

A spatial model for rare binary events

August 8, 2016

1 Introduction

The goals of spatial binary data analysis are often to estimate covariate effects while accounting for spatial dependence and to make predictions at locations without samples. A common approach to incorporate spatial dependence in the model for binary data is relating a continuous spatial process $Z(\mathbf{s}) \in \mathbb{R}$ to the binary response $Y(\mathbf{s})$ by thresholding, $Y(\mathbf{s}) = I[Z(\mathbf{s}) > c]$, where $I[\cdot]$ is an indicator function. In many spatial analyses of binary data, a Gaussian process is used to model $Z(\mathbf{s})$. This is true for both spatial probit and spatial logistic regression. In this model, spatial dependence is determined by the joint probability that two sites simultaneously exceed the threshold c . However, when c is large, and thus $Y(\mathbf{s}) = 1$ is rare, then the asymptotic theory suggests that the Gaussian process will model dependence poorly. In fact, even under strong spatial correlation for Z , it gives asymptotic independence (?), suggesting that for rare binary data, the Gaussian model will not perform very well.

We propose using a latent max-stable process (de Haan, 1984) for $Z(\mathbf{s})$ because it allows for asymptotic dependence. The max-stable process arises as the limit of the location-wise maximum of infinitely many spatial processes, and any finite-dimensional representation of a max-stable process has generalized extreme value distribution (GEV) marginal distributions. 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). Composite likelihoods have been used to model binary spatial data (Heagerty and Lele, 1998), but not using max-stable processes. 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. We chose to use this low-rank representation for our rare binary spatial regression model. Our model builds on related work by Wang and Dey (2010) who use a GEV link for non-spatial binary data. The proposed model generalizes this to have spatial dependence.

The paper proceeds as follows. In Section 2 we present the proposed latent max-stable process for rare binary data analysis. In Section 3 we give the bivariate distribution for our model. In Section 4 we show a link between a commonly used measure of dependence between binary variables and another metric for extremal dependence. The computing for our model is outlined in Section 5. Finally, we present a simulation study in Section 6 followed in Section 7 by a data analysis of two species: *Tamarix ramosissima* and *Hedysarum scoparium* in . Lastly, in Section 8 we provide some discussion and possibilities for future research.

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 $\mathcal{D} \in \mathbb{R}^2$. 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 individually 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 (over $Z(\mathbf{s})$) 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).

45 For a finite collection of locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, we denote by $\mathbf{Y} = [Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)]^T$ the vector of
 46 observations. The spatial dependence of \mathbf{Y} is determined by the joint distribution of the latent variable
 47 $\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T$. To incorporate spatial dependence, we consider the hierarchical representation of
 48 the max-stable process proposed in Reich and Shaby (2012). Consider a set of positive stable random effect
 49 $A_1, \dots, A_L \stackrel{iid}{\sim} \text{PS}(\alpha)$ associated with spatial knots $\mathbf{v}_1, \dots, \mathbf{v}_L \in \mathbb{R}^2$. The hierarchical model is given by

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

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

53 Because the latent $\mathbf{Z}(\mathbf{s})$ are independent given the random effects $\theta(\mathbf{s})$, the binary responses are also
 54 conditionally independent. This leads to the tractable likelihood

$$Y(\mathbf{s}_i) | A_1, \dots, A_L \stackrel{ind}{\sim} \text{Bern}[\pi(\mathbf{s}_i)] \quad (2)$$

55 where

$$\pi(\mathbf{s}_i) = 1 - \exp \left\{ - \sum_{l=1}^L A_l \left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)} \right)^{1/\alpha} \right\} \quad (3)$$

56 and $z(\mathbf{s}) = [1 + \xi \mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}]^{1/\xi}$. Marginally over the A_l , this gives

$$Z(\mathbf{s}) \sim \text{GEV} \left[\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}, 1, \xi \right], \quad (4)$$

57 and thus $P[Y(\mathbf{s}) = 1] = 1 - \exp \left\{ -\frac{1}{z(\mathbf{s})} \right\}$.

58 Many weight functions are possible, but the weights must be constrained so that $\sum_{l=1}^L w_l(\mathbf{s}_i) = 1$ for
 59 $i = 1, \dots, n$ to preserve the marginal GEV distribution. For example, Reich and Shaby (2012) take the
 60 weights to be scaled Gaussian kernels with knots \mathbf{v}_l ,

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

61 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
 62 spatial range of the dependence, with large ρ giving long-range dependence and vice versa.

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

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

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

66 **3 Joint distribution**

67 We give an exact expression in the case where there are only two spatial locations which is useful for
 68 constructing a pairwise composite likelihood (Padoan et al., 2010) and studying spatial dependence. When

69 $n = 2$, the probability mass function is given by

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

70 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
 71 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,
 72 as might be expected for very rare events, see Appendix A.2.

