuals. By contrast, the Dunn-Smyth residuals given by negative binomial GLLVMs are
 380 uniformly distributed around zero indicating an appropriate fit to the data.

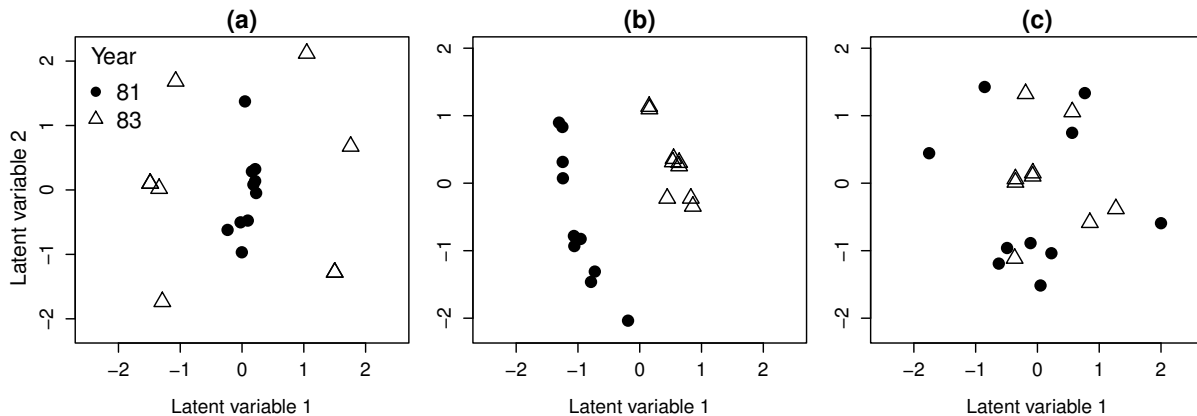
381 5.2 Coral data

382 As the second example, we consider abundances of coral reef species collected in Tikus island, Indone-
 383 sia (Warwick et al., 1990). The abundance of each reef species was measured as the length (in centimetres)

384 of a ten metre transect which intersected with the species. The data were collected during 1981-1988, but in
 385 this example we only consider measurements taken in 1981 and in 1983. The reason for this is that there was
 386 an El Niño event in 1982-1983 causing a tenfold decrease in site total abundance between the two sampling
 387 times. The aim is to study whether this event had any effect on the community structure, beyond the effect
 388 on total abundance. We consider species with more than four presences over the two years. Also one record
 389 for a site in 1983 that contained no presences was removed. The final data set thus contains $n = 19$ sites
 390 and $m = 18$ species.

391 Warwick et al. (1990) applied non-metric multidimensional scaling on this data and concluded that stress
 392 due to El Niño event increases variability in coral communities; see also Figure 4(a). Later Hui et al.
 393 (2015) applied GLLVM based ordination methods to the corresponding, converted presence-absence data,
 394 and showed that there was in fact no evidence of a difference in dispersion across the two sampling times. We
 395 now repeat their analyses using a GLLVM assuming Tweedie distributed responses. The power parameter ν
 396 was estimated using a profile likelihood approach, testing several different parameter values and selecting the
 397 one ($\nu = 1.1$) which maximised the profile likelihood. At first, the generalized linear latent variable model
 398 without site effects was fitted to produce an ordination of species abundance, i.e., including effects on total
 399 abundance as well as on relative abundance. The ordination plot in Figure 4(b) exhibits a clear location
 400 difference between coral compositions in 1981 and 1983, reflecting the El Niño event. Secondly, a GLLVM
 401 with site effects was fitted in order to study ordinations of species composition. The results in Figure 4(c)
 402 indicate that the species compositions did not change between the two sampling times.

Figure 4: The ordination plots of $n = 19$ sites based on (a) non-metric multidimensional scaling (b) Tweedie GLLVM without site effect and (c) Tweedie GLLVM with site effects. The sites in ordination plots are labeled according to the year the data was collected.



403 In Figure B5 in Appendix B the residual plots are given for the GLLVM models (b) and (c).

6 Discussion

In this paper we illustrated how generalized linear latent variable models can be used to model multivariate abundance data and biomass data, that is, data common in ecological studies. When modeling multivariate abundance data (overdispersed counts), we assumed negative binomial or zero-inflated Poisson models for responses. For biomass data (continuous but non-negative data) the Tweedie distributed responses were assumed. Notice however that these distributions just serve as examples and the method can be tailored to handle any response distribution.

Although the generalized linear latent variable models are straightforward to derive, the major challenge is the lack of computationally efficient estimation tools. In this paper, we used the Laplace approximation method for the estimation and inference. The general form for the Laplace approximation in case of exponential family is given in Huber et al. (2004), and we have extended this to the zero-inflated Poisson, negative binomial and Tweedie distributions cases, which involve additional nuisance parameters. Other case-by-case extensions may sometimes be required, e.g. to handle ordinal data, and one could argue that a disadvantage of the Laplace method is the need for case-by-case derivation of estimation algorithms. In such case, automated differentiation offers a way forward in this regard, e.g. the Template Model Builder software (Kristensen et al., 2016) can potentially simplify estimation procedures, as it requires specification of the complete likelihood only, and implementation is based on C++ code. More importantly however, such general software nevertheless employs the same Laplace approximation considered in this article as the basis for estimation and inference in GLLVMs.

