

# Joint Mean and Correlation Regression Modelling for Multivariate Data

Zhi Yang Tho, Francis Hui and Tao Zou

The Australian National University  
Research School of Finance, Actuarial Studies and Statistics

zhiyang.tho@anu.edu.au



Australian  
National  
University

**Abstract**  
 In the analysis of multivariate or multi-response data, researchers are often not only interested in studying how the mean (say) of each response evolves as a function of covariates, but also and simultaneously how the correlations between responses are related to one or more similarity/distance measures. To address such research questions, we propose a novel joint mean and correlation regression model, which is applicable to a wide variety of correlated discrete and (semi-)continuous responses. Simulations demonstrate the strong finite sample performance of the proposed estimator in terms of point estimation and inference. We apply the proposed joint mean and correlation regression model to a dataset of overdispersed counts of 38 Carabidae ground beetle species sampled throughout Scotland, with results showing in particular that beetle total length and breeding season have statistically important effects in driving the correlations between beetle species.

## Motivation

Given the data structure in Figure 1, we would like to simultaneously answer

- What is the relationship between environmental factors and different species?
- Which traits are associated with species correlation?

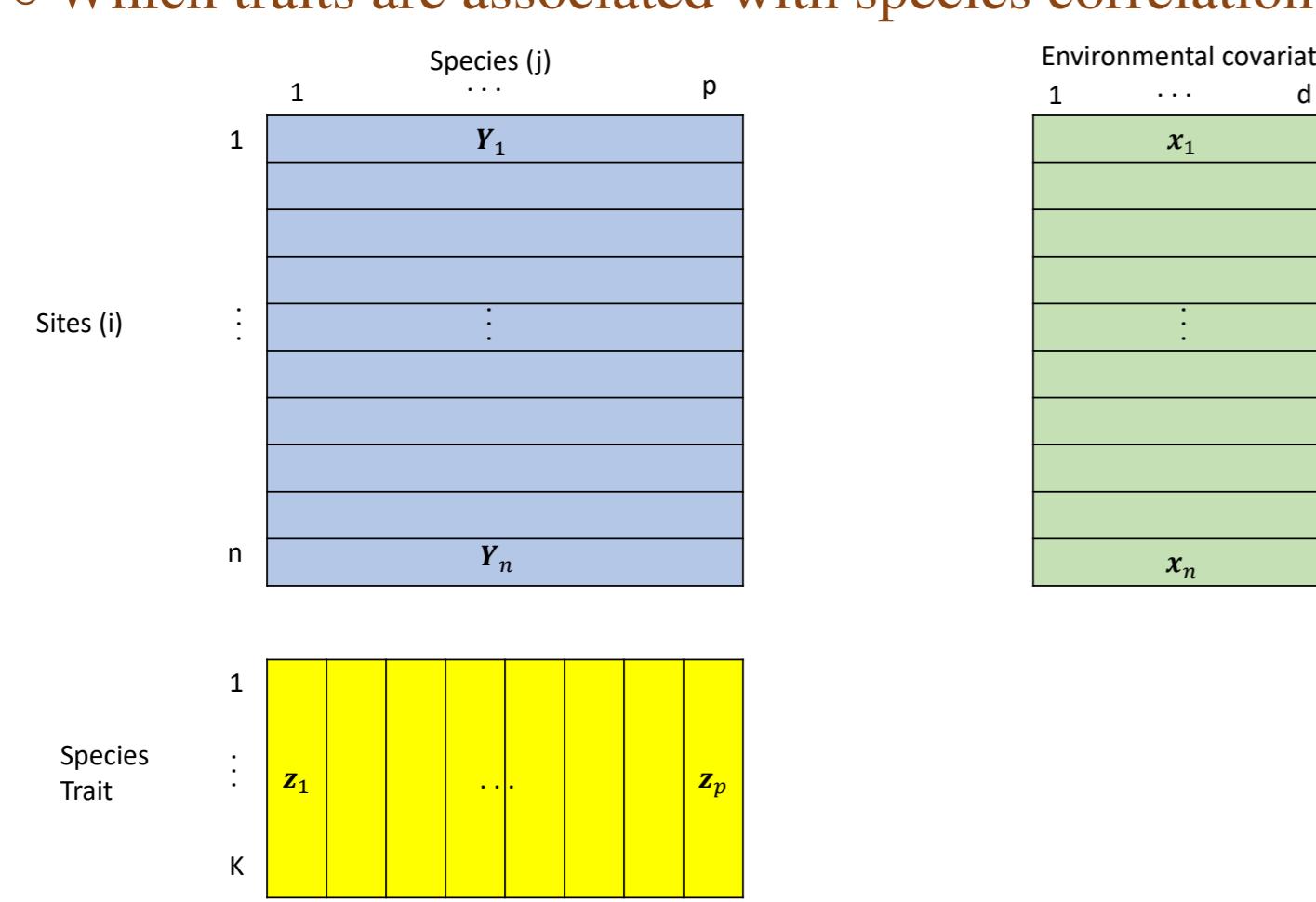


Figure 1: Data structure in community ecology that consists of response matrix, environmental covariate matrix and species trait matrix.

## Joint Mean and Correlation Regression Model

Let  $\mathbf{Y}_i$  be the  $p$ -dimensional response (discrete or continuous) vector and  $\mathbf{x}_i$  be the  $d$ -dimensional covariate vector of the  $i$ -th cluster and  $\mu_{ij}(\boldsymbol{\beta}_j) = E(Y_{ij})$ :

$$g\{\mu_{ij}(\boldsymbol{\beta}_j)\} = \mathbf{x}_i^\top \boldsymbol{\beta}_j, \text{ var}(Y_{ij}) = h\{\mu_{ij}(\boldsymbol{\beta}_j); \phi_j\}, \\ \text{corr}(Y_{ij_1}, Y_{ij_2}) = r_{j_1 j_2}(\boldsymbol{\rho}),$$

where  $g$  and  $h$  are known link and variance functions. The correlation matrix follows the correlation regression model

