



Parsimoniously Fitting Large Multivariate Random Effects in `glmmTMB`

Maeve McGillicuddy
UNSW Sydney

Gordana Popovic
UNSW Sydney

Benjamin M. Bolker
McMaster University

David I. Warton
UNSW Sydney

Abstract

Multivariate random effects with unstructured variance-covariance matrices of large dimensions, q , can be a major challenge to estimate. In this paper, we introduce a new implementation of a reduced-rank approach to fit large dimensional multivariate random effects by writing them as a linear combination of $d < q$ latent variables. By adding reduced-rank functionality to the package `glmmTMB`, we enhance the mixed models available to include random effects of dimensions that were previously not possible. We apply the reduced-rank random effect to two examples, estimating a generalized latent variable model for multivariate abundance data and a random-slopes model.

Keywords: mixed models, R.

1. Introduction

When fitting a mixed effects model, it is often necessary to use a multivariate random effect with a non-diagonal covariance matrix in order to introduce sets of correlated parameters to a model. This approach is needed, for example, when fitting random-slopes models (Bolker *et al.* 2009; Asar *et al.* 2020), to account for correlation between the slope coefficient(s) and intercept terms, and when using random effects to induce correlation in multivariate data (Coull and Agresti 2000; Pollock *et al.* 2014, for example). Without imposing structure on the variance-covariance matrix, the number of parameters that need to be estimated increases quadratically with the dimension of the random effect (specifically, $q(q+1)/2$ parameters need to be estimated for an unstructured $q \times q$ covariance matrix), and estimation quickly becomes

challenging as q gets larger.

For example, in Section 3.1 we describe a study of the effect of wind farms on fish assemblages, where we count individuals of multiple species at several sites. We wish to use a multivariate random effect to estimate correlation across species. We can do this using a mixed model fitted using the **lme4** package (Bates *et al.* 2014) in R (R Core Team 2020) as follows:

```
R> glmer(abundance ~ Zone + Year + (Species + 0 | ID), family = "poisson",
+       data = windfarm)
```

or equivalently, using the **glmmTMB** package (Brooks *et al.* 2017):

```
R> glmmTMB(abundance ~ Zone + Year + (Species + 0 | ID), family = "poisson",
+         data = windfarm)
```

We show in Appendix A that if there are only two species in the data set then this approach is reasonable (and these two lines of code produce identical answers, up to machine error), but convergence issues start to be seen when there are three or more species. The complete data set that we wish to analyse has nine species, and similar types of data frequently involve many more species, sometimes thousands (Niku *et al.* 2019).

One common way to deal with high dimensionality is to use a reduced-rank approach, making simplifying assumptions that reduce the dimension of the problem to $d < q$. Reduced-rank approaches have seen considerable use in bioinformatics (e.g., Smith *et al.* 2001; Buettner *et al.* 2015) and spatial statistics (Cressie and Johannesson 2008; Banerjee *et al.* 2008). A reduced-rank approach to fitting a multivariate random effect involves writing it as a linear combination of d latent variables, often referred to as a factor analytical model (Bartholomew *et al.* 2011); in the case of exponential family responses, it is sometimes called a generalized latent variable model (GLVM: Skrondal and Rabe-Hesketh 2004).

GLVMs have been used frequently in ecology and the social sciences, early examples being models of the presence-absence of fish species (Walker and Jackson 2011) and of polytomous party choice and rankings data (Skrondal and Rabe-Hesketh 2003). GLVMs can be technically challenging to fit, but there are a number of dedicated software solutions. The **gllamm** package (Rabe-Hesketh *et al.* 2004) in **Stata** (StataCorp 2021) uses adaptive Gaussian quadrature to integrate out the latent variables. Early software written in R, with ecologists in mind, used Bayesian Markov Chain Monte Carlo (Hui 2016; Ovaskainen *et al.* 2017). Substantially faster fits can be obtained using a Laplace (Huber *et al.* 2004; Niku *et al.* 2017) or variational approximation (Hui *et al.* 2017) to the marginal likelihood, as implemented in the R package **gllvm** (Niku *et al.* 2019). These tools were written with a particular model in mind, where a multivariate random intercept is used to induce correlation across many responses, and fixed effects are used to relate the linear predictor to measured variables, although there have been some extensions to relax this constraint in different ways (van der Veen *et al.* 2021; Niku *et al.* 2021). However, this software is unable to handle a range of common sampling designs. An example considered in Section 2.2 is when the multivariate random effect is not an intercept term.

In this paper we add reduced-rank functionality to the **glmmTMB** package (Brooks *et al.* 2017) for flexible mixed effects models that can include factor analytic terms, for multivariate random effects of large dimension. The **glmmTMB** package is built as an extension to **lme4** (Bates *et al.* 2014), widely used in the applied sciences for problems involving a range of study

designs, including multi-level and repeated measures designs. The **glmmTMB** package uses a similar interface to **lme4** but exploits automatic differentiation for faster estimation of mixed effects models (Brooks *et al.* 2017). By adding reduced-rank functionality to **glmmTMB**, we enrich the class of mixed models that can be fitted to include multivariate random effects with much larger dimension than was previously possible, such that we can now routinely fit random effects of dimension in the hundreds, or perhaps in the thousands. This new class of models includes, for example, GLVMs with functionality for any of the study designs that can be analysed using **lme4**, including multi-level or repeated measures designs.

Section 2 provides an overview of the generalized linear mixed model, the factor analytic extensions to handle multivariate random effects of high dimension, and the estimation approach used in **glmmTMB** to fit such models. We describe the usage of **glmmTMB** to fit reduced-rank multivariate random effects. We analyse two different data sets, from ecology and the social sciences, in Section 3. Section 4 concludes the paper.

2. Methods

We start by introducing a generalized linear mixed model and the factor analytic variant we use to handle multivariate random effects with large dimension. Then we discuss the estimation process for these models and the interface to fit reduced-rank multivariate random effects as implemented in **glmmTMB**.

2.1. Models

Let y_{ij} be the response for $i = 1, \dots, n_j$ observational units in cluster $j = 1, \dots, m$. A vector of p fixed effect covariates, \mathbf{x}_{ij} , and q random effect covariates, \mathbf{z}_{ij} , may also be recorded for each unit. For a generalized linear mixed model (GLMM), conditional on the vector of random effects, \mathbf{b}_j , and the vector of parameters, Ψ (defined below), the responses are assumed to come from the exponential family of distributions, $f(y_{ij}|\mathbf{b}_j, \Psi) = \exp[(y_{ij}a(\eta_{ij}) - c(\eta_{ij}))/\phi + d(y_{ij}; \phi)]$ where $a(\cdot)$, $c(\cdot)$ and $d(\cdot)$ are known functions that depend on the chosen distribution f , η_{ij} are canonical parameters, and ϕ is a dispersion parameter. Then the mean response, denoted as μ_{ij} , regressed against the fixed and random covariates can be specified as

$$\eta_{ij} = g(\mu_{ij}) = \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \mathbf{z}_{ij}^\top \mathbf{b}_j \quad (1)$$

where $g(\cdot)$ is the link function, $\boldsymbol{\beta}$ is a p -dimensional vector of regression coefficients related to the covariates, $\mathbf{x}_{ij}^\top = (1, x_{1ij}, \dots, x_{pij})$, and $\mathbf{z}_{ij}^\top = (1, z_{1ij}, \dots, z_{qij})$ is the vector of random effect covariates. The unconditional distribution of the random effects, or cluster level errors, \mathbf{b}_j , is assumed to follow a multivariate normal distribution with mean zero and a parameterized $q \times q$ variance-covariance matrix, $\boldsymbol{\Sigma}$, i.e., $\mathbf{b}_j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. The variance-covariance matrix, $\boldsymbol{\Sigma}$, controls the variances of and correlations between units in clusters. The most flexible option for $\boldsymbol{\Sigma}$ is an unstructured variance-covariance matrix, which requires $q(q+1)/2$ parameters. For models with large multivariate random effects, this flexibility becomes a problem, with the number of parameters in $\boldsymbol{\Sigma}$ increasing quadratically with the size of the random effect, q .

