



# Bayesian Inference Using the Proximal Mapping: Uncertainty Quantification Under Varying Dimensionality

Maoran Xu<sup>a</sup>, Hua Zhou<sup>b</sup>, Yujie Hu<sup>c</sup>, and Leo L. Duan<sup>d</sup>

<sup>a</sup>Department of Statistical Science, Duke University, Durham, NC; <sup>b</sup>Departments of Biostatistics and Computational Medicine, University of California, Los Angeles, CA; <sup>c</sup>Department of Geography, University of Florida, Gainesville, FL; <sup>d</sup>Department of Statistics, University of Florida, Gainesville, FL

### **ABSTRACT**

In statistical applications, it is common to encounter parameters supported on a varying or unknown dimensional space. Examples include the fused lasso regression, the matrix recovery under an unknown low rank, etc. Despite the ease of obtaining a point estimate via optimization, it is much more challenging to quantify their uncertainty. In the Bayesian framework, a major difficulty is that if assigning the prior associated with a p-dimensional measure, then there is zero posterior probability on any lower-dimensional subset with dimension d < p. To avoid this caveat, one needs to choose another dimension-selection prior on d, which often involves a highly combinatorial problem. To significantly reduce the modeling burden, we propose a new generative process for the prior: starting from a continuous random variable such as multivariate Gaussian, we transform it into a varying-dimensional space using the proximal mapping. This leads to a large class of new Bayesian models that can directly exploit the popular frequentist regularizations and their algorithms, such as the nuclear norm penalty and the alternating direction method of multipliers, while providing a principled and probabilistic uncertainty estimation. We show that this framework is well justified in the geometric measure theory, and enjoys a convenient posterior computation via the standard Hamiltonian Monte Carlo. We demonstrate its use in the analysis of the dynamic flow network data. Supplementary materials for this article are available online.

#### **ARTICLE HISTORY**

Received September 2021 Accepted May 2023

### **KEYWORDS**

Concentration of Lipschitz functions; Generalized density; Generalized projection; Hausdorff dimension; Nonexpansiveness

### 1. Introduction

Modern statistical applications often involve data that are high dimensional. To allow signal recovery under a relatively low sample size, one often needs to assume that the parameter  $\theta \in \mathbb{R}^p$  in fact lies in/near some lower dimensional space. Commonly used assumptions include sparsity (Meinshausen and Bühlmann 2006; Zhao, Rocha, and Yu 2006; Meier, Van De Geer, and Bühlmann 2008; Bickel, Ritov, and Tsybakov 2009), low rank (Shen and Huang 2008; Ji et al. 2010), geometric constraints (Goodall and Mardia 1999; Saarela and Arjas 2011), etc. In most cases, the dimensionality d is unknown. For example, we usually do not know the exact rank in the low-rank matrix factorization.

Bayesian framework provides a principled way to quantify the uncertainty on those models. A potential caveat is that if the assigned prior is associated with a p-dimensional continuous measure, then there is zero posterior probability allocated on any of the lower-dimensional subsets with dimension d < p. Instead, from a generative perspective, one should first choose a discrete prior to select d, then generate  $\theta$  within the chosen space. For example, the spike-and-slab prior (Mitchell and Beauchamp 1988) assigns a binomial distribution on d as the number of nonzero coefficients in the variable selection problem; the Bayesian adaptive regression spline uses a Poisson prior on the number of knots d, which determines the rank of the spline matrix (DiMatteo, Genovese, and Kass

2001). On the other hand, the discrete prior creates a highly combinatorial problem, and existing estimation methods such as the Reversible-jump Markov chain Monte Carlo (Green and Hastie 2009) are not very efficient to explore the high posterior probability region.

An appealing alternative is to avoid specifying any low-dimensional prior, but to induce a prior for  $\theta$  with the measure in  $\mathbb{R}^p$  yet having the mass concentrated near some low-dimensional sets. Specifically, the key is to reparameterize the parameter  $\theta$  as some transformation of a sparse vector  $\beta$ , and then assign simple continuous shrinkage prior on  $\beta$  to favor near-zero values. For example, in spline regression, one uses  $\beta$  as the sparse weights in the linear combination of some basis functions. In this category, there is a rich literature covering tasks of variable selection (Park and Casella 2008; Carvalho, Polson, and Scott 2009; Rockova and George 2018), matrix decomposition (Bhattacharya and Dunson 2011; Legramanti, Durante, and Dunson 2020), functional data analysis (Shin, Bhattacharya, and Johnson 2020), covariance estimation (Li, Craig, and Bhadra 2019; Kastner 2019), among others.