$$\mathbf{R}(\boldsymbol{\rho}) = (r_{j_1 j_2}(\boldsymbol{\rho})) = \mathbf{I}_p + \sum_{k=1}^K \rho_k \mathbf{W}_k,$$

where  $\mathbf{W}_k = (w_{ij}^{(k)})$  are similarity matrices induced from the auxiliary information vectors  $\mathbf{z}_j$ 's, e.g.,  $w_{j_1 j_2}^{(k)} = \exp(-|z_{j_1 k} - z_{j_2 k}|^2)$ . The  $\boldsymbol{\rho}$  vector possess a simple and intuitive interpretation as quantifying the impact of the similarities on the correlation between responses. For instance, a higher positive value of  $\rho_k$  implies that conditional on other trait values, two species with more similar values in their  $k$ -th trait variables are expected to have a stronger positive correlation after accounting for differences in their mean response due to the environmental covariates.

By stacking the response vectors into  $\mathbf{Y} = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_n^\top)^\top$ , we can re-express the model as

$$E(\mathbf{Y}) = \boldsymbol{\mu}(\boldsymbol{\beta}), \text{ cov}(\mathbf{Y}) = \mathbf{A}^{1/2}(\boldsymbol{\beta}) \{ \mathbf{I}_n \otimes \mathbf{R}(\boldsymbol{\rho}) \} \mathbf{A}^{1/2}(\boldsymbol{\beta}).$$

where  $\boldsymbol{\mu}(\boldsymbol{\beta})$  is the stacked vector of means,  $\mathbf{A}(\boldsymbol{\beta}) = \text{diag}\{\mathbf{A}_1(\boldsymbol{\beta}), \dots, \mathbf{A}_n(\boldsymbol{\beta})\}$  and  $\mathbf{A}_i(\boldsymbol{\beta}) = \text{diag}\{h\{\mu_{i1}(\boldsymbol{\beta}_1); \phi_1\}, \dots, h\{\mu_{ip}(\boldsymbol{\beta}_p); \phi_p\}\}$ .

## Estimation Approach

In order to obtain a valid correlation matrix estimate  $\mathbf{R}(\boldsymbol{\rho})$ , we consider a slight reparameterisation of the correlation regression model as

$$\boldsymbol{\Sigma}(\boldsymbol{\alpha}) = \alpha_0 \mathbf{I}_p + \sum_{k=1}^K \alpha_k \mathbf{W}_k = \mathbf{R}(\boldsymbol{\rho})$$

with  $\alpha_0 = 1$  and  $\alpha_k = \rho_k$  for  $k = 1, \dots, K$ . The resulting covariance equation can then be rewritten as  $\text{cov}(\mathbf{Y}) = \mathbf{A}^{1/2}(\boldsymbol{\beta}) \{ \mathbf{I}_n \otimes \boldsymbol{\Sigma}(\boldsymbol{\alpha}) \} \mathbf{A}^{1/2}(\boldsymbol{\beta})$ , and subsequently the parameters we solve for are now  $\boldsymbol{\vartheta} = (\boldsymbol{\beta}^\top, \boldsymbol{\alpha}^\top)^\top$ .