Simulation studies indicated that such estimation method performs well when modeling overdispersed counts and continuous, non-negative data. However, as shown in Joe (2008) the Laplace approximation can become less adequate when the conditional distributions of the responses are highly discrete. In such settings, such as for binary and ordinal responses, we may consider other approximations method e.g. the variational approximation approach as in Hui et al. (2016). All these choices are available in R package `gllvm`, which is associated with this article. In our two examples we illustrated how generalized linear latent variable models can be applied to produce ordination plots as well as to make inferences on environmental covariates on species communities.

The generalized latent variable model considered in this paper can be generalized in several ways. If q trait covariates \mathbf{t}_j are also recorded and one wishes to study the environmental-trait interaction, a simple way to do it is via model $g(\mu_{ij}) = \alpha_i + \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_e + \text{vec}'(\mathbf{B}'_I)(\mathbf{t}_j \otimes \mathbf{x}_i) + \mathbf{u}'_i \boldsymbol{\gamma}_j$. Here $\boldsymbol{\beta}_e$ is now a main effect for the environment, common for all species, and \mathbf{B}'_I is an interaction matrix, which tells us how well traits explain variation in the environmental response. Notice that, as compared to (1), the above model

436 includes far less parameters to be estimated and tested. In ecology, the model (without latent variables) is
 437 known as a fourth corner model (Brown et al., 2014). Another way to reduce the number of parameters is
 438 to introduce random effects into the model. For instance, using a random rather than fixed site effect might
 439 be beneficial as, based on our simulation studies, the fixed site estimates seem to be slightly biased in the
 440 case of the latter. We will consider the fourth corner latent variable model and random effect models in our
 441 future studies.

442 Acknowledgements

443 We thank the Associate Editor and the referees for their helpful comments. We also thank Dr Manoj Kumar
 444 and Dr Riitta Nissinen for providing us the plant-microbial diversity data. JN and ST were supported by
 445 the Academy of Finland grants 251965 and 283323.

446 A Proofs

447 A.1 Laplace approximations for the general exponential family

448 Assume that the responses y_{ij} come from the exponential family of distributions with mean $\mu_{ij} = E(y_{ij})$,
 449 and write $f(y_{ij}|\mathbf{u}_i, \Psi) = \exp\{y_{ij}a_j(\mu_{ij}) - b_j(\mu_{ij}) + c_j(y_{ij})\}$, where $a_j(\cdot)$, $b_j(\cdot)$ and $c_j(\cdot)$ are known functions,
 450 and Ψ includes all model parameters. The log-likelihood function (5) for parameter vector Ψ now equals

$$l(\Psi) = \sum_{i=1}^n \log \int \left[\prod_{j=1}^m \exp\{y_{ij} a_j(\mu_{ij}) - b_j(\mu_{ij}) + c_j(y_{ij})\} \right] \times (2\pi)^{-\frac{d}{2}} \exp\left(-\frac{1}{2}\mathbf{u}'_i \mathbf{u}_i\right) d\mathbf{u}_i,$$

451 and the Laplace approximation of the log-likelihood function is

$$\tilde{l}(\Psi, \hat{\mathbf{u}}_i) = \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{\Gamma(\Psi, \hat{\mathbf{u}}_i)\} + \sum_{j=1}^m \{y_{ij} a_j(\mu_{ij}) - b_j(\mu_{ij}) + c_j(y_{ij})\} - \frac{\hat{\mathbf{u}}'_i \hat{\mathbf{u}}_i}{2} \right),$$

452 where

$$\Gamma(\Psi, \hat{\mathbf{u}}_i) = \sum_{j=1}^m \frac{\partial^2 \{-y_{ij} a_j(\mu_{ij}) + b_j(\mu_{ij})\}}{\partial \mathbf{u}'_i \partial \mathbf{u}_i} \Bigg|_{\mathbf{u}_i = \hat{\mathbf{u}}_i} + \mathbf{I}_d,$$

453 and $\hat{\mathbf{u}}_i$ is the maximum of $Q(\Psi, \mathbf{u}_i) = (1/m) \left(\sum_{j=1}^m \log f(y_{ij}|\mathbf{u}_i; \Psi) - \mathbf{u}'_i \mathbf{u}_i / 2 \right)$ with respect to \mathbf{u}_i . The result
 454 has been proven in Huber et al. (2004).

455 A.2 Poisson responses

456 Species counts can be modelled as Poisson distributed responses, $y_{ij} \sim \text{Poisson}(\mu_{ij})$, and log link function.