Clearly, this strategy has its limitations—when we cannot reparameterize the low-dimensional sets of  $\theta$ , the prior specification becomes awkward. This is not uncommon. For example, the fused lasso (Tibshirani et al. 2005) is a frequentist regularization very popular in the image/signal processing, which assumes sparsity not only in the parameter  $\theta \in \mathbb{R}^p$ , but also in the

(p-1) increments between the neighboring elements  $(\theta_{i+1}-\theta_i)$ . Although we could imagine assigning some shrinkage prior on  $\beta = D\theta \in \mathbb{R}^{2p-1}$  with D the corresponding matrix, such a prior is ill-defined: as *D* is not invertible, we cannot compute  $\theta$  from  $\beta$ ; as  $\beta$  resides in the column span of D, it has a dimension p, which is less than (2p-1)—therefore, the shrinkage prior one blindly assigned would be in fact an incomplete density for a degenerate measure, making it difficult to calibrate the hyper-parameters within and assess the effects of the prior regularization.

Motivated to generalize the Bayesian approaches for handling most of the low-dimensional regularizations (including potentially complicated ones), while avoiding the caveats of having to explicitly specify a discrete prior, we consider a "projection"style approach. Starting from a continuous prior for  $\beta$  with measure in  $\mathbb{R}^p$ , we transform it into  $\theta$  using a special mapping, so that  $\theta$  has an induced prior on several low-dimensional sets. The projection idea was previously considered in several cases, such as the mixture of components with different dimensions (Petris and Tardella 2003), the isotonic regression (Dunson and Neelon 2003), monotone curve fitting (Lin and Dunson 2014) and more generally, constrained space modeling (Sen, Patra, and Dunson 2018). Nevertheless, in this article, we explore a much more general transformation known as the proximal mapping—it not only includes common Euclidean projection to a constrained set, but also useful non-projection transformation such as soft-thresholding, nuclear norm control, set expansion, etc. This mapping has been well studied in the optimization literature, with appealing properties that are convenient for canonical Bayesian inference, such as in the concentration of measure and convenient computation via the Hamiltonian Monte Carlo. We will carefully justify this prior via the geometric measure theory and demonstrate the strengths via several examples.

### 2. Method

### 2.1. Background on the Proximal Mapping

We first provide a brief review on the proximal mapping and motivate its use as a transformation tool. Let  $\theta$  be the parameter of interest in a certain space  $\Theta$ , with  $\Theta \subseteq \mathbb{R}^p$ . With another parameter  $\beta \in \boldsymbol{\beta} \subseteq \mathbb{R}^p$ , the *proximal mapping* is a transform of

$$\theta = \operatorname{prox}_{\lambda g}(\beta) := \underset{z \in \Theta}{\operatorname{arg \, min}} \left\{ \lambda g(z) + \frac{1}{2} \|z - \beta\|_2^2 \right\}, \quad (1)$$

where g is a lower semicontinuous and convex function, and  $\lambda > 0$  is a scalar as a hyper parameter. This effectively induces a parameter space:

$$\Theta_{\lambda g} = \{ \operatorname{prox}_{\lambda g}(\beta) : \beta \in \mathbb{R}^p \}.$$

For an intuitive understanding, the proximal mapping could be viewed as a generalized projection. Given a constrained set C, we can choose  $g(z) = \mathcal{X}_C(z)$ , the characteristic function of a constrained set taking value 0 if  $z \in C$ , or  $\infty$  if  $z \notin C$ . The mapping becomes  $\theta = P_C(\beta) = \arg\min_{z \in C} ||z - \beta||_2^2$ , the Euclidean projection of  $\beta$  into the set C. Furthermore, we can replace  $\mathcal{X}_C$  with other function for g, leading to a wider class of transformation.

