

1 **Assuming Independence in Spatial Latent Variable Models: Consequences and**
2 **Implications of Misspecification**

3 **Francis K.C. Hui^{1,*}, Nicole A. Hill², and A.H. Welsh¹**

¹Research School of Finance, Actuarial Studies & Statistics, Australian National University,
Acton, ACT 2601, Australia

²Institute for Marine and Antarctic Studies, University of Tasmania, Hobart, TAS 7001, Australia

**email*: francis.hui@anu.edu.au

4 **SUMMARY:** Multivariate spatial data, where multiple responses are simultaneously recorded across spatially indexed
5 observational units, are routinely collected in a wide variety of disciplines. For example, the Southern Ocean Continu-
6 ous Plankton Recorder survey collects records of zooplankton communities in the Indian sector of the Southern Ocean,
7 with the aim of identifying and quantifying spatial patterns in biodiversity in response to environmental change. One
8 increasingly popular method for modeling such data is spatial generalized linear latent variable models (GLLVMs),
9 where the correlation across sites is captured by a spatial covariance function in the latent variables. However, little is
10 known about the impact of misspecifying the latent variable correlation structure on inference of various parameters
11 in such models. To address this gap in the literature, we investigate how misspecifying and assuming independence for
12 the latent variables' correlation structure impacts estimation and inference in spatial GLLVMs. Through both theory
13 and numerical studies, we show that performance of maximum likelihood estimation and inference on regression
14 coefficients under misspecification depends on a combination of the response type, the magnitude of true regression
15 coefficient and the corresponding loadings, and, most importantly, whether the corresponding covariate is (also)
16 spatially correlated. On the other hand, estimation and inference of truly non-zero loadings and prediction of latent
17 variables is consistently not robust to misspecification of the latent variable correlation structure.

KEY WORDS: Community Ecology; Factor Analysis; Loadings; Multivariate Abundance Data; Spatial; Spatio-
temporal

1. Introduction

Multivariate spatial data, consisting of multiple responses recorded simultaneously across spatially indexed observational units, are routinely collected in a variety of disciplines. Such data are characterized by spatial correlations whose strength depends inversely on the geographic distance between units, and between response correlations for which we usually have little information *a priori* regarding their structure. This article is motivated by the Southern Ocean Continuous Plankton Recorder (SO-CPR) survey, an annual survey that collects count records of zooplankton communities using vessels traversing the Indian sector of the Southern Ocean (Hosie, 2020). Alongside species records, we have covariate information pertaining to the physical habitat and oceanographic environment, and one of the primary aims of the SO-CPR survey is to identify and quantify spatial patterns in biodiversity in response to environmental change.

In community ecology, one of the most popular emerging methods for analyzing multivariate spatial data is spatial generalized linear latent variable models (Thorson et al., 2015; Warton et al., 2016; Ovaskainen et al., 2016; Bjork et al., 2018; Shirota et al., 2019). This model is an extension of generalized linear latent variable models (GLLVMs, SkronDAL and Rabe-Hesketh, 2004) for analyzing multi-response data, where the latent variables are given additional structure to model the spatial correlations, on top of their usual role in accounting for between species correlation. For instance, the latent variables are drawn from a multivariate normal distribution with zero mean vector and a correlation matrix parameterized using a Matérn correlation function based on some distance metric between sites. Outside of community ecology, such analyses are also referred to as spatial factor analysis (e.g., Zhu et al., 2005; Lopes et al., 2011).

Compared to the standard independence assumption for latent variables i.e., when the correlation matrix is set to the identity matrix, assuming a structured correlation matrix

43 has two major implications. From a modeling perspective, it is evident that in the presence
44 of spatial correlations, assuming independence in the latent variables results in clear model
45 misspecification. There has been some research into the impacts of misspecifying the latent
46 variables for estimation and inference in GLLVMs. However, the literature so far has focused
47 on misspecification of the latent variable distribution rather than of the correlation structure,
48 and the consequences on estimation and inference of the loadings and prediction of latent
49 variables rather than the impact on regression coefficients (e.g., Ma and Genton, 2010;
50 Irincheeva et al., 2012). Furthermore, the large sample theory developed in such previous re-
51 search on misspecification in GLLVMs was motivated by non-spatial data. As such, assuming
52 independence in the latent variables across units was a reasonable foundation on which to
53 base developments. By contrast, with spatial data such theory does not carry over directly:
54 one requires alternative asymptotic frameworks e.g., increasing domain, fixed domain infill,
55 and it is not guaranteed that properties such as consistency can be necessarily achieved within
56 all these frameworks. The asymptotics of spatial and spatio-temporal modeling remains an
57 active area of research (e.g., Lu and Tjøstheim, 2014; Kurisu, 2019), and to our knowledge
58 there has been no theoretical research into the large sample effects of misspecifying the
59 correlation structure in spatial GLLVMs (although related empirical work has been done on
60 this by Shirota et al., 2019).

61 From a computational perspective, and focusing on maximum likelihood estimation, spatial
62 GLLVMs involve a high-dimensional integral and, usually, the inversion of a high-dimensional
63 matrix during the estimation process. By contrast, estimation and inference assuming inde-
64 pendence is substantially faster since the dimension of the integral in the likelihood is reduced
65 and there are no covariance parameters to estimate. While considerable progress has been
66 made into more efficient likelihood-based estimation and inference in GLLVMs and mixed
67 models in general (Niku et al., 2019), as well as on computationally faster approaches for

68 modeling spatial data in general, the substantial reduction in computation time brought
69 about by assuming independence remains.

70 Motivated by the potential tradeoff between computation time and model misspecifica-
71 tion, this article studies the consequences of assuming independence for the latent variable
72 correlation structure on estimation and inference in GLLVMs applied to multivariate spatial
73 data. We develop some large sample results for point estimation of regression coefficients
74 in GLLVMs under misspecification of the correlation structure, and complement this with
75 an extensive simulation study to examine how finite sample performance is affected by
76 response type, the magnitude of the regression coefficient and corresponding loadings, and
77 whether the measured environmental covariates are spatially correlated. The main findings
78 and contributions of this article may be summarized as follows:

- 79 • For normally distributed responses, maximum likelihood estimates of all regression co-
80 efficients excluding the intercept are estimation consistent. For non-normally distributed
81 responses, only a weaker zero consistency result can be obtained. The exception is responses
82 that are uncorrelated with all other responses, for which we (also) obtain estimation
83 consistency of all the regression coefficients (Section 3).
- 84 • For covariates that have little to no spatial correlation, empirical studies show that for
85 normal, negative binomial, and Tweedie distributed responses, point estimation of re-
86 gression coefficients is largely robust to misspecification, and sandwich-based confidence
87 intervals are close to their the nominal significance level (Sections 4.2–4.3). For binary
88 responses, point estimation and sandwich-based confidence intervals only perform well
89 under misspecification when either the magnitude of the true regression coefficient or the
90 magnitude of its corresponding loadings are not especially large (Sections 4.4).
- 91 • For covariates that have moderate to strong spatial correlation, assuming independence of
92 the latent variables almost always leads to poor estimation and inferential performance

93 for regression coefficients, irrespective of the response type and the magnitude of the
94 true coefficient. The one exception is coefficients corresponding to responses that are
95 uncorrelated with all other responses, which remain robust to misspecification when the
96 covariates are spatially correlated (Sections 4.2–4.4).

- 97 • Irrespective of the response type, estimation and inference of the non-zero loadings and
98 prediction of latent variables in the GLLVM is not robust to misspecification of the
99 latent variable correlation structure. The one exception is responses that have entirely
100 zero loadings i.e, responses that are independent of all other responses, for which a zero
101 consistency result can be obtained for the estimates of the loadings under misspecification
102 (Sections 4.2–4.4).