457 Then $a_j(\mu_{ij}) = \log(\mu_{ij})$, $b_j(\mu_{ij}) = \mu_{ij}$, and $c_j(y_{ij}) = -\log(y_{ij}!)$. Then the following Laplace approximation

458 \tilde{l} for the log-likelihood function is obtained

$$\tilde{l}(\Psi, \hat{\mathbf{u}}_i) = \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{\mathbf{\Gamma}(\Psi, \hat{\mathbf{u}}_i)\} + \sum_{j=1}^m [y_{ij} \hat{\eta}_{ij} - \exp(\hat{\eta}_{ij}) - \log(y_{ij}!)] - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right),$$

459 where $\mathbf{\Gamma}(\Psi, \hat{\mathbf{u}}_i) = \sum_{j=1}^m \exp(\hat{\eta}_{ij}) \boldsymbol{\gamma}_j \boldsymbol{\gamma}_j' + \mathbf{I}_d$, with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}_i' \boldsymbol{\beta}_j + \hat{\mathbf{u}}_i' \boldsymbol{\gamma}_j$, and $\hat{\mathbf{u}}_i$ is the maximum of

$$Q(\Psi, \mathbf{u}_i) = \frac{1}{m} \left[\sum_{j=1}^m [y_{ij} \eta_{ij} - \exp(\eta_{ij}) - \log(y_{ij}!)] - \frac{\mathbf{u}_i' \mathbf{u}_i}{2} - \frac{d}{2} \log(2\pi) \right].$$

460 A.3 Proof of Theorem 2

461 Assume that the responses y_{ij} come from the zero-inflated Poisson distribution with mean $E(y_{ij}) = (1-p_j)\mu_{ij}$

462 and density of the form (3). The log-likelihood function (5) then equals

$$\begin{aligned} l(\Psi) &= \sum_{i=1}^n \log \left(\int \prod_{j=1}^m \exp(\log [p_j + (1-p_j) \exp\{-\exp(\eta_{ij})\}]) I_{(y_{ij}=0)} \right. \\ &\quad \left. + \{\log(1-p_j) - \exp(\eta_{ij}) + y_{ij} \eta_{ij} - \log(y_{ij}!)\} I_{(y_{ij}>0)} \right) \\ &\quad \times (2\pi)^{-\frac{d}{2}} \exp\left(-\frac{1}{2} \mathbf{u}_i' \mathbf{u}_i\right) d\mathbf{u}_i. \end{aligned}$$

463 Hence, the Laplace approximation of the log-likelihood function is

$$\begin{aligned} \tilde{l}(\Psi, \hat{\mathbf{u}}_i) &= \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{\mathbf{\Gamma}(\Psi, \hat{\mathbf{u}}_i)\} + \sum_{j=1}^m \log f(y_{ij} | \hat{\mathbf{u}}_i; \Psi) - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right) \\ &= \sum_{i=1}^n \left(-\frac{1}{2} \log \det \{\mathbf{\Gamma}(\Psi, \hat{\mathbf{u}}_i)\} + \sum_{j=1}^m \left(\log (p_j + (1-p_j) \hat{A}_{ij}) I_{(y_{ij}=0)} \right. \right. \\ &\quad \left. \left. + \{\log(1-p_j) - \exp(\hat{\eta}_{ij}) + y_{ij} \hat{\eta}_{ij} - \log(y_{ij}!)\} I_{(y_{ij}>0)} \right) - \frac{\hat{\mathbf{u}}_i' \hat{\mathbf{u}}_i}{2} \right), \end{aligned}$$

464 where

$$\begin{aligned}
\Gamma(\Psi, \hat{\mathbf{u}}_i) &= \frac{\partial^2}{\partial \mathbf{u}'_i \partial \mathbf{u}_i} \left[- \sum_{j=1}^m \log f(y_{ij} | \mathbf{u}_i; \Psi) + \frac{\mathbf{u}'_i \mathbf{u}_i}{2} \right] \Big|_{\mathbf{u}_i = \hat{\mathbf{u}}_i} \\
&= \sum_{j=1}^m \frac{\partial^2 \{ \exp(\eta_{ij}) I_{(y_{ij} > 0)} - \log(p_j + (1-p_j)A_{ij}) I_{(y_{ij} = 0)} \}}{\partial \mathbf{u}'_i \partial \mathbf{u}_i} \Big|_{\mathbf{u}_i = \hat{\mathbf{u}}_i} + \mathbf{I}_d \\
&= \sum_{j=1}^m \left[\exp(\hat{\eta}_{ij}) I_{(y_{ij} > 0)} - \left(\frac{(1-p_j)\hat{A}_{ij} \exp(\hat{\eta}_{ij})(\exp(\hat{\eta}_{ij}) - 1)}{p_j + (1-p_j)\hat{A}_{ij}} \right. \right. \\
&\quad \left. \left. - \frac{(1-p_j)^2 \hat{A}_{ij}^2 \exp(2\hat{\eta}_{ij})}{(p_j + (1-p_j)\hat{A}_{ij})^2} \right) I_{(y_{ij} = 0)} \right] \gamma_j \gamma'_j + \mathbf{I}_d,
\end{aligned}$$

465 with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}'_i \beta_j + \hat{\mathbf{u}}'_i \gamma_j$ and $\hat{A}_{ij} = \exp\{-\exp(\hat{\eta}_{ij})\}$, and $\hat{\mathbf{u}}_i$ is the maximum of $Q(\Psi, \mathbf{u}_i) =$
466 $(1/m) \left(\sum_{j=1}^m \log f(y_{ij} | \mathbf{u}_i; \Psi) - \mathbf{u}'_i \mathbf{u}_i / 2 \right)$.