73 4 Quantifying spatial dependence

74 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
 75 common measure of dependence between binary variables is Cohen's Kappa (Cohen, 1960),

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

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

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

$$P_E(\beta) = 1 - 2 \exp\left\{-\frac{1}{\beta}\right\} + 2 \exp\left\{-\frac{2}{\beta}\right\},$$

78 and

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

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

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

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

82 5 Computation

83 For small K , we can evaluate the likelihood directly. When K is large, we use Markov chain Monte
 84 Carlo (MCMC) methods with the random effects model to explore the posterior distribution. To overcome
 85 challenges with evaluating the positive stable density, we follow Reich and Shaby (2012) and introduce a
 86 set of auxiliary variables B_1, \dots, B_L following the auxiliary variable technique of Stephenson (2009) (for
 87 more details, see Appendix A.3 of Reich and Shaby (2012)). So, the hierarchical model is given by

$$\begin{aligned} Y(\mathbf{s}_i) | \pi(\mathbf{s}_i) &\overset{ind}{\sim} \text{Bern}[\pi(\mathbf{s}_i)] \\ \pi(\mathbf{s}_i) &= 1 - \exp\left\{-\sum_{l=1}^L A_l \left(\frac{w_l(\mathbf{s}_i)}{z(\mathbf{s}_i)}\right)^{1/\alpha}\right\} \\ A_l &\sim \text{PS}(\alpha) \end{aligned} \quad (11)$$

with priors $\beta \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p)$, $\xi \sim N(0, \sigma_\xi^2)$, $\rho \sim \text{Unif}(\rho_l, \rho_u)$, and $\alpha \sim \text{Beta}(a_\alpha, b_\alpha)$. The model parameters are updated using Metropolis Hastings (MH) update steps, and the random effects A_1, \dots, A_L , and auxiliary variables B_1, \dots, B_L are updated using Hamiltonian Monte Carlo (HMC) update steps. The code for this is available online through <https://github.com/sammorris81/rare-binary>.

6 Simulation study

For our simulation study, we generate $n_m = 50$ datasets under 12 different simulation 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 complete datasets on a 100×100 rectangular grid of $n = 10,000$ locations. If a dataset is generated with $K < 50$, it is discarded and a new dataset is generated. This is done to guarantee that datasets have no less than 0.5% rarity. For model fitting, we select a subsample and use the remaining sites to evaluate predictive performance. 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.

6.1 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 $K = 500$ (5% rarity), on average. Because there are no covariates, we set $\xi = 0$. We then set $Y(\mathbf{s}) = I[Z(\mathbf{s}) > 0]$.

For the second process, we generate a latent variable from a spatial Gaussian process with a mean of

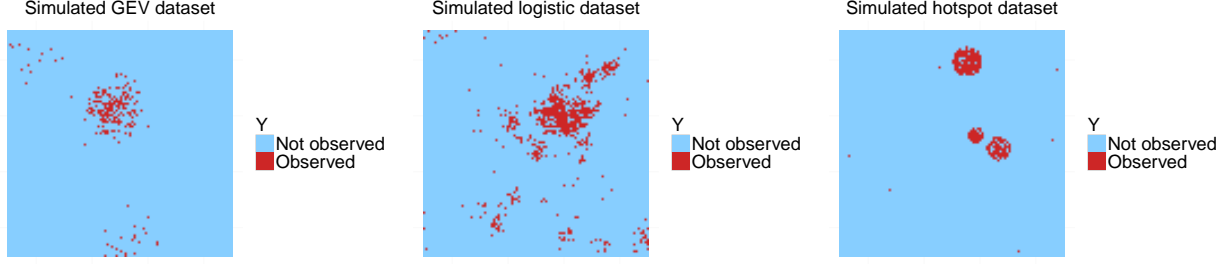


Figure 1: One simulated dataset from spatial GEV (left), spatial logistic (center), and hotspot (right) designs.

107 $\text{logit}(0.05) \approx -2.94$ and an exponential covariance given by

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

108 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 $\pi(\mathbf{s}_i) = \frac{\exp\{z(\mathbf{s})\}}{1 + \exp\{z(\mathbf{s})\}}$.

109 For the third process, we generate data using a hotspot method. For this process, we first generate
 110 hotspots throughout the space. Let n_{hs} be the number of hotspots in the space. Then $n_{\text{hs}} - 1 \sim \text{Poisson}(2)$.
 111 This generation scheme ensures that every dataset has at least one hotspot. We generate the hotspot locations
 112 $\mathbf{h}_1, \dots, \mathbf{h}_{n_{\text{hs}}} \sim \text{Unif}(0, 1)^2$. Let B_h be a circle of radius of radius r_h around hotspot $h = 1, \dots, n_{\text{hs}}$. The r_h
 113 differ for each hotspot and are generated i.i.d. from a $\text{Unif}(0.03, 0.08)$ distribution. We set $P[Y(\mathbf{s}_i) = 1] =$
 114 0.85 for all \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
 115 an average of approximately $K = 500$ for the datasets. Figure 1 gives an example dataset from each of the
 116 data settings.

