

A spatial model for rare binary events

July 27, 2016

1 Introduction

The goal in binary regression is to relate a latent variable to a response using a link function. Two common examples of binary regression include logistic regression and probit regression (Agresti, 2003). The link functions for logistic and probit regression are symmetric, so they may not be well-suited for asymmetric data. An asymmetric alternative to these link functions is the complementary log-log (cloglog) link function. More recently, Wang and Dey (2010) introduced the generalized extreme value (GEV) link function for rare binary data (a review is given in Appendix A.1). The GEV link function introduces a new shape parameter to the link function that controls the degree of asymmetry. The cloglog link is a special case of the GEV link function when the shape parameter is 0. Although this link was selected due to its ability to handle asymmetry, the GEV distribution is one of the primary distributions used for modeling extremes (Coles, 2001). Because extreme events are rare, it is therefore reasonable to use similar methods when analyzing rare binary data.

Spatial logistic and probit models commonly employ a hierarchical model assuming a latent spatial process (De Oliveira, 2000). In the hierarchical framework, spatial dependence is typically modeled with an underlying latent Gaussian process, and conditioned on this process, observations are independent. For large datasets, low-rank models can be used to ease the computational burden (Finley et al., 2015). If the latent variable is assumed to follow a GEV marginally, then a Gaussian process may not be appropriate to describe the dependence due to the fact that Gaussian processes do not demonstrate asymptotic dependence, except in the case of perfect dependence.

We propose using a latent max-stable process (de Haan, 1984) because it allows for asymptotic depen-

dence. The max-stable process arises as the limit of the location-wise maximum of infinitely many spatial processes. Max-stable processes are extremely flexible, but are often challenging to work with in high dimensions (Wadsworth and Tawn, 2014; Thibaud and Opitz, 2015). To address this challenge, methods have been proposed that implement composite likelihood techniques for max-stable processes (Padoan et al., 2010; Genton et al., 2011; Huser and Davison, 2014). As an alternative to these composite approaches, Reich and Shaby (2012) present a hierarchical model that implements a low-rank representation for a max-stable process. Although composite likelihoods have been used to model binary spatial data (Heagerty and Lele, 1998), we chose to use the low-rank representation of a max-stable process given by Reich and Shaby (2012).

Paragraph outlining the structure of the paper

2 Spatial dependence for binary regression

Let $Y(\mathbf{s})$ be the binary response at spatial location \mathbf{s} in a spatial domain of interest $\mathcal{D} \in \mathcal{R}^2$. We assume $Y(\mathbf{s}) = I[Z(\mathbf{s}) > 0]$ where $Z(\mathbf{s})$ is a latent continuous max-stable process. The marginal distribution of $Z(\mathbf{s})$ at site \mathbf{s} is GEV with location $\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}$, scale $\sigma > 0$, and shape ξ where $\mathbf{X}(\mathbf{s})$ is a p -vector of spatial covariates at site \mathbf{s} and $\boldsymbol{\beta}$ is a p -vector of regression coefficients. We set $\sigma = 1$ for identifiability because only the sign and not the scale of Z affects Y . If $\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta} = \mu$ for all \mathbf{s} , then $P(Y = 1)$ is the same for all observations, and the two parameters μ and ξ are not identifiable, so when there are no covariates, we fix $\xi = 0$. Although $\boldsymbol{\beta}$ and ξ could be permitted to vary across space, we assume that they are constant across \mathcal{D} . At spatial location \mathbf{s} , the marginal distribution is $P[Y(\mathbf{s}) = 1] = 1 - \exp\left[-\frac{1}{z(\mathbf{s})}\right]$ where $z(\mathbf{s}) = [1 - \xi \mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}]^{1/\xi}$. This is the same as the marginal distribution given by Wang and Dey (2010).

For a finite collection of locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, denote the vector of observations $\mathbf{Y} = [Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)]^T$.

The spatial dependence of \mathbf{Y} is determined by the joint distribution of $\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T$. To incor-

45 porate spatial dependence, we consider the hierarchical representation of the max-stable process proposed
 46 in Reich and Shaby (2012). Consider a set of positive stable random effect $A_1, \dots, A_L \stackrel{iid}{\sim} \text{PS}(\alpha)$ associated
 47 with spatial knots $\mathbf{v}_1, \dots, \mathbf{v}_L \in \mathcal{R}^2$. The hierarchical model is given by

$$\mathbf{Z}(\mathbf{s}_i) | A_1, \dots, A_L \stackrel{indep}{\sim} \text{GEV}[\mathbf{X}(\mathbf{s}_i)^\top \boldsymbol{\beta} + \theta(\mathbf{s}_i), \alpha \theta(\mathbf{s}_i), \xi \alpha] \quad \text{and} \quad \theta(\mathbf{s}_i) = \left[\sum_{l=1}^L A_l w_l(\mathbf{s}_i)^{1/\alpha} \right]^\alpha \quad (1)$$

48 where $w_l(\mathbf{s}_i) > 0$ are a set of L weights that vary smoothly across space and satisfy $\sum_{l=1}^L w_l(\mathbf{s}) = 1$ for all
 49 \mathbf{s} , and $\alpha \in (0, 1)$ determines the strength of dependence, with α near zero giving strong dependence and
 50 $\alpha = 1$ giving joint independence. Marginally over the A_l , this gives

$$Z(\mathbf{s}) \sim \text{GEV}(\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}, 1, \xi), \quad (2)$$

51 and thus $P[Y(\mathbf{s}) = 1] = 1 - \exp \left\{ -\frac{1}{z(\mathbf{s})} \right\}$ where $z(\mathbf{s}) = [1 - \xi \mathbf{X}(\mathbf{s}) \boldsymbol{\beta}]^{1/\xi}$.