For estimating  $\boldsymbol{\beta}$ , we solve

$$\psi_{\boldsymbol{\beta}}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \mathbf{D}^\top(\boldsymbol{\beta}) \mathbf{A}^{-1/2}(\boldsymbol{\beta}) \hat{\boldsymbol{\Sigma}}^{-1}(\boldsymbol{\alpha}) \mathbf{A}^{-1/2}(\boldsymbol{\beta}) \{ \mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}) \} = \mathbf{0}_{pd},$$

using Fisher scoring method.  
As for estimating  $\boldsymbol{\alpha}$ , we solve

$$\psi_{\boldsymbol{\alpha}}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \left( \boldsymbol{\epsilon}^\top(\boldsymbol{\beta}) \tilde{\mathbf{W}}_k \boldsymbol{\epsilon}(\boldsymbol{\beta}) \right) - (\text{tr}(\tilde{\mathbf{W}}_k \tilde{\mathbf{W}}_{k_2})) \boldsymbol{\alpha} = \mathbf{0}_{K+1}.$$

using a two-step ADMM approach [6], where

$$\begin{aligned} \circ \mathbf{D}(\boldsymbol{\beta}) &= \partial \mu(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}^\top & \circ \hat{\boldsymbol{\Sigma}}(\boldsymbol{\alpha}) &= \mathbf{I}_n \otimes \boldsymbol{\Sigma}(\boldsymbol{\alpha}) \\ \circ \tilde{\mathbf{W}}_k &= \mathbf{I}_n \otimes \mathbf{W}_k & \circ \boldsymbol{\epsilon}(\boldsymbol{\beta}) &= (\boldsymbol{\epsilon}_1^\top(\boldsymbol{\beta}), \dots, \boldsymbol{\epsilon}_p^\top(\boldsymbol{\beta}))^\top \\ \circ \boldsymbol{\epsilon}_i(\boldsymbol{\beta}) &= \mathbf{A}_i^{-1/2}(\boldsymbol{\beta}) \{ \mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\beta}) \}. \end{aligned}$$

By iteratively solving the two sets of estimating equations, we obtain estimates  $\hat{\boldsymbol{\vartheta}} = (\hat{\boldsymbol{\beta}}^\top, \hat{\boldsymbol{\alpha}}^\top)^\top$ , which can then be used to obtain the joint estimator of the mean regression coefficients and correlation regression parameters  $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}^\top, \hat{\boldsymbol{\rho}}^\top)^\top$  by setting  $\hat{\rho}_k = \hat{\alpha}_k / \hat{\alpha}_0$  for  $k = 1, \dots, K$ .

## Large Sample Theory

Assume  $d$  and  $K$  are fixed,  $n \rightarrow \infty$  and  $p$  can either be fixed or grow with increasing  $n$ . We introduce a matrix

$$\mathbf{E}^{(S)} = \begin{pmatrix} \mathbf{E}^{(S)} \otimes \mathbf{I}_d & \mathbf{0}_{qd \times (K+1)} \\ \mathbf{0}_{(K+1) \times pd} & \mathbf{I}_{K+1} \end{pmatrix} \triangleq \text{diag}\{\mathbf{E}^{(S)} \otimes \mathbf{I}_d, \mathbf{I}_{K+1}\},$$

which extracts a finite dimensional sub-vector of  $\hat{\boldsymbol{\vartheta}}$ , e.g.,  $\mathbf{E}^{(S)} \hat{\boldsymbol{\vartheta}} = (\hat{\boldsymbol{\beta}}_1^\top, \hat{\boldsymbol{\beta}}_2^\top, \hat{\boldsymbol{\alpha}}^\top)^\top$  when  $S = \{1, 2\}$ . Let  $\boldsymbol{\vartheta}^{(0)}$  and  $\boldsymbol{\theta}^{(0)}$  be the true parameter values, then we have the following asymptotic results:

**Theorem 1.** Under some regularity conditions, we have

$$\Omega^{-1/2}(\mathbf{E}^{(S)}) \begin{pmatrix} \sqrt{n} \mathbf{I}_{qd} & \mathbf{0} \\ \mathbf{0} & \sqrt{np} \mathbf{I}_{K+1} \end{pmatrix} \mathbf{E}^{(S)}(\hat{\boldsymbol{\vartheta}} - \boldsymbol{\vartheta}^{(0)}) \xrightarrow{d} N(0, \mathbf{I}_{qd+K+1}),$$

as  $n \rightarrow \infty$ , where  $p$  is either fixed or tends to infinity but satisfies  $p = o(n^{\min\{1/2, \eta/4\}})$ , and  $\eta > 0$  is related to a moment condition of the responses.