117 6.2 Sampling methods

118 We subsample the generated data using $n_s = 100, 250$ initial locations for two different sampling designs.
 119 The first is a two-stage spatially-adaptive cluster technique (CLU) taken from Pacifici et al. (2016). In this
 120 design, if an initial location is occupied, we also include the four rook neighbor (north, east, south, and west)

121 sites in the sample. For the second design, we use a simple random sample (SRS) with the same number of
 122 sites included in the cluster sample. For the GEV setting, when $n_s = 100$, there are on average 117 sites
 123 and at most 142 sites in a sample, and when $n_s = 250$, there are on average 286 sites and at most 332 sites
 124 in a sample. For the logistic setting, when $n_s = 100$, there are on average 118 sites and at most 147 sites
 125 in a sample, and when $n_s = 250$, there are on average 290 sites and at most 330 sites in a sample. For the
 126 hotspot setting, when $n_s = 100$, there are on average 110 sites and at most 128 sites in a sample, and when
 127 $n_s = 250$, there are on average 275 sites and at most 306 sites in a sample.

128 6.3 Methods

129 For each dataset, we fit the model using three different models: the proposed spatial GEV model, a spatial
 130 probit model, and a spatial logistic model. Logistic and probit methods assume the underlying process
 131 is Gaussian. In this case, we assume that $Z(\mathbf{s})$ follows a Gaussian process with mean $\mathbf{X}(\mathbf{s})^\top \boldsymbol{\beta}$ and unit
 132 variance. For the simulation study, we use an intercept only model. The marginal distributions are given by

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

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

136 logistic model. For the probit model, we use

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

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

138 6.4 Priors

139 For all models, we only include an intercept term β_0 in the model, and the prior for the intercept is
 140 $\beta_0 \sim \mathcal{N}(0, 10)$. Additionally, for all models, the prior for the bandwidth is $\rho \sim \text{Unif}(0.001, 1)$. In all
 141 methods, we place knots at all data points. For the GEV method, the prior for the spatial dependence pa-
 142 rameter is $\alpha \sim \text{Beta}(2, 5)$. We select this prior because it gives greater weight to $\alpha < 0.5$, which is the
 143 point at which spatial dependence becomes fairly weak, but also avoids values below 0.1 which can lead to
 144 numerical problems. We fix $\xi = 0$ because we do not include any covariates. For both the spatial probit
 145 and logistic models, the prior on the variance term for the random effects is $\text{IG}(0.1, 0.1)$ where $\text{IG}(\cdot)$ is an
 146 Inverse Gamma distribution.

147 6.5 Model comparisons

148 For each dataset, we fit the model using the n_s observations as a training set, and validate the model's
 149 predictive power at the remaining grid points. Let \mathbf{s}_j^* be the j th site in the validation set. From the posterior
 150 distributions of the parameters we can calculate $P[Y(\mathbf{s})_j^* = 1]$. To obtain $\hat{P}[Y(\mathbf{s}_j^*) = 1]$, we take the
 151 average of the posterior distribution for each j . We consider a few different metrics for comparing model
 152 performance. One score is the Brier scores (Gneiting and Raftery, 2007, BS). The Brier score for predicting
 153 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 s . We average the Brier scores over all test sites, and a lower score indicates a better fit. The Brier score equally penalizes false negatives and false positives, but in the case of rare data, this may not be the best metric due to the unbalanced nature of the data. Therefore, 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 AUROC across all datasets for each method and setting to obtain a single AUROC for each combination of method and setting.

6.6 Results

Overall, we find that the spatial probit model actually performs quite well in all cases. Table 1 gives the Brier scores and AUROC for each of the methods. Looking at Brier scores, we see that our model is outperformed by the probit model in all cases, by the logistic models in many settings. For AUROC, in a few of the settings, we do demonstrate a small improvement over the probit and logistic models. Because these results are somewhat surprising, we also considered the performance metrics by rareness of the data. We plot the Brier scores and AUROC for a few different sampling methods using the GEV and logistic links in Figure 2 – Figure 4 using a Loess smoother. These plots give evidence to suggest that as rareness increases, the spatial GEV method has potential to outperform the spatial probit and logistic models based on AUROC.

7 Data analysis

We compare our method to the spatial probit and logistic models for mapping the probability of the occurrence of *Tamarix ramosissima* (TR) and *Hedysarum scoparium* (HS), two 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 5. The region is split into 10-m \times 10-m grid

