

A spatial-skew model for extreme values

April 6, 2015

1 Introduction

In most climatological applications, researchers are interested in learning about the average behavior of different climate variables (e.g. ozone, temperature, rainfall). However, averages do not help regulators prepare for the unusual events that only happen once every 100 years. For example, it is important to have an idea of how much rain will come in a 100-year flood in order to construct strong enough river levees to protect lands from flooding.

Unlike multivariate normal distributions, it is challenging to model multivariate extreme value distributions (e.g. generalized extreme value and generalized Pareto distribution) because few closed-form expressions exist for the density in more than two-dimensions (Coles and Tawn, 1991). Given this limitation, pairwise composite likelihoods have been used when modeling dependent extremes (Padoan et al., 2010; Blanchet and Davison, 2011; Huser, 2013).

One way around the multi-dimensional limitation of multivariate extreme value distributions is to use skew elliptical distributions to model dependent extreme values (Genton, 2004; Zhang and El-Shaarawi, 2010; Padoan, 2011). Due to their flexibility, the skew-normal and skew- t distribution offer a flexible way to handle non-symmetric data within a framework of multivariate normal and multivariate t -distributions. As with the spatial Gaussian process, the skew-normal distribution is also asymptotically independent; however, the skew- t does demonstrate asymptotic dependence (Padoan, 2011). Although asymptotic dependence is desirable between sites that are near one another, one drawback to the skew- t is that sites remain asymptotically dependent even at far distances.

In this paper, we present a model that has marginal distributions with flexible tails, demonstrates asymp-

23 totic dependence for observations at sites that are near to one another, and has computation on the order of
24 Gaussian models for large space-time datasets. Specifically, our contribution is to incorporate thresholding
25 and random spatial partitions using a multivariate skew- t distribution. The advantage of using a thresholded
26 model as opposed to a non-thresholded model is that it allows for the tails of the distribution to inform
27 the predictions in the tails (DuMouchel, 1983). The random spatial partitions are similar to the method
28 used by Kim et al. (2005) for non-stationary Gaussian data. The partition alleviates the long-range spatial
29 dependence present in the skew- t distribution.

30 The paper is organized as follows. Section 2 is a brief review of the spatial skew- t process. In Section
31 3.3, we build upon the traditional skew- t by incorporating censoring to focus on tails, partitioning to remove
32 long-range asymptotic dependence, and extending the model to space-time data. The computing is described
33 in Section 4. In Section 5, we present a simulation study that examines the predictive capabilities of this
34 model compared with a naïve Gaussian method. We then compare our method to Gaussian and max-stable
35 methods with a data analysis of ozone measurements from throughout the US in section 6. The final section
36 provides brief discussion and direction for future research.

37 **2 Spatial skew processes**

38 Many types of data demonstrate some level of skewness and therefore should be modeled with distributions
39 that allow for asymmetry. The skew-elliptical family of distributions provides models that are mathemati-
40 cally tractable while introducing a slant parameter to account for asymmetric data (Genton, 2004). A brief
41 review of the additive process by which a skew- t process is created is given here.

2.1 Skew- t process

Let $Y(\mathbf{s})$ be the observation at spatial location $\mathbf{s} = (s_1, s_2)$. The spatial skew- t process can be written

$$Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})^T \boldsymbol{\beta} + \lambda \sigma |z| + \sigma v(\mathbf{s}) \quad (1)$$

where $\mathbf{X}(\mathbf{s})$ is a set of spatial covariates at site \mathbf{s} , $\boldsymbol{\beta}$ is the vector of regression parameters, $\lambda \in \mathcal{R}$ is a parameter controlling skew, $z \sim N(0, 1)$, $\sigma^2 \sim \text{IG}(a, b)$ is an inverse gamma random variable, and $v(\mathbf{s})$ is a spatial Gaussian process with mean zero and variance one.

Let $\mathbf{Y} = [Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)]^T$ be a set of observations at a finite collection of locations $\mathbf{s}_1, \dots, \mathbf{s}_n$. After marginalizing over both z and σ ,

$$\mathbf{Y} \sim \text{ST}_n(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Omega}, \boldsymbol{\alpha}, 2a), \quad (2)$$

that is, \mathbf{Y} follows an n -dimensional skew- t distribution with location $\mathbf{X}\boldsymbol{\beta}$, correlation matrix $\boldsymbol{\Omega}$, slant parameters $\boldsymbol{\alpha}$ and degrees of freedom $2a$, where $\mathbf{X} = [\mathbf{X}(\mathbf{s}_1)^T, \dots, \mathbf{X}(\mathbf{s}_n)^T]$, $\boldsymbol{\Omega} = \boldsymbol{\omega} \bar{\boldsymbol{\Omega}} \boldsymbol{\omega}$, $\boldsymbol{\omega} = \text{diag} \left(\frac{1}{\sqrt{ab}}, \dots, \frac{1}{\sqrt{ab}} \right)$, $\bar{\boldsymbol{\Omega}} = (\boldsymbol{\Sigma} + \lambda^2 \mathbf{1}\mathbf{1}^T)$, $\boldsymbol{\Sigma}$ is a positive definite correlation matrix, $\boldsymbol{\alpha} = \lambda(1 + \lambda^2 \mathbf{1}^T \boldsymbol{\Sigma}^{-1} \mathbf{1})^{-1/2} \mathbf{1}^T \boldsymbol{\Sigma}^{-1}$ is a vector of slant parameters. Although $\boldsymbol{\Sigma}$ can be any positive definite correlation matrix, we choose to use the stationary isotropic Matérn correlation with

$$\text{cor}[v(\mathbf{s}), v(\mathbf{t})] = \gamma I(\mathbf{s} = \mathbf{t}) + (1 - \gamma) \frac{1}{\Gamma(\nu) 2^{\nu-1}} \left(\sqrt{2\nu} \frac{h}{\rho} \right)^\nu K_\nu \left(\sqrt{2\nu} \frac{h}{\rho} \right) \quad (3)$$

where ρ is the spatial range, ν is the smoothness, γ is the proportion of variance accounted for by the spatial variation, K_ν is a modified Bessel function of the second kind, and $h = \|\mathbf{s} - \mathbf{t}\|$. This process is desirable because of its flexible tail that is controlled by the skewness parameter λ and degrees of freedom

2a. Furthermore, the marginal distributions at each location also follow a univariate skew- t distribution (Azzalini and Capitanio, 2013).

2.2 Extremal dependence

Our interest lies in spatial dependence in the tail of the skew- t process. One measure of extremal dependence is the χ statistic (Padoan, 2011). For a stationary and isotropic spatial process, the χ statistic for two location \mathbf{s} and \mathbf{t} separated by distance $h = \|\mathbf{s} - \mathbf{t}\|$ is

$$\chi(h) = \lim_{c \rightarrow \infty} \Pr[Y(\mathbf{s}) > c | Y(\mathbf{t}) > c]. \quad (4)$$

If $\chi(h) = 0$, then observations are asymptotically independent at distance h . For Gaussian processes, $\chi(h) = 0$ regardless of the distance, so they are not suitable for modeling spatially-dependent extremes. Unlike the Gaussian process, the skew- t process is asymptotically dependent. However, one problem with the spatial skew- t process is that $\lim_{h \rightarrow \infty} \chi(h) > 0$ (see Appendix A.4 for a proof). This occurs because all observations, both near and far, share the same z and σ terms. Therefore, this long-range dependence feature of the skew- t process is not ideal for spatial analysis of large geographic regions where we expect only local spatial dependence. The explicit expression for $\chi(h)$ is given in Appendix A.4.

3 Spatiotemporal skew- t model for extremes

In this section, we propose extensions to the skew- t process to model spatial extremes over a large geographic region by introducing censoring to focus on tail behavior and a random partition to remove long-range asymptotic dependence. For notational convenience, we introduce the model for a single replication, and then extend this model to the spatiotemporal setting in Section 3.3.

3.1 Censoring to focus on the tails

To avoid bias in estimating tail parameters, we model censored data. Let

$$\tilde{Y}(\mathbf{s}) = \begin{cases} Y(\mathbf{s}) & \delta(\mathbf{s}) = 1 \\ T & \delta(\mathbf{s}) = 0 \end{cases} \quad (5)$$

be the censored observation at site \mathbf{s} where $Y(\mathbf{s})$ is the uncensored observation, $\delta(\mathbf{s}) = I[Y(\mathbf{s}) > T]$, and T is a pre-specified threshold value. Then, assuming the uncensored data $Y(\mathbf{s})$ are observations from a skew- t process, we update values censored below the threshold using standard Bayesian missing data methods as described in Section 4.

