

# Spatial Bayesian hierarchical modeling of precipitation extremes over a large domain

Cameron Bracken, Balaji Rajagopalan, Linyin Cheng, Will Kleiber and Subhrendu Gangopadhyay

## Abstract

We propose a Bayesian hierarchical model for spatial extremes on a large domain. In the data layer a Gaussian elliptical copula having generalized extreme value (GEV) marginals is applied. Spatial dependence in the GEV parameters are captured with a latent spatial regression with spatially varying coefficients. Using a composite likelihood approach, we are able to efficiently incorporate a large precipitation dataset, which includes stations with missing data. The model is demonstrated by application to fall precipitation extremes at approximately 2600 stations covering the western United States,  $-125^{\circ}\text{E}$  to  $-100^{\circ}\text{E}$  longitude and  $30^{\circ}\text{N}$  to  $50^{\circ}\text{N}$  latitude. The hierarchical model provides GEV parameters on a  $1/8^{\text{th}}$  degree grid and consequently maps of return levels and associated uncertainty. The model results indicate that return levels vary coherently both spatially and across seasons, providing information about the space-time variations of risk of extreme precipitation in the western US, helpful for infrastructure planning.

## 1 Introduction

Engineering design of infrastructure such as flood protection, dams, and management of water supply and flood control require robust estimates of return levels and associated errors of precipitation extremes. Spatial modeling of precipitation extremes not only can capture spatial dependence between stations but also reduce the overall uncertainty in at-site return level estimates by borrowing strength across spatial locations [Cooley *et al.*, 2007]. Hierarchical Bayesian modeling of extremes precipitation was first introduced by [Cooley *et al.*, 2007] and since has been widely discussed in the literature [Cooley and Sain, 2010; Aryal *et al.*, 2010; Atyeo and Walshaw, 2012; Davison *et al.*, 2012; Ghosh and Mallick, 2011; Reich and Shaby, 2012; Sang and Gelfand, 2010, 2009; Apputhurai and Stephenson, 2013; Dyrddal *et al.*, 2014]. Hierarchical modeling is an alternative to regional frequency analysis providing gridded or pointwise estimates of return levels within a study region [Renard, 2011].

Bayesian hierarchical models for spatial extremes have typically been limited to small geographic regions that include on the order 100 stations covering areas on the order of ~~xxx~~ 100,000  $\text{km}^2$ . Large geographic regions with many stations present a computational challenge for hierarchical Bayesian models, specifically when computing the like-

likelihood of Gaussian processes (GPs), which for  $n$  data points requires ~~inverting an  $n \times n$  matrix~~ solving a linear system of  $n$  equations, an  $O(n^3)$  operation. Several approaches exist for speeding up GP likelihood computations such as low-rank approximations [Banerjee et al., 2008] ~~in which the GP is approximated at a small number of knots and~~ composite likelihood methods [Caragea and Smith, 2007] where the likelihood computation is broken into groups containing a small number of stations [Lindsay, 1988; Heagerty and Lele, 1998; Caragea and Smith, 2007] spectral methods [Fuentes, 2007], restricted likelihoods [Stein et al., 2004] and Laplace approximations [Rue et al., 2009]. The use of a composite likelihood approach is explored here because we not only wish to estimate covariance parameters but to also produce maps of return levels with small credible intervals.

Some attempts have been made to model extremes in large regions and with large datasets in a Bayesian hierarchical context. Reich and Shaby [2012] use a hierarchical max-stable model with climate model output in the east coast to examine spatially varying GEV parameters, with a max-stable process for the data dependence level. [Ghosh and Mallick, 2011] model gridded precipitation data over the entire US, for annual maxima at a 5x5 degree resolution (43 grid cells) and copula for data dependence, incorporating spatial dependence directly in a spatial model on the data, not parameters. [Cooley and Sain, 2010] and [Sang and Gelfand, 2009] model over 1000 grid cells of climate model output using spatial autoregressive models which take advantage of data on a regular lattice to simplify computations.

The research contributions of this study are as follows. A Bayesian hierarchical model is proposed which is capable of incorporating thousands of observation locations by utilizing a composite likelihood method. The GEV shape parameter is modeled spatially in order to capture the detailed behavior of extremes in the western US. In addition the model is capable of incorporating stations with missing data with little additional computational overhead. The model is applied to observed precipitation extremes in each season, providing estimated seasonal return levels for the western US.

In section 2 the general model structure is described. Section 3 describes details of the application to seasonal extreme precipitation in the western US. Results are discussed in Section 4 and Discussion and conclusions are given in ~~Section~~ Section 5.

## 2 Model structure

The joint distribution of the  $m$  ~~data~~ data-stations in each year is modeled as a realization from a Gaussian elliptical copula with generalized extreme value (GEV) distribution marginals. The copula is characterized by pairwise dependence matrix  $\Sigma$ . Spatial dependence is further captured through spatial processes on the location  $\mu(s)$ , scale  $\sigma(s)$  and  $\xi(s)$  parameters. We assume the parameters can be described through a latent spatial regres-

sion where the residual component  $w_\gamma(\mathbf{s})$  follows a mean 0, stationary, isotropic Gaussian process (GP) with covariance function  $C_\gamma(\mathbf{s}, \mathbf{s}')$  where  $\gamma$  represents any GEV parameter ( $\mu, \sigma, \xi$ ). The corresponding covariance matrix is  $C_\gamma(\boldsymbol{\theta}_\gamma) = [C_\gamma(\mathbf{s}_i, \mathbf{s}_j; \boldsymbol{\theta}_\gamma)]_{i,j=1}^m$  where  $\boldsymbol{\theta}_\gamma$  represents the covariance parameters. The first layer of the hierarchical model structure is:

$$(Y(\mathbf{s}_1, t), \dots, Y(\mathbf{s}_m, t)) \sim Gcop_m[\Sigma; \{\mu(\mathbf{s}), \sigma(\mathbf{s}), \xi(\mathbf{s})\}] \quad (1)$$

$$Y(\mathbf{s}, t) \sim \text{GEV}[\mu(\mathbf{s}), \sigma(\mathbf{s}), \xi(\mathbf{s})] \quad (2)$$

where  $Y(\mathbf{s}, t)$  is the response at site  $\mathbf{s}$  and time  $t$  and  $Gcop_m$  stands for “m-dimensional Gaussian elliptical copula” with dependence matrix  $\Sigma$ . The spatial data layer processes in each year are assumed independent and identically distributed. Alternatives to using a copula to construct the joint distribution are an assumption of conditional independence [Cooley *et al.*, 2007] and max-stability [Smith, 1990; Schlather, 2002; Cooley *et al.*, 2006; Shang *et al.*, 2011; Padoan *et al.*, 2010]. Marginally, observations are assumed to have a generalized extreme value (GEV) distribution.