103 The results we obtain bear some resemblance to existing research on estimation and
104 inference for spatial data as well as correlated data more generally. One article of particular
105 relevance here is Shirota et al. (2019), who empirically found that spatial GLLVMs consis-
106 tently outperform independent GLLVMs when it comes to prediction of the linear predictor
107 and estimation of the residual covariance matrix (formed from the loading matrix). The
108 findings here are concordant with these previous results, although in this article we perform
109 a broader investigation to develop large sample results as well as examine the critical issue
110 of estimation and inference of the regression coefficients (which Shirota et al., 2019, did
111 not focus on). In the more general spatial regression setting, several methods proposed for
112 variable selection (usually assuming spatially uncorrelated covariates) have demonstrated
113 that assuming independence across sites does not affect large sample selection consistency
114 results, although in finite samples the loss of efficiency and overall consequences can be
115 pronounced (Hoeting et al., 2006; Wang and Zhu, 2009; Xu et al., 2015). Also, previous
116 research into generalized linear mixed models has largely shown that estimation and inference
117 on fixed effects is fairly robust to misspecification of the random effects distribution, but

118 that random effects inference is *not* robust to the same sort of misspecification (Hui et al.,
 119 2020). This article builds on such existing research, and presents several new findings and
 120 consequences. For instance, the lack of robustness under misspecification for inference on
 121 spatially correlated covariates has major implications on how multivariate spatial data should
 122 be analyzed. In most spatially indexed observational studies in ecology, it is either known
 123 *a priori* or evidence is uncovered empirically that both the responses and one or more
 124 of the environmental predictors are spatially correlated. When analyzing such data, our
 125 results demonstrate that it is imperative that spatial correlation in the latent variables be
 126 taken into account. Otherwise, one risks inference being incorrect for regression coefficients
 127 corresponding to the spatially correlated covariates e.g., coverage probabilities of confidence
 128 intervals substantially below the nominated level, as well as poor predictive performance of
 129 the model as a whole (see also Yoon and Welsh, 2020, for similar results in the context of
 130 linear mixed models). We conclude our comparison by applying both a spatial GLLVM and a
 131 standard GLLVM assuming independence of the latent variables to the SO-CPR survey. For
 132 both models, we obtained some but not entirely similar results with regards to identifying
 133 the effects of key environmental covariates on the distributions of 24 zooplankton species
 134 across 3,900 spatial locations.

135 2. Spatial GLLVMs

136 We establish the notation of spatial GLLVMs with ecological data in mind. Let $y_j(\mathbf{s}_i)$ denote
 137 the observed response for species $j = 1, \dots, p$ at site (observational unit) $\mathbf{s}_i \in \mathcal{D}; i = 1, \dots, n$
 138 in some spatial domain \mathcal{D} . The spatial domain can be continuously (e.g., geostatistical
 139 data) or discretely (e.g., lattice data) spatially indexed, or both, and we make no explicit
 140 restrictions on this (subject to Assumption 1 discussed below). We also observe a q -vector
 141 of covariates at each site, $\mathbf{x}(\mathbf{s}_i) = (x_1(\mathbf{s}_i), \dots, x_q(\mathbf{s}_i))^\top$, representing physical environment

142 or habitat. We assume $x_1(\mathbf{s}_i) = 1$ for $i = 1, \dots, n$ to represent an intercept term, while the
 143 other terms are centered to have expectation zero, $E_X\{x_k(\mathbf{s})\} = 0$ for all $k = 2, \dots, q$.

144 The responses are assumed to be generated from a spatial GLLVM as follows: conditional
 145 on a d -vector of latent variables $\mathbf{u}(\mathbf{s})$ with $d \ll p$, as well as the covariates, the elements of the
 146 response vector $\mathbf{y}(\mathbf{s}) = (y_1(\mathbf{s}), \dots, y_p(\mathbf{s}))^\top$ are independent observations from the exponen-
 147 tial family of distributions. That is, $f\{y_j(\mathbf{s})|\mathbf{x}(\mathbf{s}), \mathbf{u}(\mathbf{s})\} = \exp\{\phi_j^{-1}[y_j(\mathbf{s})\vartheta_j(\mathbf{s}) - a\{\vartheta_j(\mathbf{s})\}] +$
 148 $c(y_j(\mathbf{s}), \phi_j)\}$ for known functions $a(\cdot)$ and $c(\cdot, \cdot)$, where $\vartheta_j(\mathbf{s})$ is the canonical parameter and
 149 ϕ_j is species-specific dispersion parameter. The mean, $E_{Y|U,X}\{y_j(\mathbf{s})|\mathbf{x}(\mathbf{s}), \mathbf{u}(\mathbf{s})\} = \mu_j(\mathbf{s}) =$
 150 $\mu_j(\mathbf{s}) = a'\{\vartheta_j(\mathbf{s})\}$, is modeled as $g\{\mu_j(\mathbf{s})\} = \eta_j(\mathbf{s}) = \mathbf{x}(\mathbf{s})^\top \boldsymbol{\beta}_j + \mathbf{u}(\mathbf{s})^\top \boldsymbol{\lambda}_j$ for a known link
 151 function $g(\cdot)$, where $\boldsymbol{\beta}_j$ is the vector of species-specific regression coefficients and $\boldsymbol{\lambda}_j$ is the
 152 vector of species-specific loadings.

153 Write the full nd -vector of latent variables generically as $\mathbf{u}(S)$, where S indexes all sampled
 154 locations. Let $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_p^\top)^\top$, $\boldsymbol{\phi} = (\phi_1, \dots, \phi_p)^\top$, and $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1^\top, \dots, \boldsymbol{\lambda}_p^\top)^\top$ denote the
 155 full vector of regression coefficients, dispersion parameters, and loadings respectively. The
 156 marginal likelihood for the spatial GLLVM is given by

$$L(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = \int \prod_{i=1}^n \prod_{j=1}^p f\{y_j(\mathbf{s}_i)|\mathbf{x}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_i)\} f(\mathbf{u}(S)|\boldsymbol{\theta}) d\mathbf{u}(S), \quad (1)$$

157 where the latent variables are assumed to come from a multivariate normal distribution
 158 with zero mean vector and a correlation matrix parameterized by a vector $\boldsymbol{\theta}$, that is,
 159 $f(\mathbf{u}(S)|\boldsymbol{\theta}) = \mathcal{N}_{nd}\{\mathbf{0}, \boldsymbol{\Sigma}(\boldsymbol{\theta})\}$. Note $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ is a correlation (as opposed to a covariance) matrix
 160 to avoid scale invariance in the GLLVM. It is the combination of covariance parameters
 161 $\boldsymbol{\theta}$ and loadings $\boldsymbol{\lambda}_j$ in (1) which models the between species and spatial correlations in
 162 the data; see the end of the following section for more discussion on the role of corre-
 163 lations between species versus spatial correlation. On the linear predictor scale, we have
 164 $\text{Cov}\{\eta_j(\mathbf{s}_i), \eta_{j'}(\mathbf{s}_{i'})\} = \boldsymbol{\lambda}_j^\top \text{Cov}\{\mathbf{u}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_{i'})\} \boldsymbol{\lambda}_{j'}$ for $i, i' = 1, \dots, n$ and $j, j' = 1, \dots, p$, where

165 $\text{Cov}\{\mathbf{u}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_{i'})\}$ depends on the precise structure of $\Sigma(\boldsymbol{\theta})$. Thus we see that the between
 166 species correlation also depends on the two spatial locations of interest. One common choice
 167 for $\Sigma(\boldsymbol{\theta})$ in community ecology is to assume the d latent variables are independent and use
 168 a spatial covariance function for each latent variable (e.g., Thorson et al., 2015; Ovaskainen
 169 et al., 2016). If we let $u_l(\mathbf{s}_i)$ denote the l -th element in $\mathbf{u}(\mathbf{s}_i)$, then for $l = 1, \dots, d$ we use
 170 the Matérn correlation function, $\text{Cov}\{u_l(\mathbf{s}_i), u_l(\mathbf{s}_{i'})\} = \rho(D(\mathbf{s}_i, \mathbf{s}_{i'}), \nu, \alpha_l)$, where $\rho(h, \nu, \alpha_l) =$
 171 $\{2^{\nu-1}\Gamma(\nu)\}^{-1}(h\alpha_l^{-1})^\nu \mathcal{K}_\nu(h\alpha_l^{-1})$ for smoothness $\nu > 0$ and spatial scale parameters $\alpha_l > 0$,
 172 where $\Gamma(\cdot)$ is the gamma function, $\mathcal{K}_\nu(\cdot)$ is the modified Bessel function of the second kind,
 173 and $D(\cdot, \cdot)$ is some distance metric. Note the scale parameter, but not the smoothness
 174 parameter, is allowed to be different across the latent variables. The vector of covariance
 175 parameters is given by $\boldsymbol{\theta} = (\nu, \alpha_1, \dots, \alpha_d)^\top$, although the smoothness may be fixed *a priori*
 176 e.g., $\nu = 1$ is a recommended choice when $\mathcal{D} = \mathbb{R}^2$ (Lindgren and Rue, 2015).