52 Because the latent $\mathbf{Z}(\mathbf{s})$ are independent given the random effects, the binary responses are also condi-
 53 tionally independent. This leads to the tractible likelihood

$$Y(\mathbf{s}_i) | A_1, \dots, A_L \stackrel{indep}{\sim} \text{Bern}[\pi(\mathbf{s}_i)] \quad (3)$$

54 where

$$\pi(\mathbf{s}_i) = 1 - \exp \left\{ - \sum_{l=1}^L A_l \left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)} \right)^{1/\alpha} \right\}. \quad (4)$$

55 Many weight functions are possible, but the weights must be constrained so that $\sum_{l=1}^L w_l(\mathbf{s}_i) = 1$ for
 56 $i = 1, \dots, n$ to preserve the marginal GEV distribution. For example, Reich and Shaby (2012) take the

57 weights to be scaled Gaussian kernels with knots \mathbf{v}_l ,

$$w_l(\mathbf{s}_i) = \frac{\exp \left[-0.5 (\|\mathbf{s}_i - \mathbf{v}_l\|/\rho)^2 \right]}{\sum_{j=1}^L \exp \left[-0.5 (\|\mathbf{s}_i - \mathbf{v}_j\|/\rho)^2 \right]} \quad (5)$$

58 where $\|\mathbf{s}_i - \mathbf{v}_l\|$ is the distance between site \mathbf{s}_i and knot \mathbf{v}_l , and the kernel bandwidth $\rho > 0$ determines the
 59 spatial range of the dependence, with large ρ giving long-range dependence and vice versa.

60 After marginalizing out the positive stable random effects, the joint distribution of \mathbf{Z} is

$$G(\mathbf{z}) = P[Z(\mathbf{s}_1) < z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n) < z(\mathbf{s}_n)] = \exp \left\{ - \sum_{l=1}^L \left[\sum_{i=1}^n \left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)} \right)^{1/\alpha} \right]^\alpha \right\}, \quad (6)$$

61 where $G(\cdot)$ is the CDF of a multivariate GEV distribution. This is a special case of the multivariate GEV
 62 distribution with asymmetric Laplace dependence function (Tawn, 1990).

63 **3 Joint distribution**

64 We give an exact expression in the case where there are only two spatial locations which is useful for
 65 constructing a pairwise composite likelihood (Padoan et al., 2010) and studying spatial dependence. When
 66 $n = 2$, the probability mass function is given by

$$P[Y(\mathbf{s}_i) = Y_i, Y(\mathbf{s}_j) = Y_j] = \begin{cases} \varphi(\mathbf{z}) & Y_i = 0, Y_j = 0 \\ \exp \left\{ -\frac{1}{z(\mathbf{s}_i)} \right\} - \varphi(\mathbf{z}), & Y_i = 1, Y_j = 0 \\ 1 - \exp \left\{ -\frac{1}{z(\mathbf{s}_i)} \right\} - \exp \left\{ -\frac{1}{z(\mathbf{s}_j)} \right\} + \varphi(\mathbf{z}), & Y_i = 1, Y_j = 1 \end{cases} \quad (7)$$

67 where $\varphi(\mathbf{z}) = \exp \left\{ - \sum_{l=1}^L \left[\left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)} \right)^{1/\alpha} + \left(\frac{w_l(\mathbf{s}_j)}{z(\mathbf{s}_j)} \right)^{1/\alpha} \right]^\alpha \right\}$. For more than two locations, we are also
68 able to compute the exact likelihood when the n is large but the number of events $K = \sum_{i=1}^n Y(\mathbf{s}_i)$ is small,
69 as might be expected for very rare events (see Appendix A.2).

70 **4 Quantifying spatial dependence**

71 Assume that Z_1 and Z_2 are both $\text{GEV}(\beta, 1, 1)$ so that $P(Y_i = 1)$ decreases to zero as β increases. A common
72 measure of dependence between binary variables is Cohen's Kappa (Cohen, 1960),

$$\kappa(\beta) = \frac{P_A - P_E}{1 - P_E} \quad (8)$$

73 where P_A is the joint probability of agreement $P(Y_1 = Y_2)$ and P_E is the joint probability of agreement
74 under an assumption of independence $P(Y_i = 1)^2 + P(Y_i = 0)^2$. For the spatial model,

$$\begin{aligned} P_A(\beta) &= 1 - 2 \exp \left\{ -\frac{1}{\beta} \right\} + 2 \exp \left\{ -\frac{\vartheta(\mathbf{s}_1, \mathbf{s}_2)}{\beta} \right\} \\ P_E(\beta) &= 1 - 2 \exp \left\{ -\frac{1}{\beta} \right\} + 2 \exp \left\{ -\frac{2}{\beta} \right\}, \end{aligned}$$

75 and

$$\kappa(\beta) = \frac{P_A(\beta) - P_E(\beta)}{1 - P_E(\beta)} = \frac{\exp \left\{ -\frac{\vartheta(\mathbf{s}_1, \mathbf{s}_2) - 1}{\beta} \right\} - \exp \left\{ -\frac{1}{\beta} \right\}}{1 - \exp \left\{ -\frac{1}{\beta} \right\}} \quad (9)$$

76 where $\vartheta(\mathbf{s}_i, \mathbf{s}_j) = \sum_{l=1}^L [w_l(\mathbf{s}_i)^{1/\alpha} + w_l(\mathbf{s}_j)^{1/\alpha}]^\alpha$ is the pairwise extremal coefficient given by Reich and
 77 Shaby (2012). To measure extremal dependence, let $\beta \rightarrow \infty$ so that events are increasingly rare. Then,

$$\kappa = \lim_{\beta \rightarrow \infty} \kappa(\beta) = 2 - \vartheta(\mathbf{s}_1, \mathbf{s}_2) \quad (10)$$

78 which is the same as the χ statistic of Coles (2001), a commonly used measure of extremal dependence.