**Theorem 2.** Under some regularity conditions, we have

$$\bar{\Omega}^{-1/2}(\mathbf{E}^{(S)}) \begin{pmatrix} \sqrt{n} \mathbf{I}_{qd} & \mathbf{0} \\ \mathbf{0} & \sqrt{np} \mathbf{I}_K \end{pmatrix} \mathbf{E}^{(S)}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{(0)}) \xrightarrow{d} N(0, \mathbf{I}_{qd+K}),$$

as  $n \rightarrow \infty$ , where  $p$  is either fixed or tends to infinity but satisfies  $p = o(n^{\min\{1/2, \eta/4\}})$ , and  $\eta > 0$  is related to a moment condition of the responses.

## Simulation Studies

- Marginal distributions of response: Bernoulli, Poisson and Gaussian;
- $n \in \{50, 100, 200, 400\}$ ,  $p \in \{10, 25, 50\}$ ,  $d = 4$  and  $K = 5$ ;
- $\boldsymbol{\beta}_{jl}^{(0)} \stackrel{i.i.d.}{\sim} N(0, 0.5)$ ,  $\rho_k^{(0)} \stackrel{i.i.d.}{\sim} U(-0.05, 0.05)$ ,  $\phi_j = 1$ ;
- $\mathbf{x}_i$  and  $\mathbf{z}_j$  vectors are obtained from standardised environmental covariates and species traits, respectively, in the real data application;
- A total of 1000 replications for each setting.