177 2.1 Assuming Independence

178 Of the various structures which can be imposed on the latent variables, the simplest one is
 179 to assume the correlation matrix is equal to the identity matrix, $\Sigma(\boldsymbol{\theta}) = \mathbf{I}_{nd}$. The latent
 180 variables are then assumed to be independently standard normally distributed, and the
 181 marginal log-likelihood for the GLLVM reduces to

$$L_{\text{Ind}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}) = \prod_{i=1}^n \left\{ \int \prod_{j=1}^p f\{y_j(\mathbf{s}_i) | \mathbf{x}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_i)\} h\{\mathbf{u}(\mathbf{s}_i)\} d\mathbf{u}(\mathbf{s}_i) \right\} = \prod_{i=1}^n L_{\text{Ind},i}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}), \quad (2)$$

182 where $h\{\mathbf{u}(\mathbf{s}_i)\} = \mathcal{N}_d(\mathbf{0}, \mathbf{I}_d)$. Equation (2) is equivalent to the standard GLLVM advocated
 183 for use by Warton et al. (2016) for community ecology data, and is a common choice for
 184 modeling multivariate data in other disciplines (Skrondal and Rabe-Hesketh, 2004). With
 185 the independence structure, the latent variables can only model between species correlation:
 186 $\text{Cov}\{\eta_j(\mathbf{s}_i), \eta_{j'}(\mathbf{s}_{i'})\} = 0$, unless $i = i'$ in which case $\text{Cov}\{\eta_j(\mathbf{s}_i), \eta_{j'}(\mathbf{s}_i)\} = \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_{j'}$.

187 It is clear that assuming spatial independence makes the likelihood function computa-
 188 tionally much easier to maximize: the product over the n observational units is pulled
 189 outside the integral and there are no covariance parameters to estimate in (2). Of course the
 190 computational gain comes with a cost, and the presence of spatial correlation means that
 191 equation (2) is misspecified in the second moment. In fact, as seen above the independence
 192 assumption means that the expected counts of two species may be correlated at any particular
 193 site, but are otherwise uncorrelated across different sites. In Appendix A, we elaborate
 194 more on the two correlation structures present in a spatial GLLVM, noting that the focus
 195 of this article is misspecification of the spatial correlation structure, and not (also) on
 196 misspecification of the between species correlation.

197 3. Estimation Under Misspecification

198 For the set of observed data, consider $\log L(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}, \boldsymbol{\theta})$ in (1) as the log-likelihood correspond-
 199 ing to the true model. Define $\ell(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = \log\{\int \prod_{j=1}^p f\{y_j(\mathbf{s})|\mathbf{x}(\mathbf{s}), \mathbf{u}(\mathbf{s})\}f(\mathbf{u}(\mathbf{s})|\boldsymbol{\theta})d\mathbf{u}(\mathbf{s})\}$,
 200 and denote the full vector of true parameters as $(\boldsymbol{\beta}^0, \boldsymbol{\phi}^0, \boldsymbol{\lambda}^0, \boldsymbol{\theta}^0)$ which satisfies
 201 $E_{X,Y}\{\nabla\ell(\boldsymbol{\beta}^0, \boldsymbol{\phi}^0, \boldsymbol{\lambda}^0, \boldsymbol{\theta}^0)\} = \mathbf{0}$. Next, consider $\log L_{\text{Ind}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda})$ in (2) as the log-likelihood
 202 corresponding to the misspecified GLLVM assuming independence, and let
 203 $\ell_{\text{Ind}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}) = \log\{\int \prod_{j=1}^p f\{y_j(\mathbf{s})|\mathbf{x}(\mathbf{s}), \mathbf{u}(\mathbf{s})\}h\{\mathbf{u}(\mathbf{s})\}d\mathbf{u}(\mathbf{s})\}$. We define the Kullback-Leibler
 204 distance $D_{\text{KL}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}) = E_{X,Y}\{\ell(\boldsymbol{\beta}^0, \boldsymbol{\phi}^0, \boldsymbol{\lambda}^0, \boldsymbol{\theta}^0) - \ell_{\text{Ind}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda})\}$, and consequently the vector
 205 of pseudo-true parameters $(\boldsymbol{\beta}^*, \boldsymbol{\phi}^*, \boldsymbol{\lambda}^*) = \arg \min_{\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}} D_{\text{KL}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda})$ which satisfy
 206 $E_{X,Y}\{\nabla\ell_{\text{Ind}}(\boldsymbol{\beta}^*, \boldsymbol{\phi}^*, \boldsymbol{\lambda}^*)\} = \mathbf{0}$.

207 Let $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) = \arg \max_{\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}} \log L_{\text{Ind}}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda})$ denote the maximum likelihood estimator
 208 of the misspecified GLLVM. We base our large sample developments under the following
 209 general regularity condition.

210 **ASSUMPTION 1:** The appropriate conditions are assumed to be satisfied to ensure weak

211 consistency of the maximum likelihood estimator under the misspecified GLLVM. That is,
 212 $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) \xrightarrow{p} (\boldsymbol{\beta}^*, \boldsymbol{\phi}^*, \boldsymbol{\lambda}^*)$ as $n \rightarrow \infty$.

213 The above assumption is rather general, and is chosen with reason. As reviewed in Section 1,
 214 large sample theory in spatial modeling remains an active area of research and is dependent
 215 both on the related asymptotic framework and the structure of the model itself (e.g., most
 216 theory has been developed assuming a univariate normal or continuous response with additive
 217 error structure, and some sort of expanding domain framework; Mardia and Marshall, 1984;
 218 Lu and Tjøstheim, 2014). Our aim is to study misspecification in spatial GLLVMs *irrespective*
 219 of the framework selected, by starting from a general position where consistency of the
 220 maximum likelihood estimator towards a set of pseudo-true parameters (true parameters in
 221 the case of the true GLLVM) can be attained. By doing so, the results we develop will apply
 222 to any framework provided estimation consistency within that framework can be achieved.
 223 In Appendix B, we elaborate on this and provide a more concrete example under which
 224 Assumption 1 holds.

225 3.1 Regression Coefficients

226 We first consider the case of conditionally normally distributed responses and the identity
 227 link, $\mu_j(\mathbf{s}) = \eta_j(\mathbf{s})$.

228 **THEOREM 1:** *For conditionally normal responses, the regression coefficients satisfy $\boldsymbol{\beta}^* =$*
 229 *$\boldsymbol{\beta}^0$. Thus under Assumption 1, the maximum likelihood estimator of the misspecified GLLVM*
 230 *satisfies $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}^0$ as $n \rightarrow \infty$, irrespective of the values of $\boldsymbol{\lambda}^0$.*

231 All proofs are provided in Appendix B. With non-normal responses, a weaker zero consis-
 232 tency result can be obtained as follows. For a nonempty subset $\mathcal{A} \subseteq \{2, \dots, q\}$, let $\mathbf{x}_{\mathcal{A}}(\mathbf{s})$ and
 233 $\mathbf{x}_{\mathcal{A}^c}(\mathbf{s})$ denote the vectors of covariates corresponding to \mathcal{A} and its complement respectively.

234 **THEOREM 2:** *Suppose there exists a subset of covariates $\mathcal{A} \subseteq \{2, \dots, q\}$ with $\mathcal{A} \neq \emptyset$, such*

235 that for all $j = 1, \dots, p$ we can write $\beta_j^0 = (\beta_{j,\mathcal{A}}^0 = \mathbf{0}, \beta_{j,\mathcal{A}^c}^0)^\top$. Assuming $\mathbf{x}_{\mathcal{A}}(\mathbf{s})$ and $\mathbf{x}_{\mathcal{A}^c}(\mathbf{s})$
 236 are independent, then the corresponding subset of β_j^* will also equal zero i.e., $\beta_{j,\mathcal{A}}^* = \mathbf{0}$. Thus
 237 under Assumption 1, the maximum likelihood estimator of the misspecified GLLVM satisfies
 238 $\hat{\beta}_{j,\mathcal{A}} \xrightarrow{p} \mathbf{0}$ for all $j = 1, \dots, p$, as $n \rightarrow \infty$, and irrespective of the values of $\boldsymbol{\lambda}^0$.

239 The above result ensures that for any covariates which are truly uninformative for all
 240 species, the corresponding estimates of the misspecified GLLVM will consistently estimate
 241 these zero coefficients. While weaker than the full consistency result for conditionally normal
 242 responses, Theorem 2 may serve as a useful basis for studying variable selection in spatial
 243 GLLVMs e.g., establishing selection consistency of information criteria or penalized likelihood
 244 methods for choosing covariates driving the entire species community. Note that for the
 245 above result to hold, we require the strong assumption of independence between the truly
 246 informative and truly uninformative covariates. Such an assumption has been made previ-
 247 ously in order to establish consistency when studying misspecification in generalized linear
 248 mixed models (e.g., Litire et al., 2007), and is similar to the partial orthogonality condition
 249 imposed in high-dimensional variable selection (e.g., Huang et al., 2008). In Section 4, we
 250 shall empirically assess the generality of Theorem 2 in settings with correlated covariates.

