

Spatial Factor Models for Multivariate Spatial Data

Jeffrey Doser¹ & Andrew Finley²

May 15, 2023

¹Department of Integrative Biology, Michigan State University.

²Department of Forestry, Michigan State University.

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.
- Examples:
 - **Environmental monitoring:** stations yield measurements on ozone, NO, CO, and PM_{2.5}.

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.
- Examples:
 - **Environmental monitoring**: stations yield measurements on ozone, NO, CO, and PM_{2.5}.
 - **Community Ecology**: assemblages/communities of species

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.
- Examples:
 - **Environmental monitoring**: stations yield measurements on ozone, NO, CO, and PM_{2.5}.
 - **Community Ecology**: assemblages/communities of species
 - **Forestry**: measurements of stand characteristics age, total biomass, and average tree diameter.

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.
- Examples:
 - **Environmental monitoring**: stations yield measurements on ozone, NO, CO, and PM_{2.5}.
 - **Community Ecology**: assemblages/communities of species
 - **Forestry**: measurements of stand characteristics age, total biomass, and average tree diameter.
 - **Atmospheric modeling**: at a given site we observe surface temperature, precipitation and wind speed

Multivariate spatial data

- Point-referenced spatial data often come as multivariate measurements at each location.
- Examples:
 - **Environmental monitoring**: stations yield measurements on ozone, NO, CO, and PM_{2.5}.
 - **Community Ecology**: assemblages/communities of species
 - **Forestry**: measurements of stand characteristics age, total biomass, and average tree diameter.
 - **Atmospheric modeling**: at a given site we observe surface temperature, precipitation and wind speed
- We anticipate dependence between measurements
 - at a particular location
 - across locations

Multivariate spatial generalized linear model

- Spatial generalized linear model for h -variate spatial data for $j = 1, 2, \dots, h$ and $i = 1, \dots, n$:

$$y_j(\mathbf{s}_i) \sim f(\mu_j(\mathbf{s}_i), \tau_j)$$

$$\mu_j(\mathbf{s}_i) = g^{-1}(\eta_j(\mathbf{s}_i)) = \mathbf{x}(\mathbf{s}_i)^\top \beta_j + \mathbf{w}_j^*(\mathbf{s}_i)$$

- We can imagine modeling $\mathbf{w}^*(\mathbf{s}_i) = (w_1^*(\mathbf{s}_i), w_2^*(\mathbf{s}_i), \dots, w_h^*(\mathbf{s}_i))'$ as an h -variate Gaussian process

Multivariate spatial generalized linear model

- Spatial generalized linear model for h -variate spatial data for $j = 1, 2, \dots, h$ and $i = 1, \dots, n$:

$$y_j(\mathbf{s}_i) \sim f(\mu_j(\mathbf{s}_i), \tau_j)$$

$$\mu_j(\mathbf{s}_i) = g^{-1}(\eta_j(\mathbf{s}_i)) = \mathbf{x}(\mathbf{s}_i)^\top \beta_j + \mathbf{w}_j^*(\mathbf{s}_i)$$

- We can imagine modeling $\mathbf{w}^*(\mathbf{s}_i) = (w_1^*(\mathbf{s}_i), w_2^*(\mathbf{s}_i), \dots, w_h^*(\mathbf{s}_i))'$ as an h -variate Gaussian process
- Could model using Multivariate NNGP as discussed previously with SVCs, works well when $h < 5$.

Multivariate spatial generalized linear model

- Spatial generalized linear model for h -variate spatial data for $j = 1, 2, \dots, h$ and $i = 1, \dots, n$:

$$y_j(\mathbf{s}_i) \sim f(\mu_j(\mathbf{s}_i), \tau_j)$$

$$\mu_j(\mathbf{s}_i) = g^{-1}(\eta_j(\mathbf{s}_i)) = \mathbf{x}(\mathbf{s}_i)^\top \beta_j + \mathbf{w}_j^*(\mathbf{s}_i)$$

- We can imagine modeling $\mathbf{w}^*(\mathbf{s}_i) = (w_1^*(\mathbf{s}_i), w_2^*(\mathbf{s}_i), \dots, w_h^*(\mathbf{s}_i))'$ as an h -variate Gaussian process
- Could model using Multivariate NNGP as discussed previously with SVCs, works well when $h < 5$.
- But what about when h is large (e.g., 10, 100)?