A reduced-rank approach to fit a multivariate random effect involves expressing it as a linear function of d latent variables:

$$\mathbf{b}_j = \boldsymbol{\Lambda} \mathbf{u}_j \quad (2)$$

where \mathbf{u}_j is a vector of d latent variables, each of which is independent and standard normal, and $\mathbf{\Lambda}$ is a $q \times d$ matrix of factor loadings. The latent variables have a zero mean and unit variance, without loss of generality. Upper triangular elements of $\mathbf{\Lambda}$ are set to zero to assist with parameter identifiability, without loss of generality. Finally, we let $\Psi = \{\beta, \mathbf{\Lambda}, \phi\}$ denote the complete set of model parameters.

Given these definitions, the multivariate random effect is distributed as

$$\mathbf{b}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}\mathbf{\Lambda}^\top) \quad (3)$$

which is a reduced-rank approximation of $\mathbf{b}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$, as often seen in factor analytic models (Bartholomew *et al.* 2011; Niku *et al.* 2017). This approach makes it possible to fit large multivariate random effects, because the number of parameters required in the variance-covariance matrix of \mathbf{b}_j is now $dq - \binom{d}{2}$, which (for fixed d) increases only linearly with q .

2.2. Estimation

Conditional on the latent variables, \mathbf{u}_j , responses are assumed to be independent, hence $f(\mathbf{y}_j | \mathbf{u}_j, \Psi) = \prod_{i=1}^n f(y_{ij} | \mathbf{u}_j, \Psi)$. As the latent variables are not observed, they are integrated out, leading to the marginal log-likelihood:

$$l(\Psi) = \sum_{j=1}^m \log(f(\mathbf{y}_j, \Psi)) = \sum_{j=1}^m \log \left(\int \prod_{i=1}^n f(y_{ij} | \Psi, \mathbf{u}_j) f(\mathbf{u}_j) d\mathbf{u}_j \right). \quad (4)$$

In some cases this expression can be explicitly solved and expressed in closed form, but for non-normal distributions it does not generally have a closed form. A number of estimation methods have been proposed to approximate the marginal likelihood including Laplace approximation, numerical integration using adaptive quadrature (Skrondal and Rabe-Hesketh 2004), Monte Carlo integration (Hui *et al.* 2015) and more recently variational approximation (Hui *et al.* 2017).

We focus on the Laplace approximation of the marginal likelihood, which is widely used for GLMMs (Raudenbush *et al.* 2000) as well as GLVMs (Huber *et al.* 2004; Niku *et al.* 2017). By writing Equation 4 in the form $l(\Psi) = \sum_{j=1}^m \log \int \exp(mQ(\mathbf{y}_j, \Psi, \mathbf{u}_j)) d\mathbf{u}_j$, where

$$Q(\mathbf{y}_j, \Psi, \mathbf{u}_j) = \frac{1}{n} \left[\sum_{i=1}^n \left\{ \frac{y_{ij} a(\eta_{ij}) - c(\eta_{ij})}{\phi} + d(y_{ij}; \phi) \right\} - \frac{\mathbf{u}_j^\top \mathbf{u}_j}{2} - \frac{q}{2} \log(2\pi) \right] \quad (5)$$

we can apply Laplace's method for integral approximation around its mode, \mathbf{u}_j . Assuming $\hat{\mathbf{u}}_j$ maximises $Q(\mathbf{y}_j, \Psi, \mathbf{u}_j)$, the approximation is derived by expanding $Q(\mathbf{y}_j, \Psi, \mathbf{u}_j)$ as a second order Taylor series around the mode, $\hat{\mathbf{u}}_j$. Following Huber *et al.* (2004), a Laplace approximation of the marginal log-likelihood function of the GLVM defined in Equation 1 can be written as

$$l(\Psi) = \sum_{j=1}^m \left(-\frac{1}{2} \log \det \{G(\Psi, \hat{\mathbf{u}}_j)\} + \sum_{i=1}^n \left\{ \frac{y_{ij} a(\eta_{ij}) - c(\eta_{ij})}{\phi} + d(y_{ij}; \phi) \right\} - \frac{1}{2} \hat{\mathbf{u}}_j^\top \hat{\mathbf{u}}_j \right) \quad (6)$$

where

$$G(\Psi, \hat{\mathbf{u}}_j) = \frac{\partial^2 - Q(\mathbf{y}_j, \Psi, \mathbf{u}_j)}{\partial \mathbf{u}_j^\top \partial \mathbf{u}_j} = \sum_{i=1}^n \frac{1}{\phi} \frac{\partial^2 \{-y_{ij} a(\eta_{ij}) + c(\eta_{ij})\}}{\partial \mathbf{u}_j^\top \partial \mathbf{u}_j} \Big|_{\mathbf{u}_j = \hat{\mathbf{u}}_j} + I_q \quad (7)$$

and $\hat{\mathbf{u}}_j$ is the maximum of $Q(\mathbf{y}_j, \Psi, \mathbf{u}_j)$.

In **glmmTMB** we maximise a Laplace approximation of the marginal log-likelihood obtained from the package Template Model Builder (**TMB**) (Kristensen *et al.* 2015) in R (or more precisely, minimise the negative log-likelihood). **TMB** evaluates the Laplace approximation and its derivatives using automatic differentiation. Any gradient-based optimization method available in R can be used to do the maximization, by default **glmmTMB** uses `nlminb()`. **TMB** then uses the generalised delta method to calculate marginal standard deviations of fixed and random effects (Kass and Steffey 1989).

2.3. Software interface

A reduced rank covariance structure is specified in **glmmTMB** using `rr`. For example, suppose \mathbf{x} is a matrix of predictors with p columns, and that we want to apply p -dimensional random coefficients that take different values for different levels of a grouping variable `group` to predictors \mathbf{x} . To specify that these random coefficients are drawn from a multivariate normal distribution whose variance-covariance matrix has rank two (that is, they can be written as a linear combination of two independent latent variables) we use the `rr` random effect structure in the formula as follows

```
rr(x | group, 2)
```

So for example to model the mean of a response (\mathbf{y}) as a function of \mathbf{x} and let the coefficients of \mathbf{x} vary randomly among groups as a function of two latent variables, the formula is

```
y ~ x + rr(x | group, 2)
```

The non-negative integer after the comma in the formula specifies the number of latent variables, or rank, of the variance-covariance matrix of the multivariate random effects, which defaults to $d = 2$. The choice of the number of latent variables can be seen as a model selection problem (Hui *et al.* 2015). Model selection tools including cross-validation, or information criteria can be used to select the rank (Bartholomew *et al.* 2011).

One issue with fitting a factor-analytic model is that the likelihood function is often multimodal, which may lead to convergence to a local maximum. To overcome this we include a data-driven method based on the work of Niku *et al.* (2019) to initialise values for parameters so that the maximising algorithm starts closer to the global maximum. The default in **glmmTMB** sets parameter values to zero, or one for fixed-effect parameters for some link functions. To control the algorithm used for initialising parameters a `start_method` argument, specified as a list with components `method` and `jitter.sd`, has been added to `glmmTMBControl()`. Setting `method = "res"` fits a generalized linear model (GLM) to the data to obtain estimates of the fixed parameters, which are then used as starting values of the fixed parameters. From the fitted model, residuals are calculated for models using the Poisson, negative binomial, and binomial families using the method due to Dunn and Smyth (1996), while the internal function `dev.residuals` is used for other families. Starting values for latent variables, \mathbf{u}_j , and factor loadings, $\mathbf{\Lambda}$, are obtained by applying a reduced-rank model to the residuals from the fitted GLM. To check stability of a solution, we suggest repeating a fit multiple times and adding variation to the starting values of latent variables when `method`

= "res" by setting `jitter.sd = 0.2`, or similar, to jitter starting values for \mathbf{u}_j by a normal variate with mean zero and standard deviation `jitter.sd`.

More examples of implementing the `rr` structure are shown in Section 3. Information on the reduced-rank and other available variance-covariance structures in **glmmTMB** can be found in `vignette("covstruct", package = "glmmTMB")`.