467 A.4 Proof of Theorem 3

468 Assume that the responses y_{ij} come from the Tweedie distribution with mean $E(y_{ij}) = \mu_{ij}$ and density of
469 the form (4). The log-likelihood function (5) then equals

$$\begin{aligned}
l(\Psi) &= \sum_{i=1}^n \log \left(\int \prod_{j=1}^m \exp \left(- \frac{\mu_{ij}^{2-\nu}}{\phi_j (2-\nu)} \right) I_{(y_{ij}=0)} + \frac{1}{y_{ij}} \tilde{W}(y_{ij}, \phi_j, \nu) \exp \left\{ \frac{1}{\phi_j} \left(\frac{y_{ij} \mu_{ij}^{1-\nu}}{1-\nu} - \frac{\mu_{ij}^{2-\nu}}{2-\nu} \right) \right\} I_{(y_{ij} > 0)} \right) \\
&\quad \times (2\pi)^{-\frac{d}{2}} \exp \left(- \frac{1}{2} \mathbf{u}'_i \mathbf{u}_i \right) d\mathbf{u}_i.
\end{aligned}$$

470 Hence, the Laplace approximation of the log-likelihood function is

$$\begin{aligned}
\tilde{l}(\Psi, \hat{\mathbf{u}}_i) &= \sum_{i=1}^n \left(- \frac{1}{2} \log \det \{ \Gamma(\Psi, \hat{\mathbf{u}}_i) \} + \sum_{j=1}^m \log f(y_{ij} | \hat{\mathbf{u}}_i; \Psi) - \frac{\hat{\mathbf{u}}'_i \hat{\mathbf{u}}_i}{2} \right) \\
&= \sum_{i=1}^n \left(- \frac{1}{2} \log \det \{ \Gamma(\Psi, \hat{\mathbf{u}}_i) \} + \sum_{j=1}^m \left[\left\{ \log \tilde{W}(y_{ij}, \phi_j, \nu) - \log(y_{ij}) \right\} I_{(y_{ij} > 0)} \right. \right. \\
&\quad \left. \left. + \frac{1}{\phi_j} \left(\frac{y_{ij} \exp\{(1-\nu)\hat{\eta}_{ij}\}}{1-\nu} - \frac{\exp\{(2-\nu)\hat{\eta}_{ij}\}}{2-\nu} \right) \right] - \frac{\hat{\mathbf{u}}'_i \hat{\mathbf{u}}_i}{2} \right),
\end{aligned}$$

471 where

$$\begin{aligned}
\Gamma(\Psi, \hat{\mathbf{u}}_i) &= \frac{\partial^2}{\partial \mathbf{u}'_i \partial \mathbf{u}_i} \left[- \sum_{j=1}^m \log f(y_{ij} | \mathbf{u}_i; \Psi) + \frac{\mathbf{u}'_i \mathbf{u}_i}{2} \right] \Bigg|_{\mathbf{u}_i = \hat{\mathbf{u}}_i} \\
&= \sum_{j=1}^m \frac{\partial^2}{\partial \mathbf{u}'_i \partial \mathbf{u}_i} \frac{1}{\phi_j} \left(- \frac{y_{ij} \exp\{(1-\nu)\eta_{ij}\}}{1-\nu} + \frac{\exp\{(2-\nu)\eta_{ij}\}}{2-\nu} \right) \Bigg|_{\mathbf{u}_i = \hat{\mathbf{u}}_i} + \mathbf{I}_d \\
&= \sum_{j=1}^m \frac{1}{\phi_j} [(2-\nu) \exp\{(2-\nu)\hat{\eta}_{ij}\} - y_{ij}(1-\nu) \exp\{(1-\nu)\hat{\eta}_{ij}\}] \boldsymbol{\gamma}_j \boldsymbol{\gamma}'_j + \mathbf{I}_d,
\end{aligned}$$

472 with $\hat{\eta}_{ij} = \alpha_i + \beta_{0j} + \mathbf{x}'_i \boldsymbol{\beta}_j + \hat{\mathbf{u}}_i' \boldsymbol{\gamma}_j$ and $\hat{A}_{ij} = \exp\{-\exp(\hat{\eta}_{ij})\}$, and $\hat{\mathbf{u}}_i$ is the maximum of $Q(\Psi, \mathbf{u}_i) =$
473 $(1/m) \left(\sum_{j=1}^m \log f(y_{ij} | \mathbf{u}_i; \Psi) - \mathbf{u}'_i \mathbf{u}_i / 2 \right)$.

474 B Additional Application Results

Figure B1: The ordination of $n = 56$ sites based on generalized linear latent variable model without any covariates assuming negative binomial distributed responses. The sites in ordination are coloured according to their (a) soil organic matter (SOM) values and (b) phosphorous (P) values, and labelled according to the sampling site.