Table 1: Brier scores ($\times 100$) [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	3.10 [0.27]	2.45 [0.19]	2.79 [0.25]	0.926 [0.009]	0.942 [0.009]	0.900 [0.020]
		SRS	2.92 [0.20]	2.54 [0.18]	2.92 [0.25]	0.938 [0.007]	0.951 [0.007]	0.879 [0.021]
	250	CLU	2.18 [0.15]	1.87 [0.13]	2.05 [0.14]	0.951 [0.008]	0.948 [0.011]	0.922 [0.017]
		SRS	2.29 [0.15]	2.06 [0.13]	2.26 [0.15]	0.949 [0.009]	0.949 [0.010]	0.908 [0.020]
Logistic	100	CLU	5.29 [0.25]	4.94 [0.23]	5.10 [0.25]	0.659 [0.012]	0.676 [0.014]	0.643 [0.013]
		SRS	5.32 [0.23]	5.09 [0.24]	5.34 [0.26]	0.690 [0.012]	0.693 [0.012]	0.613 [0.012]
	250	CLU	4.81 [0.21]	4.55 [0.21]	4.66 [0.22]	0.731 [0.010]	0.749 [0.010]	0.714 [0.014]
		SRS	4.86 [0.22]	4.63 [0.20]	5.01 [0.23]	0.742 [0.010]	0.760 [0.010]	0.698 [0.015]
Hotspot	100	CLU	2.29 [0.17]	2.01 [0.15]	1.81 [0.12]	0.841 [0.016]	0.833 [0.019]	0.824 [0.020]
		SRS	2.09 [0.13]	1.87 [0.12]	2.13 [0.15]	0.885 [0.015]	0.906 [0.013]	0.844 [0.015]
	250	CLU	1.65 [0.11]	1.25 [0.08]	1.40 [0.09]	0.934 [0.009]	0.949 [0.008]	0.939 [0.011]
		SRS	1.53 [0.10]	1.31 [0.08]	1.63 [0.11]	0.947 [0.007]	0.960 [0.005]	0.918 [0.015]

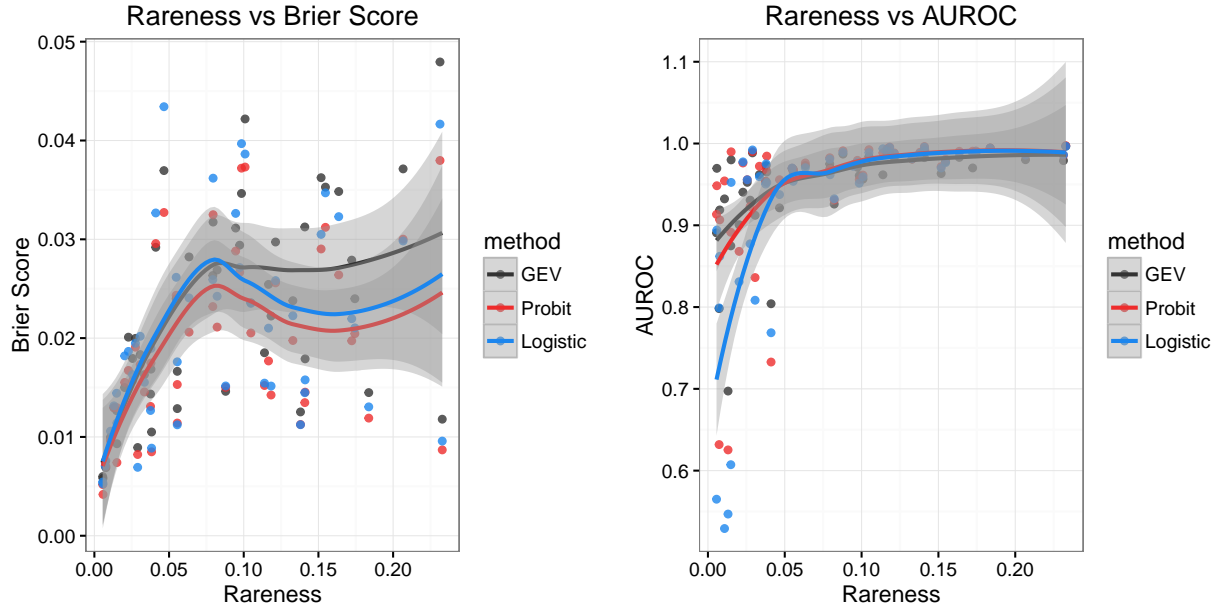


Figure 2: Smooth of BS (left) and AUROC (right) by rareness for GEV link and cluster sampling with $n_s = 250$.