3.2 Partitioning to remove long-range asymptotic dependence

We handle the problem of long-range asymptotic dependence with a random partition. As discussed in Section 2, the source of long-range dependence is the shared z and σ . Therefore, to alleviate this dependence, we allow z and σ to vary by site. The model becomes

$$Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})^T \boldsymbol{\beta} + \lambda \sigma(\mathbf{s}) |z(\mathbf{s})| + \sigma(\mathbf{s}) v(\mathbf{s}). \quad (6)$$

Let $\mathbf{w} = (w_1, w_2)$ be the location of a spatial knot. To model spatial variation, consider a set of spatial knots $\mathbf{w}_1, \dots, \mathbf{w}_K$ from a homogeneous Poisson process with intensity μ over spatial domain $\mathcal{D} \in \mathcal{R}^2$. The knots define a random partition of \mathcal{D} by subregions P_1, \dots, P_K defined as

$$P_k = \{\mathbf{s} : k = \arg \min_{\ell} \|\mathbf{s} - \mathbf{w}_{\ell}\|\}. \quad (7)$$

88 All $z(\mathbf{s})$ and $\sigma(\mathbf{s})$ for sites in subregion k are assigned common values

$$z(\mathbf{s}) = z_k \quad \text{and} \quad \sigma(\mathbf{s}) = \sigma_k \quad (8)$$

89 and the z_k and σ_k^2 are distributed as $z_k \stackrel{iid}{\sim} N(0, 1)$ and $\sigma_k^2 \stackrel{iid}{\sim} \text{IG}(a, b)$ where IG is the distribution function
 90 of an inverse gamma random variable. So, within each partition, $Y(\mathbf{s})$ follows the spatial skew- t process
 91 defined in Section 2. Across partitions, the $Y(\mathbf{s})$ remain correlated via the correlation function for $v(\mathbf{s})$
 92 because $v(\mathbf{s})$ spans all partitions.

93 When incorporating the random partition, conditional on knots $\mathbf{w}_1, \dots, \mathbf{w}_K$, the χ statistic for two sites
 94 \mathbf{s} and \mathbf{t} in partitions k_s and k_t respectively is

$$\begin{aligned} \chi(h) &= I(k_s = k_t) \chi_{\text{skew-}t}(h) + I(k_s \neq k_t) \chi_{\text{Gaus}}(h) \\ &= I(k_s = k_t) \chi_{\text{skew-}t}(h) \end{aligned} \quad (9)$$

95 where $I(\cdot)$ is an indicator function, $\chi_{\text{skew-}t}(h)$ is the χ statistic for a skew- t process, $\chi_{\text{Gaus}}(h)$ is the χ statistic
 96 for a Gaussian process, and $h = \|\mathbf{s} - \mathbf{t}\|$. Therefore, sites in different subregions are asymptotically inde-
 97 pendent because $\chi_{\text{Gaus}}(h) = 0$ for all h . Marginally, over the knots $\mathbf{w}_1, \dots, \mathbf{w}_K$, $\chi(h) = \pi(h) \chi_{\text{skew-}t}(h)$,
 98 where $\pi(h) = \Pr(k_s = k_t)$ is the probability that two sites separated by distance h are in the same partition.
 99 So, to show that $\lim_{h \rightarrow \infty} \chi(h) = 0$, we need only know that $\lim_{h \rightarrow \infty} \pi(h) = 0$. A proof of this is given in
 100 Appendix A.3.

101 In Figure 1, we give $\chi(h)$ for $K = 1, 3, 5, 10$ partitions for a skew- t distribution with $\alpha = 10$, and
 102 3 degrees of freedom. To estimate $\pi(h)$, we generate 500 sites uniformly over the unit-square. We then
 103 randomly generate 400 different sets of partitions using $K = 3, 5$, and 10. For each set of knots, we
 104 take $\pi(h)$ to be the proportion of sites in the same partition that are separated by distance h . This plot

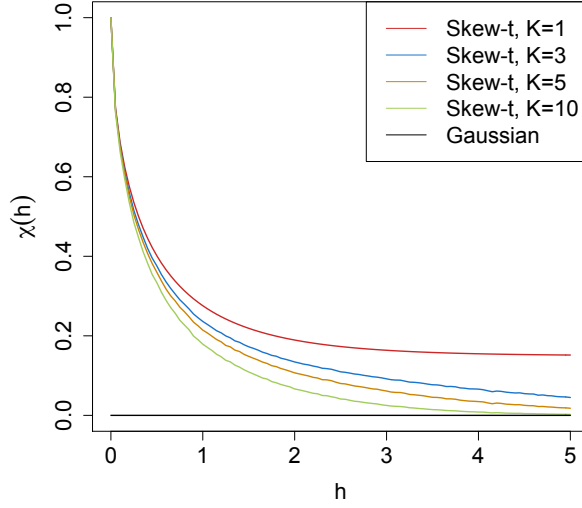


Figure 1: $\chi(h)$ for $K = 1, 3, 5$, and 10 knots as a function of distance.

demonstrates how partitioning helps to reduce extremal dependence as h increases.

3.3 Extension to space-time data

When using daily measurements, the assumption of temporal independence is inappropriate. There are several places where temporal dependence could be incorporated in the model, including the residual $v_t(\mathbf{s})$. However, we choose to allow for temporal dependence in the \mathbf{w} , z , and σ terms because these terms dictate the tail behavior which is our primary focus. In this section, we extend (6) to the spatiotemporal setting. Let

$$Y_t(\mathbf{s}) = \mathbf{X}_t(\mathbf{s})^T \boldsymbol{\beta} + \lambda \sigma_t(\mathbf{s}) |z_t(\mathbf{s})| + \sigma_t(\mathbf{s}) v_t(\mathbf{s}), \quad (10)$$

where $t \in \{1, \dots, T\}$ denotes the day of each observation. Let $\mathbf{w}_{tk} = (w_{tk1}, w_{tk2})$ be a spatial knot on day t , and let w_{t1}, \dots, w_{tK} be a collection of spatial knots on day t . As in section 3.2, these knots define a daily

113 partition P_{t1}, \dots, P_{tK} , and for $\mathbf{s} \in P_{tk}$,

$$z_t(\mathbf{s}) = z_{tk} \quad \text{and} \quad \sigma_t(\mathbf{s}) = \sigma_{tk}. \quad (11)$$

114 We use an AR(1) time series model for w_{tk} , z_{tk} , and σ_{tk} . The time series model must be specified after
 115 a transformation to preserve the skew- t process at each time point. For each time-varying parameter, we
 116 transform to obtain a standard normal marginal distribution, place a Gaussian prior with autocorrelation on
 117 the transformed parameter, and then transform back to obtain the marginal distribution required to preserve
 118 the skew- t process. We first transform the spatial knots from \mathcal{D} to \mathcal{R}^2 as follows. Let

$$w_{tki}^* = \Phi^{-1} \left[\frac{w_{tki} - \min(\mathbf{s}_i)}{\max(\mathbf{s}_i) - \min(\mathbf{s}_i)} \right], \quad i = 1, 2 \quad (12)$$

119 where Φ is a univariate standard normal density function, and $\mathbf{s}_i = [s_{1i}, \dots, s_{ni}]$. Then the transformed
 120 knots $\mathbf{w}_{tk}^* \in \mathcal{R}^2$. We use a copula on $\sigma_t^2(\mathbf{s})$ to ensure that the marginal distributions of $\sigma_t^2(\mathbf{s})$ are inverse
 121 gamma. Let

$$\sigma_t^{2*}(\mathbf{s}) = \Phi^{-1} \{ \text{IG}[\sigma_t^2(\mathbf{s})] \} \quad (13)$$

122 where IG is defined as before. We also use a copula on $z_t(\mathbf{s})$ to ensure that the marginal distributions of
 123 $z_t(\mathbf{s})$ are half-normal. Let

$$z_t^*(\mathbf{s}) = \Phi^{-1} \{ \text{HN}[z_t(\mathbf{s})] \} \quad (14)$$

124 where HN is the distribution function of a half-normal random variable. The AR(1) process for each tail

parameter is $\mathbf{w}_{1k}^* \sim N_w(0, 1)$, $z_{1k}^* \sim N(0, \sigma_{1k}^2)$, $\sigma_{1k}^{2*} \sim N(0, 1)$, and for $t > 1$ the time series is modeled as

$$\mathbf{w}_{tk}^* | \mathbf{w}_{t-1,k}^* \sim N_2 [\phi_w \mathbf{w}_{t-1,k}^*, (1 - \phi_w^2)] \quad (15)$$

$$z_{tk}^* | z_{t-1,k}^* \sim N [\phi_z z_{t-1,k}^*, \sigma_{tk}^2 (1 - \phi_z^2)] \quad (16)$$

$$\sigma_{tk}^{2*} | \sigma_{t-1,k}^{2*} \sim N [\phi_\sigma \sigma_{t-1,k}^{2*}, (1 - \phi_\sigma^2)] \quad (17)$$

where $|\phi_w|, |\phi_z|, |\phi_\sigma| < 1$. These are stationary time series models with marginal distributions $\mathbf{w}_k^* \sim N_2(0, 1)$, $z_k^* \sim N(0, \sigma_k^2)$, and $\sigma_k^{2*} \sim N(0, 1)$. After transformation back to the original space, $\mathbf{w}_{tk} \sim \text{Unif}(\mathcal{D})$, $z_{tk} \sim HN(0, \sigma_{tk}^2)$ $\sigma_{tk}^2 \sim \text{IG}(a, b)$. For each day, the model is identical to the spatial-only model in (6) by construction.

