Appendices: Joint Species distribution modeling: dimension reduction using Dirichlet processes

Daniel Taylor-Rodríguez*, Kimberly Kaufeld[†] Erin M. Schliep[‡] James S. Clark[§] and Alan E. Gelfand[¶]

1 Appendix A: Sampling algorithm

The hierarchical setup described in the previous section yields the joint distribution

$$(\sigma_{\varepsilon}^{2})^{-(\frac{nS}{2}+1)} \left\{ \prod_{i=1}^{n} \exp\left[-\frac{1}{2\sigma_{\varepsilon}^{2}} ||\mathbf{V}_{i} - \mathbf{B}\mathbf{x}_{i} - Q(\mathbf{k})\mathbf{Z}\mathbf{w}_{i}||^{2}\right] \times \exp\left[-\frac{1}{2}\mathbf{w}_{i}'\mathbf{w}_{i}\right] \right\} \times \\ |\mathbf{D}_{\mathbf{z}}|^{-\frac{1}{2}} \left\{ \prod_{j=1}^{N} \exp\left[-\frac{1}{2}Z_{j}'\mathbf{D}_{\mathbf{z}}^{-1}Z_{j}\right] \right\} \times \left\{ \prod_{l=1}^{S} \sum_{j=1}^{N} p_{j}\delta_{j}(k_{l}) \right\} \pi(\boldsymbol{p}|\mathbf{0}, \alpha) \times \\ \left\{ \mathcal{I}W\left(\mathbf{D}_{\mathbf{z}} \middle| 2 + r - 1, 4\operatorname{diag}\left\{\frac{1}{\eta_{1}}, \dots, \frac{1}{\eta_{r}}\right\}\right) \times \prod_{h=1}^{r} \mathcal{I}\mathcal{G}\left(\eta_{h} \middle| \frac{1}{2}, \frac{1}{10^{4}}\right) \right\}.$$

The full conditional densities arising from the joint density above are

- 1) $[\mathbf{Z}|\mathbf{D}_{\mathbf{z}}, \mathbf{B}, \mathbf{W}, \sigma_{\varepsilon}^2, \mathbf{V}]$: The posterior for each row of \mathbf{Z} depends on whether or not the row considered was chosen to be at least one row from \mathbf{A} . That is, for $j = 1, 2, \ldots, N$
 - 1. If $j \notin \mathbf{k}$, sample $Z_i \sim N_r(\mathbf{0}, \mathbf{D_z})$.
 - 2. Otherwise, if $j \in \mathbf{k}$, let $S_j = \{l = 1, ..., S \text{ s.t. } k_l = j\}$ and let $|S_j|$ denote the cardinality of S_j . Using these definitions the full conditional distribution for Z_j is given by

$$Z_{j}|\mathbf{D_{z}},\mathbf{B},\mathbf{W},\sigma_{\varepsilon}^{2},\mathbf{V}\sim N_{r}\left(\boldsymbol{\mu}_{\mathbf{z}_{l}},\ \Sigma_{Z_{j}}\right)$$

^{*}Postdoctoral Associate, SAMSI/Duke University, Research Triangle Park, NC 27709. taylor-rodriguez@samsi.info

[†]Postdoctoral Researcher, SAMSI/North Carolina State University, Research Triangle Park, NC 27709. kkaufeld@samsi.info

[‡]Assistant Professor, Department of Statistics, University of Missouri, Columbia, MO 65211. schliepe@missouri.edu

[§]Professor, Nicholas School of the Environment, Department of Statistical Science, Duke University, Durham, NC 27708. jimclark@duke.edu

 $[\]P$ Professor, Department of Statistical Science, Duke University, Durham, NC 27708. alan@stat.duke.edu

where $\Sigma_{Z_j} = \left(\frac{|S_j|}{\sigma_{\varepsilon}^2} \mathbf{W}^T \mathbf{W} + \mathbf{D}_{\mathbf{z}}^{-1}\right)^{-1}$, $\boldsymbol{\mu}_{Z_j} = \Sigma_{Z_j} \mathbf{W}^T \frac{1}{\sigma_{\varepsilon}^2} \sum_{l \in S_j} (\mathbf{V}^{(l)} - \mathbf{X} \boldsymbol{\beta}_l)$, and finally, $\mathbf{y}^{(l)}$ and $\boldsymbol{\beta}_l$ are the lth column of matrix \mathbf{Y} and the lth row of \mathbf{B} , respectively.