251 It is important to note that Theorem 2 demonstrates zero consistency only for completely
 252 uninformative covariates: for a partly informative covariate that is important for some but
 253 not all species, zero consistency cannot be guaranteed for species whose responses are not
 254 related to this covariate. Intuitively, this is because there remains an indirect dependence on
 255 this covariate through the combination of the direct dependence on the covariate for some
 256 species and the correlation between species induced by the latent variables. In Section 4, we
 257 shall empirically study the effects of misspecification on estimation and inference for partly
 258 informative covariates.

3.2 Loadings

Since the loadings directly relate to the correlation structures in the GLLVM, then one would expect that misspecifying the correlation structure in the latent variables should adversely affect their estimation and inference. This is largely the case, although a zero consistency result can be achieved.

THEOREM 3: *Suppose there exists a subset $\mathcal{B} \subseteq \{1, \dots, p\}$ with $\mathcal{B} \neq \emptyset$, such that $\boldsymbol{\lambda}_j^0 = \mathbf{0}$ for all $j \in \mathcal{B}$. Then the corresponding subset of $\boldsymbol{\lambda}_j^*$ will also equal zero i.e., $\boldsymbol{\lambda}_j^* = \mathbf{0}$. Thus under Assumption 1, the maximum likelihood estimator of the misspecified GLLVM satisfies $\hat{\boldsymbol{\lambda}}_j \xrightarrow{p} \mathbf{0}$ for all $j \in \mathcal{B}$, as $n \rightarrow \infty$.*

While Theorem 3 may not be important in community ecology, since we do not expect many (if any at all) species to be completely uncorrelated to other species, it may still play a useful role in other applications of GLLVMs where such independence can arise (e.g., Hirose and Konishi, 2012). Intuitively, any species that is uncorrelated with all others in the model remains uncorrelated under misspecification of the latent variable correlation structure. Therefore, we cannot borrow strength from other species to improve estimation and inference of the parameters specific to this species. This leads to the following result.

COROLLARY 1: *For any species $j \in \mathcal{B}$ that is independent of all other species on the linear predictor scale of the model, as defined in Theorem 3, the regression coefficients satisfy $\boldsymbol{\beta}_j^* = \boldsymbol{\beta}_j^0$. Thus under Assumption 1, the maximum likelihood estimator of the misspecified GLLVM satisfies $\hat{\boldsymbol{\beta}}_j \xrightarrow{p} \boldsymbol{\beta}_j^0$ for all $j \in \mathcal{B}$, as $n \rightarrow \infty$.*

4. Simulation Study

We conducted a simulation study to assess the finite sample performance of misspecified GLLVMs, relative to GLLVMs where the true spatial correlation structure for the latent

282 variables was known. Specifically, we simulated data from a spatial GLLVM with $p = 20$
 283 species, $q = 7$ covariates including the intercept, and $d = 3$ latent variables. The full details
 284 of the simulation design are provided in Appendix C, and we only provide some of the
 285 defining features below. Specifically, the n sites were arranged in a square lattice to reflect
 286 an expanding domain framework. Excluding the intercept, the covariates $\mathbf{x}(\mathbf{s}_i)$ included three
 287 spatially structured covariates and three spatially independent (but correlated to each other)
 288 covariates. The covariates were also generated in a way such that the conditions required
 289 for Theorem 2 are not satisfied, thereby allowing us to assess performance under a more
 290 realistic setting. The regression coefficients β_j varied both in magnitude and sign, and were
 291 constructed so that two covariates were informative for all species, two were informative for
 292 half of the species, and two were uninformative for all species. Similarly, the loadings λ_j in
 293 the spatial GLLVM were constructed such that species 16 to 20 were independent of all other
 294 species, while the remaining fifteen species had loadings varying in magnitude and size.

295 We generated multivariate spatial data with four possible response types i.e., four possi-
 296 ble choices for the conditional distribution $f\{y_j(\mathbf{s}_i)|\mathbf{x}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_i)\}$: 1) continuous responses
 297 from the normal distribution; 2) overdispersed counts from the negative binomial distri-
 298 bution; 3) non-negative continuous responses from the Tweedie distribution; 4) presence-
 299 absence responses from the Bernoulli distribution. We considered grids of dimension $n^{1/2} =$
 300 7, 10, 14, 22, 32, and simulated 400 datasets for each value of n . The grid sizes were chosen
 301 such that the total number of sites approximately doubled with each grid.

302 4.1 Model Fitting and Performance Assessment

303 We fitted two models to each simulated dataset: a spatial GLLVM with marginal likelihood
 304 given by (1) and assuming the true spatial correlation structure for the latent variables is
 305 known, and a misspecified independent GLLVM with marginal likelihood given by (2). In
 306 addition to point estimates, we calculated 95% Wald confidence intervals for the regression

307 coefficients and loadings using both models. For the true model, this was based on stan-
 308 dard errors obtained from the inverse of the observed information matrix, $\hat{\mathbf{I}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\theta}}) =$
 309 $-\nabla^2 \log L(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\theta}})$ where $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\theta}})^\top$ generically denotes the corresponding maximum
 310 likelihood estimates. For the misspecified model, we calculated standard errors based on the
 311 inverse of the sandwich information matrix. Let $L_{\text{Ind},i}(\boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\lambda}) =$
 312 $\int \prod_{j=1}^p f\{y_j(\mathbf{s}_i) | \mathbf{x}(\mathbf{s}_i), \mathbf{u}(\mathbf{s}_i)\} h\{\mathbf{u}(\mathbf{s}_i)\} d\mathbf{u}(\mathbf{s}_i)$ denote the likelihood for the misspecified GLLVM
 313 for site i . Then we define $\hat{\mathbf{I}}_{\text{sand}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) = \hat{\mathbf{H}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) \hat{\mathbf{J}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}})^{-1} \hat{\mathbf{H}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}})$ where

$$\begin{aligned}
 \hat{\mathbf{H}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) &= -\frac{1}{n} \sum_{i=1}^n \nabla^2 \log L_{\text{Ind},i}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) \\
 \hat{\mathbf{J}}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) &= \frac{1}{n} \sum_{i=1}^n \nabla \log L_{\text{Ind},i}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}}) \{\nabla \log L_{\text{Ind},i}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\lambda}})\}^\top.
 \end{aligned}$$

314 The sandwich information matrix is generally used as the basis for quantifying uncertainty
 315 for maximum likelihood estimates under misspecified models, and we refer the reader to
 316 (White, 1982; Verbeke and Lesaffre, 1997) on its use in misspecified models.

317 All GLLVMs were estimated using Template Model Builder (TMB, Kristensen et al., 2015),
 318 which uses a combination of automatic differentiation and the Laplace approximation to
 319 produce functions for efficiently calculating the gradient and Hessian of the (Laplace approx-
 320 imated) marginal log-likelihood. For the spatial GLLVM, we made use of R-INLA (Lindgren
 321 and Rue, 2015) to construct a stochastic partial differential equation approximation to the
 322 distribution of the spatial latent variables, which ensured the computation remained feasible
 323 for large n (see also Thorson et al., 2015). To assess performance of the regression coefficients,
 324 we calculated the empirical bias and root mean squared error (RMSE) averaged across
 325 the simulated datasets, and the empirical coverage probability and mean interval width
 326 of the 95% confidence intervals. To assess performance of the loadings and latent variables
 327 predictions, we calculated the Procrustes errors between the estimated and true loading
 328 matrix, and between the estimated and true latent variable matrix (which is formed by

stacking the n sets of vectors $\mathbf{u}(\mathbf{s}_i)$ as rows). Finally, we recorded the computation time (in minutes) used to fit the GLLVM and calculate the associated 95% confidence intervals. We present these results in Appendix C since, as to be expected, the simpler misspecified independent GLLVM was always faster and scaled better with increasing sample size across all response types.

4.2 Setting 1: Normal Responses

With regards to the regression coefficients, the biggest factor determining performance for misspecified GLLVMs was whether the corresponding covariate was spatially correlated or not (Figure 1). Thus we shall discuss our results based on these two cases separately.