4 Computation

First, we impute values below the threshold. Then, we update model parameters, Θ , using Metropolis Hastings or Gibbs sampling when appropriate. Finally, we make spatial predictions using conditional multivariate normal results and the fact that the distribution of $Y_t(\mathbf{s}) \mid \Theta, z(\mathbf{s})$ is the usual multivariate normal distribution with a Matérn spatial covariance structure.

We can use Gibbs sampling to update $Y_t(\mathbf{s})$ for censored observations that are below the threshold T . After conditioning on λ , $z_t(\mathbf{s})$ and non-censored observations, $Y_t(\mathbf{s})$ has truncated normal full conditionals. So we sample $Y_t(\mathbf{s}) \sim N_{(-\infty, T)}(\mathbf{X}_t^T(\mathbf{s})\beta + \lambda|z_t(\mathbf{s})|, \Sigma)$. After imputing the censored observations, we update the model parameters. To update the model parameters, we use standard Gibbs updates for parameters when possible. In the case Gibbs sampling is not possible, parameters are updated using a random-walk Metropolis Hastings algorithm. See Appendices A.1 and A.2 for details regarding the MCMC. The final step of the computation is to use Bayesian Kriging to generate a predictive distribution for $Y_t(\mathbf{s}^*)$ at prediction

location \mathbf{s}^* . This step is similar to the imputation for censored observations except that the full conditionals are no longer truncated at T .

4.1 Hierarchical model

Conditioned on $z_{tk}(\mathbf{s})$, $\sigma_{tk}^2(\mathbf{s})$, and P_{tk} , the marginal distributions are Gaussian and the joint distribution multivariate Gaussian. However, we do not fix the partitions, they are treated as unknown and updated in the MCMC. We model this with a Bayesian hierarchical model as follows. Let $\mathbf{w}_{t1}, \dots, \mathbf{w}_{tK}$ be a set of daily spatial knots in a spatial domain of interest, \mathcal{D} , and P_{tk} as defined in (7). Then

$$Y_t(\mathbf{s}) \mid z_t(\mathbf{s}), \sigma_t^2(\mathbf{s}), P_{tk}, \Theta = \mathbf{X}_t(\mathbf{s})^T \beta + \lambda |z_t(\mathbf{s})| + \sigma_t(\mathbf{s}) v_t(\mathbf{s}) \quad (18)$$

$$z_t(\mathbf{s}) = z_{tk} \text{ if } \mathbf{s} \in P_{tk}$$

$$\sigma_t^2(\mathbf{s}) = \sigma_{tk}^2 \text{ if } \mathbf{s} \in P_{tk}$$

$$\lambda = \lambda_1 \lambda_2$$

$$\lambda_1 = \begin{cases} +1 & \text{w.p. } 0.5 \\ -1 & \text{w.p. } 0.5 \end{cases}$$

$$\lambda_2^2 \sim IG(a, b)$$

$$v_t(\mathbf{s}) \mid \Theta \sim \text{Matérn}(0, \Sigma)$$

$$z_{tk}^* \mid z_{t-1,k}^*, \sigma_{tk}^2 \sim N(\phi_z z_{t-1,k}^*, \sigma_{tk}^2 (1 - \phi_z^2))$$

$$\sigma_{tk}^{2*} \mid \sigma_{t-1,k}^{2*} \sim N(\phi_\sigma \sigma_{t-1,k}^{2*}, (1 - \phi_\sigma^2))$$

$$\mathbf{w}_{tk}^* \mid \mathbf{w}_{t-1,k}^* \sim N_2(\phi_w \mathbf{w}_{t-1,k}^*, (1 - \phi_w^2))$$

where $\Theta = \{\rho, \nu, \gamma, \lambda, \beta\}$, and Σ is a Matérn covariance matrix as described in Section 2.1. We parameterize $\lambda = \lambda_1 \lambda_2$ to help with convergence in the MCMC.

5 Simulation study

In this section, we conduct a simulation study to investigate how the number of partitions and the level of thresholding impact the accuracy of predictions made by the model.

5.1 Design

For all simulation designs, we generate data from the model in Section 3.2 using $n_s = 144$ sites and $n_t = 50$ independent days. The sites are generated $\text{Uniform}([0, 10] \times [0, 10])$. We generate data from 5 different simulation designs:

1. Gaussian marginal, $K = 1$ knot
2. Skew- t marginal, $K = 1$ knots
3. Skew- t marginal, $K = 5$ knots
4. Max-stable
5. Transformation below $T = q(0.80)$

In the first three designs, the $v_t(\mathbf{s})$ terms are generated using a Matérn covariance with smoothness parameter $\nu = 0.5$ and spatial range $\rho = 1$. For the covariance matrices in designs 1 – 3, the proportion of the variance accounted for by the spatial variation is $\gamma = 0.9$ while the proportion of the variance accounted for by the nugget effect is 0.1. In the first design, $\sigma^2 = 2$ is used for all days which results in a Gaussian distribution. For designs 2 and 3, $\sigma_{tk}^2 \stackrel{iid}{\sim} \text{IG}(3, 8)$ to give a t distribution with 6 degrees of freedom. For designs 1, we set $\lambda = 0$. For designs 2 and 3, $\lambda = 3$ was used as to simulate moderate skewness, and the z_t are generated as described in (8). In the fourth design, we generate from a spatial max-stable distribution (Reich and Shaby, 2012). In this design, data have marginal distributions that follow a generalized extreme value distribution with parameters $\mu = 1, \sigma = 1, \xi = 0.2$. In this model, a random effect is used to induce spatial dependence using 144 spatial knots on a regular lattice in the square $[1, 9] \times [1, 9]$. For this setting, we set $\gamma = 0.5$. In

173 the final design, we generate \tilde{y} using the setting from design two, and then consider the data

$$y = \begin{cases} \tilde{y}, & \tilde{y} > T \\ T \exp\{\tilde{y} - T\}, & \tilde{y} \leq T \end{cases} \quad (19)$$

174 where $T = q(0.80)$ is the 80th sample quantile of the data. In all five designs, the mean $\mathbf{X}^T \boldsymbol{\beta} = 10$ is
 175 assumed to be constant across space.

176 $M = 50$ data sets are generated for each design. For each data set we fit the data using five models

- 177 1. Gaussian marginal, $K = 1$ knots
- 178 2. Skew- t marginal, $K = 1$ knots, $T = -\infty$
- 179 3. Symmetric- t marginal, $K = 1$ knots, $T = q(0.80)$
- 180 4. Skew- t marginal, $K = 5$ knots, $T = -\infty$
- 181 5. Symmetric- t marginal, $K = 5$ knots, $T = q(0.80)$

182 where $q(0.80)$ is the 80th sample quantile of the data. The design matrix \mathbf{X} includes an intercept with a first-
 183 order spatial trend with priors of $\beta_{\text{int}}, \beta_{\text{lat}}, \beta_{\text{long}}, \overset{iid}{\sim} \text{N}(0, 10)$. The spatial covariance parameters have priors
 184 $\log(\nu) \sim \text{N}(-1.2, 1)$, $\gamma \sim \text{Unif}(0, 1)$, $\rho \sim \text{Unif}(15)$. The skewness parameter has prior $\lambda_2 \sim \text{IG}(0.1, 0.1)$.
 185 The residual variance terms have priors $\sigma_t^2(\mathbf{s}) \sim \text{IG}(0.1, 0.1)$. The knots have priors $\mathbf{w} \sim \text{Unif}(\mathcal{D})$. We
 186 do not include skewness for the thresholded models because it cannot be identified with only one tail worth
 187 of data, and we do not fit the data using the max-stable methods from Reich and Shaby (2012) because of
 188 the computational constraints. Each chain of the MCMC ran for 20000 iterations with a burn-in period of
 189 10000 iterations. Parameters appear to converge properly; however, in the models with multiple partitions
 190 (i.e. models 4 and 5) it is hard to assess the convergence of \mathbf{w} , $z(\mathbf{s})$, and $\sigma^2(\mathbf{s})$ because of partition label
 191 switching throughout the MCMC.