The second layer of the hierarchy, also known as the process layer, involves spatial models for the GEV parameters

$$\mu(\mathbf{s}) = \beta_{\mu,0} + \mathbf{x}_\mu^T(\mathbf{s})\boldsymbol{\beta}_\mu(\mathbf{s}) + w_\mu(\mathbf{s}) \quad (3)$$

$$\sigma(\mathbf{s}) = \beta_{\sigma,0} + \mathbf{x}_\sigma^T(\mathbf{s})\boldsymbol{\beta}_\sigma(\mathbf{s}) + w_\sigma(\mathbf{s}) \quad (4)$$

$$\xi(\mathbf{s}) = \beta_{\xi,0} + \mathbf{x}_\xi^T(\mathbf{s})\boldsymbol{\beta}_\xi(\mathbf{s}) + w_\xi(\mathbf{s}) \quad (5)$$

Where  $\beta_{\gamma,0}$  are spatially independent intercept terms,  $\mathbf{x}_\gamma^T(\mathbf{s}_i)$  is a vector of  $p$  spatially varying predictors and  $\boldsymbol{\beta}_\gamma(\mathbf{s})$  is a vector of  $p$  spatially varying regression coefficients. Covariates will be discussed in Section 3.2.

The shape parameter  $\xi$  is notoriously difficult to estimate, its value determining the support of the GEV distribution. Positive values of  $\xi$  indicate a lower bound to the distribution, negative values indicate an upper bound and zero indicates no bounds. In many studies,  $\xi$  is modeled as a single value per study area or per region within the study area [Cooley *et al.*, 2007; Renard, 2011; Atyeo and Walshaw, 2012; Apputhurai and Stephenson, 2013]. As in [Cooley and Sain, 2010], we cannot assume that this parameter is constant over the large study region and so it is modeled spatially along with the other GEV parameters.

For large regions we cannot assume that a constant spatial regression holds for the entire domain and thus must introduce spatial variation in the regression coefficients. The third layer of the hierarchy involves a spatial model for these regression coefficients

$$\beta_\mu(\mathbf{s}) = \sum_{i=1}^k c_i^\mu \eta_i^\mu(\mathbf{s}) \quad (6)$$

$$\beta_\sigma(\mathbf{s}) = \sum_{i=1}^k c_i^\sigma \eta_i^\sigma(\mathbf{s}) \quad (7)$$

$$\beta_\xi(\mathbf{s}) = \sum_{i=1}^k c_i^\xi \eta_i^\xi(\mathbf{s}) \quad (8)$$

where the  $c_i$ 's are weights for  $k$  basis functions, the  $\eta_i$ 's, which are distributed throughout the domain. More details are given in section ~~xxx~~2.2.

## 2.1 Elliptical copula for data dependence

Elliptical copulas are a flexible tool for modeling multivariate data [Renard, 2011; Sang and Gelfand, 2010; Ghosh and Mallick, 2011; Renard and Lang, 2007]. This class of copulas can represent spatial data with any marginal distribution, a particularly attractive feature for extremal data. The Gaussian copula constructs the joint pdf of a random vector  $(Y_1, \dots, Y_m)$  as

$$F_{Gaussian}(y_1, \dots, y_m) = \Phi_\Sigma(u_1, \dots, u_m) \quad (9)$$

where  $\Phi_\Sigma(u_1, \dots, u_m)$  is the joint cdf of an  $m$ -dimensional multivariate normal distribution with covariance matrix  $\Sigma$ ,  $u_i = \phi^{-1}(F_i[y_i])$ ,  $\phi$  is the cdf of the standard normal distribution and  $F_i$  is the marginal GEV cdf at site  $i$ . The corresponding joint pdf is

$$f_{Gaussian}(y_1, \dots, y_m) = \frac{\prod_{i=1}^m f_i[y_i]}{\prod_{i=1}^m \psi[u_i]} \Psi_\Sigma(u_1, \dots, u_m) \quad (10)$$

where  $f_i$  is the marginal GEV pdf at site  $i$ ,  $\psi$  is the standard normal pdf and  $\Phi_\Sigma$  is the joint pdf of an  $m$ -dimensional multivariate normal distribution.

The dependence between sites is assumed to be a function of distance [Renard, 2011]. The dependence matrix is constructed with a simple exponential model

$$\Sigma(i, j) = \exp(-\|\mathbf{s}_i - \mathbf{s}_j\|/a_0) \quad (11)$$

where  $a_0$  is the copula range parameter. Note that the values in this dependence matrix are not covariances, so by analogy with the variogram, the dependence model is termed the dependogram [Renard, 2011].

## 2.2 Spatial regression model

For large regions, spatial regression relationships may not hold constant for the entire domain. In this case it is necessary to allow for spatial variation in the spatial regressions for each GEV parameter. Each regression coefficient is represented as a weighted sum of radial basis functions basis functions (Equations 6-8). The form of these basis functions are

$$\eta_i(\mathbf{s}) = \exp(-\|\mathbf{s} - \mathbf{s}_i\|^2/a_i^2) \quad (12)$$

where  $a_i^2$  is a range parameter determining the spatial extent of the basis function. These basis functions, also known as Gaussian kernels, are placed at points throughout the domain, known as knots, allowing the regression coefficients to vary smoothly in space.