Spatial Factor Model

- Approximates the dependence between multivariate (spatially-dependent) outcomes through a linear combination of a (much) lower-dimensional set of spatial factors

Spatial Factor Model

- Approximates the dependence between multivariate (spatially-dependent) outcomes through a linear combination of a (much) lower-dimensional set of spatial factors
- We represent the $h \times 1$ vector $\mathbf{w}^*(\mathbf{s}_i)$ as a linear combination of latent spatial factors and factor loadings:

$$\mathbf{w}^*(\mathbf{s}_i) = \mathbf{\Lambda} \mathbf{w}(\mathbf{s}_i)$$

- $\mathbf{\Lambda}$ is an $h \times q$ loadings matrix (tall and skinny) and $\mathbf{w}(\mathbf{s}_i)$ is a $q \times 1$ vector of realizations from q *independent* spatial GPs

Spatial Factor Model

- Approximates the dependence between multivariate (spatially-dependent) outcomes through a linear combination of a (much) lower-dimensional set of spatial factors
- We represent the $h \times 1$ vector $\mathbf{w}^*(\mathbf{s}_i)$ as a linear combination of latent spatial factors and factor loadings:

$$\mathbf{w}^*(\mathbf{s}_i) = \mathbf{\Lambda} \mathbf{w}(\mathbf{s}_i)$$

- $\mathbf{\Lambda}$ is an $h \times q$ loadings matrix (tall and skinny) and $\mathbf{w}(\mathbf{s}_i)$ is a $q \times 1$ vector of realizations from q *independent* spatial GPs
- In traditional factor analysis, $\mathbf{w}(\mathbf{s}_i)$ are realizations from independent standard normal random variables.

Spatial Factor Model

- Choosing $q \ll h$ leads to substantial computational reductions.
- Simple to code: just sample from q independent GPs as with basic univariate models.
- Yields a non-separable multivariate cross-covariance function between location \mathbf{s}_i and $\mathbf{s}_{i'}$:
$$\text{cov}(\mathbf{w}^*(\mathbf{s}_i), \mathbf{w}^*(\mathbf{s}_{i'})) = \sum_{k=1}^q \rho_k(\mathbf{s}_i, \mathbf{s}_{i'}, \phi_k) \boldsymbol{\lambda}_k \boldsymbol{\lambda}_k^\top$$
- Can simply replace the q full GPs with their corresponding NNGPs to yield a spatial factor NNGP model
- Identifiability constraints on $\boldsymbol{\Lambda}$: fix upper triangle to 0 and diagonal to 1. See Ren and Banerjee (2013) *Biometrics*

- Standard normal priors for the lower triangle of $\mathbf{\Lambda}$
- We like to model response-specific regression coefficients β_j hierarchically. For each $r = 1, \dots, p$ covariate, we model $\beta_{j,r}$ following

$$\beta_{j,r} \sim N(\mu_{\beta_r}, \tau_{\beta_r}^2)$$

- Gaussian hyperpriors for μ_{β_r} and IG or half-Cauchy priors for $\tau_{\beta_r}^2$
- Independent uniform priors for spatial decay parameters ϕ

- Full conditionals are in closed form for all parameters except ϕ for Gaussian and Binomial responses.
- Update ϕ with an Adaptive Metropolis-within-Gibbs algorithm (Roberts and Rosenthal 2009)
- See Taylor-Rodriguez et al. 2019 for Gaussian sampler, spOccupancy website for Pólya-Gamma sampler

Why we like spatial factor models

- Simple to code (don't need to deal with cross-covariance matrices).