5.2 Cross validation

Models were compared using cross validation with 100 sites used as training sites and 44 sites withheld for testing. The model was fit using the training set, and predictions were generated at the testing site locations. Because one of the primary goals of this model is to predict extreme events, we use Brier scores to select the model that best fits the data (Gneiting and Raftery, 2007). The Brier score for predicting exceedance of a threshold c is given by $[e(c) - P(c)]^2$ where $e(c) = I[y > c]$ is an indicator function indicating that a test set value, y , has exceeded the threshold, c , and $P(c)$ is the predicted probability of exceeding c . We average the Brier scores over all test sites and days. For the Brier score, a lower score indicates a better fit.

5.3 Results

We compared the Brier scores for exceeding 4 different thresholds for each dataset. The thresholds used for the Brier scores are extreme quantiles from the simulated data for $q(0.90)$, $q(0.95)$, $q(0.98)$, and $q(0.99)$. Figure 2 gives the Brier score relative to the Brier score for the Gaussian method calculated as

$$BS_{\text{rel}} = \frac{BS_{\text{method}}}{BS_{\text{Gaussian}}}. \quad (20)$$

We analyzed the results for the simulation study using a Friedman test at $\alpha = 0.05$. If the Friedman test came back with a significant results, we conducted a Wilcoxon-Nemenyi-McDonald-Thompson test to see which methods had different results. The full results for the Wilcoxon-Nemenyi-McDonald-Thompson tests are given in Appendix A.5.

Figure 2 shows that when the data come from a Gaussian process, our methods perform comparably to the Gaussian method. For data settings with skew- t marginals (settings 2 – 3), we find significant improvement over the Gaussian method. Furthermore in these data settings, we find the best performance occurs

211 when the number of knots used in the method matches the number of knots used for data generation. The
 212 non-thresholded methods tend to outperform the thresholded methods, but this is not surprising given that
 213 the data are generated directly from the model used in the method. For the max-stable data, we see that for
 214 low-extreme quantiles, the Gaussian method performs better, for more extreme quantiles, the single-partition
 215 method, both thresholded and non-thresholded, perform significantly better than the Gaussian. Finally, for
 216 setting 5, although the thresholded version of the single-partition model tends to perform the best across all
 217 of the extreme quantiles, the difference between the thresholded and non-thresholded methods is no longer
 218 significant in the more extreme quantiles.

219 **6 Data analysis**

220 To illustrate this method, we consider the daily maximum 8-hour ozone measurements for July 2005 at 1089
 221 Air Quality System (AQS) monitoring sites in the United States as the response (see Figure 3). For each
 222 site, we also have covariate information containing the estimated ozone from the Community Multi-scale
 223 Air Quality (CMAQ) modeling system. Initially, we fit a linear regression assuming a mean function of

$$\mathbf{X}_t^T(\mathbf{s})\boldsymbol{\beta} = \beta_0 + \beta_1 \cdot \text{CMAQ}_t(\mathbf{s}). \quad (21)$$

224 The data from July 10 are shown in Figure 3 along with a Q-Q plot of the residuals compared to a skew- t
 225 distribution with 10 d.f. and $\alpha = 1$. Exploratory data analysis indicates that there is dependence in the high
 226 quantile levels of the residuals beyond what we expect in the case of independence.

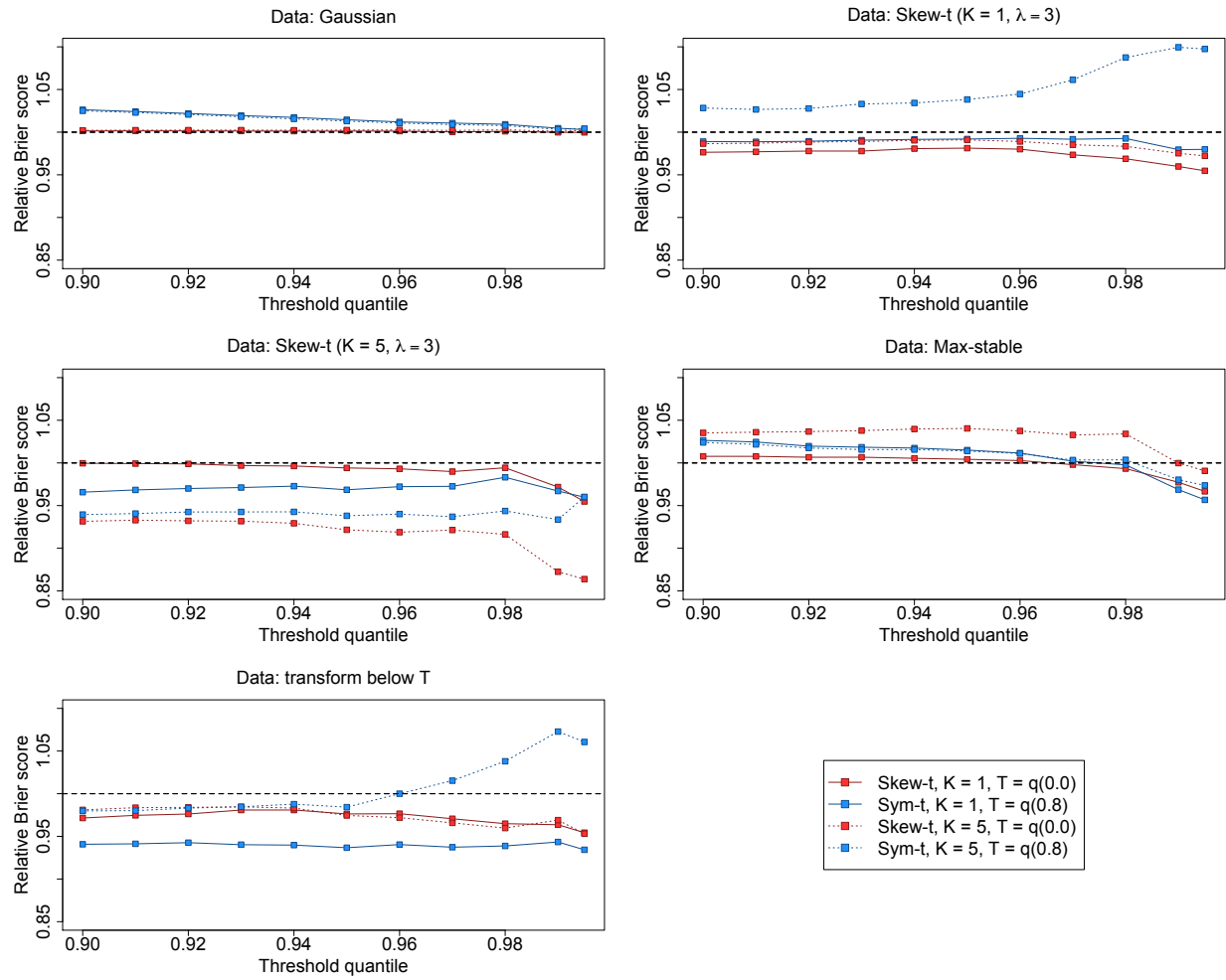


Figure 2: Brier scores relative to the Gaussian method for simulation study results. A ratio lower than 1 indicates that the method outperforms the Gaussian method.

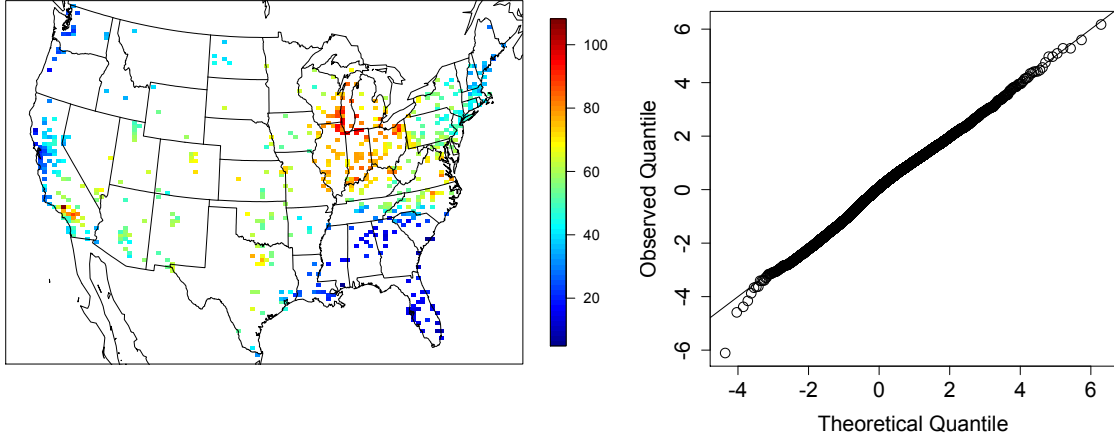


Figure 3: Ozone values on 10 July 2005 (left) Q-Q plot of the residuals (right)

6.1 Model comparisons

We fit the model using Gaussian and skew- t marginal distributions with $K = 1, 5, 6, 7, 8, 9, 10, 15$ partitions.