3. Application

We present two applications to illustrate the use of a reduced-rank covariance structure in **glmmTMB**. The first example, using wind farm data, fits a generalized latent variable model to multivariate abundance data, a form of data often gathered for ecological studies. The second example presents a random-slopes model, with many random slopes, applied to data from the Progress in International Reading Literacy Study (PIRLS).

3.1. Wind farm data

The wind farm data were gathered for the Lillgrund offshore wind farm off the southern coast of Sweden, to study the effects of the wind farm on the demersal fish community (Bergstrom *et al.* 2013). This is one of the few large-scale studies assessing the effects of offshore wind farms over a long period of time, and on more than one individual fish species. This study exemplifies a BACI (before-after-control-impact) design; abundance of fish was measured before (2003) and after construction (2010), in the wind farm zone and in two reference zones (southern and northern reference zones). Sampling occurred at fishing stations within each zone; stations remained the same throughout the study. Stations which were not sampled in both time periods were omitted. The raw abundance data (Figure 1) does not show an obvious windfarm effect, although species *Strensultra* and *Oxsimpa* may show before-after differences for the south and north zones respectively. Because we are interested in the effect of wind farms, which were established between the sampling times, our interest is in the interaction between zone and year.

A multivariate random effect is required to account for correlation in species responses within a sample – we expect correlation across species due to inter-specific interactions, or response to unobserved environmental conditions (Warton *et al.* 2015) and we wish to be able to estimate positive or negative correlations. We do not have any *a priori* structure for this correlation, and for nine species we would require 45 parameters to estimate the variance-covariance matrix among species. Thus we would expect considerable instability in the fit, especially for correlation terms involving rarer species. In comparison, a reduced-rank covariance structure of rank two would require only 17 parameters.

We model abundances, y_{ijk} , as conditionally Poisson with mean μ_{ijk} , observed at $i = 1, \dots, n$ samples, $k = 1, \dots, l$ stations, for $j = 1, \dots, p$ species such that:

$$g(\mu_{ijk}) = \mathbf{x}_i^\top \boldsymbol{\beta} + \mathbf{x}_i^\top \mathbf{b}_j^{[x]} + b_k^{[s]} + b_{ijk}^{[rr]} \quad (8)$$

where $g(\cdot)$ is the link function, \mathbf{x}_i are vectors of environmental covariates specifying the intercept, zone, year and the interaction of zone and year, and $\boldsymbol{\beta}$ is a vector of fixed coefficients. We include a multivariate random effect (with $q = 4$) on the environmental variables, $\mathbf{b}_j^{[x]} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\sigma}_x^2 \mathbf{I})$ to allow the effects of each covariate to vary across species. The random intercept

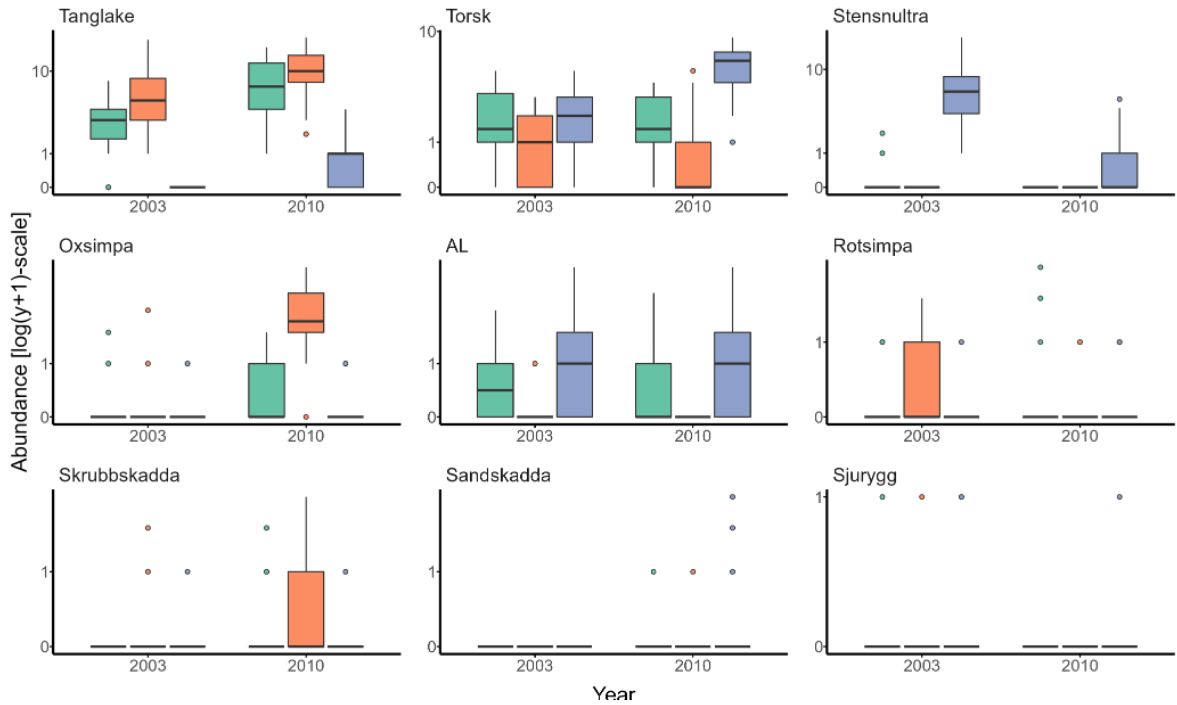


Figure 1: Boxplot of fish abundance ($\log(y + 1)$ scale) for each species at three zones, wind-farm (WF, green), north (N, orange), and south (S, lilac), before (2003) and after (2010) construction of the offshore wind farm.

for station, $b_k^{[s]} \sim \mathcal{N}(0, \sigma_s^2)$ is intended to account for paired sampling at stations, with data collected at two time points for each station. We assume each of these random effects is independent of all others and of the response (conditional on μ_{ijk}). The correlation between species is induced by the reduced-rank random effect, $b_{ijk}^{[rr]}$, assumed to satisfy

$$b_{ijk}^{[rr]} = \lambda_j \mathbf{u}_{ik} \quad (9)$$

where \mathbf{u}_{ik} is a pair (dimension $d = 2$) of latent variables, and the vector λ_j contains the corresponding factor loadings.

This model can then be fitted using the following command:

```
R> glmmTMB(abundance ~ Zone * Year + diag(Zone * Year | Species) +
+ (1 | Station) + rr(Species + 0 | ID, 2), family = "poisson",
+ control = glmmTMBControl(start_method = list(method = "res")),
+ data = windfarm)
```

The intercept was excluded from the reduced-rank random effect term (using `Species + 0` as the varying term) in order to aid interpretability of the correlation matrix discussed below. The estimated correlation matrix of the reduced-rank multivariate random effect can be obtained from the output of the `summary` method, or using the `VarCorr` method, in the same way that the estimated values for any variance-covariance matrix are returned by `glmmTMB`.