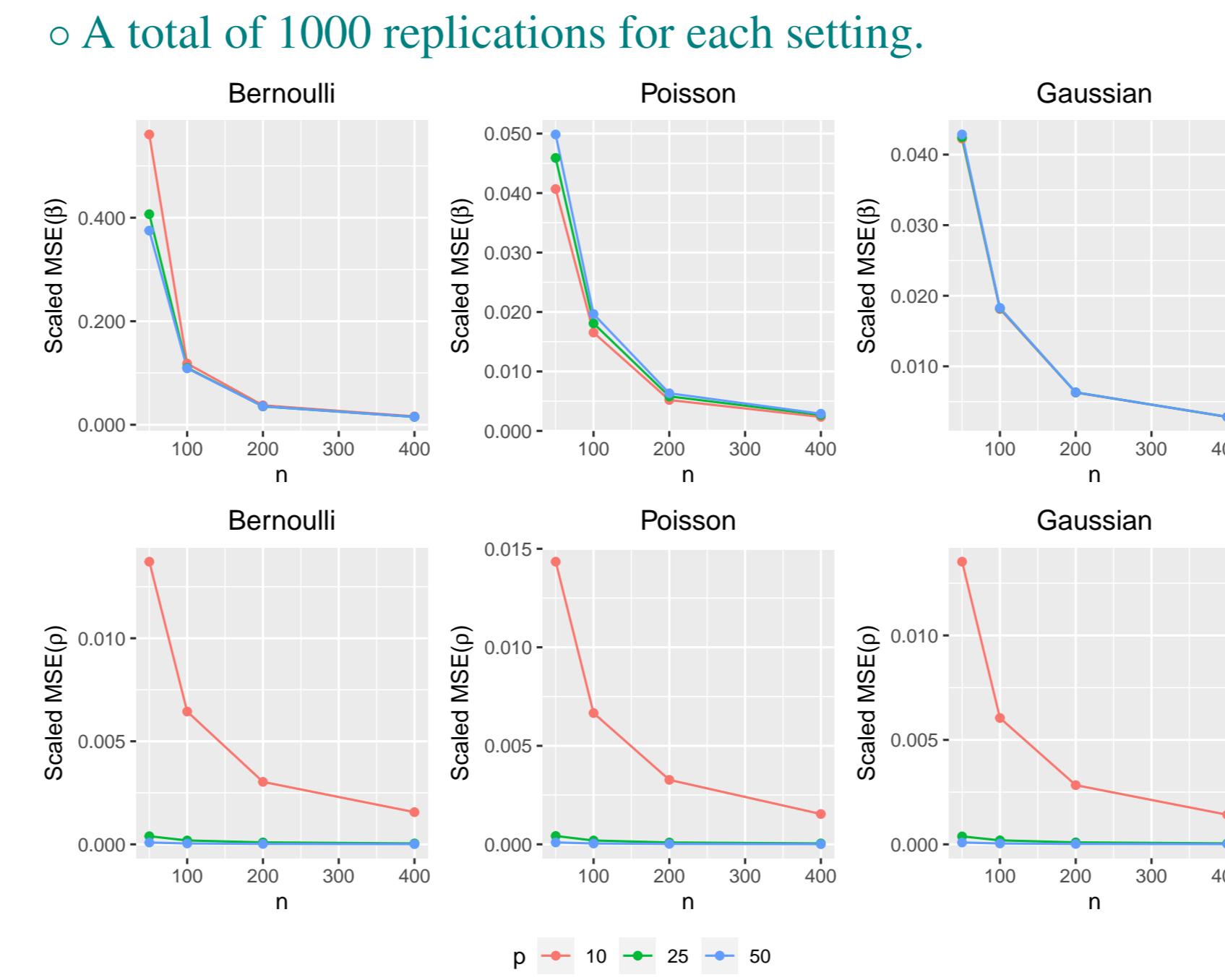


Figure 2: Scaled MSE for  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\rho}}$  by marginal distributions of responses.