Example 1 (Soft thresholding). Perhaps the most famous example is  $g(z) = ||z||_1$  from lasso (Tibshirani 1996). It has a closedform proximal mapping known as the soft-thresholding operator  $\operatorname{prox}_{\lambda g}(\beta) = \operatorname{sign}(\beta) \max(|\beta| - \lambda, 0)$ , with all operations carried out element-wise. The induced parameter space  $\Theta_{\lambda g}$ is in fact the union of multiple sets with varying dimensions:  $\{\theta \in \mathbb{R}^p : \theta_j = 0 \text{ for } j \in S, |S| = d\}$ , where S is some index set and  $d \in (0, 1, ..., p)$ , each is a Euclidean subspace of dimension (p-d) — conveniently, we do not need to explicitly specify the dimension d, since it is automatically induced through the transformation.

This suggests that the proximal mapping can be used as a convenient tool to develop priors on lower-dimension subsets. We now list a few useful proximal mappings in Table 1. In addition, the proximal mapping allows us to easily consider multiple constraints or *g* functions, since the intersection of convex sets and summation of convex functions are still convex. The general form can be computed using the alternative direction of method of multipliers algorithm Bertsekas (2014), and we will demonstrate one case in the data application. For example, consider  $\theta$  being sparse while constrained in some convex set; this would be challenging to model for conventional approaches due to the lack of reparameterization.

### 2.2. Proximal Prior

We now use the above in a Bayesian modeling framework. Suppose we have data generated from a likelihood  $L(y; \theta)$ , where we want to assign a prior on  $\theta$  in some space with dimensionality smaller or equal to p. We use the following generative

Table 1. Some useful proximal mappings.

Space of $eta$	g(z)	$prox_{\lambda g}(eta)$	Usage
$\mathbb{R}^p$	$\mathcal{X}_{\mathcal{C}}$ , $\mathcal{C}$ convex set	$P_{C}(\beta)$	Projection to a set [See Table 6.1 of Beck (2017) for an expanded list]
$\mathbb{R}^p$	<i>z</i>    <sub>1</sub>	$sign(\beta) \max( \beta  - \lambda, 0)$ , computed element-wise	Sparsity
$\{\beta \in \mathbb{R}^{k \times k}, \text{ positive semidefinite}\}$	$  Z  _*$ , nuclear norm	$U\Lambda_0V^T$ , with $\beta=U\Lambda V^T$ the singular value decomposition, $(\Lambda_0)_{ii}=\max(0,\Lambda_{ii}-\lambda)$	Low rank
$\mathbb{R}^{m \times n}$	$  Z  _{2,1} = \sum_{i} \sqrt{\sum_{j} Z_{ij}^2}$	$[\beta_i \max(1-\frac{\lambda}{\ \beta_i\ _2},0)]_{i=1}^m$ with $\beta_i$ as the <i>i</i> th row	Row / group sparsity
$\mathbb{R}^p$	$\ Dz\ _1$ with $D \in \mathbb{R}^{k \times p}$ $\operatorname{dist}_C(z) = \inf_{x \in C} \ z - x\ _2$ , $\operatorname{distance}$ to a set	Solvable via the alternating direction method of multipliers $aP_C(\beta) + (1-a)\beta$ with $a = \min\{\lambda/\text{dist}_C(\beta), 1\}$	Fused lasso, convex clustering Set expansion to <i>C</i>



process for  $\theta$ :

$$\beta \sim \Pi_{\beta}^{0},$$

$$\lambda \sim \Pi_{\lambda}^{0},$$

$$\theta = \operatorname{prox}_{\lambda\sigma}(\beta),$$
(2)

where  $\Pi^0_{\beta}$  is a continuous distribution in  $\mathbb{R}^p$ , such as the nondegenerate Gaussian  $\beta \sim N(\mu, \Sigma)$  and we use  $\Pi_1^0$  to denote a generative distribution for  $\lambda > 0$ .

Here *g* is a convex and lower-semicontinuous function such as those in Table 1. Potentially, g could be known up to some other hyper-parameter  $\gamma$ ; in that case, we denote it by  $g_{\gamma}$  and use  $\Pi^0_{\gamma}$  as the prior for  $\gamma$ . For a clear notation, we use bold subscript such as in  $\Pi_{\theta}^{0}$  as a book-keeping index to refer to the variable whose prior is being defined.