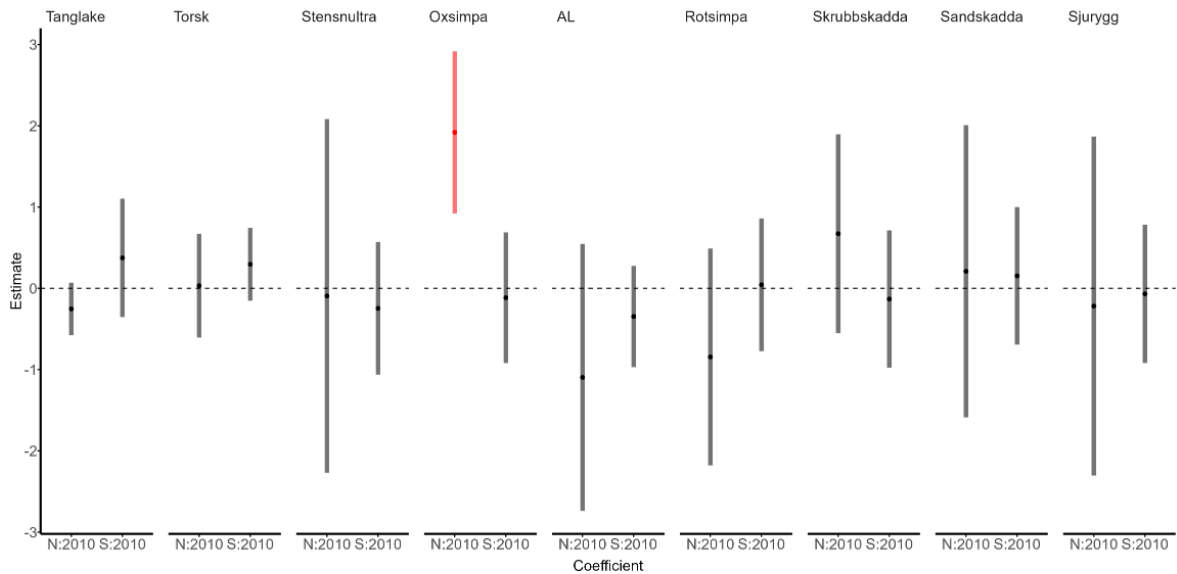


Figure 2: Conditional estimates (95% confidence interval) of the Zone by Year interaction terms for species from the diagonal random effect. The contrast between a zone (North, or South), and the Wind Farm zone in 2010 is shown.

The estimated correlation matrix of the random intercept for the wind farm data, using the rank-two model specified above, is as follows:

Conditional model:

Groups	Name	Std.Dev.	Corr							
ID	SpeciesTanglake	0.4741								
	SpeciesTorsk	0.1618	0.40							
	SpeciesStensnultra	0.6963	0.36	1.00						
	SpeciesOxsimpa	0.2092	-1.00	-0.34	-0.29					
	SpeciesAL	0.5502	-0.85	0.13	0.18	0.89				
	SpeciesRotsimpa	0.8684	0.24	0.99	0.99	-0.17	0.30			
	SpeciesSkrubbskadda	0.6541	0.43	1.00	1.00	-0.36	0.10	0.98		
	SpeciesSandskadda	1.6886	0.39	1.00	1.00	-0.32	0.14	0.99	1.00	
	SpeciesSjurygg	0.6395	0.52	0.99	0.98	-0.46	0.00	0.95	0.99	0.99

These correlations are residual correlations between species not accounted for by the covariates and other random effects in the model. For example, the correlation between species *Oxsimpa* and *AL* is 0.89 after controlling for the fixed and random covariates in the model. Note that these correlations are on the linear predictor scale (in this case the log scale), and the actual correlation observed in data is much weaker than these values, because of Poisson noise introduced by Equation 8. The marginal correlation structure is singular (for $d < p$) as it is approximated from Equation 3, where $\mathbf{\Lambda}$ is reduced rank; this is not problematic as the estimates of the factor loadings, λ_j , and latent variables, u_i , are used in further analysis. Attempting to instead fit this model using an unstructured covariance structure runs into convergence problems; the data are of insufficient quality to support estimation of 45 separate variance parameters.

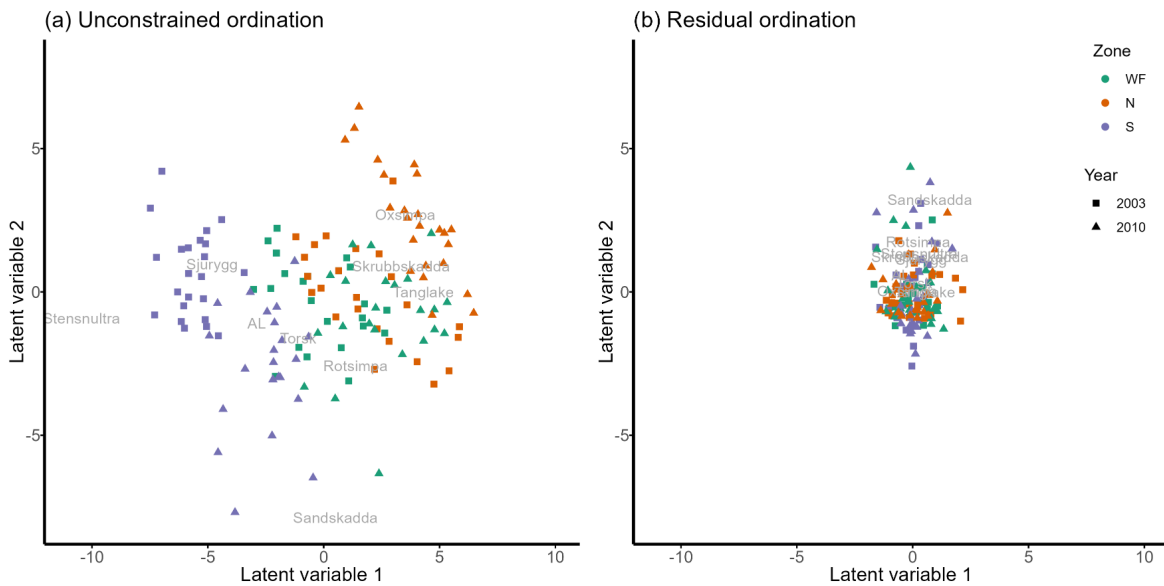


Figure 3: Ordination biplot of the wind farm data (a) for the unconstrained model, (b) after including zone, year and the interaction in the model. Zones are shown in colours, year in symbols and species factor loadings are labelled accordingly.

To test if the construction of an offshore wind farm had a significant effect on fish abundance, a parametric bootstrap analysis was conducted to test the **Zone** by **Year** interaction in both the fixed and random effect terms (code provided in Appendix B). A parametric bootstrap analysis was chosen over asymptotic tests such as Wald and likelihood ratio tests, because the asymptotic distributions of these tests are usually unknown for mixed models (Bolker *et al.* 2009). The interactions were significant ($LR = 27.35$, $p = 0.001$), indicating that for at least one of the species, mean abundance across zones has changed (in a relative sense) since construction of the wind farm (Figure 2). The estimate for Oxsimpa is clearly positive for the North-Wind Farm contrast in 2010 (Dushoff *et al.* 2019), while controlling for other covariates, although no such effect was seen for the South-Wind Farm contrast in 2010. Judging from Figure 1, this effect probably has more to do with an increase in Oxsimpa in the North Zone, rather than an effect of wind farms.

To visualise the correlations between species, an ordination biplot can be produced from the estimated latent variables and factor loadings (Figure 3). An unconstrained ordination biplot (Figure 3a), plotting the latent variables from a model without and fixed effects predictors, shows a separation in the latent variables by **Zone** and **Year**, emphasising the importance of these variables. Adding factor loadings to the plot shows us how species vary across sites, with higher abundance of a species tending to be found in sites in the same direction as the species, with respect to the origin. Oxsimpa for example could be expected to be found in high numbers in the North Zone in 2010, and Stensnultra could expect to be found in high numbers in the South Zone, especially in 2003. Figure 1 corroborates both of these results. The relative positions of the species also gives information about their correlation — because Oxsimpa and Stensnultra are negatively correlated they appear far from the origin but at opposite sides of the ordination, whereas the positively correlated Oxsimpa and Skrubbskadda are neighbours

in the ordination plot.

After fitting the model in Equation 8, which controls for the effects of **Zone**, **Year**, their interaction, and **Station**, the clustering patterns by zone and sampling time disappear (Figure 3b). The points lie much closer together, with smaller loadings, reflecting the fact that adding predictors to the model substantially reduces the magnitude of the variance-covariance terms. This verifies that the prevailing patterns seen in the unconstrained ordination, and hence in the fish communities being sampled, can be explained by where and when samples were taken.

3.2. PIRLS data

Progress in international reading literacy study (PIRLS) is a large-scale international research project measuring reading literacy in children aged nine to ten years (Martin *et al.* 2017). PIRLS has been conducted every five years since 2001, with 61 countries participating in PIRLS 2016. Published studies have proposed that school variables are more important than family background in determining academic achievement in developing countries (Heyneman and Loxley 1982). However, more recent studies report conflicting results which show that these variables are similar across countries (Marôco 2021), with authors proposing this homogenization may be due to the increase of mass schooling. Therefore, we are interested in exploring how the effect of school variables on literacy scores of students vary by country.