2) $[\mathbf{W}|\mathbf{B}, \mathbf{A}, \sigma_{\varepsilon}^2, \mathbf{V}]$: As the rows of \mathbf{W} are independent, the posterior for each of the rows is

$$\begin{aligned} \mathbf{w}_i|\mathbf{B},\mathbf{A},\sigma_{\varepsilon}^2,\mathbf{V}_i \sim N_r \left(\Sigma_{\mathbf{W}} \mathbf{A} \frac{1}{\sigma_{\varepsilon}^2} (\mathbf{V}_i - \mathbf{B} \mathbf{x}_i) \right) \end{aligned}$$
 where $\Sigma_{\mathbf{W}} = \left(\frac{1}{\sigma_{\varepsilon}^2} \mathbf{A}^T \mathbf{A} + I_r \right)^{-1}$

3) $[\mathbf{k}|\mathbf{p}, \mathbf{B}, \mathbf{Z}, \sigma_{\varepsilon}^2, \mathbf{V}]$: For the vector of labels \mathbf{k} , the full conditional is

$$\begin{split} \left[\mathbf{k}|\boldsymbol{p},\mathbf{B},\mathbf{Z},\sigma_{\varepsilon}^{2},\mathbf{V}\right] &= \prod_{l=1}^{q} \left\{\sum_{j=1}^{N} p_{lj} \delta_{j}(k_{l})\right\} \end{split}$$
 with $p_{lj} \propto p_{j} \times \exp\left[-\frac{1}{2\sigma_{\varepsilon}^{2}}||\mathbf{V}^{(l)} - \mathbf{X}\beta_{l} - \mathbf{W}Z_{j}||^{2}\right].$

4) [p|k]: The full conditional posterior for p, given conjugacy of the \mathcal{GD} distribution with multinomial sampling, the draws of p are

$$p_1 = \xi_1,$$

$$p_j = (1 - \xi_1) \cdots (1 - \xi_{j-1}) \xi_j, \text{ for } j = 2, 3, \dots, N - 1,$$

$$p_N = 1 - \sum_{j=1}^{N-1} p_j,$$

with $\xi_j \stackrel{ind}{\sim} \text{Beta}\left(\frac{\alpha}{N} + \sum_{l=1}^{S} I_{(k_l=j)}, \frac{N-j}{N}\alpha + \sum_{s=j+1}^{N} \sum_{l=1}^{S} I_{(k_l=s)}\right)$ for $j = 1, \dots, N-1$

5) $[\sigma_{\varepsilon}^2|\mathbf{A}, \mathbf{W}, \mathbf{B}, \mathbf{Y}]$: By conjugacy of the prior for σ_{ε}^2 with the normal likelihood, the full conditional is

$$\sigma_{\varepsilon}^{2}|\mathbf{A},\mathbf{W},\mathbf{B},\mathbf{V}\sim\mathcal{IG}\left(\frac{nS+\nu}{2}+1,\frac{\sum_{i=1}^{n}||\mathbf{V}_{i}-\mathbf{B}\mathbf{x}_{i}-\mathbf{A}\mathbf{w}_{i}||^{2}}{2}+\frac{\nu}{G^{2}}\right).$$

6) $[\mathbf{D_z}|\mathbf{Z}]$: By conjugacy of the prior $\mathbf{D_z}$ with the normal prior for \mathbf{Z} , the full conditional for $\mathbf{D_z}$ is

$$\mathbf{D_z}|\mathbf{Z} \sim \mathcal{IW}\left(\mathbf{D_z}\Big|2 + r + N - 1, \mathbf{Z'Z} + 4\operatorname{diag}\left\{\frac{1}{\eta_1}, \dots, \frac{1}{\eta_r}\right\}\right).$$