It is not hard to see that  $\lambda g(z) + 2^{-1} ||z - \beta||_2^2$ , as the combination of the convex g and a quadratic term, is strictly convex with a unique minimizer. Therefore, each  $\beta$  maps to a unique  $\theta$ , hence, we have a measurable mapping, which means we have a valid prior distribution for  $\theta$  using (2). We denote the conditional prior distribution for  $\theta$  as  $\Pi_{\theta}^{0}(\theta \mid \lambda, \gamma)$ , and its marginal distribution as  $\Pi^0_{\theta}(\theta)$  after integrating out  $\gamma$  and  $\lambda$ . For convenience, we will refer to either form as a "proximal prior."

We first show that, a proximal prior can produce a convenient equivalence to a hierarchical prior of first selecting a lowdimensional set and then assigning a conditional density within this set. We denote the space induced by  $\operatorname{prox}_{\lambda g}(\boldsymbol{\beta})$  as  $\Theta$ , and assume that it can be partitioned into  $\Theta = \Theta^0 \cup \Theta^1 \cup \cdots \cup \Theta^p$ , where  $\Theta^k$  denotes a k-dimensional subset of  $\Theta$ , and  $\Theta^j \cap \Theta^k = \emptyset$ if  $j \neq k$  (this can be achieved even if a higher dimensional set  $\tilde{\Theta}^k$  overlaps/contains a lower-dimensional set  $\Theta^j$ , we set  $\Theta^k = \tilde{\Theta}^k \setminus \bigcup_{j=1}^{k-1} \Theta^j$ ). Then the prior kernel (a mix of density and mass functions) evaluated at  $\theta = t$  can be written as

$$\Pi_{\boldsymbol{\theta}}^{0}(t) = \sum_{k=0}^{p} \Pi_{\boldsymbol{\theta}}^{0}(t \mid \theta \in \Theta^{k}) \mathbf{1}(t \in \Theta^{k}) \operatorname{pr}(\theta \in \Theta^{k}), \quad (3)$$

where  $\sum_{k=0}^{p} \operatorname{pr}(\theta \in \Theta^{k}) = 1$  and  $\Pi_{\theta}^{0}(t \mid \theta \in \Theta^{k})$  is a conditional density that integrates to 1 over  $t \in \Theta^k$  using an appropriate k-dimensional integral with respect to some proper measure  $\lambda^k$ , denoted by  $\int_{\Theta^k} \Pi_{\theta}^0(t \mid \theta \in \Theta^k) \lambda^k(dt) = 1$ . The integral and measure will be formally defined in the theory section.

Therefore, from a generative view, the above can be understood as first picking a set  $\Theta^k$  with probability  $pr(\theta \in \Theta^k)$ , then drawing a value t within the space of  $\Theta^k$ . This includes those corner cases where  $\text{prox}_{\lambda g}$  cannot map to some dimensional sets: that is, for some k's, we can have  $pr(\theta \in \Theta^k) = 0$ .

Accordingly, with  $L(y; \theta)$  the likelihood, the posterior of  $\theta$  can be derived as

$$\Pi(\theta = t \mid y) = \sum_{k=0}^{p} \underbrace{z_{k}^{-1} L(y; t) \Pi_{\theta}^{0}(t \mid \theta \in \Theta^{k}) \mathbf{1}(t \in \Theta^{k})}_{\Pi(\theta = t \mid \theta \in \Theta^{k}, y)}$$

$$\underbrace{z_{k} \operatorname{pr}(\theta \in \Theta^{k})}_{\operatorname{pr}(\theta \in \Theta^{k} \mid y)}, \tag{4}$$

where  $z_k = \int_{\Theta^k} L(y; \theta = t) \Pi_{\theta}^0(t \mid \theta \in \Theta^k) \lambda^k(dt)$ . We assume posterior propriety almost everywhere, such that  $z_k < \infty$  for all  $k: \operatorname{pr}(\theta \in \Theta^k) > 0.$ 