We propose that students from economically disadvantaged schools (**Eco_disad**) with a library containing more books (**Size_lib**) have higher literacy scores than students from a school with no library, but the difference in literacy scores will be less when students are from schools of higher economic background i.e., there is an interaction between the two school variables. Further, we would like to know how the interactive effect of these school-level variables will vary by country, hence we want to fit a random slopes model, with different **Eco_disad:Size_lib** effects for each country. Both **Eco_disad** and **Size_lib** are categorical variables with four levels, hence we need 15 parameters to characterise their joint effect and a 15-dimensional random effect in the model. Estimating an unstructured variance-covariance matrix Σ would require 136 parameters. In contrast, for a reduced-rank covariance structure of rank three, only 42 parameters are required – AIC was used to select the rank of three.

We consider the model for literacy score, y_{ijk} , for student $i = 1, \dots, n$, in school $j = 1, \dots, m$ and country $k = 1, \dots, p$ as follows:

$$y_{ijk} = \mathbf{x}_j^\top \boldsymbol{\beta} + b_j^{[s]} + \mathbf{x}_j^\top \mathbf{b}_k^{[rr]} + \epsilon_{ijk} \quad (10)$$

where \mathbf{x}_j are vectors of covariates relating to school factors and $\boldsymbol{\beta}$ is the vector of fixed coefficients. The random intercept for school, $b_j^{[s]} \sim \mathcal{N}(0, \sigma_s^2)$, accounts for heterogeneity of average literacy scores between schools. The random intercept and slopes of school variables $\mathbf{b}_k^{[rr]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda} \mathbf{\Lambda}^\top)$ allow the school variables to vary by country where $\mathbf{\Lambda}$ is the full matrix of factor loadings, and $\epsilon_{ijk} \sim \mathcal{N}(0, \sigma^2)$ is the residual error.

This model can be fitted in **glmmTMB** using the following command:

```
R> glmmTMB(Overall ~ Size_lib * Eco_disad + (1 | School) +
+ rr(Size_lib * Eco_disad | Country, d = 3),
+ control = glmmTMBControl(start_method = list(method = "res")),
+ family = gaussian(),
+ data = pirls)
```

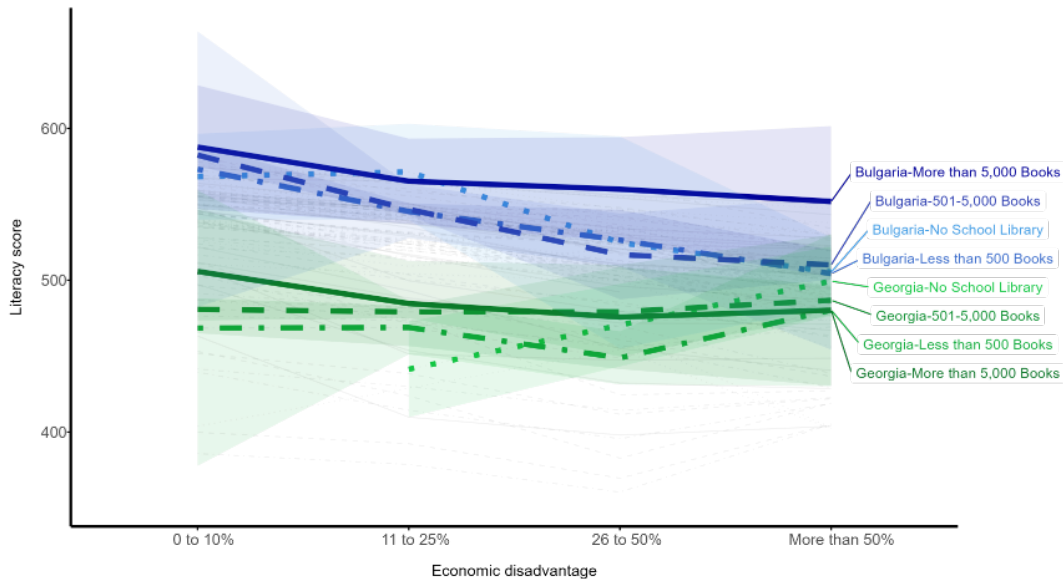


Figure 4: Plot of average literacy score (95% confidence interval shown by shaded area) of students in schools of different economic backgrounds and with varying library sizes in Georgia (green), Bulgaria (blue) and remaining countries (grey).

The starting algorithm for initialising parameters specified in the `control` argument is needed when fitting this model, otherwise there will be convergence issues.

The conditional estimates of the school random effects by country are complex (Figure 7): we will focus on a particular contrast between Bulgaria and Georgia (coloured results in Figure 7). This comparison is of interest because these two countries appear to have different patterns of library-economic disadvantage interaction (Figure 4). The most obvious pattern shown in the interaction plot (Figure 4) is that students from Bulgaria (blue lines) have higher literacy scores than students from Georgia (green lines). However, Bulgarian students' literacy scores also appear to decrease with increasing economic disadvantage; furthermore, the rate of decline appeared to be slower for better resourced-schools (i.e., those with more books). These patterns, expected from general socioeconomic principles, were generally observed across many countries. Georgia appears unusual: economic disadvantage had little overall effect on literacy scores, and the effects of library size were opposite from those expected – schools with large libraries appeared to have higher literacy scores than those with small libraries when schools were well off, but library size made little difference for strongly disadvantaged schools.

4. Discussion

In this article, we introduced a new variance-covariance structure, **rr**, in **glmmTMB**, to add reduced-rank functionality to mixed models. This feature broadens the scope of models that can be fitted by writing a large dimensional multivariate random effect as a linear combination of d latent variables, a more parsimonious structure that can be more readily estimated when the dimension of the random effect is large. In Section 1, we discussed available tools in R, such as **glvm**, which also fit latent variable models. These packages were developed with a primary focus on models for ecological data. The key advantage of our work is adding a factor analytic term to the suite of random effects structures already available in **glmmTMB**, such that generalized latent variable models can now be fitted to complex study designs, using a familiar interface.

We presented two applications to illustrate the use of a reduced-rank approach. In both examples the dimension of the random effect was moderate – 9 and 15 for these two examples – but this was already too large to be practically estimable, necessitating the use of a reduced-rank model. Reduced-rank models are capable of fitting random effects of very large dimension: for example, Niku *et al.* (2019) fitted a GLVM with a dimension of 985 to a microbial data set. When fitting models to larger data sets, difficulties can be encountered; for example, in ecology the number of species is often large compared to the number of samples. In situations like this, it may be useful to fit a model multiple times with different starting values for the parameters, and the fit with the highest log-likelihood value is considered the best fitting model.

The reduced-rank model contributes to the model-simplification toolbox, allowing for a more parsimonious random effect which may be necessary for some study designs (Matuschek *et al.* 2017). Currently, available methods for simplification are assuming a diagonal variance-covariance matrix, with homogeneous or heterogeneous variances; assuming compound symmetry; or assuming some specific form of structure (AR(1), Toeplitz, etc.). All of these structures are available in **glmmTMB**. Another alternative for controlling complexity is some form of shrinkage towards zero on the factor loadings as proposed in a Bayesian framework (Bhattacharya and Dunson 2011).

A key step in applying a reduced-rank random effect is choosing the rank d . Different strategies may be used for choosing d , depending on the analysis goal. In the wind farm application, we used a two-dimensional biplot to visualise correlations across species (Figure 3) and for this purpose $d = 2$ was appropriate. In our second application, the PIRLS study, our goal was to make inferences about correlated fixed effects, and the reduced rank approach was used to estimate this correlation. In this case we used information criteria to choose a value for d that gave us a good fit to the data. We found that estimates and confidence intervals for the fixed effect estimates were robust to different choices of rank (see supplementary material, Figure 5 and Figure 6). The extent to which fixed effects can change as covariance assumptions on random effects change Σ is a function of how much the fitted covariance structure actually changes. Factor analytical terms offer diminishing returns in terms of changes to Σ as d increases, so there is greatest capacity for changes in interpretation when d is small (in practice we have seen qualitatively important changes only for $d < 2$, as in Figure 6).