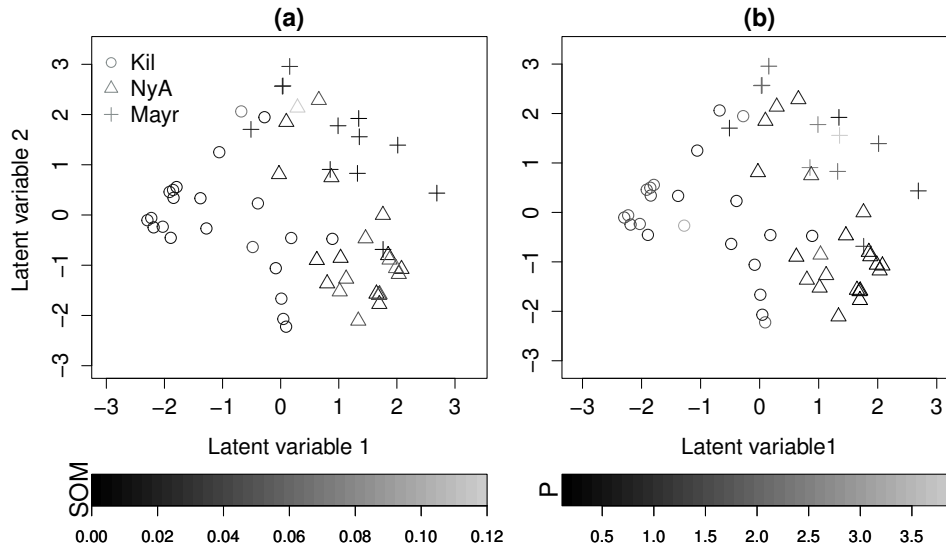


Figure B2: Ranked point estimates with 95% confidence intervals for the three environmental variables based on negative binomial GLLVM. Grey confidence intervals include the zero value.

