



This discussion paper is/has been under review for the journal Nonlinear Processes in Geophysics (NPG). Please refer to the corresponding final paper in NPG if available.

# Non-parametric Bayesian mixture of sparse regressions with application towards feature selection for statistical downscaling

D. Das<sup>1,2</sup>, J. Dy<sup>3</sup>, J. Ross<sup>3</sup>, Z. Obradovic<sup>2</sup>, and A. R. Ganguly<sup>1</sup>

<sup>1</sup>Sustainability and Data Sciences Lab, Northeastern University, Boston, USA

<sup>2</sup>Center for Data Analytics and Biomedical Informatics, Temple University, Philadelphia, USA

<sup>3</sup>Department of Electrical and Computer Engineering, Northeastern University, Boston, USA

Received: 27 February 2014 – Accepted: 14 March 2014 – Published:

Correspondence to: A. R. Ganguly (a.ganguly@neu.edu)

Published by Copernicus Publications on behalf of the European Geosciences Union & the American Geophysical Union.

## Abstract

Climate projections simulated by Global Climate Models (GCM) are often used for assessing the impacts of climate change. However, the relatively coarse resolutions of GCM outputs often precludes their application towards accurately assessing the effects of climate change on finer regional scale phenomena. Downscaling of climate variables from coarser to finer regional scales using statistical methods are often performed for regional climate projections. Statistical downscaling (SD) is based on the understanding that the regional climate is influenced by two factors – the large scale climatic state and the regional or local features. A transfer function approach of SD involves learning a regression model which relates these features (predictors) to a climatic variable of interest (predictand) based on the past observations. However, often a single regression model is not sufficient to describe complex dynamic relationships between the predictors and predictand. We focus on the covariate selection part of the transfer function approach and propose a nonparametric Bayesian mixture of sparse regression models based on Dirichlet Process (DP), for simultaneous clustering and discovery of covariates within the clusters while automatically finding the number of clusters. Sparse linear models are parsimonious and hence relatively more generalizable than non-sparse alternatives, and lends to domain relevant interpretation. Applications to synthetic data demonstrate the value of the new approach and preliminary results related to feature selection for statistical downscaling shows our method can lead to new insights.

## 1 Introduction

Climate change is one of most challenging problems facing mankind whose impacts are expected to influence policy decisions on critical infrastructures, management of natural resources, humanitarian aid, emergency preparedness along with numerous regional scale human economic and social activities. Therefore, it is imperative to accurately assess the impacts of climate change at regional scale in order to inform stakeholders for appropriate

5 decision making related to mitigation policies. Global Climate Models (GCM) are the most credible tools at present for future climate projections that accounts for the effects of greenhouse gas emissions under different socio-economic scenarios. Although GCMs perform reasonably well in projecting climate variables at a larger spatial scale ( $> 10^4 \text{ km}^2$ ), they perform poorly for regional scale climate projections. Such poor performance of the GCMs coupled with the importance of regional climate projections for impact studies have led to development of Limited Area Models (LAM) or Regional Climate Models (RCM), where finer spatial grids over a limited spatial area are embedded within a coarser GCM grid. This method is also known as dynamic downscaling. However, these models are 10 complex, computationally expensive and requires rerunning for each new region. Moreover, regional models inherit the basic gaps in understanding of climate physics that limits the performance of GCMs. A couple of recently published studies (Kumar et al., 2014; Knutti and Sedláček, 2013) rigorously compared the projections of latest generation of climate models (CMIP5) with the previous generation (CMIP3) but found no significant improvement in majority of statistical performance metrics even with higher spatial resolutions and 15 addition of new physical processes in the computational model. Uncertainties in sub-grid scale cloud-microphysics and ocean eddy processes and poor understanding of the effect of carbon cycle and other biogeochemical processes on climate systems still limits the ability of the physics-based climate models to reliably project future climate (Bader et. al., 2008), especially at regional scale. 20

A complementary approach for regional projection is statistical downscaling that uses statistical models to learn empirical statistical relationship between large scale GCM features (predictors) and regional scale climate variable(s) (predictands) to be projected. The statistical approaches of downscaling can be categorized into three broad classes – weather typing, weather generators and the transfer function approaches (Wilby et al., 25 2004). Weather typing approaches have originally been developed for weather forecasting and generally involves classifying days into similar clusters or weather states based on their synoptic similarity. Typically, weather patterns are clustered based on their similarity with nearest neighbors while the statistical models they use varies in their definition of similarity

measure. On the other hand, weather generators replicates the statistical properties of the daily predictand variable by using a stochastic model like Markov processes () that uses wet–dry and dry–wet transition probabilities as input for training while conditioning its parameters on large-scale predictors.

5 In this paper, however, we are interested in transfer function based regression models that learn a linear or nonlinear mapping between large scale predictors and regional scale predictand variables. Regression models are conceptually simplest among the three since they provide a direct mapping between the predictor and predictand values. However, the success of the regression models depends on the accurate choice of predictors. Sparse  
10 regressions based on constrained L1-norm (Tibshirani, 1994) of the coefficients became popular due to their ability to simultaneously select covariates and fit parsimonious linear models that are better generalizable and easily interpretable. Although, sparse regression models have been applied widely in many disciplines, it's application to climate has remained very limited, especially to statistical downscaling. In a recent paper (Ebtehaj  
15 et al., 2012), sparse regularization has been shown to be effective for downscaling rainfall fields for weather forecasting, whereas sparse variable selection has been used for statistical downscaling of climate variables (Phatak et al., 2011) in a separate paper. To our knowledge, there is no other published work on use of sparse regularization for statistical downscaling.

20 However, large complex climate datasets often exhibit dynamic behavior (Kannan and Ghosh, 2010) which may not be modeled well by a single regression model. Here we propose a nonparametric model for mixture of sparse regressions that can accommodate multiple sparse linear relationship inherent in the dataset. Nonparametric models are more flexible than the finite mixture models (Bishop and Svenskn, 2002) since they assume no  
25 prior knowledge about the number of distinct components in the data. We used a Dirichlet process mixture (DPM) (Antoniak, 1974) with stick-breaking construction (Ishwaran and James, 2001) to accommodate an unknown number of sparse regression models in the data. DPMs start by assuming infinite components in the data but ends up discovering a finite number of components supported by the data. We used the Bayesian version of

sparse regression (Park and Casella, 2008) to smoothly integrate the sparse regression model with the DPM, which is a nonparametric Bayesian approach where each component is represented by a set of distribution parameters specific to the corresponding component.

5 Although the number of different components may not be known, prior knowledge often exists about whether a pair of observations belong to the same component. For example, it is reasonable to assume that two observations close in time from the same location may exhibit similar behavior. We allow soft “must link” constraints between pairs of data-points that encourage the pair to belong to the same mixture component. Such constraints are incorporated in our Bayesian model with the help of a Markov random field (MRF) prior over  
10 the cluster indicator variables (Ross and Dy, 2013; Basu et al., 2006).

Variational Bayesian inferences has been shown to be much faster than stochastic alternatives for nonparametric Bayesian models (Blei and Jordan, 2006). The major contribution of this paper is to develop a fully Bayesian formulation for nonparametric mixture of sparse regression model and designing an efficient variational inference  
15 algorithm to obtain posterior distributions over the regression coefficients of potentially multiple regression components as well as the component membership probabilities of each data-point.

We have extensively demonstrated the performance of our algorithm on synthetic data. We have also applied our method for the feature selection problem for statistical  
20 downscaling of annual average rainfall over two regions in the west coast of United States. Preliminary results from the application of our algorithm to select features for regression based statistical downscaling shows that our method may lead to improved prediction and discovery of new insights.

## 2 Background

25 In this section, we provide brief descriptions of the methods in the context they were used to build our model.

## 2.1 Bayesian sparse regression

Let us assume that we are given a dataset  $D = \{\mathbf{x}_n, y_n : n = 1, \dots, N\}$  which has been generated from a linear model identified by sparse coefficients vector  $\beta$ . In a non-Bayesian setting, sparsity is enforced by a constraint on the L1-norm of the coefficients which is given by

$$\mathbf{y}_n = \beta^\top \mathbf{x}_n + \epsilon, \quad \text{subject to } \|\beta\|_1 \leq t \quad (1)$$

where  $\epsilon \sim \mathcal{N}(\mu, \tau^{-1})$ .