The reduced-rank structure has an interesting point of difference from other approaches to fitting a multivariate random effect in that it permits a singular variance-covariance matrix, or more precisely, it assumes singularity. Other methods of fitting correlated random effects

require a positive-definite variance-covariance matrix and return warnings when encountering (near-)singularity, an issue circumvented here by assuming a reduced-rank structure.

Our examples give just a few tastes of how reduced rank variance-covariance structures can be used in mixed modelling, where it has previously been technically difficult to fit models with random effects of high dimension. Some example areas where we see potential use for this approach include factor analysis with multi-level designs or repeated measures, and genotype-by-environment interaction analysis (Piepho 1997; Smith *et al.* 2001). We see a myriad of potential applications, and look forward to seeing how this new tool is used in practice.

References

- Asar Ö, Bolin D, Diggle PJ, Wallin J (2020). “Linear Mixed Effects Models for Non-Gaussian Continuous Repeated Measurement Data.” *Journal of the Royal Statistical Society C*, **69**(5), 1015–1065. doi:10.1111/rssc.12405.
- Banerjee S, Gelfand AE, Finley AO, Sang H (2008). “Gaussian Predictive Process Models for Large Spatial Data Sets.” *Journal of the Royal Statistical Society B*, **70**(4), 825–848. doi:10.1111/j.1467-9868.2008.00663.x.
- Bartholomew DJ, Knott M, Moustaki I (2011). *Latent Variable Models and Factor Analysis: A Unified Approach*. John Wiley & Sons.
- Bates D, Mächler M, Bolker BM, Walker S (2014). “Fitting Linear Mixed-Effects Models Using **lme4**.” *Journal of Statistical Software*, **67**(1), 1–48. doi:10.18637/jss.v067.i01.
- Bergstrom L, Sundqvist F, Bergstrom U (2013). “Effects of an Offshore Wind Farm on Temporal and Spatial Patterns in the Demersal Fish Community.” *Marine Ecology Progress Series*, **485**, 199–210. doi:10.3354/meps10344.
- Bhattacharya A, Dunson DB (2011). “Sparse Bayesian Infinite Factor Models.” *Biometrika*, **98**(2), 291–306. doi:10.1093/biomet/asr013.
- Bolker BM, Brooks ME, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, White JSS (2009). “Generalized Linear Mixed Models: A Practical Guide for Ecology and Evolution.” *Trends in Ecology & Evolution*, **24**(3), 127–135. doi:10.1016/j.tree.2008.10.008.
- Brooks ME, Kristensen K, van Benthem KJ, Magnusson A, Berg CW, Nielsen A, Skaug HJ, Machler M, Bolker BM (2017). “**glmmTMB** Balances Speed and Flexibility Among Packages for Zero-Inflated Generalized Linear Mixed Modeling.” *The R Journal*, **9**(2), 378–400. doi:10.3929/ethz-b-000240890.
- Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA, Marioni JC, Stegle O (2015). “Computational Analysis of Cell-to-Cell Heterogeneity in Single-Cell RNA-sequencing Data Reveals Hidden Subpopulations of Cells.” *Nature Biotechnology*, **33**(2), 155–160. doi:10.1038/nbt.3102.
- Coull BA, Agresti A (2000). “Random Effects Modeling of Multiple Binomial Responses Using the Multivariate Binomial Logit-Normal Distribution.” *Biometrics*, **56**(1), 73–80. doi:10.1111/j.0006-341X.2000.00073.x.

- Cressie N, Johannesson G (2008). “Fixed Rank Kriging for Very Large Spatial Data Sets.” *Journal of the Royal Statistical Society B*, **70**(1), 209–226. doi:10.1080/10618600.1996.10474708.
- Dunn PK, Smyth GK (1996). “Randomized Quantile Residuals.” *Journal of Computational and Graphical Statistics*, **5**(3), 236–244. doi:10.1080/10618600.1996.10474708.
- Dushoff J, Kain MP, Bolker BM (2019). “I Can See Clearly Now: Reinterpreting Statistical Significance.” *Methods in Ecology and Evolution*, **10**(6), 756–759. doi:10.1111/2041-210X.13159.
- Heyneman SP, Loxley WA (1982). “Influences on Academic Achievement Across High and Low Income Countries: A Re-Analysis of IEA Data.” *Sociology of Education*, pp. 13–21.
- Huber P, Ronchetti E, Victoria-Feser MP (2004). “Estimation of Generalized Linear Latent Variable Models.” *Journal of the Royal Statistical Society B*, **66**(4), 893–908. doi:10.1111/j.1467-9868.2004.05627.x.
- Hui FKC (2016). “**Boral** – Bayesian Ordination and Regression Analysis of Multivariate Abundance Data in R.” *Methods in Ecology and Evolution*, **7**(6), 744–750. doi:10.1111/2041-210X.12514.
- Hui FKC, Taskinen S, Pledger S, Foster SD, Warton DI (2015). “Model-Based Approaches to Unconstrained Ordination.” *Methods in Ecology and Evolution*, **6**(4), 399–411. doi:10.1111/2041-210X.12236.
- Hui FKC, Warton DI, Ormerod JT, Haapaniemi V, Taskinen S (2017). “Variational Approximations for Generalized Linear Latent Variable Models.” *Journal of Computational and Graphical Statistics*, **26**(1), 35–43. doi:10.1080/10618600.2016.1164708.
- Kass RE, Steffey D (1989). “Approximate Bayesian Inference in Conditionally Independent Hierarchical Models (Parametric Empirical Bayes Models).” *Journal of the American Statistical Association*, **84**(407), 717–726. doi:10.1080/01621459.1989.10478825.
- Kristensen K, Nielsen A, Berg CW, Skaug H, Bell B (2015). “**TMB**: Automatic Differentiation and Laplace Approximation.” *Journal of Statistical Software, Articles*, **70**(5), 1–21.
- Marôco J (2021). “What Makes a Good Reader? Worldwide Insights From PIRLS 2016.” *Reading and Writing*, **34**(1), 231–272.
- Martin MO, Mullis IVS, Hooper M (2017). “Methods and Procedures in PIRLS 2016.” url <https://timssandpirls.bc.edu/pirls2016/international-database/index.html>.
- Matuschek H, Kliegl R, Vasishth S, Baayen H, Bates D (2017). “Balancing Type I Error and Power in Linear Mixed Models.” *Journal of Memory and Language*, **94**, 305–315. doi:10.1016/j.jml.2017.01.001.
- Niku J, Hui FK, Taskinen S, Warton DI (2019). “**gllvm**: Fast Analysis of Multivariate Abundance Data With Generalized Linear Latent Variable Models in R.” *Methods in Ecology and Evolution*, **10**(12), 2173–2182. doi:10.1111/2041-210X.13303.