79 **5 Computation**

80 For small K , we can evaluate the likelihood directly. When K is large, we use Markov chain Monte
 81 Carlo (MCMC) methods with the random effects model to explore the posterior distribution. To overcome
 82 challenges with evaluating the positive stable density, we follow Reich and Shaby (2012) and introduce a
 83 set of auxiliary variables B_1, \dots, B_L following the auxiliary variable technique of Stephenson (2009). So,
 84 the hierarchical model is given by

$$\begin{aligned} Y(\mathbf{s}_i) | \pi(\mathbf{s}_i) &\stackrel{indep}{\sim} \text{Bern}[\pi(\mathbf{s}_i)] \\ \pi(\mathbf{s}_i) &= 1 - \exp \left\{ - \sum_{l=1}^L A_l \left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)} \right)^{1/\alpha} \right\} \\ Z(\mathbf{s}_i) | A_l, \dots, A_L &\stackrel{indep}{\sim} \text{GEV}[\mathbf{X}(\mathbf{s}_i)^\top \boldsymbol{\beta} + \theta(\mathbf{s}_i), \alpha \theta(\mathbf{s}_i), \xi \alpha] \\ A_l &\stackrel{iid}{\sim} \text{PS}(\alpha) \\ B_l &\stackrel{iid}{\sim} \text{Unif}(0, 1) \end{aligned} \quad (11)$$

85 with priors $\boldsymbol{\beta} \sim \text{N}(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p)$, $\xi \sim \text{N}(0, \sigma_\xi^2)$, $\rho \sim \text{Unif}(\rho_l, \rho_u)$, and $\alpha \sim \text{Beta}(a_\alpha, b_\alpha)$. The model parameters
 86 are updated using Metropolis Hastings (MH) update steps, and the random effects A_1, \dots, A_L , and auxiliary
 87 variables B_1, \dots, B_L are updated using Hamiltonian Monte Carlo (HMC) update steps.

6 Simulation study

For our simulation study, we generate $n_m = 50$ datasets under 3 different settings to explore the impact of sample size, sampling technique, and misspecification of link function. We generate data assuming three possible types of underlying process. For each of the underlying processes, we generate data on a 100×100 rectangular grid of $n = 10,000$ locations. If a dataset is generated with $K < 100$ or $K > 700$, it is discarded and a new dataset is generated.

6.1 Latent processes

The first process is a latent max-stable process that uses the GEV link described in (1) with knots on a 50×50 regularly spaced grid on $[0, 1] \times [0, 1]$. For this process, we set $\alpha = 0.35$, $\rho = 0.1$, and $\beta_0 \approx 2.97$ which gives $K = 500$, on average. Because there are no covariates, we set $\xi = 0$. We then set $Y(\mathbf{s}) = I[z(\mathbf{s}) > 0]$ where $I[\cdot]$ is an indicator function.

For the second process, we generate a latent variable from a spatial Gaussian process with a mean of $\text{logit}(0.05) \approx -2.94$ and an exponential covariance given by

$$\text{cov}(\mathbf{s}_1, \mathbf{s}_2) = \tau_{\text{Gau}}^2 \exp \left\{ -\frac{\|\mathbf{s}_1 - \mathbf{s}_2\|}{\rho_{\text{Gau}}} \right\} \quad (12)$$

where $\tau_{\text{Gau}} = 10$ and $\rho_{\text{Gau}} = 0.1$. Finally, we generate $Y(\mathbf{s}_i) \stackrel{\text{ind}}{\sim} \text{Bern}[\pi(\mathbf{s}_i)]$ where $\pi(\mathbf{s}_i) = \exp \left\{ \frac{z(\mathbf{s})}{1 + z(\mathbf{s})} \right\}$

For the third process, we generate data using a hotspot method. For this process, we first generate hotspots throughout the space. Let n_{hs} be the number of hotspots in the space. Then $n_{\text{hs}} - 1 \sim \text{Poisson}(2)$. This generation scheme ensures that every dataset has at least one hotspot. We generate the hotspot locations $\mathbf{h}_1, \dots, \mathbf{h}_{n_{\text{hs}}}$. Let B_h be a circle of radius of radius r_h around hotspot $h = 1, \dots, n_{\text{hs}}$. The r_h differ for each hotspot and are generated i.i.d. from a $\text{Unif}(0.03, 0.08)$ distribution. We set $P[Y(\mathbf{s}_i) = 1] = 0.85$ for all

107 \mathbf{s}_i in B_h , and $P[Y(\mathbf{s}_i)] = 0.0005$ for all \mathbf{s}_i outside of B_h . These settings are selected to give an average of
 108 close $K = 500$ for the datasets.

109 6.2 Methods

110 For each dataset, we fit the model using three different models, the proposed spatial GEV model, a spatial
 111 probit model, and a spatial logistic model. Because logistic and probit methods represent two of the more
 112 common spatial techniques for binary data, we chose to compare our method to them. One way these meth-
 113 ods differ from our proposed method is that they assume the underlying process is Gaussian. In this case,
 114 we assume that $Z(\mathbf{s})$ follows a Gaussian process with mean $\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}$ and exponential covariance function.
 115 The marginal distributions are given by

$$P[Y(\mathbf{s}) = 1] = \begin{cases} \frac{\exp[\mathbf{X}^\top(\mathbf{s})\boldsymbol{\beta} + \mathbf{W}(\mathbf{s})\boldsymbol{\epsilon}]}{1 + \exp[\mathbf{X}^\top(\mathbf{s})\boldsymbol{\beta} + \mathbf{W}(\mathbf{s})\boldsymbol{\epsilon}]}, & \text{logistic} \\ \Phi[\mathbf{X}^\top\boldsymbol{\beta}(\mathbf{s}) + \mathbf{W}(\mathbf{s})\boldsymbol{\epsilon}], & \text{probit} \end{cases} \quad (13)$$

116 where $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \tau_L^2 \mathbf{I}_L)$ are Gaussian random effects at the knot locations, and $\mathbf{W}(\mathbf{s})$ are a set of L basis
 117 functions given to recreate the Gaussian process at all sites. We use our own code for the spatial probit
 118 model, but we use the `spGLM` function in the `spBayes` package (Finley et al., 2015) to fit the spatial
 119 logistic model. For the probit model, we use

$$\mathbf{W}_l(\mathbf{s}_i) = \frac{\exp[-(\|\mathbf{s}_i - \mathbf{v}_l\|/\rho)^2]}{\sqrt{\sum_{j=1}^L \exp[-(\|\mathbf{s}_i - \mathbf{v}_j\|/\rho)^2]^2}}. \quad (14)$$

120 For the logistic model, the $\mathbf{W}_l(\mathbf{s}_i)$ are the default implementation from the `spGLM`.

6.3 Sampling technique

We subsample the generated data using $n_s = 100$ and $n_s = 250$ initial locations for two different sampling designs. The first is a two-stage spatially-adaptive cluster technique (CLU) taken from Pacifici et al. (2016). In this design, if an initial location is occupied, we also include the four rook neighbor (north, east, south, and west) sites in the sample. For the second design, we use a simple random sample (SRS) with the same number of sites included in the cluster sample.