However, in a Bayesian setting, the sparsity can be imposed by a Laplace prior (also known as double exponential distribution) on  $\beta$  which is given by Park and Casella (2008)

$$p(\beta | \gamma, \tau) = \prod_{j=1}^D \frac{\sqrt{\gamma_j \tau}}{2} \exp(-\sqrt{\gamma_j \tau} |\beta_j|) \quad (2)$$

However, due to the analytical intractability of the Laplace prior, it is often represented in the following scale-mixture (of Gaussians) form using an additional random variable  $\alpha$ .

$$p(\beta | \tau, \gamma) = \prod_{j=1}^D \frac{\sqrt{\gamma_j \tau}}{2} \exp(-\sqrt{\gamma_j \tau} |\beta_j|) = \prod_{j=1}^D \int \mathcal{N}(\beta_j; 0, \tau^{-1} \alpha_j^{-1}) \text{InvGa}(\alpha_j; 1, \frac{\gamma_j}{2}) d\alpha_j$$

For a fully hierarchical Bayesian setting, Gamma prior is imposed on parameter  $\tau$  as well as on individual penalty parameters  $\gamma_j$ . So the joint distribution over all the parameters can be given by

$$p(\beta, \tau, \alpha, \gamma) = \text{Ga}(\tau; c_0, d_0) \prod_{j=1}^D \left\{ \mathcal{N}(\beta_j; 0, \tau^{-1} \alpha_j^{-1}) \times \text{InvGa}(\alpha_j; 1, \frac{\gamma_j}{2}) \text{Ga}(\gamma_j; a_0, b_0) \right\} \quad (3)$$

## 2.2 Markov random fields

A Markov random field (MRF) is represented by an undirected graphical model in which the nodes represent variables or groups of variables and the edges indicate dependence relationships. An important property of MRFs is that a collection of variables is conditionally independent of all others in the field given the variables in their Markov blanket. The Hammersley–Clifford theorem states that the distribution,  $p(\mathbf{Z})$ , over the variables in a MRF factorizes according to

$$p(\mathbf{Z}) = \frac{1}{\mathcal{Z}} \exp \left( - \sum_{c \in C} H_c(z_c) \right) \quad (4)$$

where  $\mathcal{Z}$  is a normalization constant called the *partition function*,  $C$  is the set of all cliques in the MRF,  $z_c$  are the variables in clique  $c$ , and  $H_c$  is the *energy function* over clique  $c$  (Geman and Geman, 1984). Cliques are sets of variables or nodes in the graphical model that are fully connected and the smallest clique is an edge. The energy function captures the desired configuration of local variables. *Partition function*  $\mathcal{Z}$  normalizes the probability measure and it is computed by summing the exponentiated energy functions of all possible configurations.

## 2.3 Dirichlet process mixture

The Dirichlet Process (DP) was first introduced in statistics literature as a measure on measures (Ferguson, 1973). It is parameterized by a base measure,  $G_0$ , and a positive scaling parameter  $\lambda$ :

$$G|\{G_0, \lambda\} \sim \text{DP}(G_0, \lambda) \quad (5)$$

The notion of a Dirichlet process mixture (DPM) arises if we treat the  $k$ th draw from  $G$  as a parameter of the distribution over some observation (Antoniak, 1974) representing a particular mixture component. DPMs can be interpreted as mixture models with an

infinite number of mixture components in the sense that data exhibits a finite number of components but previously unseen components represented by new data can still be accommodated. More recently, a variational inference algorithm for DPMs was introduced (Blei and Jordan, 2006) using the stick-breaking construction (Sethuraman, 1994) which

5 uses two infinite collections of random variables  $V_k \sim \text{Beta}(1, \lambda)$  and  $\boldsymbol{\eta}_k^* \sim G_0$  to construct  $G$  as

$$\theta_k = V_k \prod_{j=1}^{k-1} (1 - V_j)$$

$$G(\boldsymbol{\eta}) \sim \sum_{k=1}^{\infty} \theta_k \delta(\boldsymbol{\eta}, \boldsymbol{\eta}_k^*).$$
(6)

For a mixture of sparse regression models, if the parameters for each components are given

10 by  $\boldsymbol{\eta}_k$ , the subsequent data generation process for such a mixture model can be described in the following steps using a stick-breaking construction:

1. Draw  $v_k \sim \text{Beta}(1, \lambda)$   $k = \{1, 2, \dots, \infty\}$ .

2. Draw  $\boldsymbol{\eta}_k \sim G_0$ ,  $k = \{1, 2, \dots, \infty\}$ .

3. Generate  $\theta_k = v_k \prod_{m=1}^{k-1} (1 - v_m)$ .

15 4. For each data point  $n$ :

(a) Draw  $z_n \sim \text{Mult}(\boldsymbol{\theta})$ .

(b) Draw  $y_n \sim \mathcal{N}(y_n; \mathbf{x}_n, \boldsymbol{\eta}_{z_n})$ .

We can truncate the construction process at  $k = K$  by enforcing  $V_{k-1} = 0$  which forces all  $\theta_k$  for  $k > K$  to be zero (see step 3). The resulting construction is called a truncated Dirichlet process (TDP) which can be shown to approximate the true Dirichlet process quite

20 well given  $K$  is large relative to the number of the data-points (Ishwaran and James, 2001).



### 3 Methodology

Now, let us assume that we are given a dataset  $D = \{\mathbf{x}_n, y_n : n = 1, \dots, N\}$  which has been generated from a mixture of  $K$  different sparse models identified by sparse coefficients  $\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(K)}$ . Let us also assume that the number of components  $K$  is unknown.

- 5 We use a Bayesian formulation of the sparse regression model for each component  $\beta^{(k)}$ , with  $k = 1, 2, \dots, K$ . Let us first state the Bayesian version of the  $k$ th sparse model. The linear regression model of the  $k$ th component can be represented by the following Gaussian distribution.

$$p\left(y_n | \mathbf{x}_n, \beta^{(k)}\right) \sim \mathcal{N}\left(y_n; \beta^{(k)\top} \mathbf{x}_n, \tau_k^{-1}\right) \quad (7)$$

#### 10 3.1 Mixture of sparse regressions

We introduce  $K$ -dimensional latent indicator variables  $\{z_n : n = 1, \dots, N\}$  to represent the component membership of each data-point  $\{\mathbf{x}_n, y_n\}$ . If the data-point belongs to the  $k$ th component, then  $z_{nk}$  will be 1 and all other elements of  $z_n$  will be 0. We further denote  $\mathbf{Z} = [z_1 \quad z_2 \quad \dots \quad z_n]$ . We can now rewrite Eq. (7) in terms of  $z_n$  as

$$15 \quad p\left(y_n | \mathbf{x}_n, \{\beta^{(k)}\}\right) \sim \prod_{k=1}^K \left\{ \mathcal{N}\left(y_n; \beta^{(k)\top} \mathbf{x}_n, \tau_k^{-1}\right) \right\}^{z_{nk}} \quad (8)$$