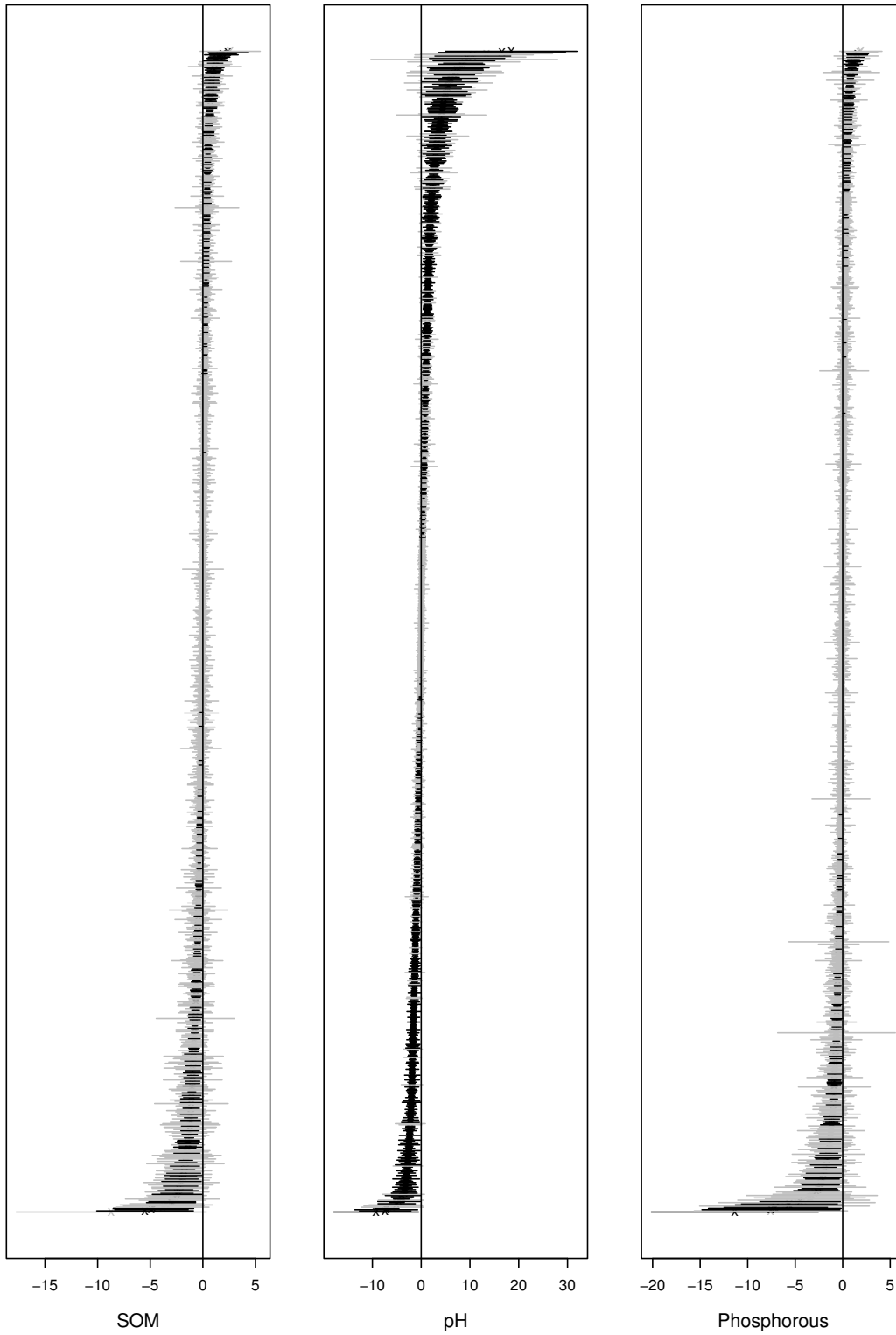


Figure B3: The ordination of $n = 56$ sites based on generalized linear latent variable model with pH, soil organic matter and phosphorous as covariates, and assuming negative binomial distributed responses. The sites in ordination are coloured according to their (a) pH values, (b) soil organic matter (SOM) values and (c) phosphorous (P) values, and labeled according to the sampling site. The effect of environmental variables vanishes, but the ordination is affected by the sampling location few Kilpisjärvi sites being different from the others what comes to species composition.

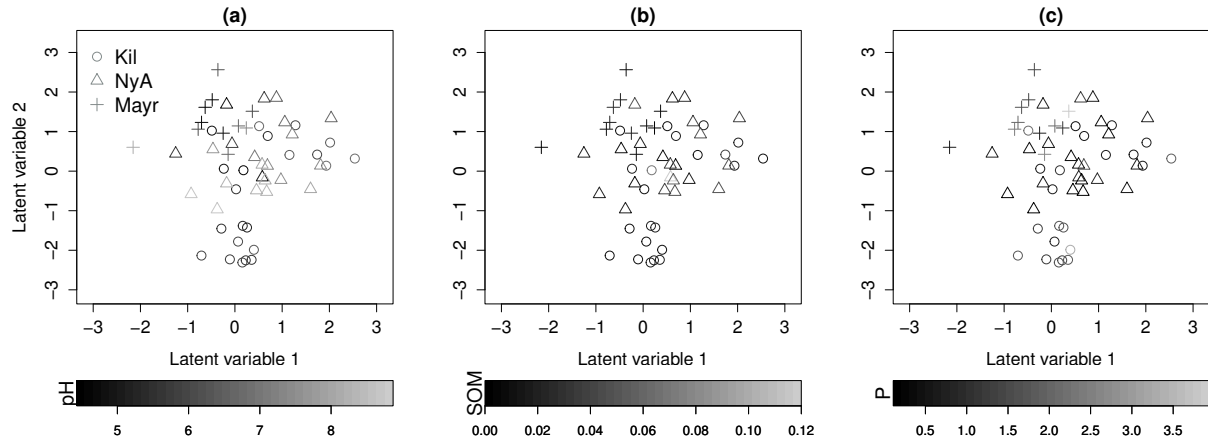


Figure B4: Dunn-Smyth residuals against linear predictors for the (a) Poisson, (b) zero inflated Poisson and (c) negative binomial GLLVM models with pH, soil organic matter, phosphorous and categorical site as covariates. Lowess curves are included in the plots.

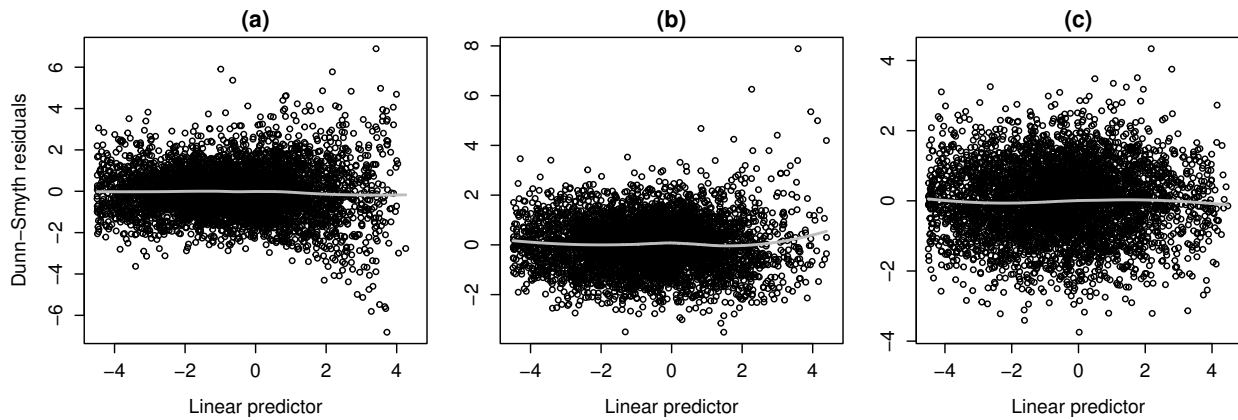
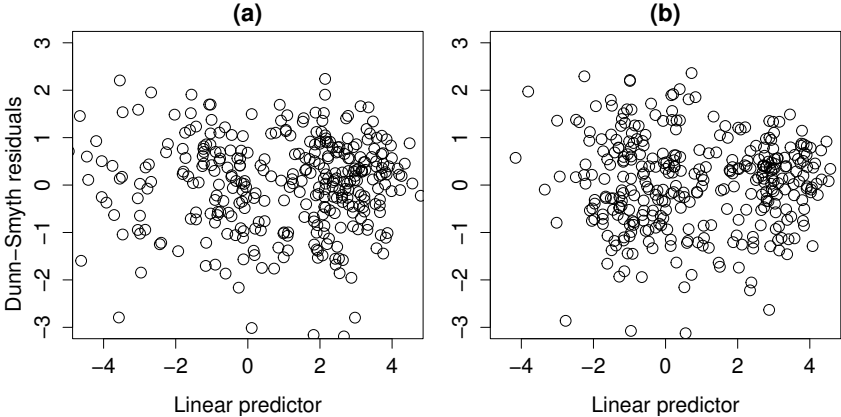