We choose to censor $Y(\mathbf{s})$ at $T = 0, 50$ (0.42 sample quantile), and 75 (0.92 sample quantile) ppb in order to compare results from no, moderate, and high censoring. The upper threshold of 75 ppb was used because the current air quality standard is based on exceedance of 75 ppb. As with the simulation study, for models with a threshold of $T = 75$, we use a symmetric- t marginal distribution. We also compare models with no time series to models that include the time series. Finally, as a comparison to max-stable methods, we fit the model using the hierarchical max-stable model of Reich and Shaby (2012). All methods assume the mean function given in (21). To ensure that the max-stable method runs in a reasonable amount of time, we take a stratified sample of the sites to get 800 sites and consider this our new dataset. We conduct two-fold cross validation using 400 training sites and 400 validation sites as described in Section 5.2

Each chain of the MCMC ran for 30000 iterations with a burn-in period of 25000 iterations. Parameters appear to converge properly; however, as before, for models with multiple partitions it is hard to assess the convergence of \mathbf{w} , $z(\mathbf{s})$, and $\sigma^2(\mathbf{s})$ because of partition label switching throughout the MCMC. For each

model, Brier scores were averaged over all sites and days to obtain a single Brier score for each dataset. At a particular threshold or quantile level, the model that fits the best is the one with the lowest score. We then compute the relative (to Gaussian) Brier scores (see Section 5.3) to compare each model.

6.2 Results

The results suggest that the skew- t , thresholded, partitioned, and time series models all give an improvement in predictions over the Gaussian model, whereas the max-stable method results in relative Brier scores between 1.07 and 1.15 indicating poorer performance than the Gaussian model. The plots in Figure 4 show the relative Brier scores for time-series and non-time-series models, using $K = 1, 7$, and 15 knots at thresholds $T = 0, 50$, and 75 ppb. Most of the models perform similarly across all the Brier scores; however, for single-partition models without thresholding, performance tends to diminish in the extreme quantiles. The results also suggest that thresholding improves performance for estimates in the extreme quantiles. Both plots have similar features suggesting that most settings do reasonably well. In particular, for all extreme quantiles, selecting a moderate number of knots (e.g. $K = 5, \dots, 10$) tends to give the best results. Table 1 shows the best two models for selected extreme quantiles.

We illustrate the predictive capability of our model in Figure 5 by plotting the 99th quantile of the posterior predictive density for July in South Carolina and Georgia. We fit the model using four methods, two reference and two that performed better. These four methods are

[1.]Gaussian (reference) Skew- t , $K = 1$ knot, $T = 0$, no time series (reference) Skew- t , $K = 5$ knots, $T = 50$, no time series (comparison) Symmetric- t , $K = 10$ knots, $T = 75$, time series (comparison)

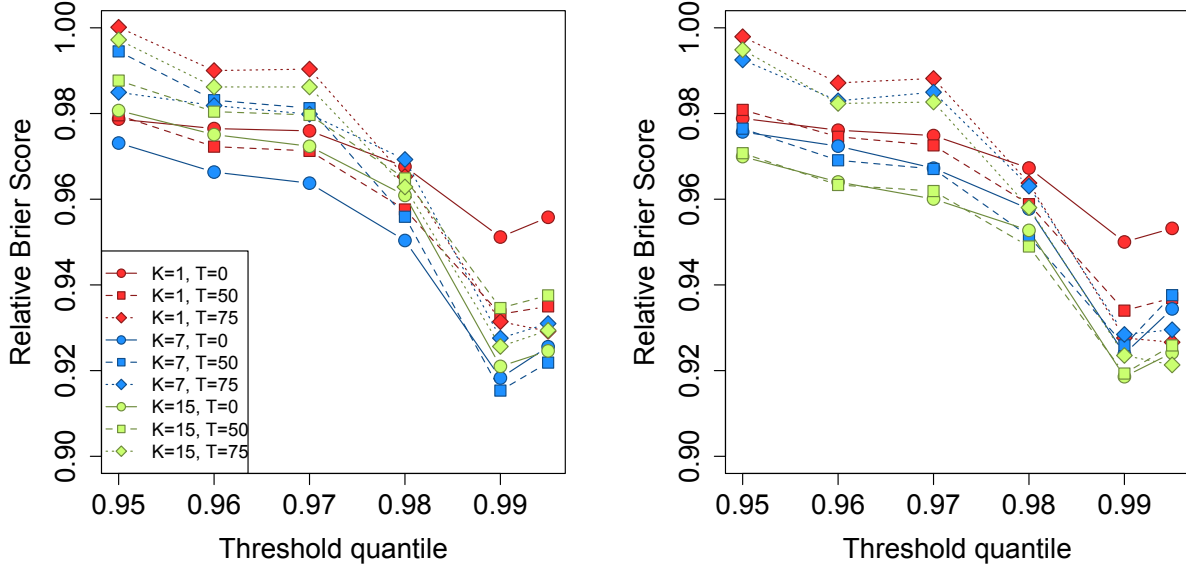


Figure 4: Relative Brier scores for time-series models (left) and non-time-series models (right). Relative brier score for the max-stable model is between 1.07 and 1.15

Table 1: Top two performing models for ozone analysis at extreme quantiles with Relative Brier score

	1st				2nd			
$q(0.90)$	No time series	$K = 7$	$T = 0$	BS: 0.980	No time series	$K = 9$	$T = 0$	BS: 0.980
$q(0.95)$	No time series	$K = 15$	$T = 50$	BS: 0.970	No time series	$K = 9$	$T = 50$	BS: 0.970
$q(0.98)$	No time series	$K = 5$	$T = 50$	BS: 0.945	No time series	$K = 10$	$T = 50$	BS: 0.946
$q(0.99)$	Time series	$K = 10$	$T = 75$	BS: 0.912	Time series	$K = 6$	$T = 75$	BS: 0.913
$q(0.995)$	Time series	$K = 6$	$T = 75$	BS: 0.917	Time series	$K = 10$	$T = 75$	BS: 0.918