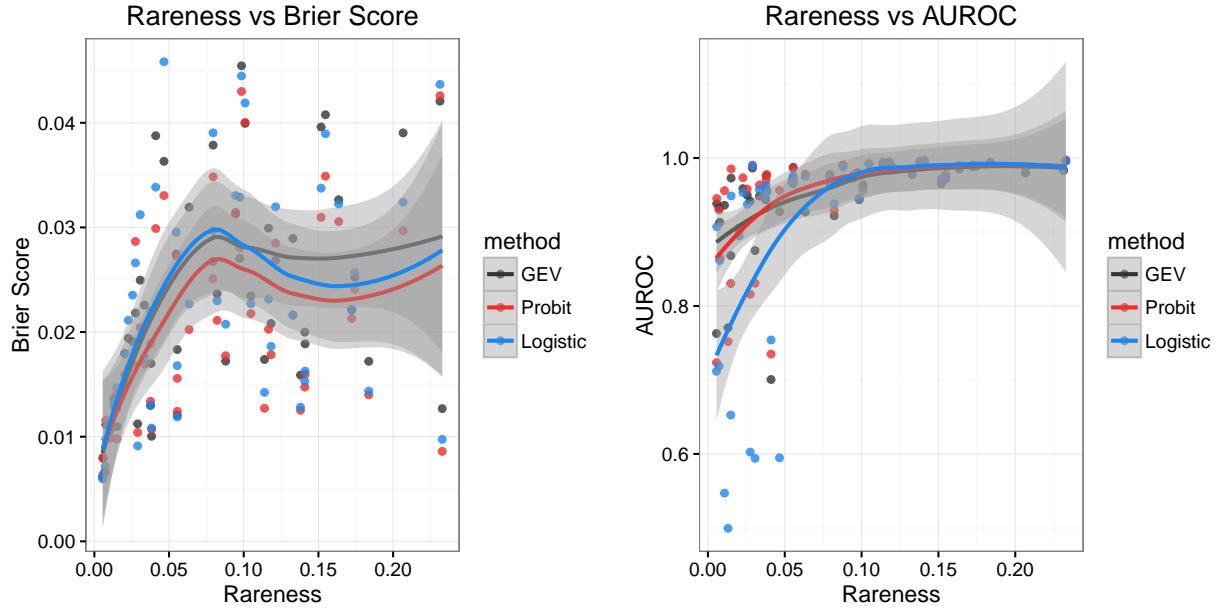


Figure 3: Smooth of BS (left) and AUROC (right) by rareness for GEV link and simple random sampling with $n_s = 250$.

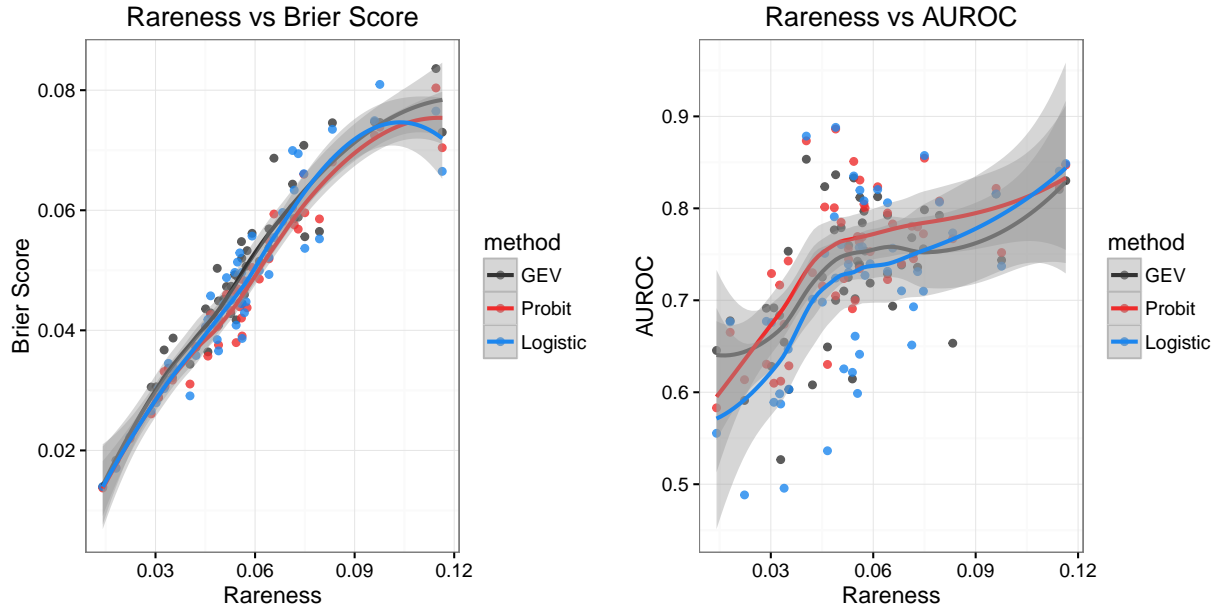


Figure 4: Smooth of BS (left) and AUROC (right) by rareness for logistic link and cluster sampling with $n_s = 250$.