- Niku J, Hui FK, Taskinen S, Warton DI (2021). “Analyzing Environmental-Trait Interactions in Ecological Communities With Fourth-Corner Latent Variable Models.” *Environmetrics*, p. e2683. doi:10.1002/env.2683.
- Niku J, Warton DI, Hui FK, Taskinen S (2017). “Generalized Linear Latent Variable Models for Multivariate Count and Biomass Data in Ecology.” *Journal of Agricultural, Biological and Environmental Statistics*, **22**(4), 498–522. doi:10.1007/s13253-017-0304-7.
- Ovaskainen O, Tikhonov G, Norberg A, Guillaume Blanchet F, Duan L, Dunson D, Roslin T, Abrego N (2017). “How to Make More Out of Community Data? A Conceptual Framework and Its Implementation as Models and Software.” *Ecology Letters*, **20**(5), 561–576. ISSN 1461-0248. doi:10.1111/ele.12757.
- Piepho HP (1997). “Analyzing Genotype-Environment Data by Mixed Models With Multiplicative Terms.” *Biometrics*, pp. 761–766.
- Pollock LJ, Tingley R, Morris WK, Golding N, O’Hara RB, Parris KM, Vesik PA, McCarthy MA (2014). “Understanding Co-Occurrence by Modelling Species Simultaneously With a Joint Species Distribution Model (JSDM).” *Methods in Ecology and Evolution*, **5**(5), 397–406. doi:10.1111/2041-210X.12180.
- Rabe-Hesketh S, Skrondal A, Pickles A (2004). “Generalized Multilevel Structural Equation Modeling.” *Psychometrika*, **69**(2), 167–190. doi:10.1007/BF02295939.
- Raudenbush SW, Yang ML, Yosef M (2000). “Maximum Likelihood for Generalized Linear Models With Nested Random Effects via High-Order, Multivariate Laplace Approximation.” *Journal of Computational and Graphical Statistics*, **9**(1), 141–157. doi:10.1080/10618600.2000.10474870.
- R Core Team (2020). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Skrondal A, Rabe-Hesketh S (2003). “Multilevel Logistic Regression for Polytomous Data and Rankings.” *Psychometrika*, **68**(2), 267–287. doi:10.1007/BF02294801.
- Skrondal A, Rabe-Hesketh S (2004). *Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models*. Chapman and Hall/CRC.
- Smith A, Cullis B, Thompson R (2001). “Analyzing Variety by Environment Data Using Multiplicative Mixed Models and Adjustments for Spatial Field Trend.” *Biometrics*, **57**(4), 1138–1147. doi:10.1111/j.0006-341X.2001.01138.x.
- StataCorp (2021). *Stata Statistical Software: Release 17*. StataCorp LLC, College Station, TX.
- van der Veen B, Hui FKC, Hovstad KA, Solbu EB, O’Hara RB (2021). “Model-Based Ordination for Species With Unequal Niche Widths.” *Methods in Ecology and Evolution*, **12**(7), 1288–1300. doi:10.1111/2041-210X.13595.
- Walker SC, Jackson DA (2011). “Random-Effects Ordination: Describing and Predicting Multivariate Correlations and Co-Occurrences.” *Ecological Monographs*, **81**(4), 635–663.

Warton DI, Blanchet FG, O'Hara RB, Ovaskainen O, Taskinen S, Walker SC, Hui FK (2015). "So Many Variables: Joint Modeling in Community Ecology." *Trends in Ecology & Evolution*, **30**(12), 766–779. doi:[10.1016/j.tree.2015.09.007](https://doi.org/10.1016/j.tree.2015.09.007).

A. Model summaries for the wind farm example

Summary of the model for the wind farm example in Section 1 fitted using `lme4` is as follows:

```
R> summary(glmer(abundance ~ Zone + Year + (Species + 0 | ID),
+   family = "poisson", data = wf.ex))
```

Generalized linear mixed model fit by maximum likelihood (Laplace Approximation) [`'glmerMod'`]

Family: poisson (log)

Formula: abundance ~ Zone + Year + (Species + 0 | ID)

Data: wf.ex

AIC	BIC	logLik	deviance	df.resid
1310.5	1336.1	-648.3	1296.5	277

Scaled residuals:

Min	1Q	Median	3Q	Max
-1.38779	-0.67079	0.06129	0.43556	1.58605

Random effects:

Groups	Name	Variance	Std.Dev.	Corr
ID	SpeciesTorsk	0.7603	0.8719	
	SpeciesTanglake	0.9839	0.9919	-0.74

Number of obs: 284, groups: ID, 142

Fixed effects:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.63001	0.10942	5.758	8.52e-09 ***
ZoneN	0.07242	0.12807	0.565	0.57174
ZoneS	-0.57875	0.19527	-2.964	0.00304 **
Year2010	0.57457	0.10360	5.546	2.92e-08 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Correlation of Fixed Effects:

	(Intr)	ZoneN	ZoneS
ZoneN		-0.343	
ZoneS		-0.525	-0.014
Year2010		-0.493	0.057 -0.130

The summary of the model fitted using **glmmTMB** is as follows:

```
R> summary(glmmTMB(abundance ~ Zone + Year + (Species + 0 | ID),
+   family = "poisson", data = wf.ex))
```

```
Family: poisson ( log )
Formula:      abundance ~ Zone + Year + (Species + 0 | ID)
Data: wf.ex
```

AIC	BIC	logLik	deviance	df.resid
1310.5	1336.1	-648.3	1296.5	277

Random effects:

Conditional model:

Groups Name	Variance	Std.Dev.	Corr
ID SpeciesTorsk	0.7603	0.8719	
SpeciesTanglake	0.9839	0.9919	-0.74

Number of obs: 284, groups: ID, 142

Conditional model:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	0.63000	0.10942	5.758	8.53e-09	***
ZoneN	0.07243	0.12807	0.566	0.57170	
ZoneS	-0.57876	0.19528	-2.964	0.00304	**
Year2010	0.57457	0.10360	5.546	2.92e-08	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

B. Parametric bootstrap analysis for the wind farm example

A parametric bootstrap to test the interaction terms of zone and year for the wind farm data. The P-value is estimated by comparing the observed likelihood ratio statistic to the simulated distribution of the test statistic under the null hypothesis. Bootstrap replications which failed to converge are ignored.

```
R> LRobs <- 2 * logLik(wf.glmm) - 2 * logLik(wf.glmm.null)
R> library(boot)
R> lrt.fun <- function(data) {
+   library(glmmTMB)
+   null <- try(glmmTMB(abundance ~ Zone + Year + diag(Zone + Year|Species) +
+     (1|Station) + rr(Species + 0 | ID, d = 2),
+     family = "poisson", data = data))
+   alt <- try(glmmTMB(abundance ~ Zone * Year + diag(Zone * Year|Species) +
+     (1|Station) + rr(Species + 0 | ID, d = 2),
+     family = "poisson", data = data))
+   LR <- tryCatch({2*logLik(alt) - 2*logLik(null)}, error = function(e) {NA})
+   return(LR)
+ }
R> sim.abund <- function(data, mle) {
+   library(glmmTMB)
+   out <- data
+   out$abundance <- simulate(mle)$sim_1 #simulate data under the null
+   out
+ }
R> wf.boot <- boot(data = windfarm,
+   ran.gen = sim.abund,
+   lrt.fun,
+   mle = wf.glmm.null,
+   sim = "parametric",
+   R = 1000,
+   parallel= "snow", #for Windows
+   ncpus = 4)
R> p <- ( sum(wf.boot$t[,1] >= LRobs, na.rm=TRUE) +1 )/(1000 + 1)
```