- For this mixture of sparse regression model, each component has separate parameter set  $\{\beta^{(k)}, \tau_k\}$ . Moreover, after adding the parameters related to the scale-mixture representation of the Laplace prior on  $\beta^{(k)}$  (refer Sect. 2.1), the set of parameters is finally given by  $\eta_k = \{\beta^{(k)}, \tau_k, \alpha_k, \gamma_k\}$ . The prior distribution  $G_0$  from which these parameters can be drawn jointly is given in Eq. (3). We can now use the stick-breaking construction described in Sect. 2.3 to formulate our mixture model. The overall generative process is
- 20

then:

$$\begin{aligned}
 p(\mathbf{y}, \mathbf{Z}, \mathbf{v}, \{\boldsymbol{\beta}^{(k)}\}, \boldsymbol{\tau}, \{\boldsymbol{\alpha}^{(k)}\}, \{\boldsymbol{\gamma}^{(k)}\}, \lambda | \mathbf{X}) &= p(\mathbf{y} | \mathbf{X}, \{\boldsymbol{\beta}^{(k)}\}, \boldsymbol{\tau}) p(\mathbf{Z} | \mathbf{v}) p(\mathbf{v} | \lambda) p(\lambda | m_0) \\
 &\times p(\{\boldsymbol{\beta}^{(k)}\} | \boldsymbol{\tau}, \{\boldsymbol{\alpha}^{(k)}\}) p(\{\boldsymbol{\alpha}^{(k)}\} | \{\boldsymbol{\gamma}^{(k)}\}) \\
 &\times p(\{\boldsymbol{\gamma}^{(k)}\} | a_0, b_0) p(\boldsymbol{\tau} | c_0, d_0)
 \end{aligned} \tag{9}$$

The graphical model that represents the dependence relationships between all the parameters involved in this current mixture model is shown in Fig. 1. The shaded circles denote observed variables whereas the unshaded circles denote unobserved variables. We have used a Gamma prior on  $\lambda$  having a hyper-parameter  $m_0$ . We have omitted the hyper-parameters  $a_0, b_0, c_0, d_0$ , and  $m_0$  from the list of conditioning variables in the left side to avoid clutter. The individual distributions in Eq. (9) are given below

$$\mathbf{y} | \mathbf{X}, \{\boldsymbol{\beta}^{(k)}\}, \boldsymbol{\tau} \sim \prod_{n=1}^N \prod_{k=1}^K \left\{ \mathcal{N} \left( y_n; \mathbf{x}_n^\top \boldsymbol{\beta}^{(k)}, \tau_k^{-1} \right) \right\}^{z_{nk}} \tag{10a}$$

$$\mathbf{Z} | \mathbf{v} \sim \prod_{n=1}^N \prod_{k=1}^K \left\{ v_k \prod_{j=1}^{k-1} (1 - v_j) \right\}^{z_{nk}} \tag{10b}$$

$$\mathbf{v} | \lambda \sim \prod_{k=1}^K \text{Beta}(v_k; 1, \lambda) \tag{10c}$$

$$\lambda \sim \text{Ga}(\lambda; m_0, 1) \tag{10d}$$

$$\{\boldsymbol{\beta}^{(k)}\} | \boldsymbol{\tau}_k, \{\boldsymbol{\alpha}^{(k)}\} \sim \prod_{k=1}^K \prod_{j=1}^D \mathcal{N} \left( \beta_j^{(k)}; 0, \left( \tau_k \alpha_j^{(k)} \right)^{-1} \right) \tag{10e}$$

$$\boldsymbol{\tau} \sim \prod_{k=1}^K \text{Ga}(\tau_k; c_0, d_0) \tag{10f}$$

$$(\{\boldsymbol{\alpha}^{(k)}\}, \{\boldsymbol{\gamma}^{(k)}\}) \sim \prod_{k=1}^K \prod_{j=1}^D \text{InvGa} \left( \alpha_j^{(k)}; 1, \frac{\gamma_j^{(k)}}{2} \right) \times \text{Ga} \left( \gamma_j^{(k)}; a_0, b_0 \right) \quad (10g)$$

### 3.2 Accommodating “must link” constraints

Prior knowledge about must link constraints between pairs of data-points can be enforced via a MRF prior on the indicator variables  $z_n$  where each data point is considered a node and each constraint between a pair of data point is regarded as an edge between the respective nodes. We denote the collection of edges by  $\mathcal{C}$  and the MRF prior is given by Eq. (4). We define the energy function as:

$$H(\mathbf{z}_i, \mathbf{z}_j) = \begin{cases} -1, & \mathbf{z}_i^\top \mathbf{z}_j = 1 \text{ and } (i, j) \text{ is ML} \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

Here ML means must link. This prior encourages similar values of indicator variables  $z_i$  and  $z_j$  if they happen to share a “must link” edge. Since the MRF prior is assigned only on the indicator variables  $\mathbf{Z}$ , it only alters Eq. (10b) and the new prior on  $\mathbf{Z}$  is given by

$$\mathbf{Z} | \mathbf{v} \sim \frac{1}{\mathcal{Z}} \exp \left( - \sum_{(i,j) \in \mathcal{C}} H(\mathbf{z}_i, \mathbf{z}_j) \right) \times \prod_{n=1}^N \prod_{k=1}^K \left\{ v_k \prod_{j=1}^{k-1} (1 - v_j) \right\}^{z_{nk}} \quad (12)$$

### 3.3 Variational inference

Let us consider all the unknown parameters in our model as latent variables and denote all the latent variables by  $\mathbf{H} = \{\mathbf{Z}, \mathbf{v}, \{\boldsymbol{\beta}^{(k)}\}, \boldsymbol{\tau}, \{\boldsymbol{\alpha}^{(k)}\}, \{\boldsymbol{\gamma}^{(k)}\}, \boldsymbol{\lambda}\}$ . Moreover, from now on, we will ignore feature variables  $\mathbf{X}$  from the list of conditioning variables as they are observed. Using Jensen’s inequality, we can find a lower-bound of the log-marginal  $\ln p(\mathbf{y})$  which is given as

$$\ln p(\mathbf{y}) > \int q(\mathbf{H}) \ln \left\{ \frac{p(\mathbf{y}, \mathbf{H})}{q(\mathbf{H})} \right\} d\mathbf{H} \quad (13)$$

For any arbitrary distribution  $q(\mathbf{H})$ . The variational inference is performed by restricting  $q(\mathbf{H})$  within a parametric family so that the maximization of the lower bound given in Eq. (13) is tractable. We consider only those  $q(\mathbf{H})$  which factorize over some disjoint groups of the component random variables of  $\mathbf{H}$  in the following way

$$q(\mathbf{H}) = \prod_{j=1}^L q_j(\mathbf{h}_j) \quad (14)$$

We can now maximize the lower bound given in Eq. (13) with respect to each components  $q_j(\mathbf{h}_j)$  in Eq. (14) and obtain the parametric form of  $q_j(\mathbf{h}_j)$  given by

$$q_j^*(\mathbf{h}_j) = \frac{\exp(E_{i \neq j}[\ln p(\mathbf{y}, \mathbf{H})])}{\int \exp(E_{i \neq j}[\ln p(\mathbf{y}, \mathbf{H})]) d\mathbf{h}_j} \quad (15)$$

where the expectation is taken with respect to all the other factors  $\{q_i\}$  for  $i \neq j$ . It can be shown that the  $q(\mathbf{H})$  obtained this way is the closest approximation of the actual posterior  $p(\mathbf{H}|\mathbf{y})$  in terms of KL-divergence out of all possible alternatives of the form Eq. (14). Therefore this is a deterministic but approximate posterior inference method unlike stochastic inference methods like MCMC which samples from the actual posterior. However, variational inference is much faster and approximates the true posterior reasonably well for practical purposes.

Once we apply Eq. (15) to the joint distribution described in Eqs. (9) and (10), we can get the update equations for the approximate posterior distributions for each of the latent variables involved.