Figure B5: Dunn-Smyth residuals against linear predictors for the Tweedie models (a) without site effect and (b) with site effect.



References

- 475
- 476 Araújo, M. B. and Luoto, M. (2007). The importance of biotic interactions for modelling species distributions
477 under climate change. *Global Ecology and Biogeography*, 16:743–753.
- 478 Bartholomew, D. J., Knott, M., and Moustaki, I. (2011). *Latent variable models and factor analysis: A*
479 *unified approach*. Wiley: New York.
- 480 Bianconcini, S. and Cagnone, S. (2012). Estimation of generalized linear latent variable models via fully
481 exponential Laplace approximation. *Journal of Multivariate Analysis*, 112:183–193.
- 482 Blanchet, F. (2014). *HMSC: Hierarchical modelling of species community*. R package version 0.6-2.
- 483 Brown, A. M., Warton, D. I., Andrew, N. R., Binns, M., Cassis, G., and Gibb, H. (2014). The fourth-corner
484 solution - using predictive models to understand how species traits interact with the environment. *Methods*
485 *in Ecology and Evolution*, 5:344–352.
- 486 Burnham, K. and Anderson, D. (2002). *Model selection and multimodel inference: A practical information-*
487 *theoretic approach*. Springer.
- 488 Chu, H., Fierer, N., Lauber, C. L., Caporaso, J. G., Knight, R., and Grogan, P. (2010). Soil bacterial
489 diversity in the arctic is not fundamentally different from that found in other biomes. *Environmental*
490 *Microbiology*, 12:2998–3006.
- 491 Cressie, N., Calder, C. A., Clark, J. S., Hoef, J. M. V., and Wikle, C. K. (2009). Accounting for uncer-
492 tainty in ecological analysis: the strengths and limitations of hierarchical statistical modeling. *Ecological*
493 *Applications*, 19(3):553–570.
- 494 Dunn, P. K. and Smyth, G. K. (1996). Randomized quantile residuals. *Journal of Computational and*
495 *Graphical Statistics*, 5:236–244.
- 496 Dunn, P. K. and Smyth, G. K. (2005). Series evaluation of tweedie exponential dispersion model densities.
497 *Statistics and Computing*, 15:267–280.
- 498 Dunstan, P. K., Foster, S. D., Hui, F., and Warton, D. I. (2013). Finite mixture of regression modeling for
499 high-dimensional count and biomass data in ecology. *Journal of Agricultural, Biological and Environmental*
500 *Sciences*, 18:357–375.
- 501 Foster, S. D. and Bravington, M. V. (2013). A Poisson–Gamma model for analysis of ecological non-negative
502 continuous data. *Environmental and ecological statistics*, 20:533–552.

503 Hall, P., Ormerod, J. T., and Wand, M. (2011a). Theory of gaussian variational approximation for a poisson
504 mixed model. *Statistica Sinica*, 21:369–389.

505 Hall, P., Pham, T., Wand, M. P., Wang, S. S., et al. (2011b). Asymptotic normality and valid inference for
506 Gaussian variational approximation. *The Annals of Statistics*, 39:2502–2532.

507 Huber, P. and Ronchetti, E. (2009). *Robust Statistics*. Wiley: New York.

508 Huber, P., Ronchetti, E., and Victoria-Feser, M. (2004). Estimation of generalized linear latent variable
509 models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66:893–908.

510 Hui, F. K. C. (2016). boral–Bayesian Ordination and Regression Analysis of Multivariate Abundance Data
511 in R. *Methods in Ecology and Evolution*, 7:744–750.

512 Hui, F. K. C., Taskinen, S., Pledger, S., Foster, S. D., and Warton, D. I. (2015). Model-Based Approaches
513 to Unconstrained Ordination. *Methods in Ecology and Evolution*, 6:399–411.

514 Hui, F. K. C., Warton, D., Ormerod, J., Haapaniemi, V., and Taskinen, S. (2016). Variational Approxima-
515 tions for Generalized Linear Latent Variable Models. *Journal of Computational and Graphical Statistics*.
516 In press.

517 Joe, H. (2008). Accuracy of laplace approximation for discrete response mixed models. *Computational*
518 *Statistics & Data Analysis*, 5066–5074:52.

519 Jorgensen, B. (1997). *The Theory of Dispersion Models*. Chapman & Hall.

520 Kendal, W. S. (2004). Taylor’s ecological power law as a consequence of scale invariant exponential dispersion
521 models. *Ecological Complexity*, 1(3):193–209.