C. Sensitivity analysis

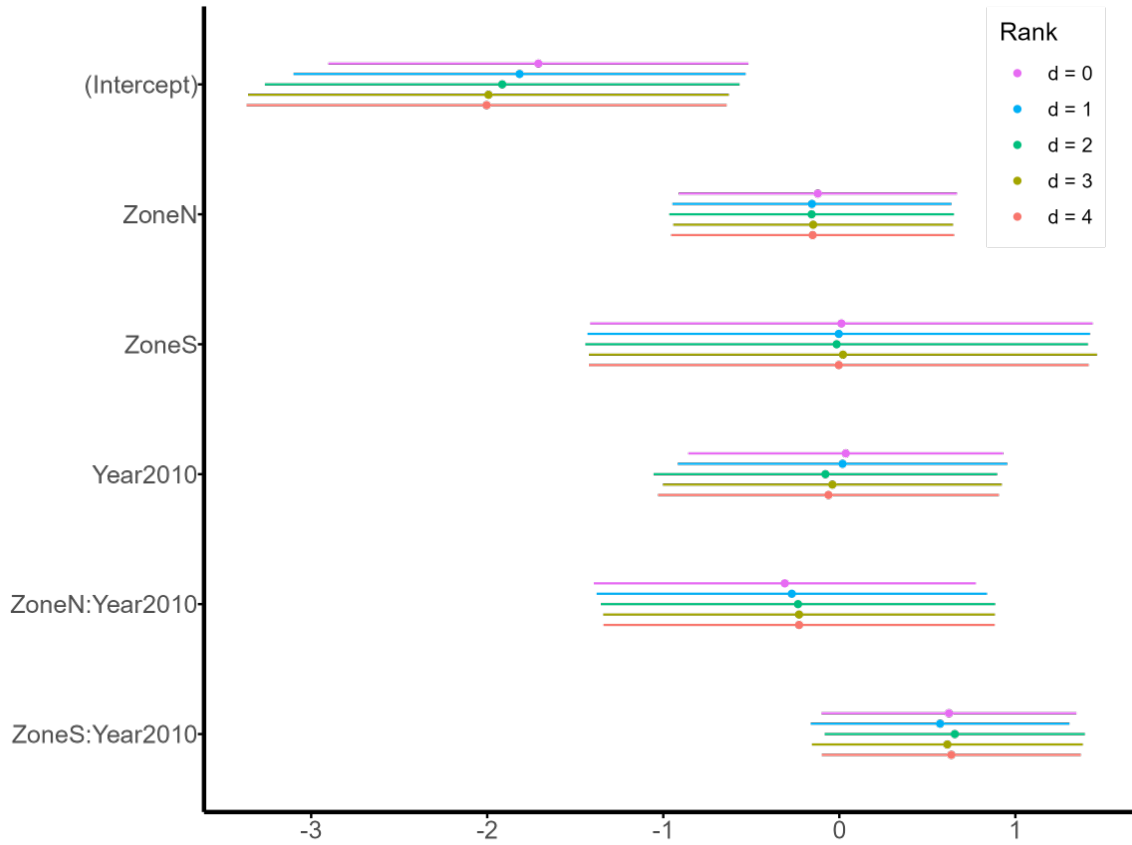


Figure 5: The fixed effect estimates and 95% confidence intervals for the wind farm model are similar when the rank (d) of the reduced-rank random effect varies from zero to four.

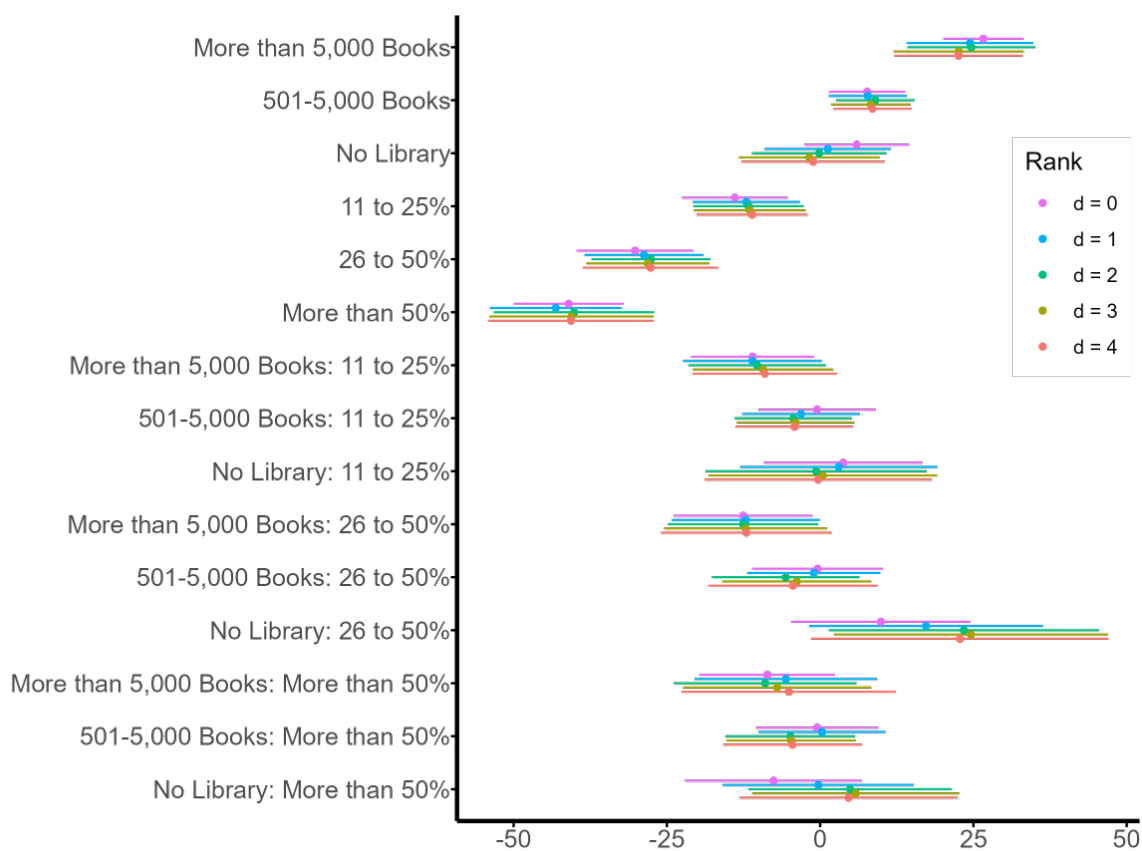


Figure 6: The fixed effect estimates and 95% confidence intervals for the PIRLS model are similar when the rank (d) of the reduced-rank random effect varies from one to four. When the reduced-rank random effect is replaced by a random intercept ($d = 0$), the estimates of the fixed effects are less similar and the standard errors may be smaller.

D. Reduced-rank random effect estimates from the PIRLS model

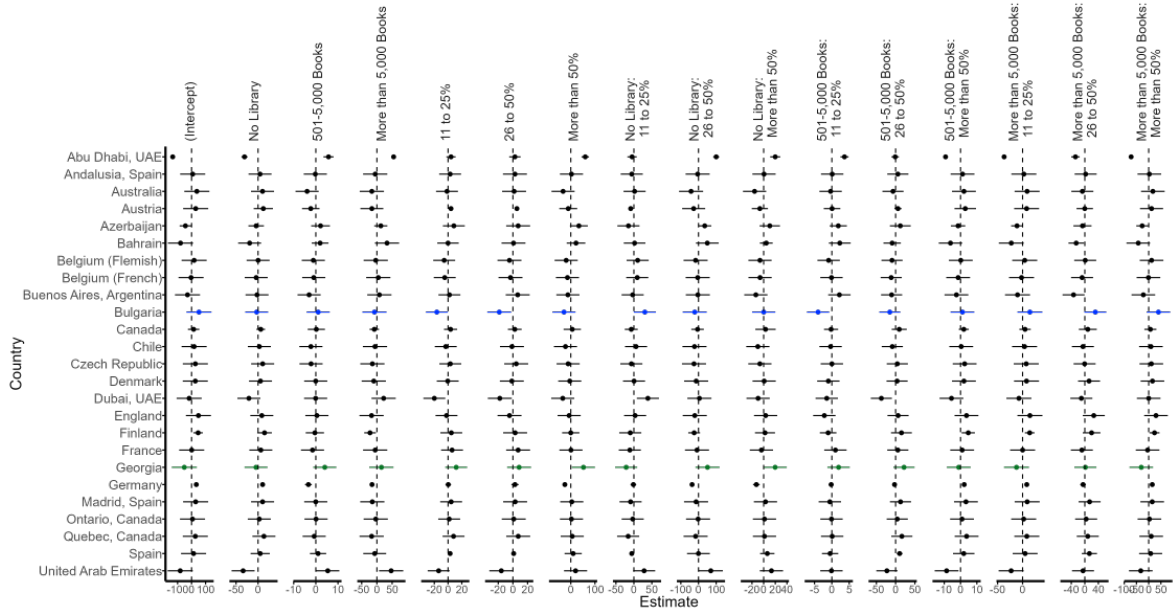


Figure 7: Conditional estimates of school variables by country from the reduced-rank random effect in the PIRLS model.

Affiliation:

Maeve McGillicuddy
School of Mathematics and Statistics
Evolution & Ecology Research Centre
UNSW Sydney
NSW 2052, Australia
E-mail: m.mcgillicuddy@unsw.edu.au

Gordana Popovic
Stats Central, Mark Wainwright Analytical Centre
UNSW Sydney
NSW 2052, Australia

Benjamin M. Bolker
Department of Mathematics & Statistics
McMaster University
1280 Main Street West
Hamilton, Ontario L8S 4K1, Canada

David I. Warton
School of Mathematics and Statistics
Evolution & Ecology Research Centre
UNSW Sydney
NSW 2052, Australia