The knots are placed according to a space-filling design [Johnson *et al.*, 1990; Nychka and Saltzman, 1998]. For each GEV parameter, we use 10 knots (Figure 1) since based on the author's experience, regression relationships in the western US region tend to hold for regions of a few square degrees. For simplicity, the same knot locations were used for each GEV parameter and the copula but this is not required.

## 2.3 Missing Data

Stations with missing data can be easily incorporated in the model. When the GEV likelihood is computed, years with missing data are simply skipped. With at least 30 years of data at each station, the GEV parameters can be estimated adequately based on only the available data. For simplicity, the copula was fit to only stations with complete data, though missing data could be incorporated by varying the size of the covariance matrix for each year.

## 2.4 Likelihood and priors

The marginal distribution of  $Y(\mathbf{s}_i)$  is  $\text{GEV}(\mu(\mathbf{s}_i), \sigma(\mathbf{s}_i), \xi(\mathbf{s}_i))$ .  $Y(\mathbf{s}_i, t)$  is  $\text{GEV}(y(\mathbf{s}_i, t) | \mu(\mathbf{s}_i), \sigma(\mathbf{s}_i), \xi(\mathbf{s}_i))$  where the log-likelihood for some data point  $y$  is:

$$\log \text{GEV}(y | \mu, \sigma, \xi) = -\log(\sigma) - (1 + 1/\xi) \log(b) - b^{-1/\xi} \quad (13)$$

where  $b = 1 + \xi(y - \mu)/\sigma$ .

Let  $\gamma$  represent any of the GEV parameters  $(\mu, \sigma, \xi)$ . The residual Gaussian processes likelihood  $p(\mathbf{w}_\gamma | \boldsymbol{\theta}_\gamma)$  is obtained from the multivariate normal density function  $\mathbf{w}_\gamma | \boldsymbol{\theta}_\gamma \sim \text{MVN}(\mathbf{0}, \Sigma_\gamma)$ , where  $\Sigma_\gamma = C(\boldsymbol{\theta}_\gamma)$ . We use an exponential covariance function with parameters  $\delta_\gamma^2$  (the partial sill or marginal variance),  $a_\gamma$  (the range) and  $\tau_\gamma^2$  (the nugget), so  $\boldsymbol{\theta}_\gamma = (\delta_\gamma^2, a_\gamma, \tau_\gamma^2)$ . The parametric form of the covariance function is

$$C(\mathbf{s}_i, \mathbf{s}_j; \boldsymbol{\theta}_\gamma) = \begin{cases} \delta_\gamma^2 \exp(-\|\mathbf{s}_i - \mathbf{s}_j\|/a_\gamma) & i \neq j \\ \delta_\gamma^2 + \tau_\gamma^2 & i = j \end{cases}$$

We use weakly informative normal priors centered at 0, with a standard deviations as follows:  $0.1 (\delta_\xi^2, \tau_\xi^2)$ ,  $1 (\delta_\mu^2, \delta_\sigma^2, \tau_\mu^2, \tau_\sigma^2, \beta_0^\xi, c_i^\mu, c_i^\sigma, c_i^\xi; i = 1, \dots, n)$ ,  $10 (\beta_0^\mu, \beta_0^\sigma)$ ,  $1000 (a_\mu, a_\sigma, a_\xi, a_0, a_i; i = 1, \dots, n)$ . For  $\xi$  we restrict values to the range  $[-0.5, 0.5]$ , motivated by the typical ranges seen in precipitation data [Cooley and Sain, 2010].

## 2.5 Composite likelihood

When using Gaussian processes for large datasets, inversion (or Cholesky decomposition) of the covariance matrix is the main computational bottleneck. We use a composite likelihood approach to approximate the true likelihood [Lindsay, 1988][Lindsay, 1988; Varin et al., 2011]. In our approach, the data is broken up into  $G$  groups each with  $n_g$  stations. The composite likelihood estimate of the true likelihood is a product of the likelihood in each group.

$$L_{cl} = \prod_{g=1}^G \mathcal{N}(\mathbf{0}, \Sigma_g(\boldsymbol{\theta})) \quad (14)$$

Approximating the likelihood in this way requires  $O(Gn_g^3)$  computations as opposed to  $O(n^3)$ . This approximation is applied to the copula as well as each of the GEV parameter residuals.

What remains in the model are a few application specific details: selection of the knot locations and the selection of covariates. These are described in the next sections.

## 2.6 Composite likelihood group size and distribution

In order to use a composite likelihood approach we must decide how many stations to use in each group ( $n_g$ ). The number of stations in each group should be small enough so as not to incur substantial computational cost but large enough so that the covariance parameters can be adequately estimated. ~~In We use~~ We used 30 stations per group or approximately 1% of the total number of stations. The consequences of this choice are explored in Section ~~xxx~~4.2.

We must also choose how stations are to be grouped. Several approaches come to mind such as selecting groups based on climatological regions or elevation bands. We group stations randomly, expecting that groups will have a mixture of stations with a range of proximities, allowing for proper estimation of both small and large scale behavior.

## 3 Application to the Western US

### 3.1 Precipitation Data

Daily precipitation data was obtained from the Global Historical Climatology Network (GHCN). We use all available stations in the western US which contain more than 30 years of data from 1950-2013. 3-day maxima were computed fall (SON). For a season to be included for a particular year, we require no more than 25% of the days be missing. The number of stations included (with the number of complete stations in parentheses) was 2618 (848). Figure 1 shows the station locations, with solid black points indicating stations with complete data and filled grey points indicating stations with incomplete data. Red asterisks indicate the centers (knots) for the radial basis functions.