### 1. Distribution of $\mathbf{z}$ :

$$q_{\mathbf{Z}}(\mathbf{Z}) = \prod_{V \in \mathcal{V}} \left[ \frac{1}{Z_V} \exp \left( - \sum_{\substack{(i,j) \in C \\ i,j \in V}} H(\mathbf{z}_i, \mathbf{z}_j) \right) \prod_{n \in V} \prod_{k=1}^K \rho_{nk}^{\mathbf{z}_{nk}} \right] \quad (16)$$

with

$$\rho_{nk} = \frac{r_{nk}}{\sum_k r_{nk}} \quad (17)$$

$$\begin{aligned} \ln r_{nk} = & \frac{1}{2} \langle \ln \tau_k \rangle - \frac{1}{2} \ln 2\pi - \frac{\langle \tau_k \rangle}{2} \left( y_n^2 - 2 \langle \boldsymbol{\beta}^{(k)} \rangle^\top \mathbf{x}_n y_n + \mathbf{x}_n^\top \langle \boldsymbol{\beta}^{(k)} \rangle \langle \boldsymbol{\beta}^{(k)} \rangle^\top \mathbf{x}_n \right) \\ & + \langle \ln v_k \rangle + \sum_{j=1}^{k-1} \langle \ln(1 - v_j) \rangle \end{aligned} \quad (18)$$

5

2. Distribution of  $\{\boldsymbol{\beta}^{(k)}\}$ :

$$q_{\boldsymbol{\beta}}(\{\boldsymbol{\beta}^{(k)}\}) = \prod_{k=1}^K \mathcal{N}(\{\boldsymbol{\beta}^{(k)}\}; \boldsymbol{\mu}_k, \Sigma^{(k)}) \quad (19)$$

with

$$\Sigma^{(k)} = \left( \langle \tau_k \rangle \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \mathbb{E}[\mathbf{Z}]_{nk} + \langle \tau_k \rangle \text{diag}(\langle \boldsymbol{\alpha}^{(k)} \rangle) \right)^{-1} \quad (20)$$

$$\boldsymbol{\mu}_k = \Sigma^{(k)} \left( \sum_{n=1}^N \mathbf{x}_n y_n \mathbb{E}[\mathbf{Z}]_{nk} \right) \langle \tau_k \rangle \quad (21)$$

10

Here  $\text{diag}(\langle \boldsymbol{\alpha}^{(k)} \rangle)$  corresponds to the LASSO (Tibshirani, 1994) shrinkage. The moments are given by<sup>1</sup>

$$\langle \boldsymbol{\beta}^{(k)} \rangle = \boldsymbol{\mu}_k; \quad \left\langle \left( \boldsymbol{\beta}_p^{(k)} \right)^2 \right\rangle = \Sigma_{pp}^{(k)} + \mu_{kp}^2$$

$$\langle \boldsymbol{\beta}^{(k)} \rangle \langle \boldsymbol{\beta}^{(k)} \rangle^\top = \Sigma^{(k)} + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top$$

15

<sup>1</sup> $\langle f(s) \rangle$  means expected value of  $f(s)$  with respect to the distribution of  $s$ .

### 3. Distribution of $\boldsymbol{\tau}$ :

$$q_{\boldsymbol{\tau}}(\boldsymbol{\tau}) = \prod_{k=1}^K \text{Ga}(\tau_k; c_k, d_k) \quad (22)$$

with

$$c_k = c_0 + \frac{1}{2} \left( \sum_{n=1}^N \text{E}[\mathbf{Z}]_{nk} + p \right) \quad (23)$$

$$d = d_0 + \frac{I}{2} + \frac{J}{2} \quad (24)$$

where

$$I = \sum_{n=1}^N \left( y_n^2 \text{E}[\mathbf{Z}]_{nk} - 2 \text{E}[\mathbf{Z}]_{nk} \mathbf{x}_n^\top y_n \langle \boldsymbol{\beta}^{(k)} \rangle + \text{E}[\mathbf{Z}]_{nk} \mathbf{x}_n^\top \langle \boldsymbol{\beta}^{(k)} \rangle (\boldsymbol{\beta}^{(k)})^\top \mathbf{x}_n \right)$$

$$J = \sum_{p=1}^D \left\langle \alpha_p^{(k)} \right\rangle \left\langle \left( \beta_p^{(k)} \right)^2 \right\rangle$$

The relevant moments are

$$\langle \tau_k \rangle = c_k / d_k \quad \text{and} \quad \langle \ln \tau_k \rangle = \psi(c_k) - \ln(d_k)$$

### 4. Distribution of $\boldsymbol{v}$

$$q_{\boldsymbol{v}}(\boldsymbol{v}) = \prod_{k=1}^K \text{Beta}(v_k; \xi_k, \kappa_k) \quad (25)$$

with

$$\xi_k = 1 + \sum_{n=1}^N \mathbb{E}[\mathbf{Z}]_{nk} \quad \text{and} \quad \kappa_k = \langle \lambda \rangle + \sum_{j=k+1}^K \sum_{n=1}^N \mathbb{E}[\mathbf{Z}]_{nj}$$

Relevant moments are given by  $\langle \ln v_k \rangle = \psi(\xi_k) - \psi(\xi_k + \kappa_k)$  and  $\langle \ln(1 - v_k) \rangle = \psi(\kappa_k) - \psi(\xi_k + \kappa_k)$ .

5 5. Distribution of  $\{\alpha^{(k)}\}$ :

$$q_{\alpha}(\{\alpha^{(k)}\}) = \prod_{k=1}^K \prod_{p=1}^D \text{InvGaussian}(\alpha_p^{(k)}; g_p^k, h_p^k) \quad (26)$$

with

$$g_j^k = \sqrt{\frac{\langle \gamma_j^{(k)} \rangle}{\langle \tau_k \rangle \langle (\beta_j^{(k)})^2 \rangle}}$$

$$h_j^k = \langle \gamma_j^{(k)} \rangle$$

10 where  $\text{InvGaussian}(\alpha_j^{(k)}; g_j^k, h_j^k)$  denotes inverse Gaussian distribution with mean  $g_j^k$  and shape parameter  $h_j^k$  having the following density function.

$$p_{\text{IG}}(\alpha_j^{(k)}; g_j^k, h_j^k) = \sqrt{\frac{h_j^k}{2\pi (\alpha_j^{(k)})^3}} \times \exp\left(-\frac{h_j^k (\alpha_j^{(k)} - g_j^k)^2}{2 (g_j^k)^2 \alpha_j^{(k)}}\right) (\alpha_j^{(k)} > 0)$$

The relevant moments are given by

$$\langle \alpha_j^{(k)} \rangle = g_j^k \quad \text{and} \quad \langle (\alpha_j^{(k)})^{-1} \rangle = (g_j^k)^{-1} + (h_j^k)^{-1}$$

15

6. Distribution of  $\{\gamma^{(k)}\}$ :

$$q_{\gamma}(\{\gamma^{(k)}\}) = \prod_{p=1}^D \text{Ga}(\gamma_j^{(k)}; a_j^k, b_j^k) \quad (27)$$

with

$$a_j^k = a_0 + 1$$

$$b_j^k = b_0 + \frac{1}{2} \left\langle \left( \alpha_j^{(k)} \right)^{-1} \right\rangle$$

and the relevant moment is  $\langle \gamma_j^{(k)} \rangle = a_j^k / b_j^k$

7. Distribution of  $\lambda$ :

$$q_{\lambda}(\lambda) = \text{Ga}(\lambda; u, w) \quad (28)$$

where

$$u = m_0 + K; \quad w = - \sum_{k=1}^K \langle \ln(1 - v_k) \rangle$$

Relevant moment is  $\langle \lambda \rangle = \frac{u}{w}$ .

The first part of the variational posterior of  $q_Z(\mathbf{Z})$  in Eq. (16) arises from the MRF prior and contributes towards enforcing “must link” constraints. Note that  $\mathcal{V}$  in Eq. (16) is a set of sets and  $V$  is a component set of connected nodes within  $\mathcal{V}$ . Basically,  $\mathcal{V}$  denotes the set of connected components within the constraint graph described in Sect. 3.2. Therefore the partition function  $\mathcal{Z}_{\mathcal{V}}$  needs to be computed only for the connected components, not for the



entire graph. Computing  $\mathcal{Z}_V$  becomes tractable if the connected components are small, i.e. the constraint set is sparse.

In order to automatically generate a sparse constraints set, we first implemented all the constraints in form of edges and then used a graph partitioning algorithm (Hespanha, 2004) to partition the constraint graph in such a way that none of the partitions are left with more than a predefined number of nodes. At the time of inference we used a “backtracking” algorithm (Tarjan, 1972) to find the strongly connected components within the graph. To compute the expectation  $E[\mathbf{z}]$ , we first computed the multinomial probabilities  $\rho_{nk}$  and then did an MRF update on each connected component by computing the probabilities of each possible state combination and summing the probability-weighted state matrices. The partition function is computed by summing the exponentiated sum of energy function of each state matrix. Note that, isolated nodes (not part of the any connected components) will not need their  $\rho_{nk}$  updated.