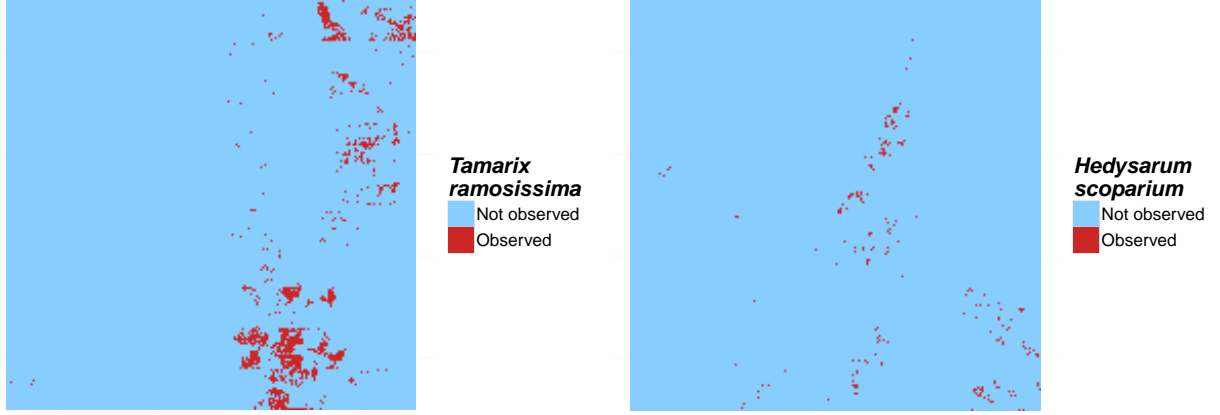


Figure 5: True occupancy of *Tamarix ramosissima*(left) and *Hedysarum scoparium*(right) from a 1-km² study region of PR China.

cells. *Tamarix ramosissima* can be found in approximately 6% of the grid cells, and *Hedysarum scoparium* can be found in approximately 0.54% of the grid cells.

7.1 Methods

For the data analysis, we generate 50 subsamples using the CLU and SRS sampling methods with $n_s = 100$, 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.

7.2 Results

As with the simulation study, we find that in most cases the spatial probit model gives the best performance. Table 2 give summary Brier scores ($\times 100$) and AUROC for the *Tamarix ramosissima* and *Hedysarum scoparium* analysis. Figure 6 gives the vertically averaged ROC curves for each method and sampling setting for *Tamarix ramosissima* and Figure 7 gives the vertically averaged ROC curves for *Hedysarum scoparium*. These results appear to support the suggestion from the simulation study that spatial GEV gives an advan-

Table 2: Brier scores ($\times 100$) [SE] and AUROC [SE] for GEV, Probit, and Logistic methods for *Tamarix ramosissima* and *Hedysarum scoparium*.

Species	n	Sample	BS			AUROC		
			GEV	Probit	Logistic	GEV	Probit	Logistic
TR	100	CLU	5.120 [0.050]	5.039 [0.049]	5.382 [0.029]	0.732 [0.014]	0.731 [0.014]	0.699 [0.012]
		SRS	4.997 [0.045]	4.938 [0.055]	5.500 [0.027]	0.798 [0.008]	0.802 [0.009]	0.636 [0.012]
	250	CLU	4.779 [0.049]	4.657 [0.045]	4.950 [0.051]	0.771 [0.013]	0.784 [0.013]	0.798 [0.011]
		SRS	4.823 [0.053]	4.735 [0.048]	5.120 [0.071]	0.827 [0.011]	0.851 [0.007]	0.717 [0.019]
HS	100	CLU	1.765 [0.018]	1.831 [0.029]	1.679 [0.002]	0.642 [0.010]	0.617 [0.012]	0.573 [0.008]
		SRS	1.914 [0.066]	1.996 [0.083]	1.685 [0.002]	0.686 [0.009]	0.683 [0.011]	0.587 [0.008]
	250	CLU	1.667 [0.005]	1.657 [0.006]	1.679 [0.001]	0.659 [0.009]	0.648 [0.011]	0.566 [0.005]
		SRS	1.691 [0.017]	1.666 [0.010]	1.684 [0.001]	0.691 [0.010]	0.709 [0.012]	0.574 [0.007]

tage as rareness increases. For *Tamarix ramosissima*, when $n = 100$, there is a small distinction between the spatial GEV and probit models. The results suggest that the GEV model has a small improvement over probit in the case of cluster sampling, but that using a probit model demonstrates a small improvement over the GEV model in the simple random setting. However, when $n = 250$, both logistic and probit models appear to outperform the GEV model. For the rarer species, *Hedysarum scoparium*, we find more conclusive evidence that the GEV model provides an improvement for cluster sampling when $n = 100$. At this sample size, there is also some evidence to suggest that the GEV model gives some improvement over the probit model for simple random sampling. For $n = 250$, we still have evidence that using a cluster sampling strategy, the GEV model gives the best performance, but for simple random sampling, the probit model performs the best.

8 Discussion and future research

In this paper, we present a max-stable spatial method for rare binary data. The principal finding in this paper is that the spatial probit model tends to outperform the proposed model. This finding is surprising given that the max-stable process is the theoretically justified spatial process for extreme value distributions, and it leads to possible research questions in the future.