### 3.2 Covariates

For all GEV parameters the same covariates are used, i.e.,  $\mathbf{x}_\mu(\mathbf{s}) = \mathbf{x}_\sigma(\mathbf{s}) = \mathbf{x}_\xi(\mathbf{s}) = \mathbf{x}(\mathbf{s})$ . The covariates are elevation and mean seasonal precipitation. Typically, latitude and longitude are used as well but the spatially variation of the regression coefficients precludes this. Covariates were obtained at knot locations, station locations and at a 1/8th degree grid throughout the study area. Elevation data was obtained from the NASA Land Data

Assimilation Systems (NLDAS) website<sup>1</sup> [Xia *et al.*, 2012a, b]. Mean seasonal precipitation was computed from the Maurer dataset [Maurer *et al.*, 2002].

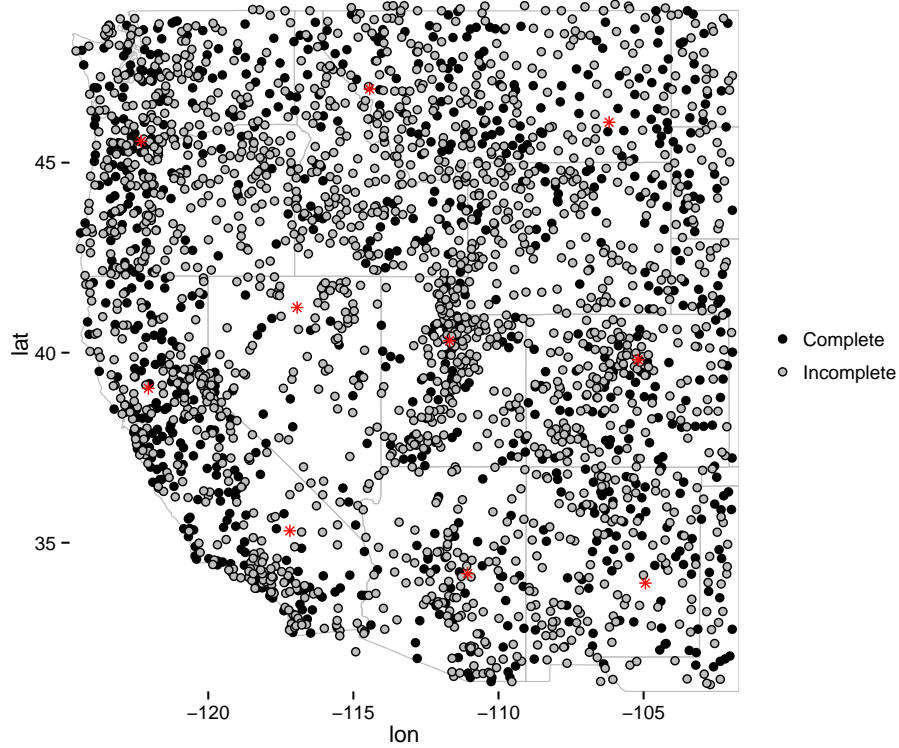


Figure 1: Station locations with complete data (black solid dots) and station locations with incomplete data (grey filled dots). Red asterisks are knot locations for the spatially varying regression coefficients.

### 3.3 Implementation

The model was implemented in the Stan modeling language [Stan Development Team, 2015b] using the RStan interface [Stan Development Team, 2015a]. Stan uses the No-U-

<sup>1</sup><http://ldas.gsfc.nasa.gov/nldas/NLDASelevation.php>



Turn Sampler (NUTS), an implementation of Hamiltonian Monte Carlo (HMC) [Betancourt, 2013; Hoffman and Gelman, 2014]. The NUTS sampler deals well with highly correlated parameters, tends to need very few warmup iterations and typically produces nearly uncorrelated samples. For these reasons, very long chains are usually not needed, nor is thinning. The tradeoff in using the NUTS sampler in this application was much longer computation time per sample compared to a traditional Metropolis-Hastings or Gibbs sampler.

Three chains of length 3,000 were run, with the first 1000 iterations discarded as warmup, resulting in 6,000 samples for each parameter in each season. To assess convergence, we compute the  $\hat{R}$  statistic to ensure it is below 1.1, as well as visually inspect trace plots.

### 3.4 Computation of gridded return levels

After computing  $\mu = [\mu_i]_{i=1}^n$ ,  $\sigma = [\sigma_i]_{i=1}^n$  and  $\xi = [\xi_i]_{i=1}^n$  distributions of each GEV parameter are obtained at each 1/8th degree grid cell via conditional simulation. The gridded parameter values are used to compute return levels at each grid cell using the GEV return level formula

$$z_i(r) = \mu_i + \sigma_i((-\log(1 - 1/r))^{-\xi_i} - 1)/\xi_i,$$

where  $r$  is the return period in years (100 years for example).

## 4 Results

### 4.1 Testing the validity of the Gaussian copula

An implication of the Gaussian copula is that marginal distributions are asymptotically independent, or  $P(F_x(X) > p | F_y(Y) > p) \rightarrow 0$  as  $p \rightarrow 1$ . To test this we conducted asymptotic independence tests (—ref—) [Reiss and Thomas, 2007] for all pairs of stations. The null hypothesis of this test is dependence, so setting a significance level of 99% ensures that stations passing the test exhibit strong asymptotic dependence. At the 99% significance level, 0.15% of pairwise stations exhibited dependence, less than the 1% expected from chance. In addition we examined plots of the station locations when dependence was indicated by the test. These plots did not show any discernible spatial pattern of dependence, for example dependent stations did not tend to fall near each other.

## 4.2 Group size selection

To demonstrate that the selection of group size has little effect on return levels, a small experiment is conducted. We run the model for a region encompassing most of the state of Oregon, using 4 knots. The group size is set to be 2, 5, 10, 15, 20 and 30 stations representing approximately 1%, 2%, 4%, 6%, 8% and 13% of the total number of stations respectively. The same 240 stations (60 complete, 180 incomplete) are used in each model run.

Figure 2 shows the median return level for each model run. The results are nearly identical for this range of group sizes, indicating that median return levels are not sensitive to the choice of group size. Credible intervals of return levels (not shown) were quite similar as well, with credible intervals decreasing as group size is increased indicating that a larger group size yields more accurate results, as expected. In light of this we chose a group size of 30 for the large domain which provides both a diversity in the distribution of stations within a group but is small enough to not significantly hinder computation.