The parameters of each of the distributions has dependency on moments of one or more of the other variables. We therefore find a locally optimum solution via an iterative process that starts with random initial values of the relevant moments and stops when the indicator variables  $\mathbf{Z}$  stop changing any more. Note that, once the approximate solution is reached, we can compute the marginal distributions over coefficients  $\beta_p^{(k)}$  which is a Gaussian with mean  $\mu_p^{(k)}$  and variance  $\Sigma_{pp}^{(k)}$  for each  $k$ . We can thereby perform a  $t$  test to determine whether the corresponding feature has a non-zero coefficient.

### 3.4 Computational considerations

One computational bottleneck of the proposed VB algorithm is the inversion of the  $D \times D$  matrix in Eq. (20). If  $D < N$ , then faster matrix inversion can be achieved by first applying a Cholesky decomposition and then inverting the resulting upper triangular matrix. However if  $D > N$ , we can first apply a fast (approximate) singular value decomposition on  $\Sigma^{(k)^{-1}}$  and then use Woodbury matrix inversion identity so that we now have to invert a  $N \times N$  matrix instead.

We have truncated the infinite Dirichlet Process at  $K = 20$  for most of our experiments. The speed of the algorithm can be further improved by parallelizing the updates for each of  $K$  components which is straightforward as they are updated independent of each other. Another major computational challenge was the MRF updates. Apart from controlling the maximum size of the connected components, we parallelized the MRF updates over each subgraph by making the state generation independent of the previous state.

## 4 Experiments

We have evaluated our method both on synthetic and climate datasets. Typical values used for the hyper-parameters were  $a_0 = b_0 = c_0 = d_0 = 0.01$  and  $\lambda = 1$ . Selecting these values within a reasonable range does not affect the results significantly. We made sure that the cardinality of the largest connected component in the constraints graph never exceeds 8.

### 4.1 Synthetic dataset

We compared the performance of both constrained and unconstrained version of our method with non-parametric mixture of linear regression (NPMLR) model without any regularization. We set-up three experiments: (1) to test whether or not our algorithm can learn the number of clusters; (2) to evaluate the effect of constraints; and (3) to check the sensitivity of our approach to noise.

For all our experiments involving synthetic data, we used  $N = 1000$  data points and  $D = 30$  features. In our first set of experiments we tested our method for  $K = 2 \dots 5$  actual clusters. Each column of the  $N \times D$  input matrix  $\mathbf{X}$  is generated from a uniform distribution. For each value of  $K$ , we partitioned the input matrix  $\mathbf{X}$  in  $K$  equal parts  $\mathbf{X}_1 \dots \mathbf{X}_K$ . Then for each partition  $\mathbf{X}_k$  ( $k = 1 \dots K$ ), we generate sparse coefficients  $\beta_k$  by randomly selecting 10 out of 30 components to be non-zero. We assign a value of  $5k$  (where  $k$  is the index of the cluster,  $k = 1, \dots, K$ ) to the non-zero components within the  $k$ -th cluster so that two clusters are distinctly identifiable in case the indices of non-zero components of the clusters are the same. We then generate the output  $\mathbf{y}_k$  for the  $k$ th cluster using the linear regression model

of (1). The fixed noise variance  $\tau_k^{-1}$  for the first experiment was generated by randomly choosing a number between 0 and 0.1 to introduce diversity. A final dataset was obtained by merging  $\{\mathbf{X}_k, \mathbf{y}_k\}$  for all  $k = 1 \dots K$ . The process is repeated 30 times and mean and variance of the evaluation metrics were reported in the form of errorbars for each value of  $K$  in Fig. 2. For all these experiments, the total number of constraints were kept at 20 per cluster while the size of the largest subgraph was kept below 7.

The second experiment was performed to evaluate the effect of number of “must link” constraints on the performance of the constrained version of the algorithm. Here, the actual number of clusters were fixed at  $K = 3$  along with the base noise variance (0.1) and the number of constraints per clusters were varied from 0 to 30 incremented by 5, although the actual number of constraints may be less since we removed some constraints to achieve sparsity in the constraint graph. The result is reported in Fig. 3.

In our third experiment, we evaluated the effect of noise on the performance of our algorithm. Again, we kept the number of clusters fixed at  $K = 3$  and number of constraints fixed at 20 per cluster (for the constrained version). We varied the base noise level in each cluster from 0 to 0.5 and added a randomly generated value between 0 to 0.1 with the base noise level for each cluster to maintain diversity among the clusters. Average and variance of 30 repetitions are reported in Fig. 4.

#### 4.1.1 Evaluation metrics

We measured two aspects of the performance of our algorithm. First, whether it can cluster the datapoints correctly. We put a data point into one of the possible 20 components (since we truncated the infinite Dirichlet process at  $K = 20$  for all experiments) depending on the value of the row  $\mathbf{E}[\mathbf{Z}]_n$  (a vector) in the  $N \times 20$  matrix  $\mathbf{E}[\mathbf{Z}]$  estimated by the variational inference algorithm. The estimated cluster membership  $\hat{c}_n$  (a scalar) is given by  $\hat{c}_n = \underset{k}{\operatorname{argmax}} \mathbf{E}[\mathbf{Z}]_{nk}$ . We retain all the valid components out of 20 possible, which have at least one member initially. Then we run an update algorithm to merge very small clusters with the closest larger ones. Note that, the estimated cluster indices (a value between 1 and

20) may not correspond directly to the actual cluster indices (a value between 1 to actual value of  $K$ ) since the variational inference algorithm is not aware of the actual order of the cluster indices (e.g. actual cluster index 1 may correspond to estimated cluster index 9). So we use a metric called normalized mutual information (NMI) that evaluates the match between estimated cluster memberships  $\hat{c}$  and actual ones  $c$  without needing there be direct correspondence. NMI is given by  $\text{NMI}(c, \hat{c}) = \frac{H(c) - H(c|\hat{c})}{\sqrt{H(c)H(\hat{c})}}$ , where  $H(\cdot)$  is the entropy. Higher NMI values mean that the clustering results are more similar to ground-truth. The metric reaches its maximum value of one when there is perfect agreement.

A second metric is used to evaluate the quality of the sparse regression model estimated within each discovered cluster. Here we are only interested in finding whether our algorithm picks the non-zero coefficients correctly. We use  $F$  score to measure the match between actual and estimated non-zero coefficients within each cluster.  $F$  score for the  $k$ th component is given by  $F_k = \frac{2P_kR_k}{P_k+R_k}$  where  $P_k$  is the precision and  $R_k$  is the recall of the estimated coefficients for the  $k$ th component. We reported the average of  $F_k$  values over all components discovered by our algorithm. Unlike the previous metric, here we need to know the direct correspondence between the cluster indices so that we can match the actual and estimated coefficient vectors. We developed an algorithm to find such a correspondence.

#### 4.1.2 Discussion of results

We can see the performance of all three algorithms are comparable in terms of identifying the clusters correctly, although the NMI value of NPMLR degrades significantly for  $K = 5$ . However, as desired, our method outperforms NPMLR in terms of correctly retrieving the sparse structure of regression coefficients within each cluster. There is general downward trend of performance for all algorithms with increasing number of actual components in the data. This is an inherent problem with the DPM models as it tends to attach each new data-point to the largest current component and thereby favoring models with less components. Also, as the number of actual components grow, the probability of two components to be similar increases.

The increased flexibility of non-parametric methods comes at a cost of hitting local optima more likely and finding solutions that are not interpretable. Adding more constraints may decrease this probability but at the same time restricts the variational method from finding solutions leading to a larger lower bound, especially in the presence of more components in the data. Therefore increasing the number of constraints may result more interpretable solutions, but not improved accuracy. It is also encouraging to see that our method is relatively robust to added noise, a major challenge with the real datasets, especially in terms of correctly identifying the sparse structure.