For spatially uncorrelated covariates, both the bias and RMSE of the regression coefficients tended to zero for misspecified GLLVMs as sample size increased (Figure 1 left and middle columns, solid circles). However the RMSE for misspecified GLLVMs was consistently higher than for true GLLVMs, suggesting more variability in the estimates under misspecification. There was also little difference between the true and misspecified GLLVMs in terms of coverage probability for coefficients corresponding to spatially uncorrelated covariates, with the intervals for both models tending to the nominal level irrespective of the size of the regression coefficient and the norm of the corresponding loadings (Figure 1 right column). The sandwich based confidence intervals from the misspecified GLLVM tended to be wider than the intervals based on the true GLLVM, especially for small to medium sized regression coefficients (see Appendix C).

It is startling to see how much the performance deteriorated for coefficients corresponding to spatially correlated covariates under misspecification: the empirical bias and RMSE were substantially larger compared to both estimates from the true GLLVM, as well as estimates corresponding to spatially uncorrelated covariates (Figure 1 left and middle columns, empty circles). The sandwich based confidence interval suffered severe undercoverage, with

354 performance becoming worse with increasing sample size. This poor performance occurred
355 irrespective of the true magnitude of the regression coefficient, but did depend on the norm
356 of the corresponding loadings.

357 The difference in performance between spatially correlated versus uncorrelated covariates,
358 under misspecification, can be attributed to the former exhibiting a strong correlation to
359 the spatially structured latent variables. Specifically, a non-negligible proportion of the
360 variation in each spatially correlated covariate can be explained by a linear combination of
361 the eigenvectors of the spatial covariance matrix $\Sigma(\boldsymbol{\theta})$. Under misspecification, because the
362 latent variables are assumed to be independent, the independent GLLVM then erroneously
363 attributes part of this spatial correlation in the latent variables to the covariates. This causes
364 estimation and inference for the corresponding elements of $\boldsymbol{\beta}_j$ to deteriorate as the estimated
365 coefficients become affected by the values of $\boldsymbol{\lambda}_j$. By contrast, spatially uncorrelated covariates
366 (in general) exhibit little correlation to spatially structured latent variables, and so the
367 erroneous attribution does not occur and estimation and inference for the corresponding
368 coefficients are largely unaffected. This difference in performance has close connections to
369 the issue of spatial confounding, and we expand upon this issue in Section 6 and in Appendix
370 F.

371 [Figure 1 about here.]

372 Neither estimation of the loadings nor prediction of the latent variables was robust to
373 misspecification of the correlation structure, with the Procrustes errors for the misspecified
374 GLLVM being consistently higher than those from the true GLLVM (Figure 2 top row).
375 Additional analyses (not shown) of the results for loadings reveal that, consistent with the
376 zero consistency result of Theorem 3, it was the non-zero loadings (from species 1 to 15) that
377 were poorly estimated under misspecification, while the truly zero loadings (from species 16
378 to 20) were estimated well under misspecification of the latent variable correlation structure.

379 Of course, given the above results a natural question to ask is whether performance of the
 380 misspecified independent GLLVM would improve if we were to fit a model with a larger
 381 number of latent variables. In Appendix C, we perform such additional simulations where
 382 we fit independent GLLVMs with $d = 5, 7, 9$, noting that true number of latent variables in
 383 the simulation is $d = 3$. In summary, results for normal responses show that estimation and
 384 inference on the regression coefficients is hardly affected by increasing the number of latent
 385 variables, while prediction performance actually deteriorates (as quantified by the Frobenius
 386 norm between the matrix of estimated and true linear predictor component $\mathbf{u}(\mathbf{s}_i)^\top \boldsymbol{\lambda}_j$).

387 [Figure 2 about here.]

388 4.3 *Setting 2: Negative Binomial Counts*

389 The overall patterns seen for the normal response case largely carry over to this setting. For
 390 spatially uncorrelated covariates, the empirical bias and RMSE of the regression coefficients
 391 tended to zero for misspecified GLLVMs (Figure 3 left and middle columns, solid circles). This
 392 finding expands the zero consistency result of Theorem 2 to the case of covariates that are
 393 correlated with each other, and to truly zero coefficients coming from partly informative as
 394 well as completely uninformative covariates. There was also little difference between the true
 395 and misspecified GLLVMs in terms of the coverage probability of the confidence intervals for
 396 coefficients corresponding to spatially uncorrelated covariates, with intervals for both models
 397 tending to the nominal level irrespective of the size of the coefficient and the corresponding
 398 norm of the loadings (Figure 3 right column). By contrast, for spatially correlated covariates,
 399 estimation and inference of regression coefficients suffered severely under misspecification,
 400 with considerably higher RMSE and major undercoverage especially for coefficients that
 401 also had a large corresponding norm of the loadings. Again, the poor performance can be
 402 explained by the spatially correlated covariates exhibiting a strong correlation to the spatially
 403 structured latent variables, and so the misspecified GLLVM erroneously attributes part of

404 the latent variable component of the model to the covariates, leading to poor estimation and
405 inference on the corresponding coefficients.

406 [Figure 3 about here.]

407 Estimation of the loadings and prediction of the latent variables both suffered under
408 misspecification of the correlation structure for negative binomial GLLVMs (Figure 2 bottom
409 row), with additional analyses again revealing that it was the non-zero loadings that were
410 particularly poorly estimated under misspecification, while the truly zero loadings were
411 estimated relatively well.

412 Simulation results for the Tweedie response GLLVM presented similar trends to those for
413 the negative binomial counts and normal responses seen above, and the results for these
414 are presented in Appendix C. In particular, estimation and inference on the regression
415 coefficients was largely robust to misspecification for spatially uncorrelated variables, but
416 suffered severely misspecification for spatially correlated variables. Estimation of the loadings
417 and prediction of the latent variables performed consistently poorly under misspecification
418 of the latent variable correlation structure.

419 4.4 *Setting 3: Binary Responses*

420 In contrast to previous settings where estimation and inference of coefficients correspond-
421 ing to spatially uncorrelated covariates was largely robust to misspecification, for binary
422 responses we observe that the degree of robustness also depended heavily on whether the
423 true value of the coefficient was itself close to or exactly equal to zero. In the left and
424 middle columns of Figure 4, we observe eight notably outlying regression coefficients (four
425 corresponding to spatially uncorrelated covariates and four corresponding to spatially corre-
426 lated covariates) that were both large in their true magnitude and had large corresponding
427 norm of the loadings, which performed especially poorly under misspecification. These eight
428 coefficients also had coverage probabilities for the corresponding sandwich based confidence

429 intervals that were notably below the nominal level (Figure 4 right column). This is despite
430 the fact that the sandwich based confidence intervals tended to be relatively wide in this
431 setting overall (see Appendix C).

432 [Figure 4 about here.]

433 Comparing against results from previous simulating settings, this suggests that for binary
434 responses, misspecification of the latent variable correlation structure can have a major
435 impact on inference for non-zero β_j 's (even for spatially uncorrelated covariates), *if* the
436 corresponding norm of the loadings is also relatively large. For coefficients corresponding
437 to spatially correlated covariates, misspecification again led to poor performance overall
438 (Figure 4 empty circles).

439 The results for estimation of the loadings and prediction of the latent variables were similar
440 to those seen with the other response types in Sections 4.2 and 4.3, and so for brevity we
441 present them in the Appendix C.

442 5. Application to SO-CPR Survey

443 We fitted negative binomial GLLVMs to data collected in season 2007–2008 of the SO-
444 CPR survey, with the aim being to quantify the relationship between the species responses
445 and select environmental covariates. The data consisted of $n = 3,875$ sampling locations
446 irregularly spaced across the Southern Ocean. Furthermore, we focused on $p = 24$ species
447 detected at more than 5% of the sites (see Appendix D for a figure of the sampling locations
448 and the species included in the analyses). We included two environmental covariates that were
449 *a priori* considered to be important in influencing one or more of the species distributions:
450 salinity (which is highly correlated with temperature) and photosynthetically active radiation
451 (PAR). The latter is a measure of amount of light available for photosynthesis. Using
452 orthogonal polynomials, we entered both covariates into the GLLVMs as linear and quadratic

453 terms, and also considered a pairwise interaction term between them. Along with an intercept
454 term, this produced a total of $q = 6$ covariates in $\mathbf{x}(\mathbf{s}_i)$ and $pq = 144$ coefficients in $\boldsymbol{\beta}$. We
455 fitted GLLVMs with $d = 3$ latent variables and considered two possible spatial correlation
456 structures: 1) a Matérn covariance function with $\nu = 1$; 2) an independence correlation
457 structure.