## 4.3 Gridded return levels

Figure 4 shows the median of the GEV parameters after interpolation by conditional simulation. The location and shape fields are highly correlated; locations with higher average extreme precipitation tend to have more variability in these extremes. Values of  $\xi$  are always positive, indicating a heavy upper tail. The southern coastal region in California in the summer indicates a very heavy upper tail. Figure shows the ratio of the median return level to the width of the 90% CI indicating the largest relative uncertainties actually occur mostly in southern California, where the GEV tail is the fattest.

## 4.4 Validation

Cross validation was conducted by dropping 885 stations or approximately 35% of the total stations. Gridded return levels were computed for this subset of data. Figure 6 shows the difference between the median return level for the full data and subset data. The difference map shows some spatial coherence but none that indicates any strong bias in a single region (states for example). The largest differences occur in areas in the northwest where influential stations were dropped randomly.

## 4.5 A case for composite likelihood

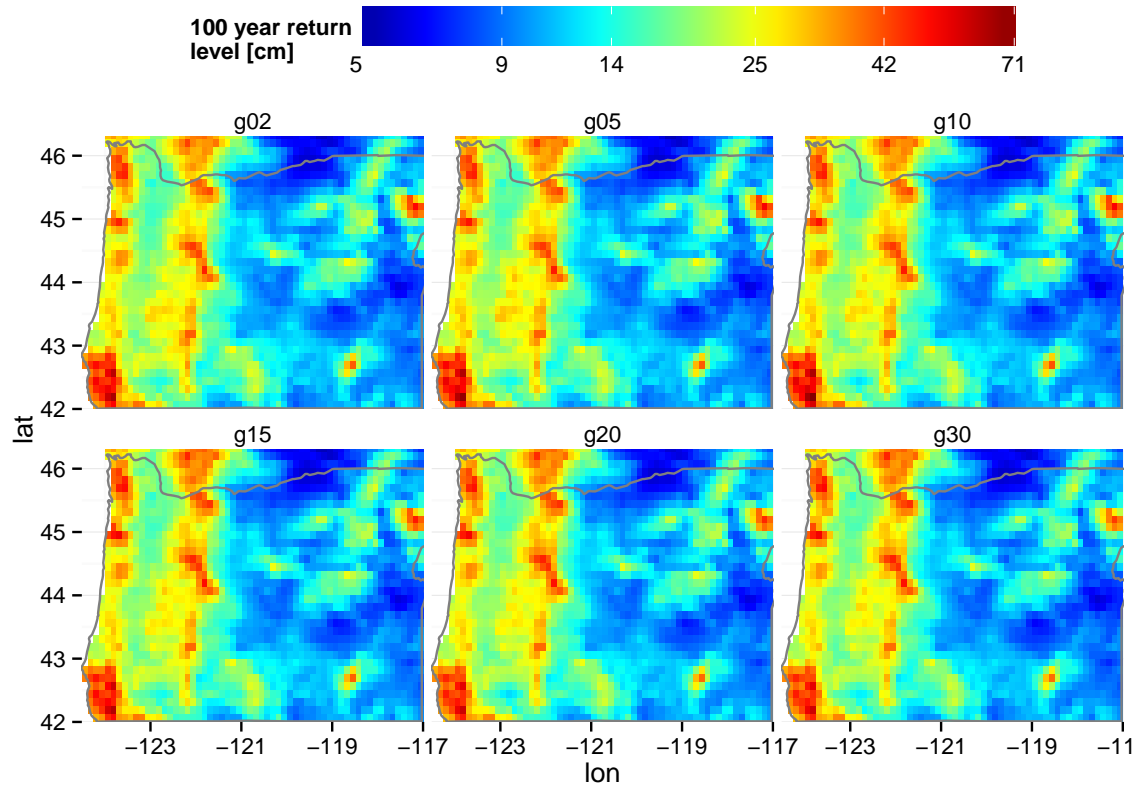


Figure 2: Median return levels using a group sizes of 2, 5, 10, 15, 20 and 30. Note the logarithmic color scale.

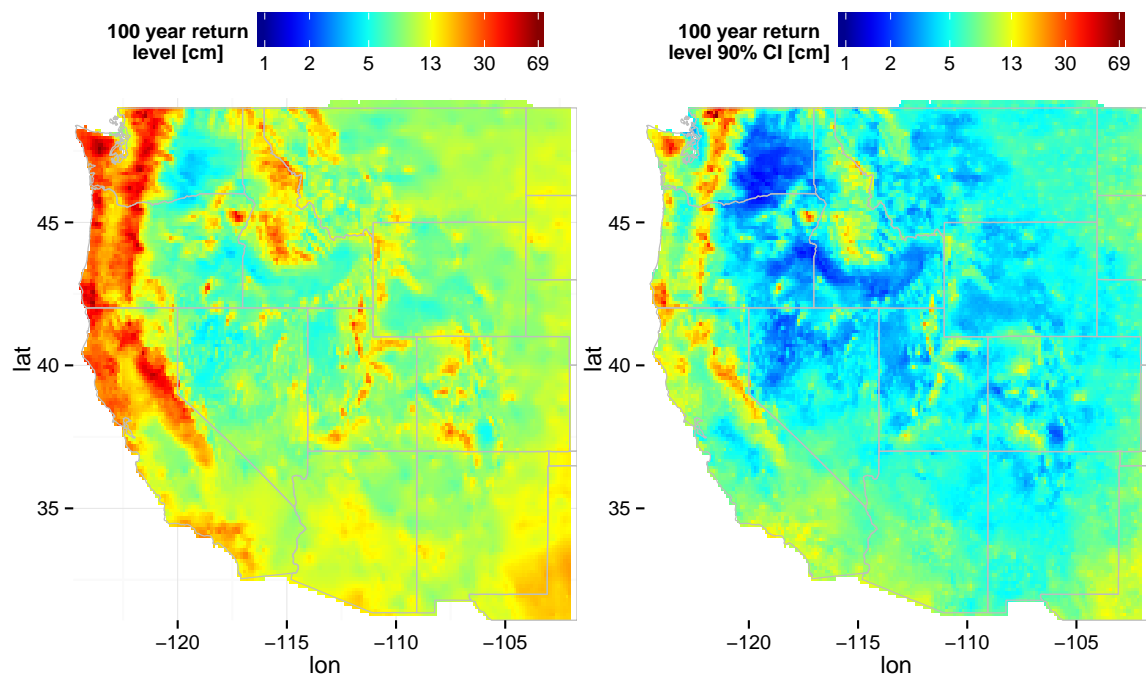


Figure 3: Median 100-year return levels for fall (left) and width of corresponding 95% credible interval (right). Note the logarithmic color scale.