The proximal priors simplify these procedures. Using the transformation  $\theta = \operatorname{prox}_{\lambda g}(\beta)$ , for any measurable set  $A \in \Theta$ ,

$$\begin{split} &\operatorname{pr}(\theta \in \mathcal{A} \mid y) \\ &= \sum_{k=0}^{p} z_{k}^{-1} \left\{ \int_{\mathcal{A} \cap \Theta^{k}} L(y;t) \Pi_{\theta}^{0}(t \mid \theta \in \Theta^{k}) \mathrm{d}t \right\} \frac{z_{k} \operatorname{pr}(\theta \in \Theta^{k})}{\sum_{k=0}^{d} z_{k} \operatorname{pr}(\theta \in \Theta^{k})} \\ &\stackrel{(a)}{=} \sum_{k=0}^{p} z_{k}^{-1} \left[ \int_{\operatorname{prox}_{\lambda_{g}}^{-1}(\mathcal{A} \cap \Theta^{k})} L\left\{ y; \operatorname{prox}_{\lambda_{g}}(b) \right\} \frac{\Pi_{\beta}^{0}(b) \mathbf{1}[\operatorname{prox}_{\lambda_{g}}(\beta) \in \Theta^{k}]}{\operatorname{pr}[\operatorname{prox}_{\lambda_{g}}(\beta) \in \Theta^{k}]} \mathrm{d}b \right] \\ &\frac{z_{k} \operatorname{pr}(\theta \in \Theta^{k})}{\sum_{k=0}^{d} z_{k} \operatorname{pr}(\theta \in \Theta^{k})} \\ &= \frac{1}{\sum_{k=0}^{d} z_{k} \operatorname{pr}(\theta \in \Theta^{k})} \sum_{k=0}^{p} \int_{\operatorname{prox}_{\lambda_{g}}^{-1}(\mathcal{A} \cap \Theta^{k})} L\left\{ y; \operatorname{prox}_{\lambda_{g}}(b) \right\} \Pi_{\beta}^{0}(b) \mathrm{d}b, \end{split}$$

where in (a) we mean that  $\Pi_{\theta}^{0}(t \mid \theta \in \Theta^{k})$  contains the Jacobian term in the change-of-variables  $t = \text{prox}_{\lambda g}(b)$  (details of the Jacobian calculation provided in the theory section). At any given  $\beta = b$ , we can omit the integral and summation, and obtain a remarkably simple posterior density of  $\beta$ :

$$\Pi(\beta = b \mid y) \propto L\left\{y; \operatorname{prox}_{\lambda g}(b)\right\} \Pi_{\beta}^{0}(b). \tag{5}$$

Remark 1. To clarify, although the hierarchical form of  $\Pi_{\mathbf{q}}^{0}$ provides a nice interpretation to our proximal prior, such an equivalence is not strictly necessary for the proximal modeling framework to work. To be rigorous, the above equivalence requires a few regularity conditions, to be formalized in the theory section.

Therefore, compared to (4), the posterior density (5) is much easier for Bayesian applications. This also suggests a new strategy of "data augmentation using optimization" [instead of marginalization as in Tanner and Wong (1987)]—if we can write the parameter  $\theta$  as some proximal mapping from  $\beta$ , then we can sample  $\beta$  first as an augmented variable; after sampling, we compute  $\theta = \text{prox}_{\lambda g}(\beta)$  and discard the information from  $\beta$ .

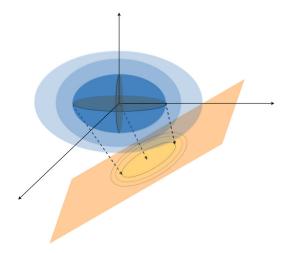
We now use one example to illustrate the equivalence.

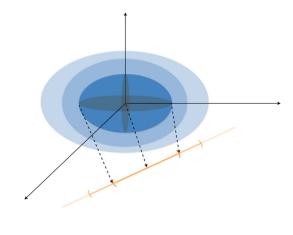
Example 2 (Affinely constrained prior under varying rank). Suppose we want to assign a prior for  $\theta$  in a set of affine constraints  $C = \{\theta \in \mathbb{R}^p : A^T \theta = b\}$ , where  $A^T \in \mathbb{R}^{m \times p}$  is another parameter, with m < p and  $b \in Col(A^T)$  the column space of  $A^{\mathrm{T}}$  (so that C is not empty). Since A is not fixed, we do not know the rank of A, hence, not the dimensionality of C. Using the proximal prior with  $g_A(z) = 0$  if  $A^T \theta = b$ ,  $g_A(z) = \infty$  otherwise (hence,  $\operatorname{prox}_{\lambda g_A}$  is invariant to any finite value of  $\lambda > 0$ ), and  $\beta \sim N(\mu, \Sigma)$ , we have a closed-form proximal mapping

$$\theta = \operatorname{prox}_{\lambda g}(\beta) = \beta - A(A^{\mathrm{T}}A)^{-}(A^{\mathrm{T}}\beta - b),$$

where  $(\cdot)^-$  is the Moore-Penrose inverse. We illustration this mapping in Figure 1.