458 The results are presented in Appendix D, from which we observe that many of the estimated
459 coefficients from the two fitted models were similar in magnitude and sign, revealing that
460 salinity and PAR, along with the quadratic and pairwise interactions terms, were important
461 factors in driving many of the species responses. There were however some notable differences:
462 focusing on whether the 95% confidence intervals contained zero or not, each covariate
463 had several species that produced differing conclusions across the two models (5 for the
464 linear effect of salinity, 3 for the quadratic effect of salinity, 3 for the linear effect of PAR,
465 6 for the quadratic effect of PAR, and 5 for the interaction between salinity and PAR).
466 The majority of species had a negative coefficient of the quadratic effect of salinity, thus
467 supporting the idea of a “niche” in this environment space and consequently of particular
468 water masses between Tasmania and Antarctica. By contrast, many species had significant
469 negative linear effects and significant positive quadratic effects of PAR, suggesting almost all
470 marine species for analysis tended to prefer low light environments. Regarding computation
471 time, the independent GLLVM took 3.70 hours to complete (on an Intel Xeon E5-2680 V3
472 at 2.50 GHz with 5 CPUs), while the spatial GLLVM took 15.5 hours to fit.

473 6. Discussion

474 We explored the consequences of misspecifying the spatial correlation structure of the latent
475 variables in GLLVMs. If the aim of the analysis centers on estimation and inference for the
476 regression coefficients, our findings show that by far the strongest determinant of performance
477 is whether the covariate is spatially correlated or not. For spatially uncorrelated covariates,

478 estimation and inference for the coefficients is largely robust to misspecification, with the
479 exception of binary responses. But for spatially correlated covariates, which arise quite
480 frequently in observational studies in ecology, assuming independence of the latent variables
481 leads to poor performance for the corresponding regression coefficients. It is also important
482 to account for potential spatial correlation if the aim of the analysis centers on estimation
483 and inference for loadings, and for prediction from the GLLVM as a whole.

484 In our simulation study, we only generated covariates that were either spatially independent
485 or covariates with a fixed spatial scale. This raises the question of what happens if we were
486 to consider covariates with varying degrees of spatial correlation. To answer this question,
487 we conducted a new set of simulations where multivariate spatial data were generated from
488 a spatial GLLVM where the covariates possessed varying strengths of spatial strengths.
489 The full details are presented in Appendix E, but in brief, results exhibited a general
490 and consistent trend across all four response types, such that estimation and inference
491 on the regression under the misspecified independent GLLVM gradually worsened as the
492 corresponding covariate became more strongly spatially correlated. By contrast, estimation
493 and inference from fitting the true spatial GLLVM were largely unaffected by the strength
494 of spatial correlation for the covariates. These results provide further evidence that, if there
495 is *a priori* information covariates are spatially correlated, then it is imperative to account
496 for the spatial correlation when specifying the latent variable structure.

497 Comparing across empirical and theoretical results, there may appear to be some disagree-
498 ment in terms of point estimation of the regression coefficients: Theorems 1 and 2 state that
499 consistency is achievable under misspecification for all or some of the coefficients respectively,
500 yet our simulations suggest issues with consistency for coefficients corresponding to spatially
501 correlated covariates. While not definitive, and we encourage future research into the large
502 sample properties of maximum likelihood estimation for misspecified spatial GLLVMs, this

503 discordance may be because spatially correlated covariates possibly do not entirely satisfy
504 the conditions underlying Assumption 1. As discussed in Section 4.2, there may be potential
505 issues surrounding the identifiability of coefficients for spatially correlated covariates when
506 a misspecified independent GLLVM is fitted, and the problems bears similarity to the issue
507 of spatial confounding. We expand upon this connection in Appendix F, noting that while
508 spatial confounding has been investigated in some detail for univariate spatial regression,
509 there is very little research on its impact for spatial GLLVMs (see Shirota et al., 2019, for
510 a recent exception that presented an interesting discussion on this topic), let alone what it
511 means to “correct” for spatial confounding in misspecified independent GLLVM where the
512 latent variables are, by definition, spatially independent.

513 While the above results have focused on the most extreme type of misspecification in the
514 latent variable correlation structure (by assuming independence), a natural avenue of research
515 is to examine cases of “slight misspecification” e.g., we include spatially correlated latent
516 variables but the form of the spatial covariance is misspecified, or the true spatial correlation
517 structures differ across the latent variables and but we misspecify and assume the same
518 structure. Of course, we must be mindful that there is a limitless number of ways one can
519 “mix and match” different types of spatially structured latent variables (e.g., we could have
520 latent variables included at different scales to reflect a nested design, similar to Ovaskainen
521 et al., 2016, and misspecify the correlation structure at one scale but not another), let alone
522 the issue of what happens if we misspecify any combination or set of combinations of these.
523 In Appendix G we present results from an additional simulation study where we generate
524 data from a spatial GLLVM, and compare the performance of an independent GLLVM versus
525 a “slightly misspecified” spatial GLLVM where the wrong smoothness parameter ν in the
526 Matérn correlation function is assumed. Overall, and not surprisingly, the performance of
527 the slightly misspecified spatial GLLVM is generally better than the independent GLLVM.

528 However there is effectively zero computational gain from fitting this slightly misspecified
529 spatial GLLVM, and so in some sense is against the motivation and spirit of this article i.e.,
530 our aim is to examine the consequences on estimation and inference when we are far away
531 from the true spatial structure, where this is traded off against substantial reductions in
532 computation time.

533 Finally, in our simulations we assumed that the true number of latent variables is known,
534 when in practice this also has been chosen (e.g., Hui et al., 2018). How misspecifying the
535 latent variable correlation structures affects the number of latent variables chosen is an avenue
536 of future research to pursue, along with the more general topic of how misspecification of
537 the latent variable correlation structure affects variable selection as well as other aspects of
538 inference for GLLVMs as a whole (see Hoeting et al., 2006; Xu et al., 2015, for examples of
539 related research on variable selection in geostatistical models when the spatial correlation is
540 misspecified).

541 ACKNOWLEDGEMENTS

542 Both FKCH and AHW were supported by Australia Research Council Discovery Grants.
543 Thanks to Noel Cressie, Anders Nielsen, and Graham Hosie for useful discussions.

544 DATA AVAILABILITY STATEMENT

545 The data that support the findings in this paper are openly available as part of the “Aus-
546 tralian Antarctic Data Centre, at <http://dx.doi.org/doi:10.26179/5ee84f77cc4ec>, cited
547 as Hosie (2020).

548 REFERENCES

549 Bjork, J. R., Hui, F. K. C., O’Hara, R. B., and Montoya, J. M. (2018). Uncovering the
550 drivers of host-associated microbiota with joint species distribution modeling. *Molecular*

- 551 *ecology* **27**, 2714–2724.
- 552 Hirose, K. and Konishi, S. (2012). Variable selection via the weighted group lasso for factor
553 analysis models. *Canadian Journal of Statistics* **40**, 345–361.
- 554 Hoeting, J. A., Davis, R. A., Merton, A. A., and Thompson, S. E. (2006). Model selection
555 for geostatistical models. *Ecological Applications* **16**, 87–98.
- 556 Hosie, G. (2020). Southern Ocean Continuous Plankton Recorder Zooplankton Records.
557 *Australian Antarctic Data Centre*. Version 8. 10.26179/5ee84f77cc4ec .
- 558 Huang, J., Horowitz, J. L., and Ma, S. (2008). Asymptotic properties of bridge estimators
559 in sparse high-dimensional regression models. *The Annals of Statistics* **36**, 587–613.
- 560 Hui, F. K. C., Mueller, S., and Welsh, A. H. (2020). Random effects misspecification can have
561 severe consequences for random effects inference in linear mixed models. *International*
562 *Statistical Review* **In Press**,.
- 563 Hui, F. K. C., Tanaka, E., and Warton, D. I. (2018). Order selection and sparsity in latent
564 variable models via the ordered factor LASSO. *Biometrics* **74**, 1311–1319.
- 565 Irincheeva, I., Cantoni, E., and Genton, M. G. (2012). Generalized linear latent variable
566 models with flexible distribution of latent variables. *Scandinavian Journal of Statistics*
567 **39**, 663–680.
- 568 Kristensen, K., Nielsen, A., Berg, C., Skaug, H., and Bell, B. (2015). TMB: Automatic
569 Differentiation and Laplace Approximation. *Journal of Statistical Software* **70**, 1–21.
- 570 Kurisu, D. (2019). On nonparametric inference for spatial regression models under domain
571 expanding and infill asymptotics. *Statistics & Probability Letters* **154**, 108543.
- 572 Lindgren, F. and Rue, H. (2015). Bayesian spatial modelling with R-INLA. *Journal of*
573 *Statistical Software* **63**, 1–25.
- 574 Litire, S., Alonso, A., and Molenberghs, G. (2007). Type I and Type II Error Under Random-
575 Effects Misspecification in Generalized Linear Mixed Models. *Biometrics* **63**, 1038–1044.