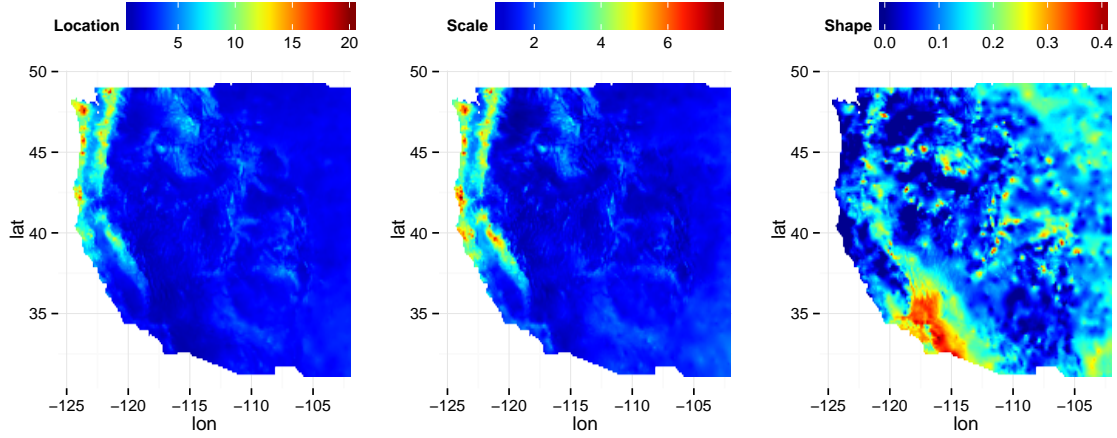


Figure 4: Median of underlying GEV parameters, location ( $\mu$ ), scale ( $\sigma$ ) and shape ( $\xi$ ).

To highlight the usefulness of the composite likelihood approach for this application we present results using a Gaussian predictive process (GPP) model [Banerjee et al., 2008] for the latent GEV parameter processes (Figure 7). A Gaussian predictive process model approximates the likelihood at a small set of knots to reduce the dimension of the covariance matrix and the computational burden of inverting it. We originally set out using GPPs for this application but switched to a composite likelihood approach when we realized the uncertainty was unacceptably large away from knot locations.

The median return levels with the GPP approach were nearly identical to those from the composite likelihood method (Figure 3) but large differences are apparent when looking at the credible intervals of the return levels. Clear artifacts are present at the locations of knots, where uncertainty is greatly reduced. Uncertainty away from knot locations is typically large, rendering this method much less useful than the composite likelihood approach.

## 5 Discussion and conclusions

We describe a general hierarchical model for extreme data observed over space and time. The data is assumed to originate from a Gaussian elliptical copula having generalized extreme value (GEV) marginal distributions. Spatial dependence is further captured by

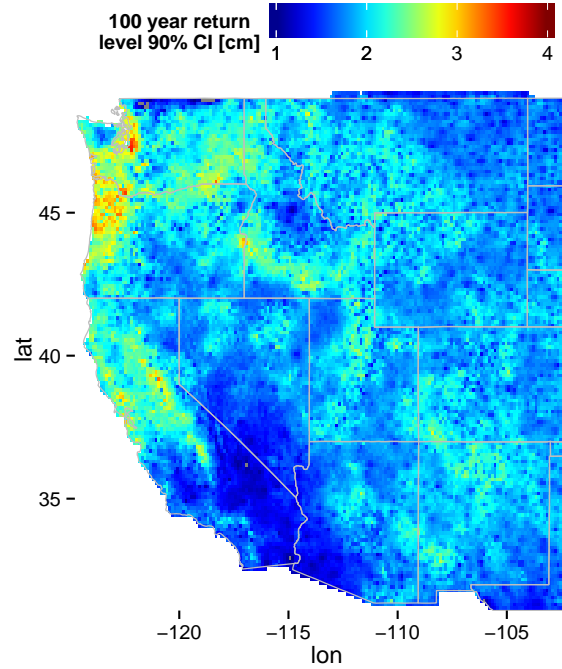


Figure 5: Ratio of 50th percentile return level and 90% credible interval width. ~~Higher values~~

Gaussian processes on the three GEV parameters (location, scale and shape). Using a composite likelihood approach, we are able to incorporate 2595 observation locations with 54 years of data. With spatially varying regression coefficients, the model can be applied to arbitrarily large regions. The model was applied to extreme 3-day precipitation in fall in the western United States, a climatically and geographically diverse region. The model was fit using a standard Bayesian methodology, implicitly capturing uncertainty in the parameter estimates and spatial predictions.

In Section 4.5 we briefly examine results for the same region using a Gaussian predictive process (GPP) model for the latent GEV parameters. In this application, GPPs produced unreasonably large posterior credible intervals when moving away from knot locations. In light of this we recommend a composite likelihood approach for regions of equal or larger size than the western US.