## 4.2 Feature selection for downscaling rainfall

A grand challenge in climate science relevant for adaptation and policy remains our inability to provide credible stakeholder-relevant “statistical downscaling”, or developing statistical techniques for more accurate, precise and interpretable high-resolution projections with lower-resolution climate model data (Benestad et al., 2008). Regression models of statistical downscaling (Benestad et al., 2008; Ghosh, 2010) works by first selecting a set of climate variables that have information about the target variable, and then fitting a regression model to predict the target variable at higher resolution. In this application, selecting the right set of predictors are as important as building a prediction model since even a good prediction with a model that is physically not interpretable is less desirable as it may not generalize well. We focus on the feature selection problem for statistical downscaling of annual average rainfall. The use of annual averages reduce the amount of noise in the observed rainfall data, which enables us to examine the robustness of our methods with less ambiguity.

Existence of multiple states or patterns is acknowledged in regression-based statistical downscaling literature for rainfall (e.g. Kannan and Ghosh, 2010) where parametric methods like  $K$ -means was used to find distinct clusters. Here we used our model to simultaneously find clusters, if any, and select features for the purpose of statistical downscaling of station observed annual average rainfall over two climatologically homogeneous regions over the

continental US. Figure 5 shows the climatologically homogeneous regions over United States.

Since rainfall follows a log-normal distribution (), the target variable we used is logarithm of annual average rainfall. In Fig. 6, we have shown the distribution of average rainfall over all sites in western US before and after taking logarithm.

Potential features used can fall under one of two broad categories – local atmospheric variables and large-scale climate indices. Local covariates originate from each station and exhibits both spatial and temporal variability. Annual and seasonal averages of maximum temperature falls in this category along with elevation, sea level pressure (SLP) and convective available potential energy (CAPE). A dependence on any of these variables roughly indicates dominance of local convective rainfall in the region. Daily rainfall station data obtained from United States Historical Climatology Network (USHCN) (Easterling et al., 1996). All other features are described in Table 1.

Climate indices are global variables that represent large-scale signals in climate variables. A list of covariates used for each category is given in Table 1. A dependence on any of these variables roughly indicates rainfall due to large-scale circulation.

We could use the covariates between 1979 and 2011 as SLP and CAPE is available only for that period. Also, if more than 50% of the daily observations out of a year are found to be missing for any of variables at a specific location, we simply discarded all variables for that year and for that specific location. We averaged monthly climate indices and daily local variables over a year. Finally the annual/seasonal average time-series of predictors for each station were merged for a homogeneous region under consideration. West (CA,NV) and Northwest (WA, OR, ID) regions are shown by gray shaded areas over the US map in Fig. 7 (left and right panels, respectively).

#### 4.2.1 Results and discussion

We applied spatial “must-link” constraints among pairs of data-points belonging from the same location. Ideally, if there are  $n$  points in a cluster, we will be required to put  $\binom{n}{2}$  constraints to cover all pairs of data-points. To reduce complexity, initially we kept only

those constraints that connect data-points from consecutive years. However, this reduced set of constraints proved to be too restrictive and all data-points tended to merge into a single cluster. So, we kept removing the constraints in an intuitive manner until more than one cluster emerged for a region. We found more than one cluster for all regions except the southern region. We stopped removing constraints until new clusters stopped emerging for a region. Here we show only the clusters in western and northwestern regions, since majority of stations were mostly split into obtained clusters in these regions. In other regions, almost all stations had mixed membership. We assign a station to a cluster if more than 80 % of its data-points belong to that cluster.

A quick look at the histogram of target variable (right panel in Fig. 6) also supports the possibility of two distinct rainfall modes in the region. As mentioned earlier, we obtained one sparse linear model for each of the discovered components within a region. Since a non-zero coefficient in the sparse model implies dependence on the corresponding covariate, we can obtain interesting insights about the dependence of average rainfall on various atmospheric and climate indices from the coefficients of the individual sparse models within each cluster. Interestingly in the north-west region there is only a single member station in the first component which exhibits dependence on the local temperature variables and SLP whereas the larger cluster shows dependence on a larger number of climate indices. In the western region, the first cluster shows dependence on local temperature variables and the second cluster shows more dependence on large scale variables. Both clusters show dependence on elevation. While dependence on large scale indices are not surprising for both these regions due to the known effect of westerlies in these coastal regions, dependence of smaller clusters (especially in the northwest) on local variables may hint towards existence of some regional small scale atmospheric mechanisms. While spatially coherent clusters are more likely to occur in nature, geographical features such as mountains and lakes and even man-made structures such as large dams and reservoirs may abruptly disturb the spatial smoothness of clusters, since their presence may alter the climate pattern of the nearby areas with respect to the surrounding regions. However, before we can build statistical downscaling models, more rigorous statistical and physical analysis

is required based on these preliminary insights obtained using our method. The clusters discovered here, and the corresponding covariates, can be utilized to develop individual non-linear prediction models per cluster.

DPMs automatically find the number of clusters  $K$  and adapts to varying values of  $K$ . However, DPMs prevent the model to “learn” an unnecessarily large value of  $K$ , if a smaller  $K$  is sufficient to describe the model, thus managing complexity. Based on the results of experiments on the synthetic dataset shown in figure 2, we found that the performance of the method degrades as the number of components  $K$  grows larger. We believe it is reasonable to expect that there will only be a limited number of distinct relationships between average rainfall and their covariates when we apply our method at the regional scale. However, even in situations where a large number of relationships exist within a particular region, our method may not be able to identify all of the distinct methods, but it can nevertheless be expected to outperform the use of a single model. The single model will attempt to learn a relationship that is the average of all distinct relations, which our approach will still attempt to distinguish among major categories of relationships even though some of them may be lumped together.

## 5 Conclusions

In this paper, we proposed a nonparametric Bayesian mixture of sparse regression models for simultaneous clustering and discovery of covariates within each cluster using Dirichlet process mixture model. Moreover, our model can accommodate prior knowledge about “must link” constraints between the pair of data-points using a Markov Random Field prior on the cluster membership variables. Our major contribution is to develop an efficient and scalable variational inference algorithm for inference on the fully Bayesian model. We applied our method both on synthetic and real climate data and successfully discovered multiple underlying behaviors in the data. Preliminary results of applying our method towards feature selection for statistical downscaling of rainfall shows promise towards finding new climate insights with appropriate caveats. Going forward, we would like to



incorporate priors for diversity among the clusters in order to discourage merging of close but dissimilar clusters. We intend to further extend our model for predictive analysis and build a full-scale statistical downscaling method using the features selected by the current model.

*Acknowledgements.* This work was funded by the NSF Expeditions in Computing grant “Understanding Climate Change: A Data Driven Approach”, award number 1029166.

## References

- Antoniak, C.: Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems, *Ann. Stat.*, 2, 1152–1174, 1974. 4, 7
- 5 Basu, S., Bilenko, M., Banerjee, A., and Mooney, R.: Probabilistic semi-supervised clustering with constraints, *Journal Machine Learning Research*, 71–98, 2006. 5
- Benestad, R., Hanssen-Bauer, I., and Chen, D.: *Empirical-Statistical Downscaling*, World Scientific Publishing Company, New Jersey, London, 2008. 21
- 10 Bishop, C. and Svenskn, M.: Bayesian hierarchical mixtures of experts, in: *Uncertainty in Artificial Intelligence*, 57–64, Morgan Kaufman, San Francisco, CA, 2002. 4
- Blei, D. M. and Jordan, M. I.: Variational inference for Dirichlet process mixtures, *Bayesian Analysis*, 1, 121–143, 2006. 5, 8
- Climate Indices: Monthly Atmospheric and Ocean Time Series, available at: <http://www.esrl.noaa.gov/psd/data/climateindices/list/>, retrieved: 2 April 2014. 29
- 15 Easterling, D. R, Karl, T. R., Mason, E. H., Hughes, P. Y. and Bowman, D. P. United States Historical Climatology Network (USHCN) Monthly Temperature and Precipitation Data, Tech. rep., Oak Ridge National Laboratory, US Department of Energy, Oak Ridge, Tennessee, available at: <http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn.html>, retrieved: 2 April 2014. 22, 29
- 20 Ebtehaj, A. M., Fofoula-Georgiou, E., and Lerman, G.: Sparse regularization for precipitation downscaling, *J. Geophys. Res.*, 117, 1–12, 10.1029/2011JD017057, 2012. 4
- Ferguson, T.: A Bayesian analysis of some nonparametric problems, *Ann. Stat.*, 1, 209–230, 1973. 7
- Geman, S. and Geman, D.: Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images, *IEEE T. Pattern Anal.*, PAMI-6, 721–741, 1984. 7
- 25 Ghosh, S.: SVM-PGSL coupled approach for statistical downscaling to predict rainfall from GCM output, *J. Geophys. Res.*, 115, D22102, 10.1029/2009JD013548, 2010. 21
- Hespanha, J. P.: An efficient MATLAB Algorithm for Graph Partitioning, Tech. rep., University of California, Santa Barbara, available at: <http://www.ece.ucsb.edu/~hespanha/techrep.html>,  
30 retrieved: 2 April 2014. 17