- 576 Lopes, H. F., Gamerman, D., and Salazar, E. (2011). Generalized spatial dynamic factor
577 models. *Computational Statistics & Data Analysis* **55**, 1319–1330.
- 578 Lu, Z. and Tjostheim, D. (2014). Nonparametric estimation of probability density functions
579 for irregularly observed spatial data. *Journal of the American Statistical Association*
580 **109**, 1546–1564.
- 581 Ma, Y. and Genton, M. G. (2010). Explicit estimating equations for semiparametric
582 generalized linear latent variable models. *Journal of the Royal Statistical Society Series*
583 *B* **72**, 475–495.
- 584 Mardia, K. V. and Marshall, R. J. (1984). Maximum likelihood estimation of models for
585 residual covariance in spatial regression. *Biometrika* **71**, 135–146.
- 586 Niku, J., Brooks, W., Herliansyah, R., Hui, F. K. C., Taskinen, S., and Warton, D. I. (2019).
587 Efficient estimation of generalized linear latent variable models. *PloS one* **14**, e0216129.
- 588 Ovaskainen, O., Roy, D. B., Fox, R., and Anderson, B. J. (2016). Uncovering hidden spatial
589 structure in species communities with spatially explicit joint species distribution models.
590 *Methods in Ecology and Evolution* **7**, 428–436.
- 591 Shirota, S., Gelfand, A. E., and Banerjee, S. (2019). Spatial joint species distribution
592 modeling using Dirichlet processes. *Statistica Sinica* **29**, 1127–1154.
- 593 Skrondal, A. and Rabe-Hesketh, S. (2004). *Generalized Latent Variable Modeling: Multilevel,*
594 *Longitudinal, and Structural Equation Models*. CRC Press.
- 595 Thorson, J. T., Scheuerell, M. D., Shelton, A. O., See, K. E., Skaug, H. J., and Kristensen,
596 K. (2015). Spatial factor analysis: a new tool for estimating joint species distributions
597 and correlations in species range. *Methods in Ecology and Evolution* **6**, 627–637.
- 598 Verbeke, G. and Lesaffre, E. (1997). The effect of misspecifying the random-effects distri-
599 bution in linear mixed models for longitudinal data. *Computational Statistics & Data*
600 *Analysis* **23**, 541–556.

- 601 Wang, H. and Zhu, J. (2009). Variable selection in spatial regression via penalized least
602 squares. *Canadian Journal of Statistics* **37**, 607–624.
- 603 Warton, D. I., Blanchet, F. G., OHara, R., Ovaskainen, O., Taskinen, S., Walker, S. C., and
604 Hui, F. K. C. (2016). Extending Joint Models in Community Ecology: A Response to
605 Beissinger et al. *Trends in Ecology & Evolution* **31**, 737–738.
- 606 White, H. (1982). Maximum Likelihood Estimation of Misspecified Models. *Econometrica*
607 **50**, 1–26.
- 608 Xu, L., Wang, Y.-G., Zheng, S., and Shi, N.-Z. (2015). Model selection with misspecified
609 spatial covariance structure. *Journal of Statistical Computation and Simulation* **85**,
610 2276–2294.
- 611 Yoon, H.-J. and Welsh, A. H. (2020). On the effect of ignoring correlation in the covariates
612 when fitting linear mixed models. *Journal of Statistical Planning and Inference* **204**,
613 18–34.
- 614 Zhu, J., Eickhoff, J., and Yan, P. (2005). Generalized linear latent variable models for
615 repeated measures of spatially correlated multivariate data. *Biometrics* **61**, 674–683.

SUPPORTING INFORMATION

616

617 Web Appendices, figures, and additional details and discussion referenced in Sections 3 to 5
618 are available with this paper at the Biometrics website on Wiley Online Library. Additionally,
619 template R scripts for estimating GLLVMs and for performing all simulations are provided.

Figure 1. Simulation results for empirical bias (left column), root mean squared error (middle column), and coverage probability (right column) of regression coefficients β_j , for normal response GLLVMs. Points are differentiated by shape ('•' for misspecified GLLVM coefficients corresponding to spatially uncorrelated covariates; 'o' for misspecified GLLVM coefficients corresponding to spatially correlated covariates; '△' for true GLLVM coefficients) and shading (darker shades are coefficients with larger corresponding norms of the loadings).

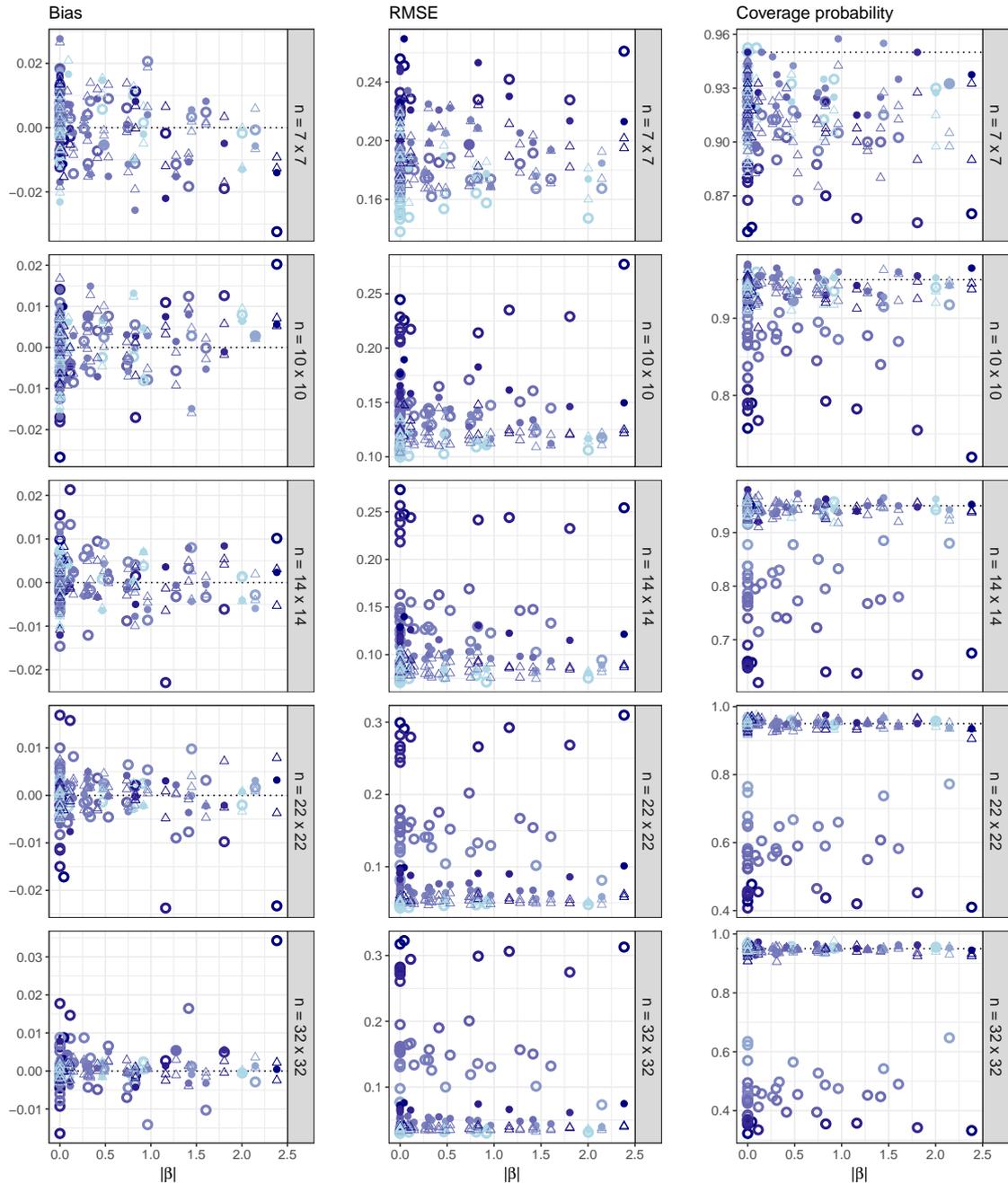


Figure 2. Simulation results showing comparative boxplots (left and red for misspecified GLLVM, right and blue for true GLLVM) for the Procrustes error of the estimated loadings Λ (left column) and Procrustes error of the predicted latent variables $\mathbf{u}(s_i)$ (middle right column), for GLLVMs with normal response (top row) and negative binomial response (bottom row).