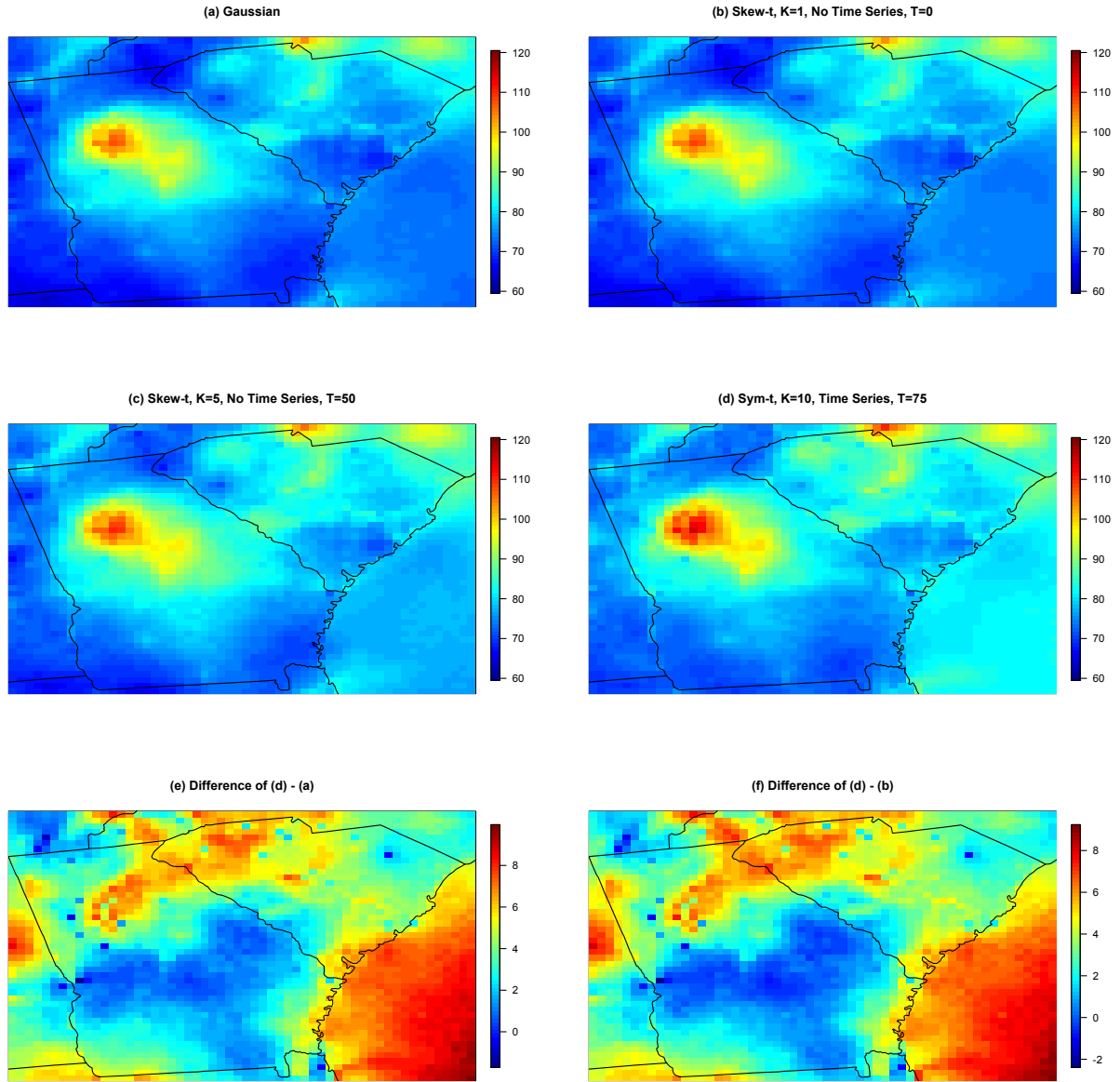


Figure 5: (a) – (d) give give the posterior predictive $\hat{q}(0.99)$ for the month of July under four different models, (e) gives the difference between $\hat{q}(0.99)$ in plots (d) and (a), (f) gives the difference between $\hat{q}(0.99)$ in plots (d) and (b).

7 Discussion

In this paper we propose a new approach for spatitemporal modeling of extreme values. The proposed model gives flexible tail behavior, demonstrates asymptotic dependence for observations at sites that are near to one another, and has computation on the order of Gaussian models for large space-time datasets. In the simulation study, we demonstrate that this model shows statistically significant improvements over a naïve Gaussian approach. In both the simulation study, and the application to ozone data, we find that incorporating a partition in the model improves extreme prediction. Furthermore the results from the data analysis suggest that thresholding can improve performance when predicting in the extreme tails of the data.

This model presents new avenues for future research. One possibility is the implementation of a different partition structure. We choose to define the random effects for a site by using an indicator function based on closeness to a knot. However, this indicator function could be replaced by kernel function that would allow for multiple knots to impact each site, with the weight of each knot to be determined by some characteristic such as distance. Another area that should be explored is the temporal dependence in the model. Instead of implementing a time series on the random effects, a three-dimensional covariance structure on the residuals could be implemented to address temporal dependence. Finally, we acknowledge that by specifying the number of knots, we may be underestimating the uncertainty in the model. This could be incorporated by treating the number of knots as a model parameter instead of fixing it to be a specific value.

Acknowledgments

A Appendices

A.1 MCMC details

The MCMC sampling for the model 4.1 is done using R (<http://www.r-project.org>). Whenever possible, we select conjugate priors (see Appendix A.2); however, for some of the parameters, no conjugate prior distributions exist. When no conjugate prior distribution exists, we use a random walk Metropolis Hastings update step. In each Metropolis Hastings update, we tune the algorithm to give acceptance rates near 0.40.

Spatial knot locations

For each day, we update the spatial knot locations, $\mathbf{w}_1, \dots, \mathbf{w}_K$, using a Metropolis Hastings block update. Because the spatial domain is bounded, we generate candidate knots using the transformed knots $\mathbf{w}_1^*, \dots, \mathbf{w}_K^*$ (see section 3.3) and a random walk bivariate Gaussian candidate distribution

$$\mathbf{w}_k^{*(c)} \sim N(\mathbf{w}_k^{*(r-1)}, s^2 I_2)$$

where $\mathbf{w}_k^{*(r-1)}$ is the location for the transformed knot at MCMC iteration $r - 1$, s is a tuning parameter, and I_2 is an identity matrix. After candidates have been generated for all K knots, the acceptance ratio is

$$R = \left\{ \frac{l[Y_t(\mathbf{s}|\mathbf{w}_1^{(c)}, \dots, \mathbf{w}_K^{(c)}, \dots)]}{l[Y_t(\mathbf{s}|\mathbf{w}_1^{(r-1)}, \dots, \mathbf{w}_K^{(r-1)}, \dots)]} \right\} \times \left\{ \frac{\prod_{k=1}^K \phi(\mathbf{w}_k^{(c)})}{\prod_{k=1}^K \phi(\mathbf{w}_k^{(r-1)})} \right\} \times \left\{ \frac{\prod_{k=1}^K p(\mathbf{w}_k^{*(c)})}{\prod_{k=1}^K p(\mathbf{w}_k^{*(r-1)})} \right\}$$

where l is the likelihood given in (18), and $p(\cdot)$ is the prior either taken from the time series given in (3.3) or assumed to be uniform over \mathcal{D} . The candidate knots are accepted with probability $\min\{R, 1\}$.

Spatial random effects

If there is no temporal dependence amongst the observations, we use a Gibbs update for z_{tk} , and the posterior distribution is given in A.2. If there is temporal dependence amongst the observations, then we update z_{tk} using a Metropolis Hastings update. Because this model uses $|z_{tk}|$, we generate candidate random effects using the z_{tk}^* (see Section 3.3) and a random walk Gaussian candidate distribution

$$z_{tk}^{*(c)} \sim \mathbf{N}(z_{tk}^{*(r-1)}, s^2)$$

where $z_{tk}^{*(r-1)}$ is the value at MCMC iteration $r - 1$, and s is a tuning parameter. The acceptance ratio is

$$R = \left\{ \frac{l[Y_t(\mathbf{s})|z_{tk}^{(c)}, \dots]}{l[Y_t(\mathbf{s})|z_{tk}^{(r-1)}]} \right\} \times \left\{ \frac{p[z_{tk}^{(c)}]}{p[z_{tk}^{(r-1)}]} \right\}$$

where $p[\cdot]$ is the prior taken from the time series given in Section 3.3. The candidate is accepted with probability $\min\{R, 1\}$.

Variance terms

When there is more than one site in a partition, then we update σ_{tk}^2 using a Metropolis Hastings update. First, we generate a candidate for σ_{tk}^2 using an $\text{IG}(a^*/s, b^*/s)$ candidate distribution in an independence Metropolis Hastings update where $a^* = (n_{tk} + 1)/2 + a$, $b^* = [Y_{tk}^T \Sigma_{tk}^{-1} Y_{tk} + z_{tk}^2]/2 + b$, n_{tk} is the number of sites in partition k on day t , and Y_{tk} and Σ_{tk}^{-1} are the observations and precision matrix for partition k on day t . The acceptance ratio is

$$R = \left\{ \frac{l[Y_t(\mathbf{s})|\sigma_{tk}^{2(c)}, \dots]}{l[Y_t(\mathbf{s})|\sigma_{tk}^{2(r-1)}]} \right\} \times \left\{ \frac{l[z_{tk}|\sigma_{tk}^{2(c)}, \dots]}{l[z_{tk}|\sigma_{tk}^{2(r-1)}, \dots]} \right\} \times \left\{ \frac{p[\sigma_{tk}^{2(c)}]}{p[\sigma_{tk}^{2(r-1)}]} \right\} \times \left\{ \frac{c[\sigma_{tk}^{2(r-1)}]}{c[\sigma_{tk}^{2(c)}]} \right\}$$

306 where $p[\cdot]$ is the prior either taken from the time series given in Section 3.3 or assumed to be $\text{IG}(a, b)$, and
 307 $c[\cdot]$ is the candidate distribution. The candidate is accepted with probability $\min\{R, 1\}$.