For all models, we place knots on a 15×15 regularly spaced grid over the domain, and we also place knots at all sites where $Y(\mathbf{s}) = 1$. This knot arrangement is selected for two reasons. First, the regular grid provides a natural cutoff for the lower bound of the Uniform prior on ρ . This lower bound is important because if ρ is too small, relative to the knot placement, it is possible to end up with predictions at locations that are independent from all data points. There is a trade-off in selecting the number of knots to use for the random effects. If the knot spacing is too far apart, we risk a negative bias when estimating $P(Y = 1)$. One way to address this challenge is to provide a finer grid, but this can quickly become computationally burdensome. Furthermore, with a finer grid, we often end up increasing the resolution of the grid in areas that may not be important. Therefore, we choose to place additional knots at sites where $Y = 1$ as a balance between grid size and detail in the important areas.

6.4 Priors

For all models, we only include an intercept term β_0 in the model, and the prior for the intercept is $\beta_0 \sim N(0, 10)$. Additionally, for all models, the prior for the bandwidth is $\rho \sim \text{Unif}(\frac{1}{30}, 1)$. This lower bound is selected because it is half of the distance between the rook neighbors of the knots. For the GEV method, the prior for the spatial dependence parameter is $\alpha \sim \text{Beta}(2, 5)$. We select this prior because it gives greater weight to $\alpha < 0.5$, which is the point at which spatial dependence becomes fairly weak, but

also avoids values below 0.1 which can lead to numerical problems. We fix $\xi = 0$ because we do not include any covariates. For both the spatial probit and logistic models, the prior on the variance term for the random effects is $IG(0.1, 0.1)$ where $IG(\cdot)$ is an Inverse Gamma distribution. Both the spatial probit and logit models assume an exponential covariance structure. For all models, we run the MCMC sampler for 25,000 iterations with a burn-in period of 20,000 iterations. Convergence is assessed through visual inspection of traceplots.

6.5 Model comparisons

For each dataset, we fit the model using the n_s observations as a training set, and validate the model's predictive power at the remaining grid points. Let \mathbf{s}_j^* be the j th site in the validation set. To obtain the posterior predictive distribution, at each iteration of the MCMC, we generate a spatial field of zeros and ones at the validation locations. Then to obtain $\hat{P}[Y(\mathbf{s}_j^*) = 1]$, we take the average of the posterior distribution for each j . We consider a few different metrics for comparing model performance. One score is the Brier scores (Gneiting and Raftery, 2007, BS). The Brier score for predicting an occurrence at site \mathbf{s} is given by $\{I[Y(\mathbf{s}) = 1] - \hat{P}[Y(\mathbf{s}) = 1]\}^2$ where $I[Y(\mathbf{s}) = 1]$ is an indicator function indicating that an event occurred at site \mathbf{s} . We average the Brier scores over all test sites, and a lower score indicates a better fit. We also consider the receiver operating characteristic (ROC) curve, and the area under the ROC curve (AUROC) for the different methods and settings. The ROC curve and AUROC are obtained via the `ROCR` (Sing et al., 2005) package in R (R Core Team, 2016). We then average AUCs across all datasets for each method and setting to obtain a single AUC for each combination of method and setting.

6.6 Results

Needs updating

Table 1: Brier scores (SE) and AUROC (SE) for GEV, Probit, and Logistic methods from the simulation study.

Setting	n	Sample Type	BS			AUROC		
			GEV	Probit	Logistic	GEV	Probit	Logistic
GEV	100	CLU						
		SRS						
	250	CLU						
		SRS						
Probit	100	CLU						
		SRS						
	250	CLU						
		SRS						
Hotspot	100	CLU						
		SRS						
	250	CLU						
		SRS						

Table 2 gives the Brier scores and AUC for each of the methods. In Figure 5 – Figure 3, for each setting we present the vertically averaged ROC curve for each method.

7 Data analysis

Needs updating

We compare our method to the spatial probit and logit for mapping the probability of the occurrence of *Tamarix ramosissima*, a plant species, for a 1-km² study region of PR China Smith et al. (2012). The Chinese Academy of Forestry conducted a full census of the area, and the true occupancy of the species are plotted in Figure 4. The region is split into 10-m \times 10-m grid cells, and *Tamarix ramosissima* can be found in approximately 6% of the grid cells.