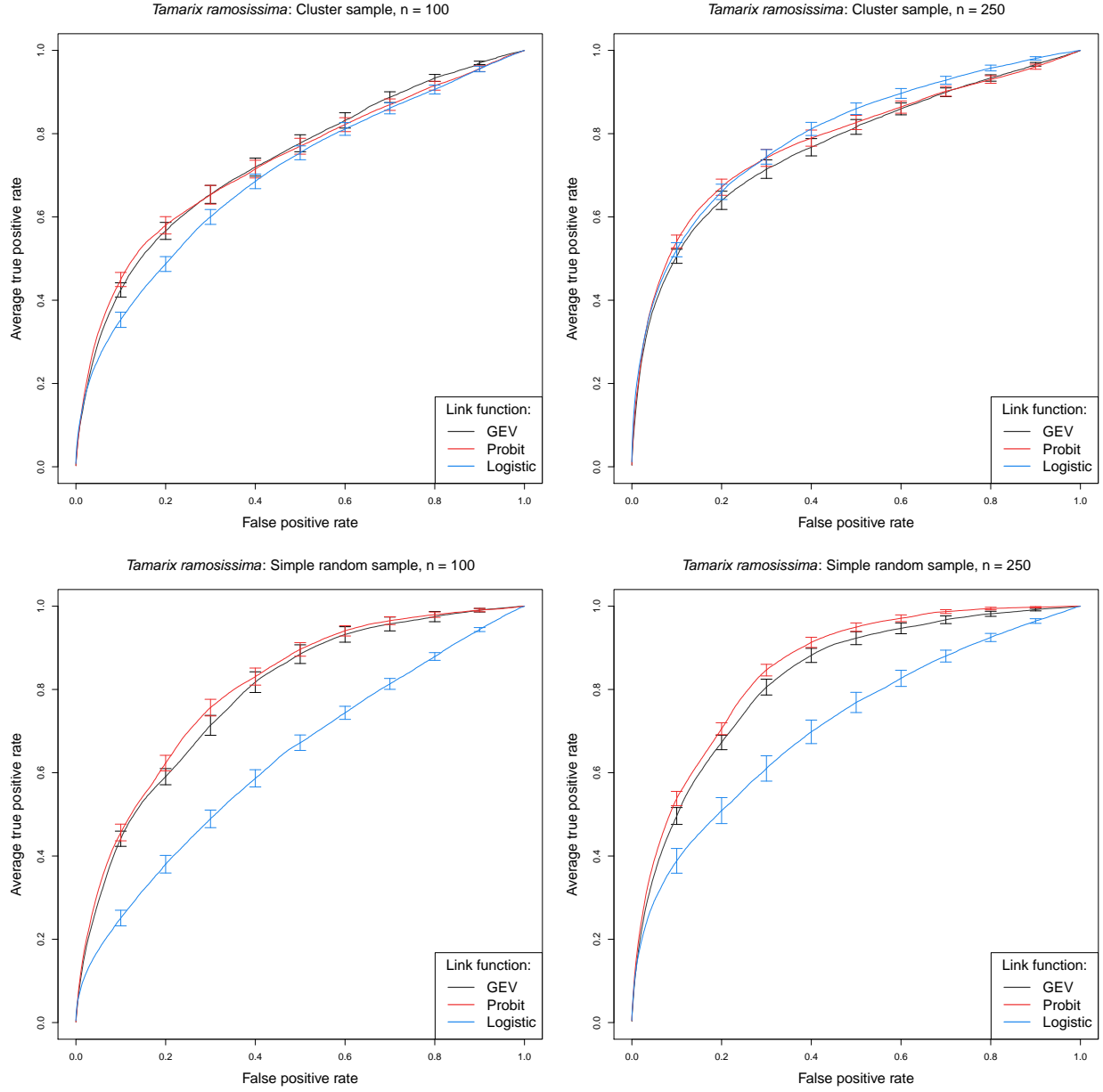


Figure 6: Vertically averaged ROC curves for *Tamarix ramosissima*.