308 **Spatial covariance parameters**

309 We update the three spatial covariance parameters, $\log(\rho)$, $\log(\nu)$, γ , using a Metropolis Hastings block
 310 update step. First, we generate a candidate using a random walk Gaussian candidate distribution

$$\log(\rho)^{(c)} \sim \text{N}(\log(\rho)^{(r-1)}, s^2)$$

311 where $\log(\rho)^{(r-1)}$ is the value at MCMC iteration $r - 1$, and s is a tuning parameter. Candidates are
 312 generated for $\log(\nu)$ and γ in a similar fashion. The acceptance ratio is

$$R = \left\{ \frac{\prod_{t=1}^T l[Y_t(\mathbf{s}) | \rho^{(c)}, \nu^{(c)}, \gamma^{(c)}, \dots]}{\prod_{t=1}^T l[Y_t(\mathbf{s}) | \rho^{(r-1)}, \nu^{(r-1)}, \gamma^{(r-1)}, \dots]} \right\} \times \left\{ \frac{p[\rho^{(c)}]}{p[\rho^{(r-1)}]} \right\} \times \left\{ \frac{p[\nu^{(c)}]}{p[\nu^{(r-1)}]} \right\} \times \left\{ \frac{p[\gamma^{(c)}]}{p[\gamma^{(r-1)}]} \right\}.$$

313 All three candidates are accepted with probability $\min\{R, 1\}$.

314 **A.2 Posterior distributions**

315 **Conditional posterior of $z_{tk} \mid \dots$**

316 If knots are independent over days, then the conditional posterior distribution of $|z_{tk}|$ is conjugate. For
 317 simplicity, drop the subscript t , let $\tilde{z}_{tk} = |z_{tk}|$, and define

$$R(\mathbf{s}) = \begin{cases} Y(\mathbf{s}) - X(\mathbf{s})\beta & s \in P_l \\ Y(\mathbf{s}) - X(\mathbf{s})\beta - \lambda \tilde{z}(\mathbf{s}) & s \notin P_l \end{cases}$$

318 Let

$R_1 = \text{the vector of } R(\mathbf{s}) \text{ for } s \in P_l$

$R_2 = \text{the vector of } R(\mathbf{s}) \text{ for } s \notin P_l$

$$\Omega = \Sigma^{-1}.$$

319 Then

$$\begin{aligned} \pi(z_l | \dots) &\propto \exp \left\{ -\frac{1}{2} \left[\begin{pmatrix} R_1 - \lambda \tilde{z}_l \mathbf{1} \\ R_2 \end{pmatrix}^T \begin{pmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{21} & \Omega_{22} \end{pmatrix} \begin{pmatrix} R_1 - \lambda \tilde{z}_l \mathbf{1} \\ R_2 \end{pmatrix} + \frac{\tilde{z}_l^2}{\sigma_l^2} \right] \right\} I(z_l > 0) \\ &\propto \exp \left\{ -\frac{1}{2} [\Lambda_l \tilde{z}_l^2 - 2\mu_l \tilde{z}_l] \right\} \end{aligned}$$

320 where

$$\mu_l = \lambda(R_1^T \Omega_{11} + R_2^T \Omega_{21}) \mathbf{1}$$

$$\Lambda_l = \lambda^2 \mathbf{1}^T \Omega_{11} \mathbf{1} + \frac{1}{\sigma_l^2}.$$

321 Then $\tilde{Z}_l | \dots \sim N_{(0,\infty)}(\Lambda_l^{-1} \mu_l, \Lambda_l^{-1})$

322 **Conditional posterior of β | ...**

323 Let $\beta \sim N_p(0, \Lambda_0)$ where Λ_0 is a precision matrix. Then

$$\begin{aligned}\pi(\beta \mid \dots) &\propto \exp \left\{ -\frac{1}{2} \beta^T \Lambda_0 \beta - \frac{1}{2} \sum_{t=1}^T [\mathbf{Y}_t - X_t \beta - \lambda |z_t|]^T \Omega [\mathbf{Y}_t - X_t \beta - \lambda |z_t|] \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \left[\beta^T \Lambda_\beta \beta - 2 \sum_{t=1}^T [\beta^T X_t^T \Omega (\mathbf{Y}_t - \lambda |z_t|)] \right] \right\} \\ &\propto N(\Lambda_\beta^{-1} \mu_\beta, \Lambda_\beta^{-1})\end{aligned}$$

324 where

$$\begin{aligned}\mu_\beta &= \sum_{t=1}^T [X_t^T \Omega (\mathbf{Y}_t - \lambda |z_t|)] \\ \Lambda_\beta &= \Lambda_0 + \sum_{t=1}^T X_t^T \Omega X_t.\end{aligned}$$

325 **Conditional posterior of σ^2 | ...**

326 In the case where $L = 1$ and temporal dependence is negligible, then σ^2 has a conjugate posterior distribu-

327 tion. Let $\sigma_t^2 \stackrel{iid}{\sim} \text{IG}(\alpha_0, \beta_0)$. For simplicity, drop the subscript t . Then

$$\begin{aligned}\pi(\sigma^2 \mid \dots) &\propto (\sigma^2)^{-\alpha_0 - 1/2 - n/2 - 1} \exp \left\{ -\frac{\beta_0}{\sigma^2} - \frac{|z|^2}{2\sigma^2} - \frac{(\mathbf{Y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\mu})}{2\sigma^2} \right\} \\ &\propto (\sigma^2)^{-\alpha_0 - 1/2 - n/2 - 1} \exp \left\{ -\frac{1}{\sigma^2} \left[\beta_0 + \frac{|z|^2}{2} + \frac{1}{2} (\mathbf{Y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\mu}) \right] \right\} \\ &\propto \text{IG}(\alpha^*, \beta^*)\end{aligned}$$

328 where

$$\alpha^* = \alpha_0 + \frac{1}{2} + \frac{n}{2}$$

$$\beta^* = \beta_0 + \frac{|z|^2}{2} + \frac{1}{2}(\mathbf{Y} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{Y} - \boldsymbol{\mu}).$$

329 In the case that $L > 1$, a random walk Metropolis Hastings step will be used to update σ_{lt}^2 .

330 **Conditional posterior of $\lambda \mid \dots$**

331 For convergence purposes we model $\lambda = \lambda_1 \lambda_2$ where

$$\lambda_1 = \begin{cases} +1 & \text{w.p.0.5} \\ -1 & \text{w.p.0.5} \end{cases} \quad (22)$$

$$\lambda_2^2 \sim IG(\alpha_\lambda, \beta_\lambda). \quad (23)$$

$$(24)$$

332 Then

$$\pi(\lambda_2 \mid \dots) \propto \lambda_2^{2(-\alpha_\lambda - 1)} \exp \left\{ -\frac{\beta_\lambda}{\lambda_2^2} \right\} \prod_{t=1}^T \prod_{k=1}^K \frac{1}{\lambda_2} \exp \left\{ -\frac{z_{tk}^2}{2\lambda_2^2 \sigma_{tk}^2} \right\}$$

$$\propto \lambda_2^{2(-\alpha_\lambda - kt - 1)} \exp \left\{ -\frac{1}{\lambda_2^2} \left[\beta_\lambda + \frac{z^2}{2\sigma_{tk}^2} \right] \right\}$$

333 Then $\lambda_2 \mid \dots \sim IG(\alpha_\lambda + kt, \beta_\lambda + \frac{z^2}{2\sigma_{tk}^2})$

334 **A.3 Proof that $\lim_{h \rightarrow \infty} \pi(h) = 0$**

335 Let $N(A)$ be the number of knots in A , the area between sites \mathbf{s}_1 and \mathbf{s}_2 . Consider a spatial Poisson process
 336 with intensity $\mu(A)$. So,

$$P[N(A) = k] = \frac{\mu(A)^k \exp\{-\mu(A)\}}{k!}.$$

337 Then for any finite k , $\lim_{h \rightarrow \infty} P[N(A) = k] = 0$ because $\lim_{h \rightarrow \infty} \mu(A) = \infty$. With each additional knot
 338 in A , the chance that \mathbf{s}_1 and \mathbf{s}_2 will be in the same partition will decrease, because partition membership
 339 is defined by the closest knot to a site. Therefore, $\lim_{h \rightarrow \infty} \pi(h) = 0$.

340 **A.4 Skew- t distribution**

341 **Univariate extended skew- t distribution**

342 We say that Y follow a univariate extended skew- t distribution with location $\xi \in \mathcal{R}$, scale $\omega > 0$, skew
 343 parameter $\alpha \in \mathcal{R}$, extended parameter $\tau \in \mathcal{R}$, and degrees of freedom ν if has distribution function

$$f_{\text{EST}}(y) = \omega^{-1} \frac{f_T(z; \nu)}{F_T(\tau/\sqrt{1 + \alpha^2}; \nu)} F_T \left[(\alpha z + \tau) \sqrt{\frac{\nu + 1}{\nu + z^2}}; 0, 1, \nu + 1 \right] \quad (25)$$

344 where $f_T(t; \nu)$ is a univariate Student's t with ν degrees of freedom, $F_T(t; \nu) = P(T < t)$, and $z = (y - \xi)/\omega$.