A crux of this model is the use of appropriate spatial covariates. Mean seasonal precipitation (MSP) had a correlation of 95% with the MLE estimates of  $\mu$  and 75% with the MLE estimates of  $\sigma$ . ~~This covariate went a long way~~ Even with spatially varying regression coefficients, appropriate covariates are key. The covariates here helped in generating re-

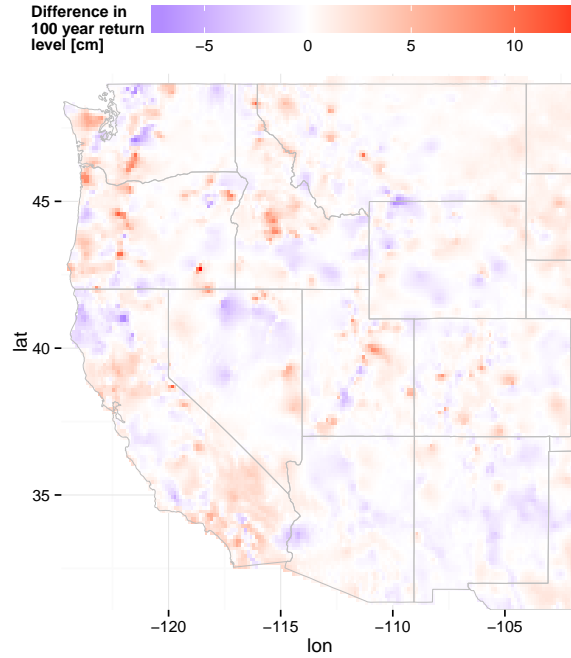


Figure 6: Difference between 50th percentile return levels from the full model and the validation model dropping 35% of the data.

alistic spatial variability ~~-.The covariates also help~~ and helped to reveal a complex spatial pattern for the shape parameter,  $\xi$ . The strongest covariate for  $\xi$  was elevation. The spatial variability in  $\xi$  shows that it is inappropriate to model without spatial variation for anything but the smallest regions.

A number of extensions can be made to this framework. The most obvious extension is to allow temporal variation in the GEV parameters by including temporal covariates. While this extension remains infeasible for the size of the current study region, it may be feasible for smaller regions, say a single state or moderate sized river basin. Additional spatial covariates could be included; for example, seasonal temperature, winds or evapotranspiration. A model such as the one presented here can be used to investigate changes in risk under specific climate regimes (i.e. ENSO); one would simply include the mean seasonal precipitation field from strong El Niño or La Niña years. Because we incorporate a data layer, this model could be used to simulate realistic fields of extremes under specific climate regimes. Finally, we plan to explore the linking of streamflow data into the hierarchy, so that streamflow extremes can be simultaneously estimated.

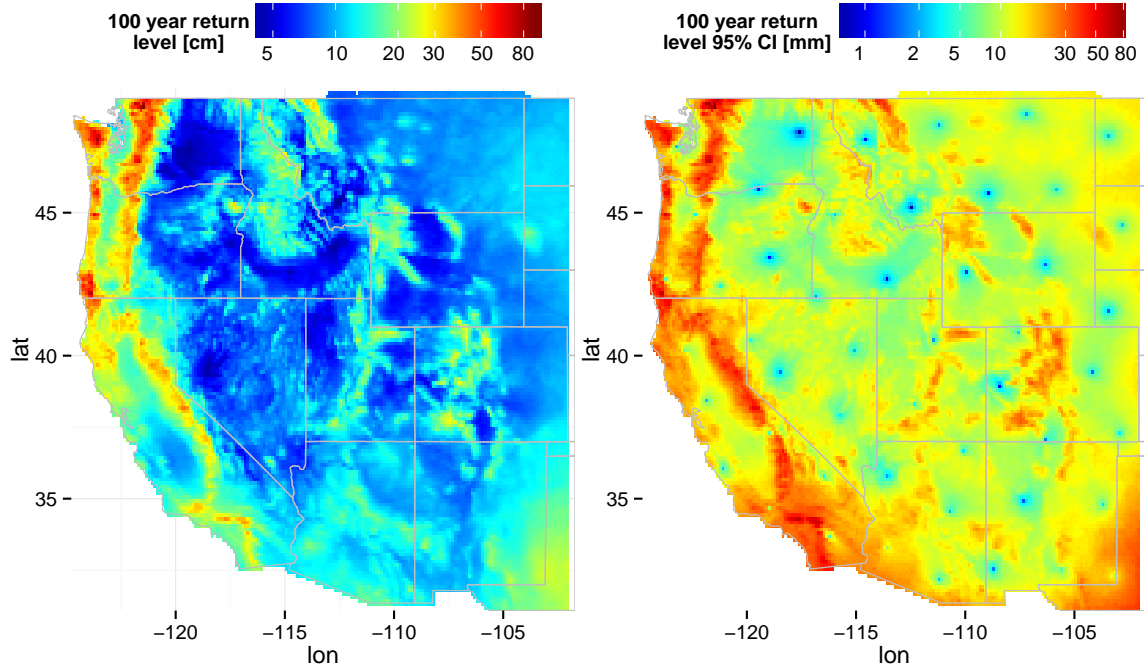


Figure 7: [Return levels maps produced using latent gaussian predictive processes.](#)

## 6 Acknowledgments

Funding for this research by a Science and Technology grant from Bureau of Reclamation is gratefully acknowledged. This work utilized the Janus supercomputer, which is supported by the National Science Foundation (award number CNS-0821794) and the University of Colorado Boulder. The Janus supercomputer is a joint effort of the University of Colorado Boulder, the University of Colorado Denver and the National Center for Atmospheric Research. The authors are thankful for support from the Janus supercomputer staff at the University of Colorado.

Pre- and postprocesseing analysis was conducted using the R language [R Core Team, 2014].

Data is available at: [http://bechtel.colorado.edu/~bracken/spatial\\_extremes/](http://bechtel.colorado.edu/~bracken/spatial_extremes/).