522 Kristensen, K., Nielsen, A., Berg, C., Skaug, H., and Bell, B. (2016). Tmb: Automatic differentiation and
523 laplace approximation. *Journal of Statistical Software, Articles*, 70(5):1–21.

524 Letten, A. D., Keith, D. A., Tozer, M. G., and Hui, F. K. (2015). Fine-scale hydrological niche differentiation
525 through the lens of multi-species co-occurrence models. *Journal of Ecology*, 103:1264–1275.

526 Männistö, M. K., Tirola, M., and Häggblom, M. M. (2007). Bacterial communities in arctic fjelds of finnish
527 lapland are stable but highly ph-dependent. *FEMS Microbiology Ecology*, 59:452–465.

528 Martin, T. G., Wintle, B. A., Rhodes, J. R., Kuhnert, P. M., Field, S. A., Low-Choy, S. J., Tyre, A. J., and
529 Possingham, H. P. (2005). Zero tolerance ecology: improving ecological inference by modelling the source
530 of zero observations. *Ecology letters*, 8:1235–1246.

- 531 Morales-Castilla, I., Matias, M. G., Gravel, D., and Araújo, M. B. (2015). Inferring biotic interactions from
532 proxies. *Trends in ecology & evolution*, 30(6):347–356.
- 533 Moustaki, I. (1996). A latent trait and a latent class model for mixed observed variables. *British Journal of*
534 *Mathematical and Statistical Psychology*, 49:313–334.
- 535 Moustaki, I. and Knott, M. (2000). Generalized latent trait models. *Psychometrika*, 65:391–411.
- 536 Nakagawa, S. and Schielzeth, H. (2013). A general and simple method for obtaining R² from generalized
537 linear mixed-effects models. *Methods In Ecology And Evolution*, 4:133–142.
- 538 Nissinen, R., Männistö, M., and van Elsas, J. (2012). Endophytic bacterial communities in three arctic
539 plants from low arctic fell tundra are cold-adapted and host-plant specific. *FEMS Microbiology Ecology*,
540 82:510–522.
- 541 Ovaskainen, O., Abrego, N., Halme, P., and Dunson, D. (2016a). Using latent variable models to identify large
542 networks of species-to-species associations at different spatial scales. *Methods in Ecology and Evolution*,
543 7:549–555.
- 544 Ovaskainen, O., de Knecht, H. J., and Delgado Sanchez, M. d. M. (2016b). *Quantitative Ecology and Evolu-*
545 *tionary Biology: Integrating Models with Data*. Oxford: Oxford University Press.
- 546 Rabe-Hesketh, S., Skrondal, A., and Pickles, A. (2002). Reliable estimation of generalized linear mixed
547 models using adaptive quadrature. *Stata Journal*, 2:1–21.
- 548 Rodrigues-Motta, M., Pinheiro, H. P., Martins, E. G., Araujo, M. S., and dos Reis, S. F. (2013). Multivariate
549 models for correlated count data. *Journal of Applied Statistics*, 40:1586–1596.
- 550 Sammel, M. D., Ryan, L. M., and Legler, J. M. (1997). Latent variable models for mixed discrete and
551 continuous outcomes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59:667–
552 678.
- 553 Skrondal, A. and Rabe-Hesketh, S. (2004). *Generalized Latent Variable Modeling: Multilevel, Longitudinal*
554 *and Structural Equation Models*. Chapman & Hall, Boca Raton.
- 555 Taylor, L. R. (1961). Aggregation, variance and the mean. *Nature*, 189:732 – 735.
- 556 Warton, D. I. (2005). Many zeros does not mean zero inflation: comparing the goodness-of-fit of parametric
557 models to multivariate abundance data. *Environmetrics*, 16:275–289.

- 558 Warton, D. I., Blanchet, F. G., O'Hara, R., Ovaskainen, O., Taskinen, S., Walker, S. C., and Hui, F. K.
559 (2016). Extending Joint Models in Community Ecology: A Response to Beissinger et al. *Trends in Ecology*
560 *& Evolution*, 31:737–738.
- 561 Warton, D. I., Blanchet, F. G., O'Hara, R., Ovaskainen, O., Taskinen, S., Walker, S. C., and Hui, F. K. C.
562 (2015). So many variables: Joint modeling in community ecology. *Trends in Ecology and Evolution*,
563 30:766–779.
- 564 Warwick, R., Clarke, K., and Suharsono (1990). A statistical analysis of coral community responses to the
565 1982–83 el niño in the thousand islands, indonesia. *Coral Reefs*, 8:171–179.
- 566 Welsh, A. H., Cunningham, R. B., Donnelly, C., and Lindenmayer, D. B. (1996). Modelling the abundance
567 of rare species: statistical models for counts with extra zeros. *Ecological Modelling*, 88:297–308.
- 568 Yu, D. W., Ji, Y., Emerson, B. C., Wang, X., Ye, C., Yang, C., and Ding, Z. (2012). Biodiversity soup:
569 metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring. *Methods in Ecology*
570 *and Evolution*, 3:613–623.