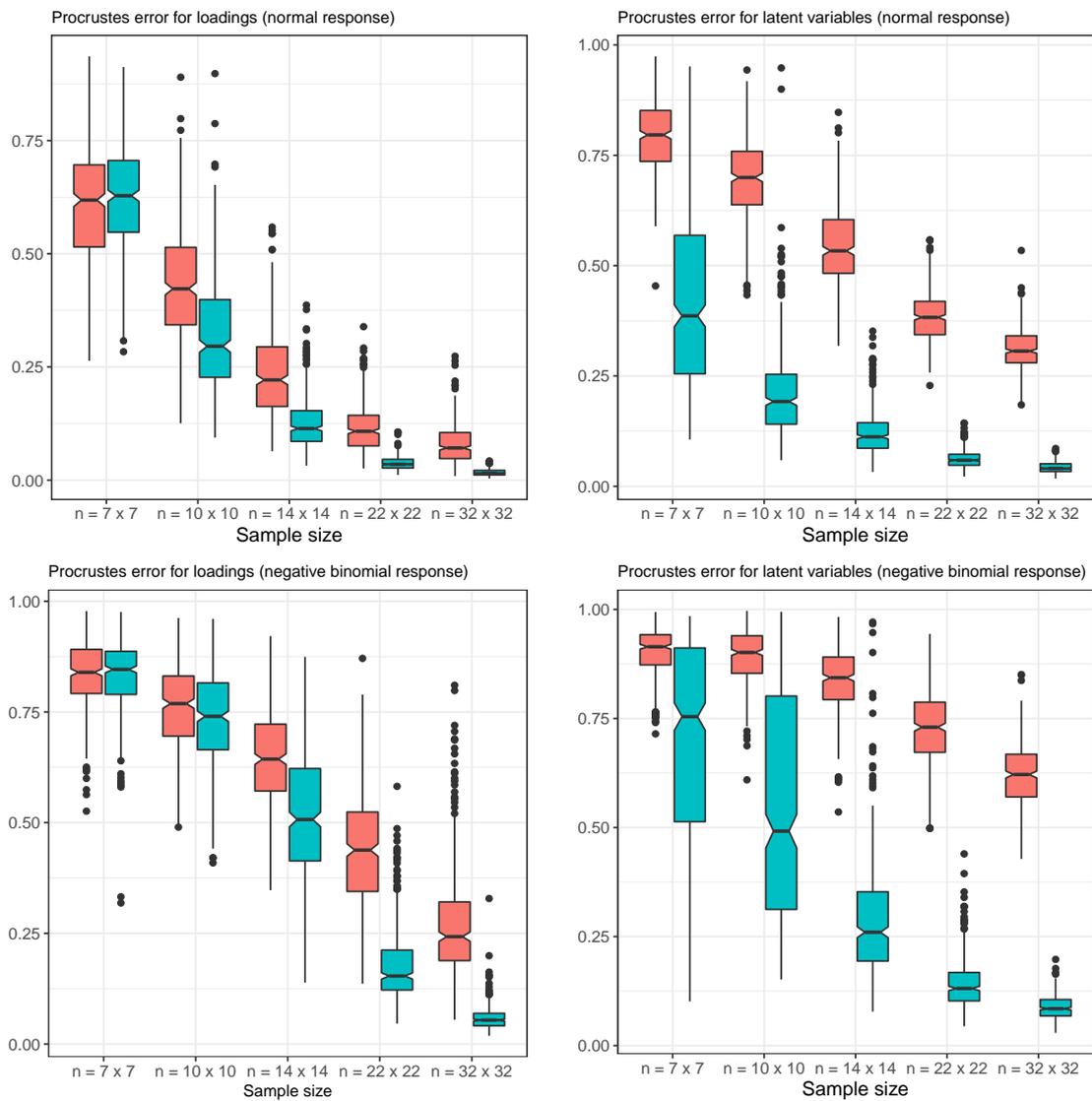


Figure 3. Simulation results for empirical bias (left column), root mean squared error (middle column), and coverage probability (right column) of regression coefficients β_j , for negative binomial GLLVMs. Points are differentiated by shape ('•' for misspecified GLLVM coefficients corresponding to spatially uncorrelated covariates; 'o' for misspecified GLLVM coefficients corresponding to spatially correlated covariates; '△' for true GLLVM coefficients) and shading (darker shades are coefficients with larger corresponding norms of the loadings).

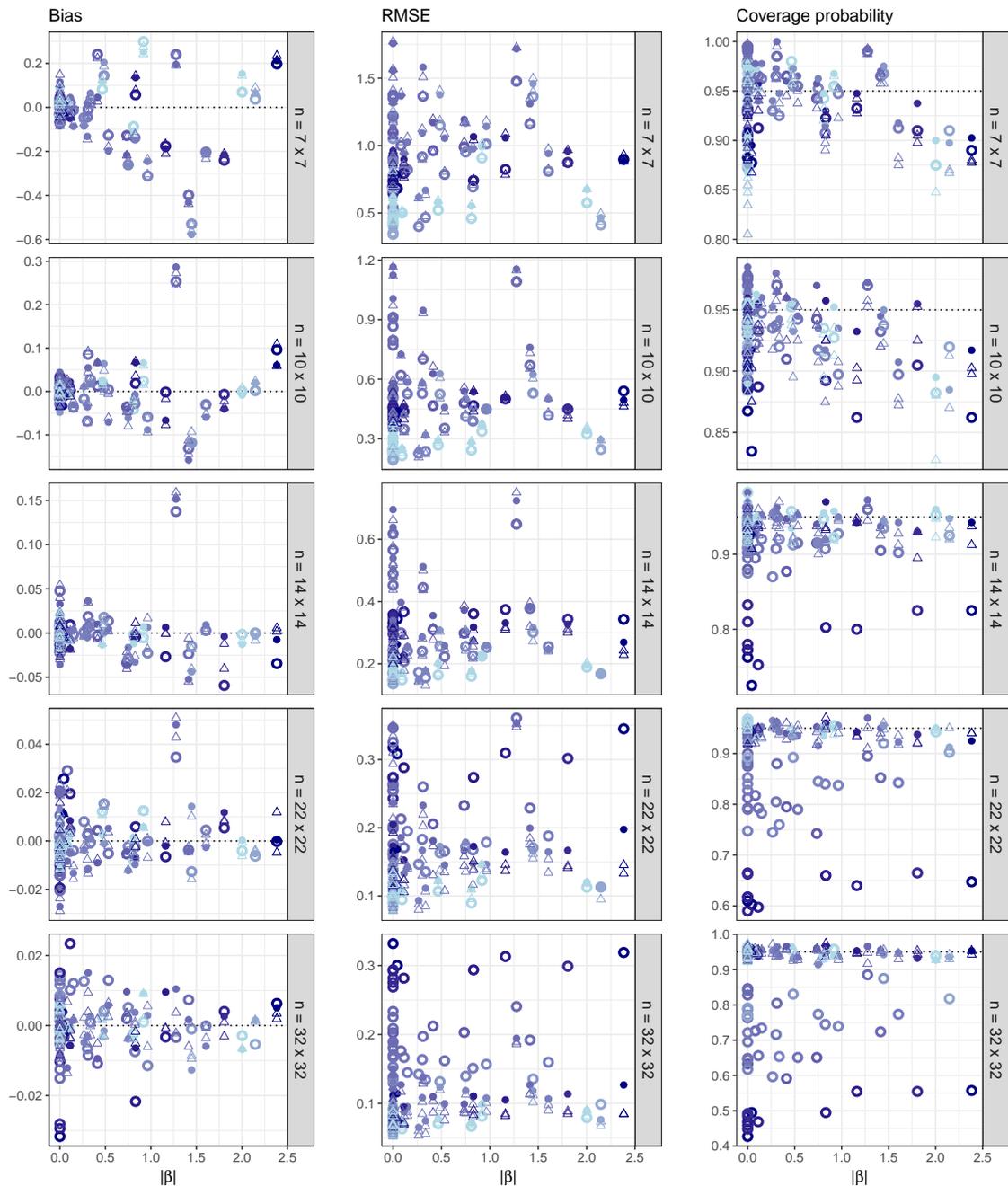


Figure 4. Simulation results for empirical bias (left column), root mean squared error (middle column), and coverage probability (right column) of regression coefficients β_j , for Setting 3 with binary GLLVMs. Points are differentiated by shape ('•' for misspecified GLLVM coefficients corresponding to spatially uncorrelated covariates; 'o' for misspecified GLLVM coefficients corresponding to spatially correlated covariates; '△' for true GLLVM coefficients) and shading (darker shades are coefficients with larger corresponding norms of the loadings).