Why we like spatial factor models

- Simple to code (don't need to deal with cross-covariance matrices).
- Relatively fast and efficient (well, at least for Gaussian and Binomial).

Why we like spatial factor models

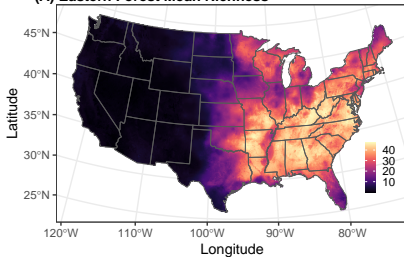
- Simple to code (don't need to deal with cross-covariance matrices).
- Relatively fast and efficient (well, at least for Gaussian and Binomial).
- Factors and factor loadings can be used for model-based ordination.

Why we like spatial factor models

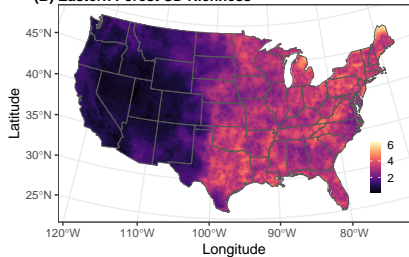
- Simple to code (don't need to deal with cross-covariance matrices).
- Relatively fast and efficient (well, at least for Gaussian and Binomial).
- Factors and factor loadings can be used for model-based ordination.
- Straightforward extensions to spatially-varying coefficient models.

Example: bird communities across the continental US

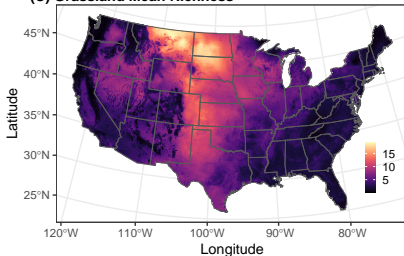
(A) Eastern Forest Mean Richness



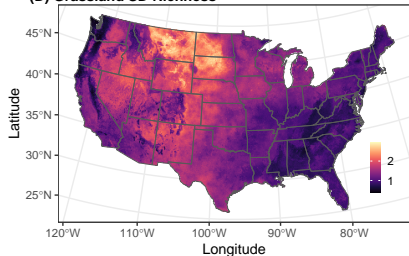
(B) Eastern Forest SD Richness



(C) Grassland Mean Richness

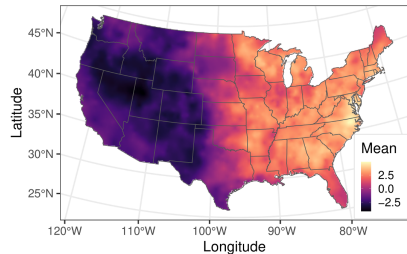
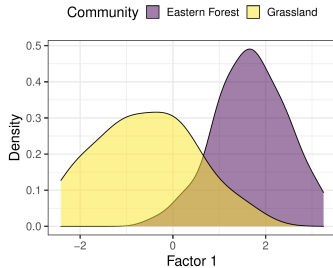


(D) Grassland SD Richness



Example: bird communities across the continental US

Visualization of the first spatial factor and corresponding factor loadings



Some downsides to spatial factor models

- Convergence assessment is not always straightforward
- Sensitivity to initial values
- Order of the first q responses has important implications for convergence and mixing.
- Assume a multivariate stochastic process can be represented as a linear combination of independent univariate processes

- `spOccupancy`: spatial NNGP and non-spatial factor models for binary data
- `spAbundance`: Gaussian, Poisson, and NB spatial NNGP and non-spatial factor models.
- `boral`: many distributions for non-spatial and spatial factor models (Hui 2015 *MEE*; spatial use full GPs fit in JAGS)
- `Hmsc`: spatial models using NNGPs (Tikhonov et al. 2019; *MEE*)
- `spBFA`: a variety of spatial models with some nifty priors (Berchuck et al. 2022 *Bayesian Analysis*)

Modeling the distribution of 10 tree species across Vermont