## References

- Apputhurai, P., and A. G. Stephenson, Spatiotemporal hierarchical modelling of extreme precipitation in Western Australia using anisotropic Gaussian random fields, *Environmental and Ecological Statistics*, 20(4), 667–677, 2013.
- Aryal, S. K., B. C. Bates, E. P. Campbell, Y. Li, M. J. Palmer, and N. R. Viney, Characterizing and Modeling Temporal and Spatial Trends in Rainfall Extremes, *dx.doi.org*, 10(1), 241–253, 2010.
- Atyeo, J., and D. Walshaw, A region-based hierarchical model for extreme rainfall over the UK, incorporating spatial dependence and temporal trend, *Environmetrics*, 23(6), 509–521, 2012.
- Banerjee, S., A. E. Gelfand, A. O. Finley, and H. Sang, Gaussian predictive process models for large spatial data sets, *Journal of the Royal Statistical Society*, 2008.
- Betancourt, M. J., Generalizing the No-U-Turn Sampler to Riemannian Manifolds, *arXiv*, 1304(1920), 2013.
- Caragea, P. C., and R. L. Smith, Asymptotic properties of computationally efficient alternative estimators for a class of multivariate normal models, *Journal of Multivariate Analysis*, 98(7), 1417–1440, 2007.
- Cooley, D., and S. R. Sain, Spatial Hierarchical Modeling of Precipitation Extremes From a Regional Climate Model, *Journal of Agricultural, Biological, and Environmental Statistics*, 15(3), 381–402, 2010.
- Cooley, D., P. Naveau, and P. Poncet, Variograms for spatial max-stable random fields, *Dependence in probability and statistics*, 2006.
- Cooley, D., D. Nychka, and P. Naveau, Bayesian spatial modeling of extreme precipitation return levels, *Journal of the American Statistical Association*, 2007.
- Davison, A. C., S. A. Padoan, and M. Ribatet, Statistical Modeling of Spatial Extremes, *Statistical Science*, 27(2), 161–186, 2012.
- Dyrrdal, A. V., A. Lenkoski, T. L. Thorarinsdottir, and F. Stordal, Bayesian hierarchical modeling of extreme hourly precipitation in Norway, *Environmetrics*, 2014.
- Fuentes, M., Approximate Likelihood for Large Irregularly Spaced Spatial Data, *Journal of the American Statistical Association*, 102(477), 321–331, 2007.
- Ghosh, S., and B. K. Mallick, A hierarchical Bayesian spatio-temporal model for extreme precipitation events, *Environmetrics*, 22(2), 192–204, 2011.

- Heagerty, P. J., and S. R. Lele, A Composite Likelihood Approach to Binary Spatial Data, *Journal of the American Statistical Association*, 93(443), 1099, 1998.
- Hoffman, M. D., and A. Gelman, The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo, *Journal of Machine Learning Research*, 15(Apr), 1593–1623, 2014.
- Johnson, M. E., L. M. Moore, and D. Ylvisaker, Minimax and maximin distance designs, *Journal of Statistical Planning and Inference*, 26(2), 131–148, 1990.
- Lindsay, B. G., Composite Likelihood Methods, *Contemporary Mathematics*, 80, 221–239, 1988.
- Maurer, E. P., A. W. Wood, J. C. Adam, D. P. Lettenmaier, and B. Nijssen, A Long-Term Hydrologically Based Dataset of Land Surface Fluxes and States for the Conterminous United States, *Journal of Climate*, 15(22)(22), 3237–3251, 2002.
- Nychka, D., and N. Saltzman, Design of Air-Quality Monitoring Networks, in *Case Studies in Environmental Statistics*, pp. 51–76, Springer US, New York, NY, 1998.
- Padoan, S. A., M. Ribatet, and S. A. Sisson, Likelihood-Based Inference for Max-Stable Processes, *dx.doi.org*, 105(489), 263–277, 2010.
- R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2014.
- Reich, B. J., and B. Shaby, A hierarchical max-stable spatial model for extreme precipitation, *The annals of applied statistics*, 6(4), 1430–1451, 2012.
- Reiss, R.-D., and M. Thomas, *Statistical Analysis of Extreme Values: with Applications to Insurance, Finance, Hydrology and Other Fields*, 3rd edition ed., Birkhäuser, 2007.
- Renard, B., A Bayesian hierarchical approach to regional frequency analysis, *Water Resources Research*, 2011.
- Renard, B., and M. Lang, Use of a Gaussian copula for multivariate extreme value analysis: Some case studies in hydrology, *Advances in Water Resources*, 30(4), 897–912, 2007.
- Rue, H., S. Martino, and N. Chopin, Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(2), 319–392, 2009.
- Sang, H., and A. E. Gelfand, Hierarchical modeling for extreme values observed over space and time, *Environmental and Ecological Statistics*, 16(3), 407–426, 2009.
- Sang, H., and A. E. Gelfand, Continuous Spatial Process Models for Spatial Extreme Values, *Journal of Agricultural, Biological, and Environmental Statistics*, 15(1), 49–65, 2010.

- Schlather, M., Models for Stationary Max-Stable Random Fields, *Extremes*, 5(1), 33–44, 2002.
- Shang, H., J. Yan, and X. Zhang, El Niño–Southern Oscillation influence on winter maximum daily precipitation in California in a spatial model, *Water Resources Research*, 47(11), n/a–n/a, 2011.
- Smith, R. L., Max-stable processes and spatial extremes, *Unpublished manuscript*, 1990.
- Stan Development Team, *RStan: the R interface to Stan*, Version 2.7.0, 2015a.
- Stan Development Team, Stan: A C++ Library for Probability and Sampling, Version 2.7.0, 2015b.
- Stein, M. L., Z. Chi, and L. J. Welty, Approximating likelihoods for large spatial data sets, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(2), 275–296, 2004.
- Varin, C., N. Reid, and D. Firth, An overview of composite likelihood methods, *Statist. Sinica*, pp. 5–42, 2011.
- Xia, Y., et al., Continental-scale water and energy flux analysis and validation for North American Land Data Assimilation System project phase 2 (NLDAS-2): 2. Validation of model-simulated streamflow, *J. Geophys. Res.*, 117(D3), D03,110, 2012a.
- Xia, Y., et al., Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of model products, *J. Geophys. Res.*, 117(D3), D03,109, 2012b.