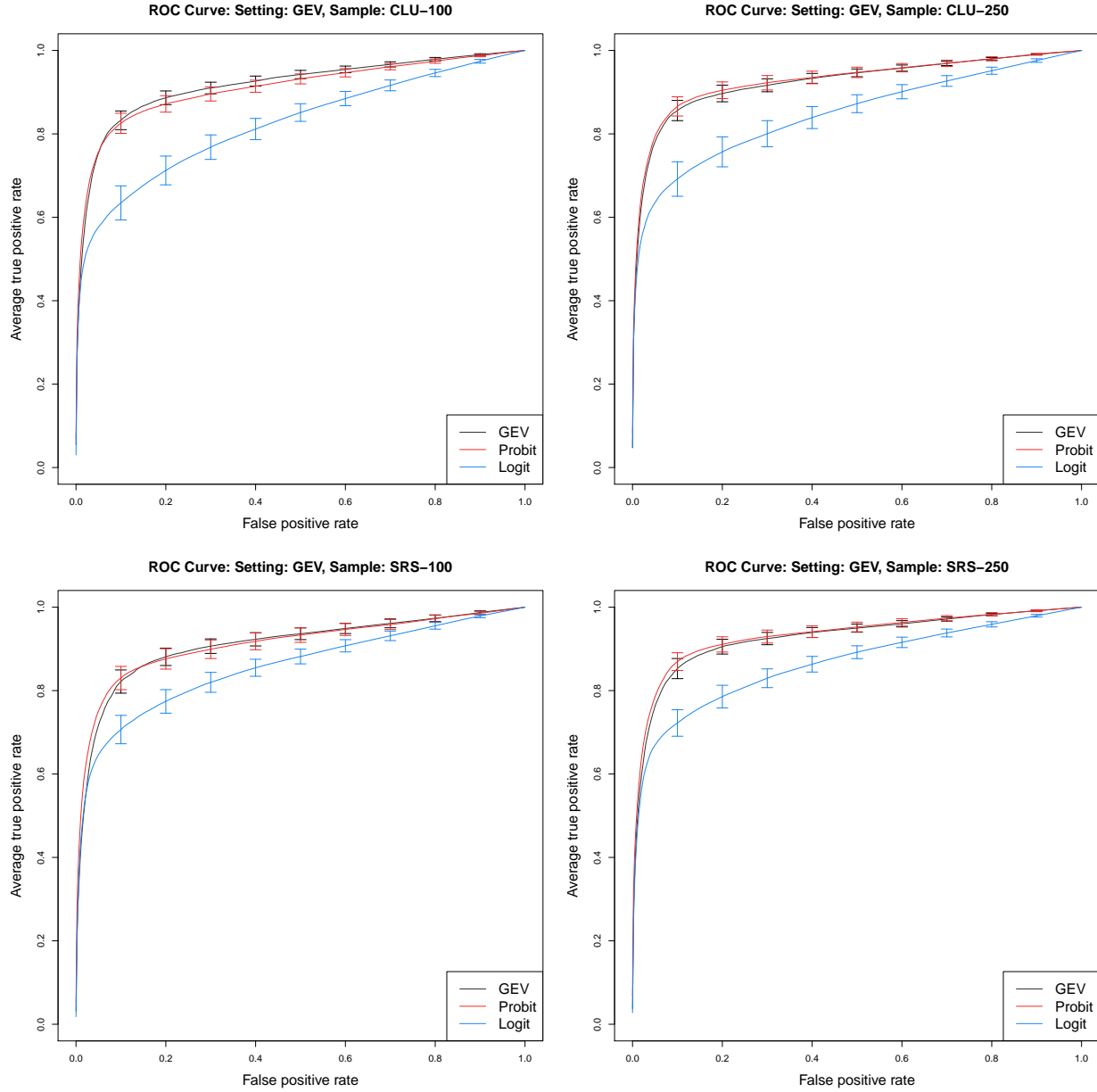


Figure 1: Vertically averaged ROC curves for GEV setting.

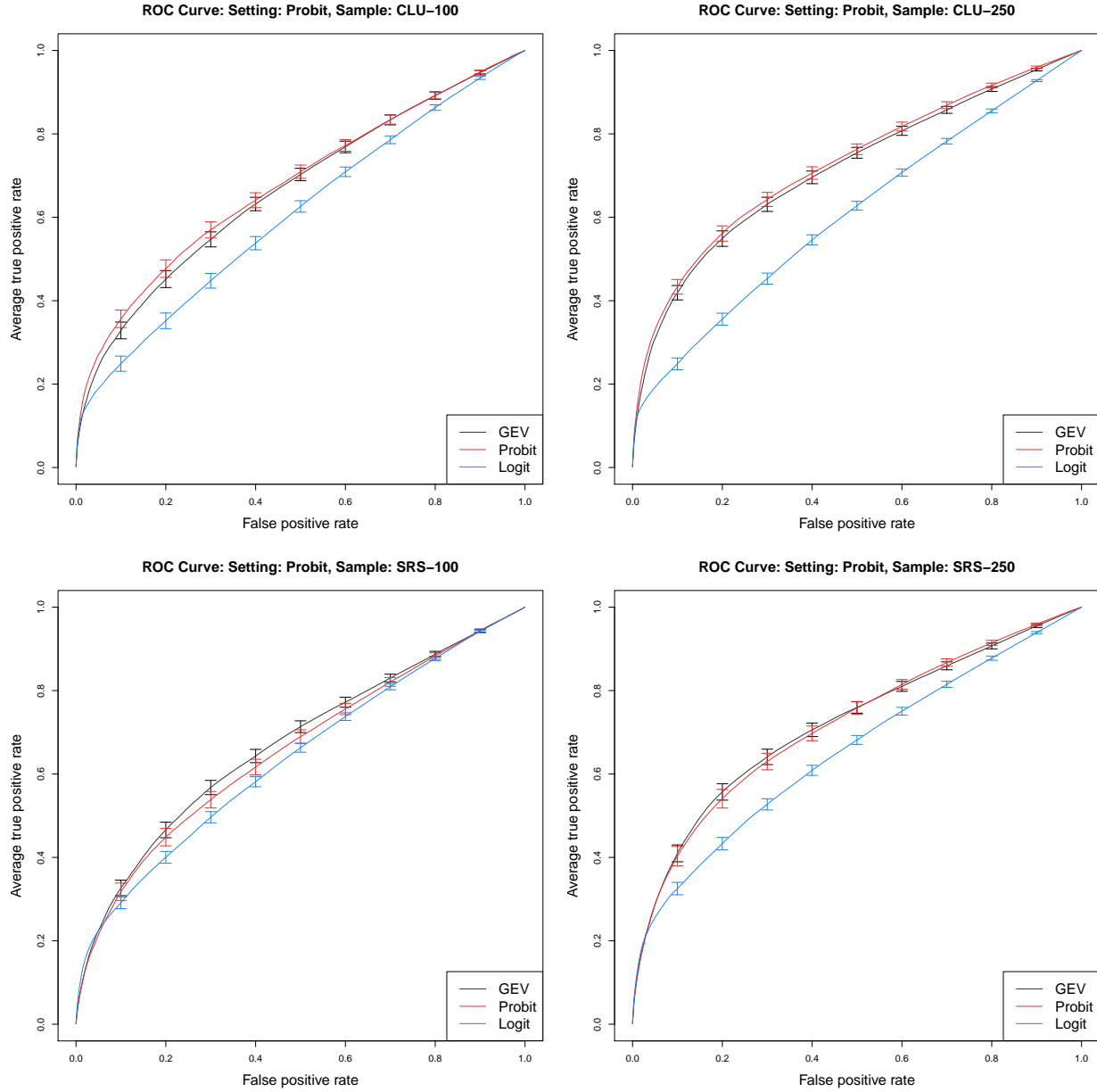


Figure 2: Vertically averaged ROC curves for Probit setting.