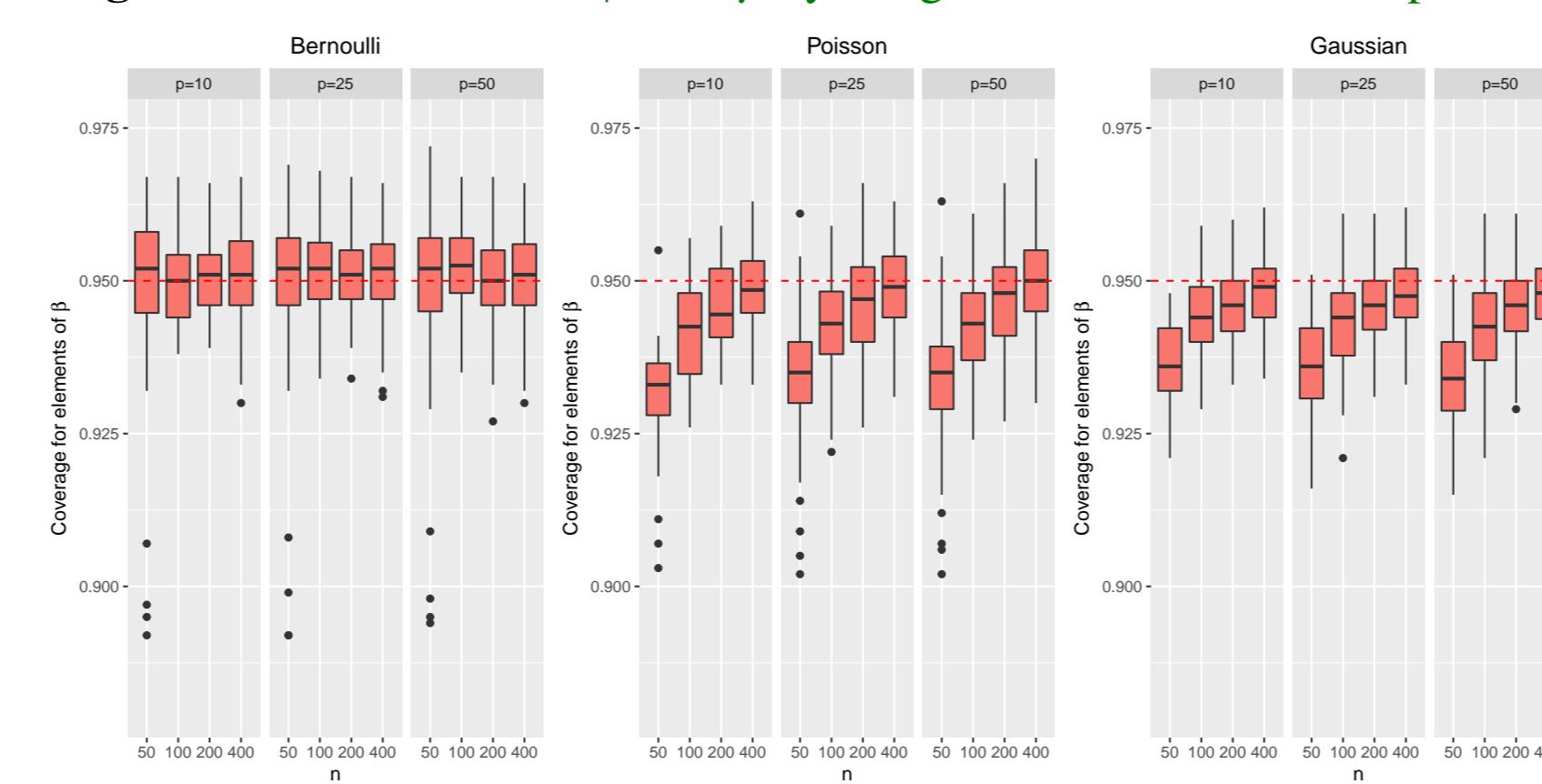


Figure 3: Boxplots summarising coverage of 95% approximate confidence intervals for each of the mean regression coefficients in the simulation studies where the total number of coefficients depends on the dimension of responses  $p$ .

## Real Data Application

- Data comes from [4], which consists of the counts of  $p = 38$  carabid ground beetle species:  $\mathbf{Y}_i$ , for  $i = 1, \dots, n$ , on  $n = 87$  sites spread across nine main areas in Scotland.
- Study the effects of environmental processes on different ground beetle species, and quantify how different traits influence residual correlation between species after accounting for species-environmental responses.

Species $j$	Estimation of $\beta_j$			
	Intercept	Soil pH	Elevation	Land management
<i>A.Muelleri</i>	2.444 (2.015, 2.872)	1.382 (0.631, 2.132)	-0.730 (-1.396, -0.064)	-0.022 (-0.725, 0.682)
<i>A.Apricaria</i>	-0.697 (-1.528, 0.135)	-0.151 (-1.186, 0.883)	-1.649 (-2.904, -0.393)	1.615 (0.692, 2.538)
<i>A.Bifrons</i>	-1.349 (-2.462, -0.236)	1.415 (0.134, 2.696)	-2.145 (-3.720, -0.570)	0.314 (-0.775, 1.404)
<i>A.Communis</i>	0.604 (0.145, 1.062)	-0.080 (-0.876, 0.715)	-0.439 (-1.058, 0.181)	-0.970 (-1.741, -0.198)
<i>A.Familiaris</i>	0.221 (-0.382, 0.823)	0.454 (-0.593, 1.500)	-0.371 (-1.287, 0.544)	0.474 (-0.509, 1.458)
<i>A.Lunicollis</i>	-0.016 (-0.533, 0.502)	-0.714 (-1.644, 0.217)	-0.841 (-1.600, -0.082)	-0.542 (-1.401, 0.317)
<i>A.Plebeja</i>	3.074 (2.667, 3.480)	1.576 (0.860, 2.291)	-0.933 (-1.567, -0.298)	-0.390 (-1.063, 0.283)
<i>A.Dorsalis</i>	0.106 (-0.513, 0.725)	0.305 (-0.544, 1.153)	-0.144 (-1.064, 0.776)	2.557 (1.720, 3.394)
<i>B.Aeneum</i>	2.010 (1.265, 2.756)	1.467 (0.262, 2.673)	-2.993 (-4.314, -1.672)	-0.236 (-1.345, 0.873)
<i>B.Guttula</i>	1.948 (1.560, 2.336)	0.921 (0.264, 1.577)	-1.024 (-1.672, -0.376)	0.448 (-0.159, 1.055)

Total length	Estimation of $\rho_k$				
	Traits $k$	Leg colour	Wing development	Overwintering	Breeding season
0.061 (0.023, 0.098)	-0.002 (-0.022, 0.018)	0.021 (-0.001, 0.044)	0.004 (-0.018, 0.026)	0.036 (0.008, 0.064)	

Table 1: Point estimates and 95% confidence intervals (in parentheses) for mean regression coefficients of the first 10 species, and correlation regression parameters for all five different trait variables.

- Carabid ground beetle species possess differing relationship with the environmental factors e.g., *A.Muelleri* prefers higher levels of soil pH, while *A.Communis* presents no clear evidence of being influenced by this habitat factor.
- The correlation regression parameter estimate for total length presented the strongest magnitude and significance, suggesting that conditional on other trait values, the abundances of carabid species with similar total lengths are more positively associated after accounting for differences in their mean abundance due to environmental processes.



Figure 4: <i