- Ishwaran, H. and James, L. F.: Gibbs sampling methods for stick-breaking priors, *J. Am. Stat. Assoc.*, 96, 161–173, 2001. 4, 8
- Kannan, S. and Ghosh, S.: Prediction of daily rainfall state in a river basin using statistical downscaling from GCM output, *Stoch. Env. Res. Risk. A.*, 25, 457–474, 10.1007/s00477-010-0415-y, 2010. 4, 21
- 5 Knutti, R. and Sedláček, J.: Robustness and uncertainties in the new CMIP5 climate model projections, *Nature Climate Change*, 3, 369–373, 10.1038/nclimate1716, 2013. 3
- Kumar, D., Kodra, E., and Ganguly, A. R.: Regional and seasonal intercomparison of CMIP3 and CMIP5 climate model ensembles for temperature and precipitation, *Clim. Dynam.*, online first, 10.1007/s00382-014-2070-3, 2014. 3
- 10 Mesinger, F., DiMego, G., Kalnay, E., Mitchell, K., Shafran, P. C., Ebisuzaki, W., Jović, D., Woollen, J., Rogers, E., Berbery, E. H., Ek, M. B., Fan, Y., Grumbine, R., Higgins, W., Li, H., Lin, Y., Manikin, G., Parrish, D., and Shi, W.: North American regional reanalysis, *Bull. Am. Meteorol. Soc.*, 87, 343–360, 10.1175/BAMS-87-3-343, 2006. 29
- 15 Park, T. and Casella, G.: The Bayesian LASSO, *J. Am. Stat. Assoc.*, 103, 681–686, 10.1198/016214508000000337, 2008. 5, 6
- Phatak, A., Bates, B., and Charles, S.: Statistical downscaling of rainfall data using sparse variable selection methods, *Environ. Modell. Softw.*, 26, 1363–1371, 10.1016/j.envsoft.2011.05.007, 2011. 4
- 20 Ross, J. and Dy, J.: Nonparametric Mixture of Gaussian Processes with Constraints, *International Conference of Machine Learning*, 28, 2013. 5
- Sethuraman, J.: A constructive definition of Dirichlet priors, *Stat. Sinica*, 4, 639–650, 1994. 8
- Tarjan, R.: Depth-first search and linear graph algorithms, *SIAM J. Comput.*, 1, 146–160, 1972. 17
- Tibshirani, R.: Regression shrinkage and selection via the LASSO, *J. Roy. Stat. Soc. B*, 58, 267–288, 1994. 4, 13
- 25 Wilby, R., Charles, S., Zorita, E., and Timbal, B.: Guidelines for use of climate scenarios developed from statistical downscaling methods, *Tech. Rep. August*, Inter-Governmental Panel for Climate Change, available at: <http://www.narccap.ucar.edu/doc/tgica-guidance-2004.pdf>, retrieved: 2 April 2014. 3
- 30 Bader, D.C., C. Covey, W.J. Gutkowski, Jr., I.M. Held, K.E. Kunkel, R.L. Miller, R.T. Tokmakian, and M.H. Zhang, 2008: *Climate Models: An Assessment of Strengths and Limitations*. U.S. Climate Change Science Program Synthesis and Assessment Product 3.1. Department of Energy, Office

of Biological and Environmental Research, 124 pp. available at: [http://pubs.giss.nasa.gov/docs/2008/2008\\_Bader\\_etal\\_1.pdf](http://pubs.giss.nasa.gov/docs/2008/2008_Bader_etal_1.pdf) retrieved: July 20, 2014. 3

Greene, A. M., A. W. Robertson, P. Smyth, and S. Triglia, 2011: Downscaling projections of Indian monsoon rainfall using a nonhomogeneous hidden Markov model. *Quart. J. Royal Meteor. Soc.*, 137, 347-359 4

Kedem, B., and L. S. Chiu, On the lognormality of rain rate, *Proc. Natl. Acad. Sci. U.S.A.*, 84, 901-905, 1987

22

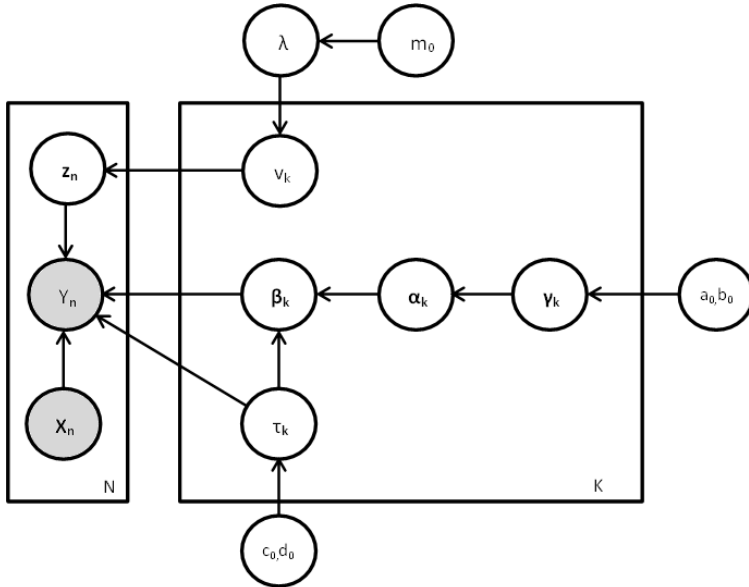
**Table 1.** Potential features used for statistical downscaling of rainfall.

---

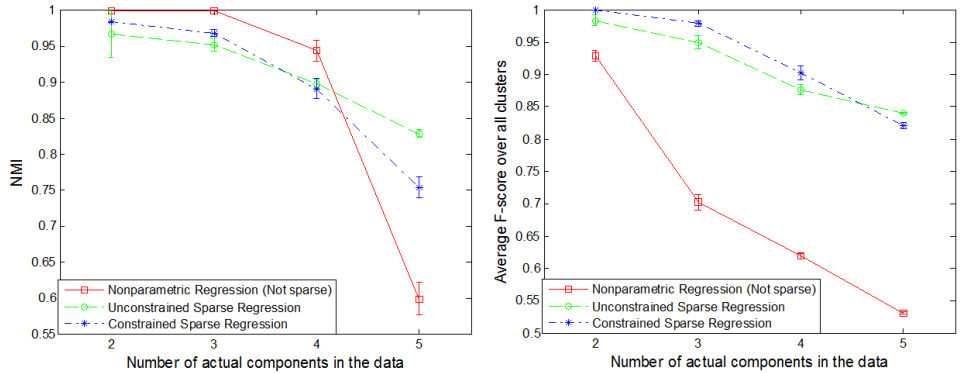
**Atmospheric: (Easterling et al., 1996; Mesinger et al., 2006)** Mean Annual Maximum Temperature (MATmax), Mean Winter Maximum Temperature (DJFTmax), Mean Spring Maximum Temperature (MAMTmax), Mean Summer Maximum Temperature (JJATmax), Mean Autumn Maximum Temperature (SONTmax) (Easterling et al., 1996), Sea Level Pressure (SLP), Convective Available Potential Energy (CAPE) (Mesinger et al., 2006)