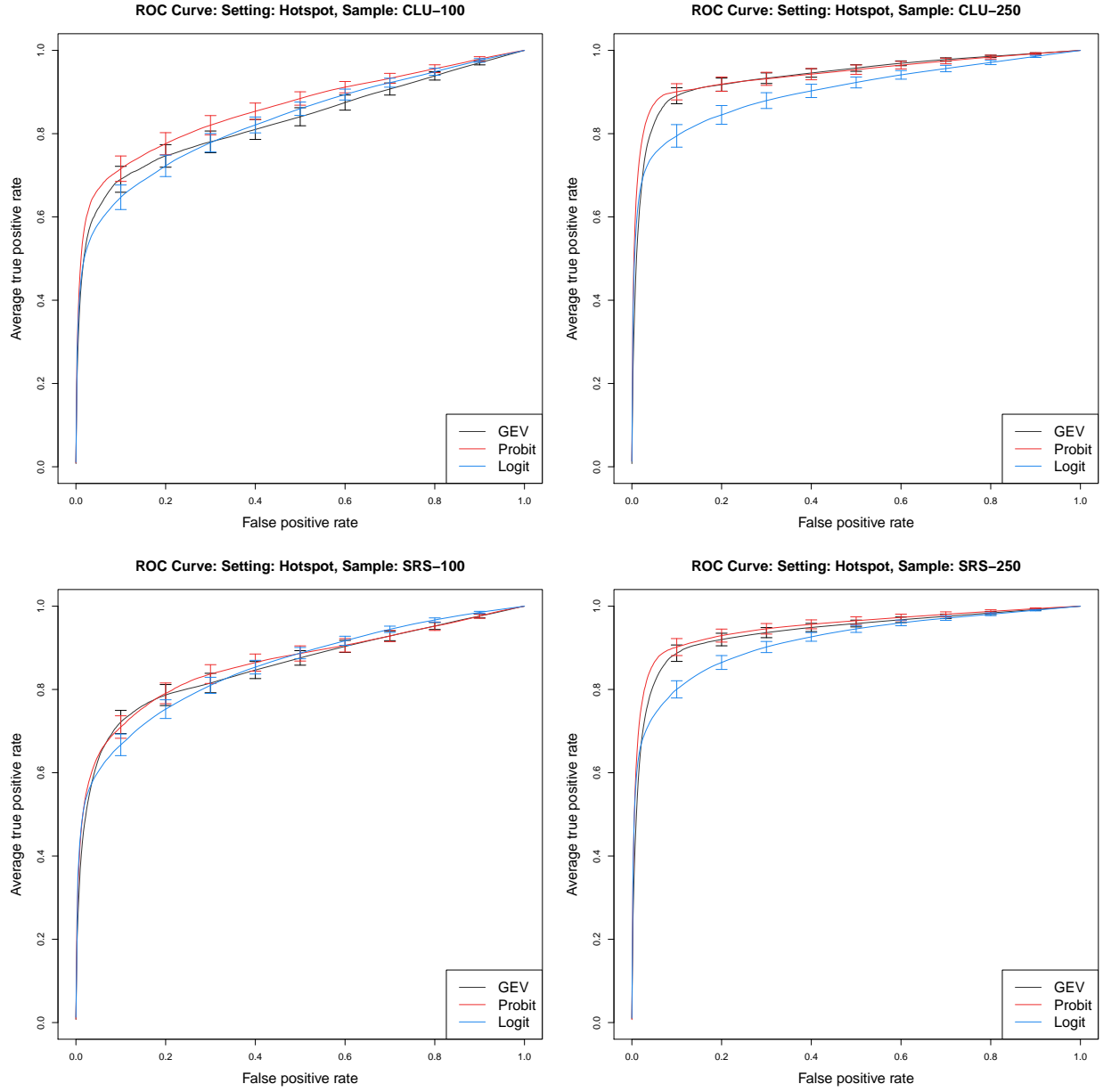


Figure 3: Vertically averaged ROC curves for Hotpost setting.

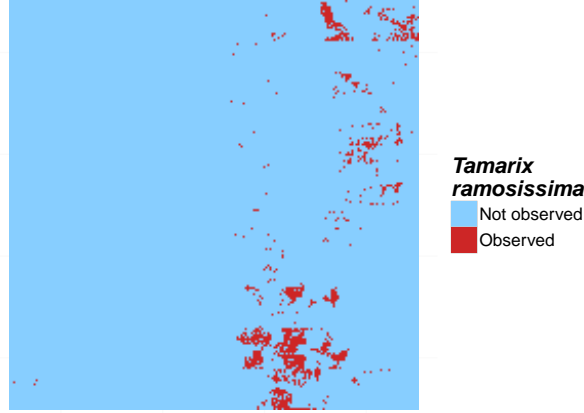


Figure 4: True occupancy of *Tamarix ramosissima* from a 1-km² study region of PR China.

Table 2: Brier scores (SE) and AUROC (SE) for GEV, Probit, and Logistic methods for *Tamarix ramosissima*.

n	Sample Type	BS			AUROC		
		GEV	Probit	Logistic	GEV	Probit	Logistic
100	CLU						
	SRS						
250	CLU						
	SRS						

7.1 Methods

For the data analysis, we generate 100 subsamples using the CLU and SRS sampling methods with $n_s = 100$ and $n_s = 250$ initial locations. For each subsample, we fit the spatial GEV, spatial probit, and spatial logistic models. Knot placement, prior distributions, and MCMC details for the data analysis are the same as the simulation study. To compare models, we use similar metrics as in the simulation study, but we average the metrics over subsamples.