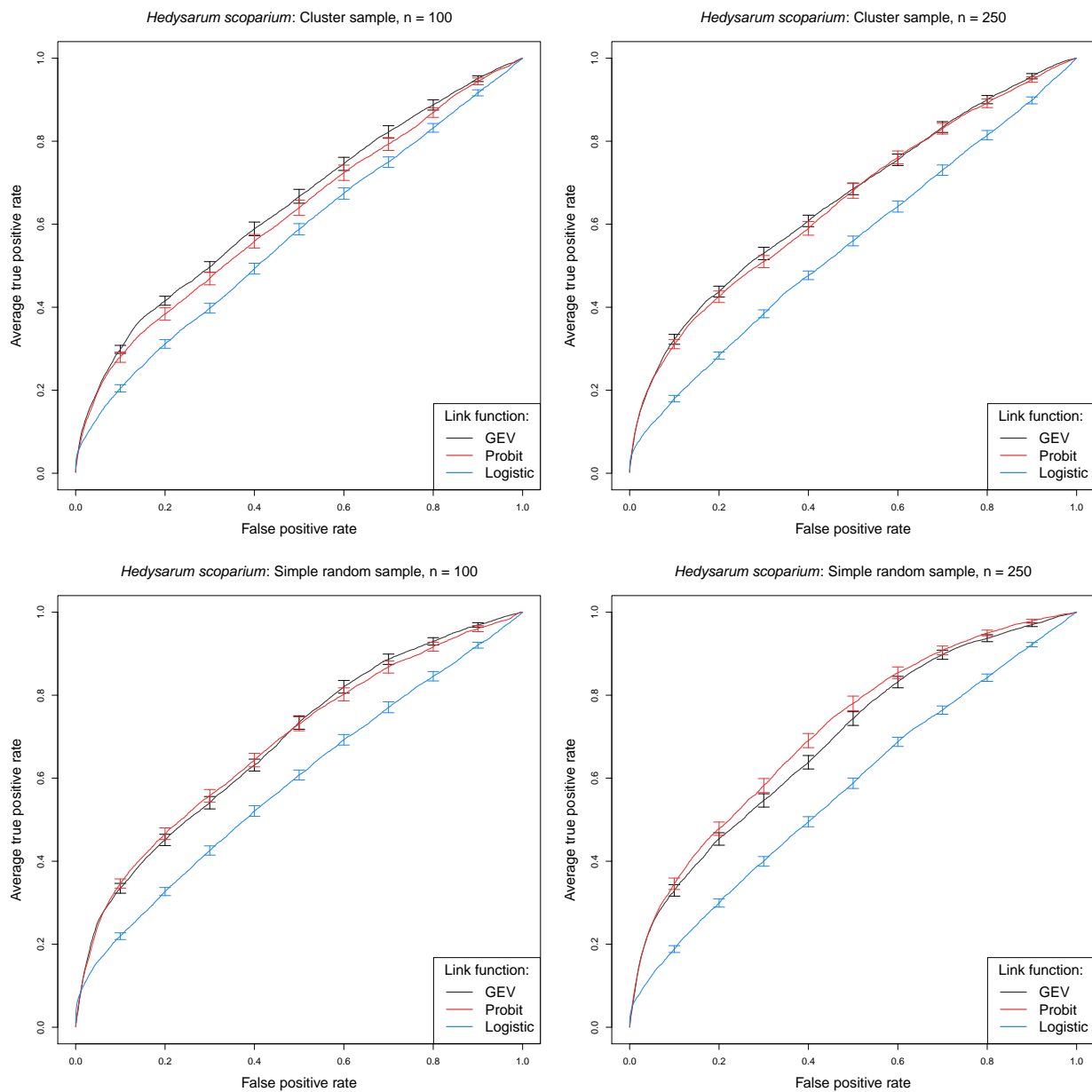


Figure 7: Vertically averaged ROC curves for *Hedysarum scoparium*.

It is unusual that the spatial probit model should outperform the proposed model, particularly when the data are generated directly from the proposed model. One possible explanation is that for the simulated data, there is a wide range of rarity in the data (GEV: 0.5% – 35.9%, Logistic: 1.4% – 14.4%, and Hotspot: 0.5% – 6.8%). Given that for both the GEV and logistic data settings, we have a number of datasets with a relatively high rate of occurrence, it is possible that probit is competitive partly due to the fact that the data are not rare. Both the simulation study and data analysis appear to support the idea that the GEV method will perform better on rarer datasets. Therefore, it may be useful to conduct more research on rare datasets or through simulation with a slightly more restrictive data generation strategy (i.e. restrict datasets to $K < 500$).

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

224 $Y_{K+1} = \dots = Y_n = 0$. Then the joint likelihood conditional on the random effect θ is

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

225 Finally marginalizing over the random effect, we obtain

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

226 Consider the first term in the summation,

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

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

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

229 where

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

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

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

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

References

- 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.
- 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.
- Padoan, S. A., Ribatet, M. and Sisson, S. A. (2010) Likelihood-Based Inference for Max-Stable Processes. *Journal of the American Statistical Association*, **105**, 263–277.
- R Core Team (2016) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Reich, B. J. and Shaby, B. A. (2012) A hierarchical max-stable spatial model for extreme precipitation. *The Annals of Applied Statistics*, **6**, 1430–1451.
- Sing, T., Sander, O., Beerenwinkel, N. and Lengauer, T. (2005) ROCr: visualizing classifier performance in R. *Bioinformatics*, **21**, 3940–3941.
- Smith, D. R., Yuancai, L., Walter, C. A. and Young, J. A. (2012) Incorporating predicted species distribution in adaptive and conventional sampling designs. In *Design and Analysis of Long-term Ecological Monitoring Studies* (eds. R. A. Gitzen, J. J. Millspaugh, A. B. Cooper and D. S. Licht), chap. 17, 381–396. Cambridge University Press.
- Stephenson, A. G. (2009) High-Dimensional Parametric Modelling of Multivariate Extreme Events. *Australian & New Zealand Journal of Statistics*, **51**, 77–88.
- Tawn, J. A. (1990) Modelling multivariate extreme value distributions. *Biometrika*, **77**, 245–253.

- 265 Thibaud, E. and Opitz, T. (2015) Efficient inference and simulation for elliptical Pareto processes.
266 *Biometrika*, **102**, 855–870.
- 267 Wadsworth, J. L. and Tawn, J. A. (2014) Efficient inference for spatial extreme value processes associated
268 to log-Gaussian random functions. *Biometrika*, **101**, 1–15.
- 269 Wang, X. and Dey, D. K. (2010) Generalized extreme value regression for binary response data: An appli-
270 cation to B2B electronic payments system adoption. *The Annals of Applied Statistics*, **4**, 2000–2023.