2 Appendix B: Additional simulation plots from Section 6

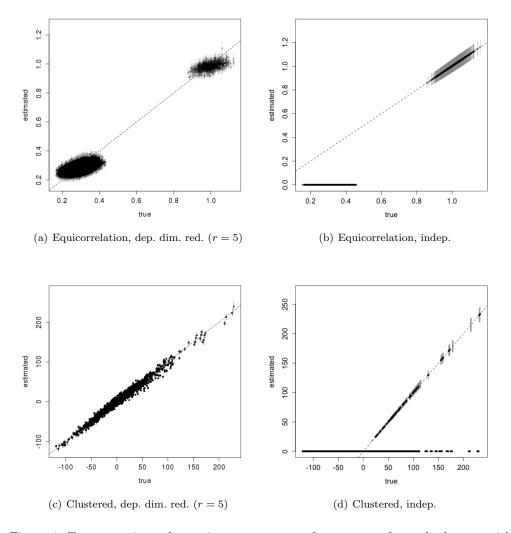
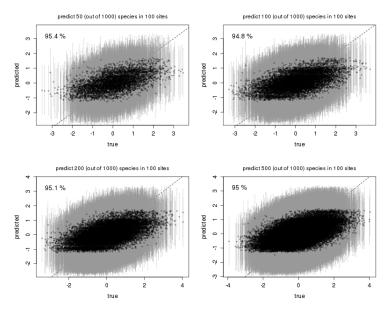
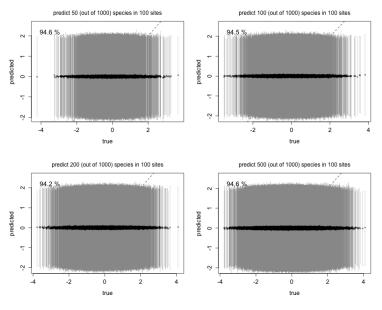


Figure 1: True vs estimated covariance parameters for one out-of-sample dataset with the equicorrelated and clustered covariance structures for n=3000, S=1000 comparing the dimension reduced dependent approach (with r=5) and the independence model



(a) Dependence with dimension red. (r=5), equicorrelated covariance structure



(b) Independence, equicorrelated covariance structure

Figure 2: Mean and 95% credible set for the conditional prediction for 50, 100, 200 and 500 species conditional on 950, 900, 800 and 500 species, respectively, in 100 plots with true data generated from equicorrelated covariance structure for n=3000, S=1000. Comparison between the dimension reduced dependent approach (with r=5) and the independence model. On each of the plots the empirical coverage for the prediction intervals is provided.

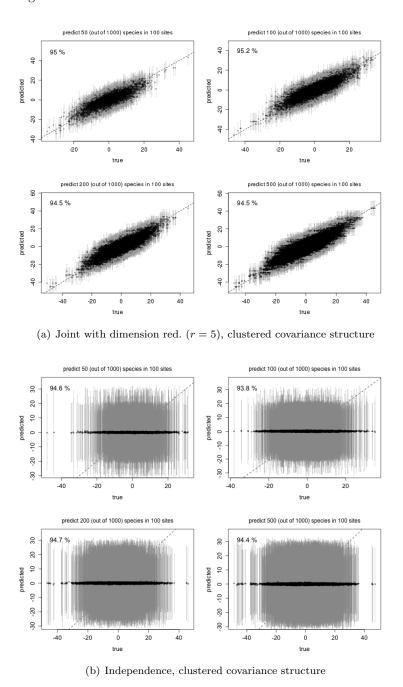
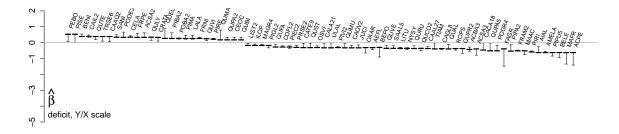
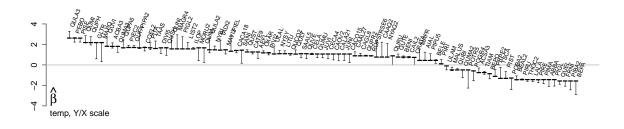


Figure 3: Mean and 95% credible set for the prediction of 50, 100, 200 and 500 species conditional on 950, 900, 800 and 500 species, respectively, in 100 plots with true data generated from clustered covariance structure for n=3000, S=1000. Comparison between the dimension reduced dependent approach (with r=5) and the independence model. On each of the plots the empirical coverage for the prediction intervals is provided.

3 Appendix C: Fitted regression coefficients for the FIA data



(a) deficit



(b) temperature

Figure 4: Mean and 95% credible sets for the "significant" fitted regression coefficients in the FIA data.