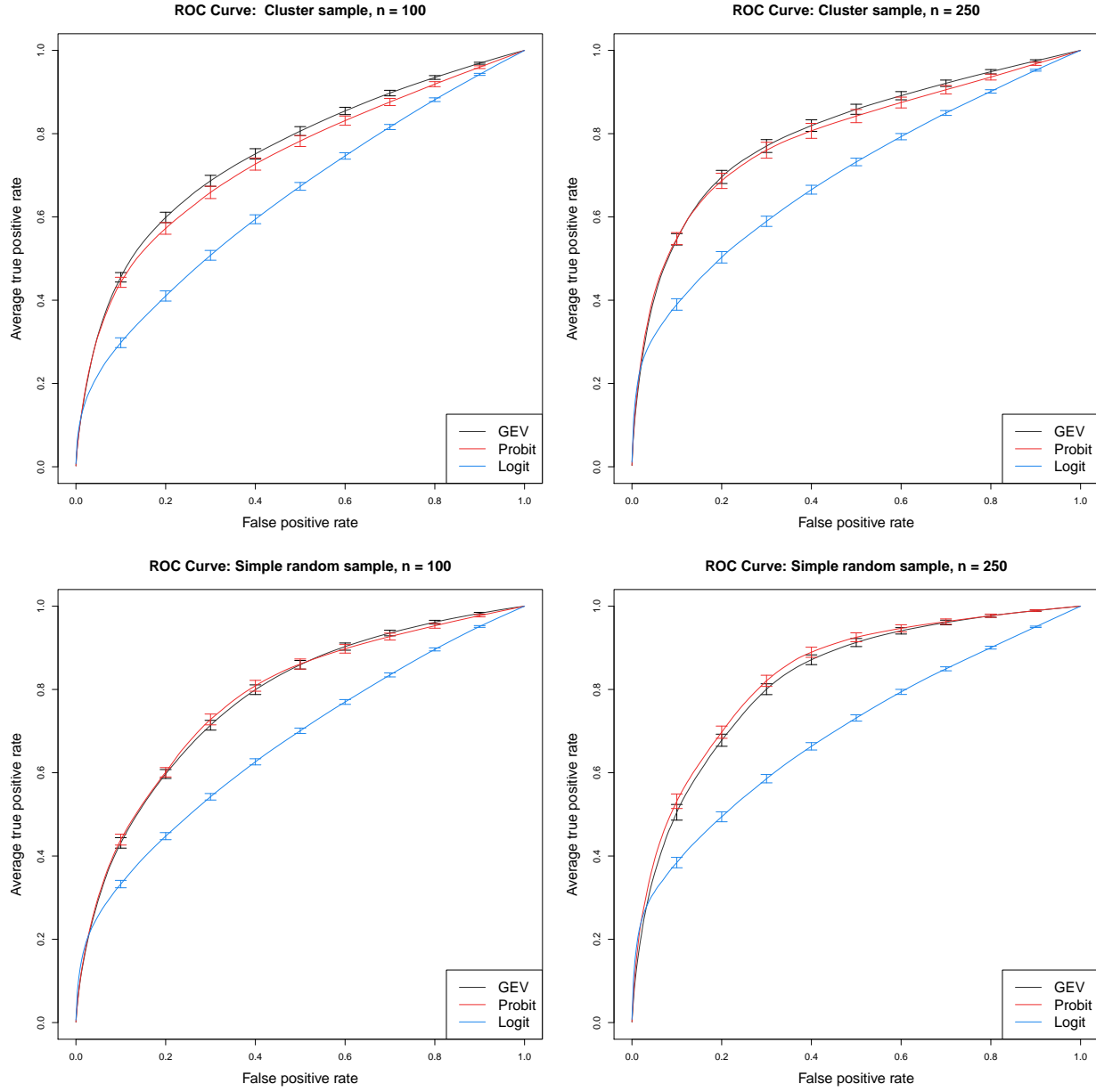


Figure 5: Vertically averaged ROC curves for *Tamarix ramosissima*.

7.2 Results

8 Conclusions

Acknowledgments

A Appendices

A.1 Binary regression using the GEV link

Here, we provide a brief review of the the GEV link of Wang and Dey (2010). Let $Y_i \in \{0, 1\}, i = 1, \dots, n$ be a collection of i.i.d. binary responses. It is assumed that $Y_i = I(z_i > 0)$ where $I(\cdot)$ is an indicator function, $z_i = [1 - \xi \mathbf{X}_i \boldsymbol{\beta}]^{1/\xi}$ is a latent variable following a $\text{GEV}(1, 1, 1)$ distribution, \mathbf{X}_i is the associated p -vector of covariates with first element equal to one for the intercept, and $\boldsymbol{\beta}$ is a p -vector of regression coefficients. Then, $Y_i \stackrel{\text{ind}}{\sim} \text{Bern}(\pi_i)$ where $\pi_i = 1 - \exp\left(-\frac{1}{z_i}\right)$.

A.2 Derivation of the likelihood

We use the hierarchical max-stable spatial model given by Reich and Shaby (2012). If at each margin, $Z_i \sim \text{GEV}(1, 1, 1)$, then $Z_i | \theta_i \stackrel{\text{indep}}{\sim} \text{GEV}(\theta, \alpha\theta, \alpha)$. We reorder the data such that $Y_1 = \dots = Y_K = 1$, and $Y_{K+1} = \dots = Y_n = 0$. Then the joint likelihood conditional on the random effect θ is

$$\begin{aligned}
P(Y_1 = y_1, \dots, Y_n = y_n) &= \prod_{i \leq K} \left\{ 1 - \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \right\} \prod_{i > K} \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \\
&= \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] - \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \sum_{i=1}^K \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \\
&\quad + \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \sum_{1 < i < j \leq K} \left\{ \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} - \left(\frac{\theta_j}{z_j} \right)^{1/\alpha} \right] \right\} \\
&\quad + \dots + (-1)^K \exp \left[- \sum_{i=1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right]
\end{aligned} \tag{15}$$

193 Finally marginalizing over the random effect, we obtain

$$\begin{aligned}
P(Y_1 = y_1, \dots, Y_n = y_n) &= \int G(\mathbf{z}|\mathbf{A})p(\mathbf{A}|\alpha)d\mathbf{A}. \\
&= \int \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] - \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \sum_{i=1}^K \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \\
&\quad + \exp \left[- \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] \sum_{1 < i < j \leq K} \left\{ \exp \left[- \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} - \left(\frac{\theta_j}{z_j} \right)^{1/\alpha} \right] \right\} \\
&\quad + \dots + (-1)^K \exp \left[- \sum_{i=1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right] p(\mathbf{A}|\alpha)d\mathbf{A}.
\end{aligned} \tag{16}$$