345 In the case that $\tau = 0$, then Y follows a univariate skew- t distribution.

346 Multivariate skew- t distribution

347 If $\mathbf{Z} \sim \text{ST}_d(0, \bar{\boldsymbol{\Omega}}, \boldsymbol{\alpha}, \eta)$ is a d -dimensional skew- t distribution, and $\mathbf{Y} = \boldsymbol{\xi} + \boldsymbol{\omega}\mathbf{Z}$, where $\boldsymbol{\omega} = \text{diag}(\omega_1, \dots, \omega_d)$,

348 then the density of \mathbf{Y} at \mathbf{y} is

$$f_{\mathbf{y}}(\mathbf{y}) = \det(\boldsymbol{\omega})^{-1} f_{\mathbf{z}}(\mathbf{z}) \quad (26)$$

349 where

$$f_{\mathbf{z}}(\mathbf{z}) = 2t_d(\mathbf{z}; \bar{\boldsymbol{\Omega}}, \eta) T \left[\boldsymbol{\alpha}^T \mathbf{z} \sqrt{\frac{\eta + d}{\nu + Q(\mathbf{z})}}; \eta + d \right] \quad (27)$$

$$\mathbf{z} = \boldsymbol{\omega}^{-1}(\mathbf{y} - \boldsymbol{\xi}) \quad (28)$$

350 where $t_d(\mathbf{z}; \bar{\boldsymbol{\Omega}}, \eta)$ is a d -dimensional Student's t -distribution with scale matrix $\bar{\boldsymbol{\Omega}}$ and degrees of freedom

351 η , $Q(\mathbf{z}) = \mathbf{z}^T \bar{\boldsymbol{\Omega}}^{-1} \mathbf{z}$ and $T(\cdot; \eta)$ denotes the univariate Student's t distribution function with η degrees of

352 freedom (Azzalini and Capitanio, 2013).

353 Extremal dependence

354 For a bivariate skew- t random variable $\mathbf{Y} = [Y(\mathbf{s}), Y(\mathbf{t})]^T$, the $\chi(h)$ statistic (Padoan, 2011) is given by

$$\chi(h) = \bar{F}_{\text{EST}} \left\{ \frac{[x_1^{1/\eta} - \varrho(h)]\sqrt{\eta + 1}}{\sqrt{1 - \varrho(h)^2}}; 0, 1, \alpha_1, \tau_1, \eta + 1 \right\} + \bar{F}_{\text{EST}} \left\{ \frac{[x_2^{1/\eta} - \varrho(h)]\sqrt{\eta + 1}}{\sqrt{1 - \varrho(h)^2}}; 0, 1, \alpha_2, \tau_2, \eta + 1 \right\}, \quad (29)$$

355 where \bar{F}_{EST} is the univariate survival extended skew- t function with zero location and unit scale, $\varrho(h) = \text{cor}(y_1, y_2)$,

356 $\alpha_j = \alpha_i \sqrt{1 - \varrho^2}$, $\tau_j = \sqrt{\eta + 1}(\alpha_j + \alpha_i \varrho)$, and $x_j = F_T(\bar{\alpha}_i \sqrt{\eta + 1}; 0, 1, \eta) / F_T(\bar{\alpha}_j \sqrt{\eta + 1}; 0, 1, \eta)$ with

357 $j = 1, 2$ and $i = 2, 1$ and where $\bar{\alpha}_j = (\alpha_j + \alpha_i \varrho) / \sqrt{1 + \alpha_i^2 [1 - \varrho(h)^2]}$.

358 **Proof that** $\lim_{h \rightarrow \infty} \chi(h) > 0$

359 Consider the bivariate distribution of $\mathbf{Y} = [Y(\mathbf{s}), Y(\mathbf{t})]^T$, with $\varrho(h)$ given by (3). So, $\lim_{h \rightarrow \infty} \varrho(h) = 0$.

360 Then

$$\lim_{h \rightarrow \infty} \chi(h) = \bar{F}_{\text{EST}} \left\{ \sqrt{\eta + 1}; 0, 1, \alpha_1, \tau_1, \eta + 1 \right\} + \bar{F}_{\text{EST}} \left\{ \sqrt{\eta + 1}; 0, 1, \alpha_2, \tau_2, \eta + 1 \right\}. \quad (30)$$

361 Because the extended skew- t distribution is not bounded above, for all $\bar{F}_{\text{EST}}(x) = 1 - F_{\text{EST}} > 0$ for all
 362 $x < \infty$. Therefore, for a skew- t distribution, $\lim_{h \rightarrow \infty} \chi(h) > 0$.

363 A.5 Simulation study pairwise difference results

364 The following tables show the methods that have significantly different Brier scores when using a Wilcoxon-
 365 Nemenyi-McDonald-Thompson test. In each column, different letters signify that the methods have signifi-
 366 cantly different Brier scores. For example, there is significant evidence to suggest that method 1 and method
 367 4 have different Brier scores at $q(0.90)$, whereas there is not significant evidence to suggest that method 1
 368 and method 2 have different Brier scores at $q(0.90)$. In each table group A represents the group with the
 369 lowest Brier scores. Groups are significant with a familywise error rate of $\alpha = 0.05$.

Table 2: Setting 1 – Gaussian marginal, $K = 1$ knot

	$q(0.90)$	$q(0.95)$	$q(0.98)$	$q(0.99)$
Method 1	A	A	A	A
Method 2	A B	A B	A	A
Method 3	C	C	C	B
Method 4	B	B	B	B
Method 5	C	C	C	B

Table 3: Setting 2 – Skew- t marginal, $K = 1$ knot

	$q(0.90)$	$q(0.95)$	$q(0.98)$	$q(0.99)$
Method 1	B	B	B	B
Method 2	A	A	A	A
Method 3	B	B	B	A B
Method 4	B	B	B	B
Method 5	C	C	C	C

Table 4: Setting 3 – Skew- t marginal, $K = 5$ knots

	$q(0.90)$	$q(0.95)$	$q(0.98)$	$q(0.99)$
Method 1	B	B	B	C
Method 2	B	B	B	B C
Method 3	A	A	B	C
Method 4	A	A	A	A B
Method 5	A	A	A	A

Table 5: Setting 4 – Max-stable

	$q(0.90)$	$q(0.95)$	$q(0.98)$	$q(0.99)$
Method 1	A	A	A B	C
Method 2	B	A B	A	A B
Method 3	C	B C	A B	A
Method 4	D	D	C	C
Method 5	C D	C	B	B C

Table 6: Setting 5 – Transformation below $T = q(0.80)$

	$q(0.90)$	$q(0.95)$	$q(0.98)$	$q(0.99)$
Method 1	C	B	C	C
Method 2	B	B	A B	A B
Method 3	A	A	A	A
Method 4	B C	B	B	B C
Method 5	B C	B	C	C

References

- Azzalini, A. and Capitanio, A. (2013) *The Skew-Normal and Related Families*. Institute of Mathematical Statistics Monographs. Cambridge University Press.
- Blanchet, J. and Davison, A. C. (2011) Spatial modeling of extreme snow depth. *The Annals of Applied Statistics*, **5**, 1699–1725.
- Coles, S. G. and Tawn, J. A. (1991) Modelling Extreme Multivariate Events. *Journal of the Royal Statistical Society: Series B (Methodological)*, **53**, 377–392.
- DuMouchel, W. H. (1983) Estimating the stable index α in order to measure tail thickness: a critique. *The Annals of Statistics*, **11**, 1019–1031.
- Genton, M. G. (2004) *Skew-Elliptical Distributions and Their Applications: A Journey Beyond Normality*. Statistics (Chapman & Hall/CRC). Taylor & Francis.
- Gneiting, T. and Raftery, A. E. (2007) Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, **102**, 359–378.
- Huser, R. (2013) *Statistical Modeling and Inference for Spatio-Temporal Extremes*. Ph.D. thesis.
- Kim, H.-M., Mallick, B. K. and Holmes, C. C. (2005) Analyzing Nonstationary Spatial Data Using Piecewise Gaussian Processes. *Journal of the American Statistical Association*, **100**, 653–668.
- Padoan, S. A. (2011) Multivariate extreme models based on underlying skew- and skew-normal distributions. *Journal of Multivariate Analysis*, **102**, 977–991.
- 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.

- 390 Reich, B. J. and Shaby, B. A. (2012) A hierarchical max-stable spatial model for extreme precipitation. *The*
391 *Annals of Applied Statistics*, **6**, 1430–1451.
- 392 Zhang, H. and El-Shaarawi, A. (2010) On spatial skewGaussian processes and applications. *Environmetrics*,
393 **21**, 33–47.