**Climate Indices: (NOAA)** North Atlantic Oscillation (NAO), East Atlantic Pattern (EA), West Pacific Pattern (WP), East Pacific/North Pacific Pattern (EPNP), Pacific/North American Pattern (PNA), East Atlantic/West Russia Pattern (EAWR), Scandinavia Pattern (SCA), Tropical/Northern Hemisphere Pattern (TNH), Polar/Eurasia Pattern (POL), Pacific Transition Pattern (PT), Nino 1+2, Nino 3, Nino 3.4, Nino 4, Southern Oscillation Index (SOI), Pacific Decadal Oscillation (PDO), Northern Pacific Oscillation (NP), Tropical/Northern Atlantic Index (TNA), Tropical/Southern Atlantic Index (TSA), Western Hemisphere Warm Pool (WHWP), Global Mean Temperature Anomaly (GlobalMeanTemp) (NOAA)

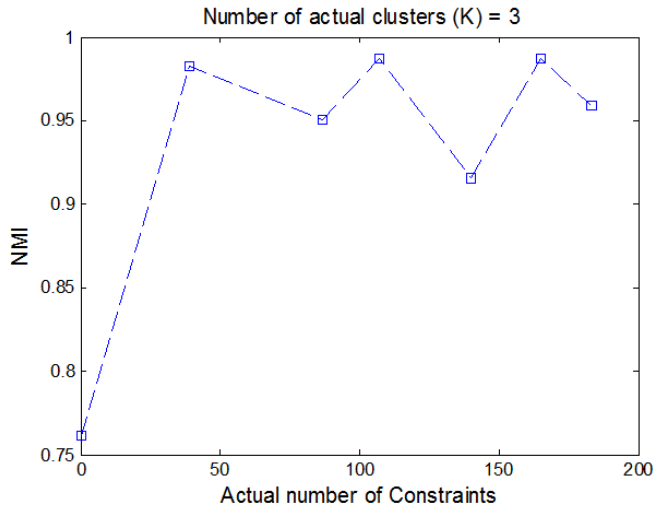
---



**Figure 1.** Graphical representation of the of complete Bayesian hierarchical model.

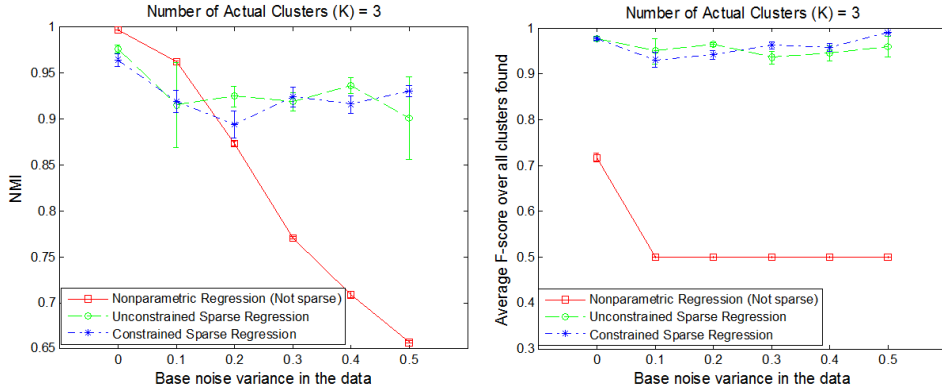


**Figure 2.** Left: ability of nonparametric unregularized and sparse regressions (unconstrained and constrained) to correctly identify clusters in presence of increased number of actual components in the data. Right: ability of nonparametric unregularized and sparse regressions (unconstrained and constrained) to correctly retrieve the sparse structure within each cluster.

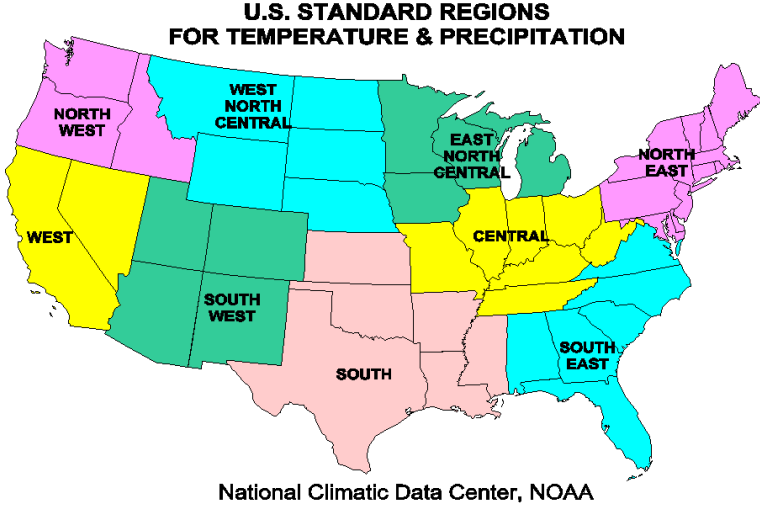


**Figure 3.** Performance of the constrained version of the algorithm (in terms of NMI (more the better)) with number of “must link” constraints.

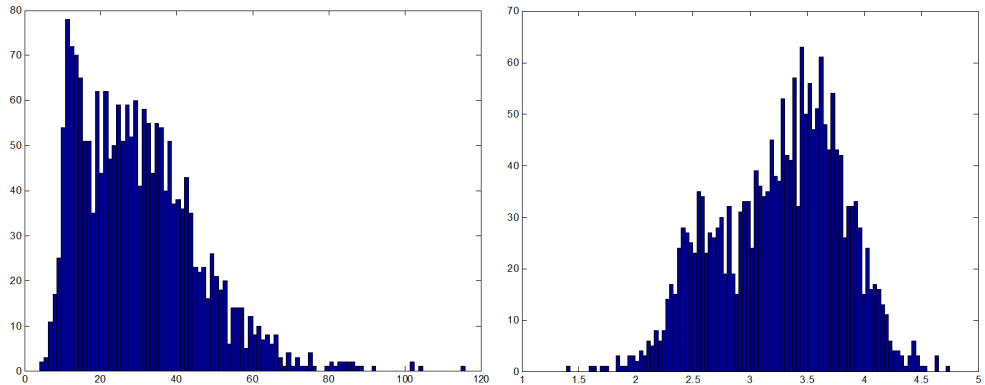




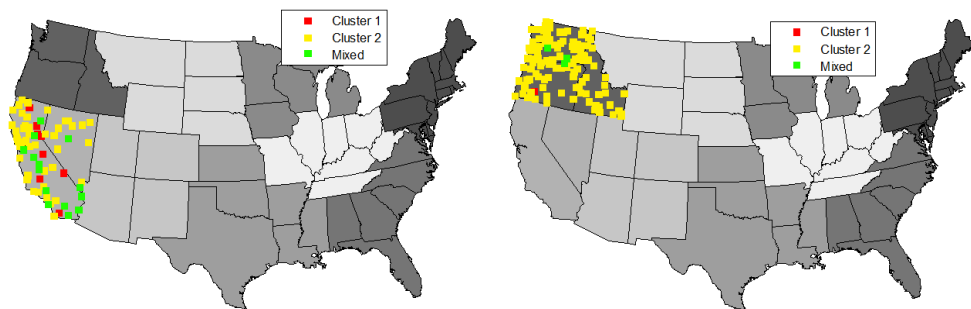
**Figure 4.** Left: ability of nonparametric unregularized and sparse regressions (unconstrained and constrained) to correctly identify clusters (indicated by NMI) with increasing noise. Right: ability of nonparametric unregularized and sparse regressions (unconstrained and constrained) to correctly retrieve the sparse structure within each cluster (indicated by average  $F$  score).



**Figure 5.** Map showing climatologically homogeneous regions over continental United States.



**Figure 6.** (Left) Distribution of average rainfall over all sites in the western US. (Right) Distribution of average rainfall after transformation.



**Figure 7.** Left: location of stations and their cluster membership in the western region. Right: location of stations and their cluster membership in the northwestern region.