194 Consider the first term in the summation,

$$\begin{aligned}
\int \exp \left\{ - \sum_{i=K+1}^n \left(\frac{\theta_i}{z_i} \right)^{1/\alpha} \right\} p(\mathbf{A}|\alpha) d\mathbf{A} &= \int \exp \left\{ - \sum_{i=K+1}^n \left(\frac{\left[\sum_{l=1}^L A_l w_l(\mathbf{s}_i)^{1/\alpha} \right]^\alpha}{z_i} \right)^{1/\alpha} \right\} p(\mathbf{A}|\alpha) d\mathbf{A} \\
&= \int \exp \left\{ - \sum_{i=K+1}^n \sum_{l=1}^L A_l \left(\frac{w_l(\mathbf{s}_i)}{z_i} \right)^{1/\alpha} \right\} p(\mathbf{A}|\alpha) d\mathbf{A} \\
&= \exp \left\{ - \sum_{l=1}^L \left[\sum_{i=K+1}^n \left(\frac{w_l(\mathbf{s}_i)}{z_i} \right)^{1/\alpha} \right]^\alpha \right\}. \tag{17}
\end{aligned}$$

195 The remaining terms in equation (16) are straightforward to obtain, and after integrating out the random
196 effect, the joint density for $K = 0, 1, 2$ is given by

$$P(Y_1 = y_1, \dots, Y_n = y_n) = \begin{cases} G(\mathbf{z}) & K = 0 \\ G(\mathbf{z}_{(1)}) - G(\mathbf{z}) & K = 1 \\ G(\mathbf{z}_{(12)}) - G(\mathbf{z}_{(1)}) - G(\mathbf{z}_{(2)}) + G(\mathbf{z}) & K = 2 \end{cases} \tag{18}$$

197 where

$$G[\mathbf{z}_{(1)}] = P[Z(\mathbf{s}_2) < z(\mathbf{s}_2), \dots, Z(\mathbf{s}_n) < z(\mathbf{s}_n)]$$

$$G[\mathbf{z}_{(2)}] = P[Z(\mathbf{s}_1) < z(\mathbf{s}_1), Z(\mathbf{s}_3) < z(\mathbf{s}_3), \dots, Z(\mathbf{s}_n) < z(\mathbf{s}_n)]$$

$$G[\mathbf{z}_{(12)}] = P[Z(\mathbf{s}_3) < z(\mathbf{s}_3), \dots, Z(\mathbf{s}_n) < z(\mathbf{s}_n)].$$

198 Similar expressions can be derived for all K , but become cumbersome for large K .

199 A.3 Simulation study pairwise difference results

200 Needs updating

The following tables show the methods that have significantly different Brier scores when using a Wilcoxon-Nemenyi-McDonald-Thompson test. In each column, different letters signify that the methods have significantly different Brier scores.

Table 3: Pairwise BS comparisons

	Setting 1	Setting 2	Setting 3	Setting 4	Setting 5	Setting 6
Method 1	A	A	A	C	B	B
Method 2	A B	B	A	B	A	A
Method 3	B	B	A	A	A B	A

References

- Agresti, A. (2003) *Categorical Data Analysis*. Wiley Series in Probability and Statistics. Wiley, 2nd edn.
- Cohen, J. (1960) A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, **20**, 37–46.
- Coles, S. (2001) *An Introduction to Statistical Modeling of Extreme Values*. Lecture Notes in Control and Information Sciences. London: Springer.
- De Oliveira, V. (2000) Bayesian prediction of clipped Gaussian random fields. *Computational Statistics and Data Analysis*, **34**, 299–314.
- Finley, A. O., Banerjee, S. and Gelfand, A. E. (2015) spBayes for Large Univariate and Multivariate Point-Referenced Spatio-Temporal Data Models. *Journal of Statistical Software*, **63**.
- Genton, M. G., Ma, Y. and Sang, H. (2011) On the likelihood function of Gaussian max-stable processes. *Biometrika*, **98**, 481–488.
- Gneiting, T. and Raftery, A. E. (2007) Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, **102**, 359–378.
- de Haan, L. (1984) A Spectral Representation for Max-stable Processes. *The Annals of Probability*, **12**, 1194–1204.
- Heagerty, P. and Lele, S. (1998) A Composite Likelihood Approach to Binary Spatial Data. *Journal of the American Statistical Association*, **1459**, 1099–1111.
- Huser, R. and Davison, A. C. (2014) Space-time modelling of extreme events. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **76**, 439–461.
- Pacifici, K., Reich, B. J., Dorazio, R. M. and Conroy, M. J. (2016) Occupancy estimation for rare species using a spatially-adaptive sampling design. *Methods in Ecology and Evolution*, **7**, 285–293.

- 226 Padoan, S. A., Ribatet, M. and Sisson, S. A. (2010) Likelihood-Based Inference for Max-Stable Processes.
227 *Journal of the American Statistical Association*, **105**, 263–277.
- 228 R Core Team (2016) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
229 Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 230 Reich, B. J. and Shaby, B. A. (2012) A hierarchical max-stable spatial model for extreme precipitation. *The*
231 *Annals of Applied Statistics*, **6**, 1430–1451.
- 232 Sing, T., Sander, O., Beerenwinkel, N. and Lengauer, T. (2005) ROCR: visualizing classifier performance
233 in R. *Bioinformatics*, **21**, 3940–3941.
- 234 Smith, D. R., Yuancai, L., Walter, C. A. and Young, J. A. (2012) Incorporating predicted species distribution
235 in adaptive and conventional sampling designs. In *Design and Analysis of Long-term Ecological Mon-*
236 *itoring Studies* (eds. R. A. Gitzen, J. J. Millspaugh, A. B. Cooper and D. S. Licht), chap. 17, 381–396.
237 Cambridge University Press.
- 238 Stephenson, A. G. (2009) High-Dimensional Parametric Modelling of Multivariate Extreme Events. *Aus-*
239 *tralian & New Zealand Journal of Statistics*, **51**, 77–88.
- 240 Tawn, J. A. (1990) Modelling multivariate extreme value distributions. *Biometrika*, **77**, 245–253.
- 241 Thibaud, E. and Opitz, T. (2015) Efficient inference and simulation for elliptical Pareto processes.
242 *Biometrika*, **102**, 855–870.
- 243 Wadsworth, J. L. and Tawn, J. A. (2014) Efficient inference for spatial extreme value processes associated
244 to log-Gaussian random functions. *Biometrika*, **101**, 1–15.
- 245 Wang, X. and Dey, D. K. (2010) Generalized extreme value regression for binary response data: An appli-
246 cation to B2B electronic payments system adoption. *The Annals of Applied Statistics*, **4**, 2000–2023.