The  $\theta$ -marginal proximal prior is a discrete mixture over different rank of *A*:





- (a) Given  $\operatorname{rank}(A)=1$ , the mapping  $\theta=\operatorname{prox}_{\lambda g}(\beta)$  creates a two-dimensional prior on the hyperplane (orange).
- (b) Given  $\operatorname{rank}(A) = 2$ , the mapping  $\theta = \operatorname{prox}_{\lambda g}(\beta)$  creates a one-dimensional prior on the line (orange).

Figure 1. Illustrative example of constructing a prior on an affinely constrained set  $C = \{\theta : A^T\theta = b\}$ . A challenge arises when the rank(A) is unknown, the dimensionality of the prior is unknown. The proximal prior bypasses this hurdle by transforming a continuous prior (blue) into the constrained space (orange), without the need to explicitly specify the dimensionality.

$$\begin{split} \Pi_{\pmb{\theta}}^0(\theta) &= \int \Pi_{\pmb{\theta}}^0\left(\theta \mid A\right) \Pi_{\pmb{A}}^0(A) \mathrm{d}A \\ &= \sum_{d=0}^p \mathrm{pr}^0(\mathrm{rank}(A) = d) \underbrace{\frac{\int_{A:\mathrm{rank}(A) = d} \Pi_{\pmb{\theta}}^0\left(\theta \mid A\right) \Pi_{\pmb{A}}^0(A) \mathrm{d}A}{\int_{A:\mathrm{rank}(A) = d} \Pi_{\pmb{\theta}}^0\left[\theta \mid \mathrm{rank}(A) = d\right]} \end{split}$$

and  $\Pi_{\theta}^{0}(\theta \mid A)$  is the degenerate Gaussian density with mean  $A\left(A^{T}A\right)^{-1}b+P_{A^{\perp}}\mu$  and covariance  $P_{A^{\perp}}\Sigma P_{A^{\perp}}$ , where  $P_{A^{\perp}}=I-A\left(A^{T}A\right)^{-1}A^{T}$ . Although the summation may not have a closed-form, the weights and conditional density can be tractable in applications.

For illustration, we consider a Bayesian envelope linear regression for multivariate response  $Y_i \in \mathbb{R}^p$ :

$$Y_i = \mu + \tilde{\Theta}X_i + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \text{N}(0_p, \ \Gamma\Omega\Gamma^{\text{T}} + \Gamma_*\Omega_*\Gamma_*^{\text{T}})$$
  
 $\theta = \Gamma \eta,$ 

for  $i=1,\ldots,n$ , each covariate  $X_i\in\mathbb{R}^m$  (with  $p\leq m$ ), and the noise  $\epsilon_i\in\mathbb{R}^p$ ;  $[\Gamma\ \Gamma_*]$  together from a  $p\times p$  orthonormal matrix, with  $\Gamma$  a  $p\times u$  submatrix,  $\eta\in\mathbb{R}^{u\times m}$  full rank, and both  $\Omega$  and  $\Omega_*$  positive definite matrices. The regression coefficient matrix  $\widetilde{\Theta}\in\mathbb{R}^{p\times m}$  is of rank u, with u unknown. The motivation is that by making  $\widetilde{\Theta}X_i$  in the subspace spanned by the leading eigenvectors of the covariance,  $\Gamma_*^TY_i$  is small in magnitude and independent from  $X_i$ , leading to a sufficient dimension reduction (Cook, Li, and Chiaromonte 2010). For Bayesian inference, Khare, Pal, and Su (2017) proposed to use a matrix-Bingham prior on  $\Gamma$  with a pre-specified u, so that it has conjugate forms in a Gibbs sampler for posterior computation. On the other hand, since  $\Gamma$  is in an

orthogonal and low-rank space, it is difficult to generalize to other forms of prior such as letting u vary.

Using the affine constraint proximal mapping, we can bypass these challenges. We reparameterize  $\Gamma_*\Omega_*\Gamma_*^{\rm T}=AA^{\rm T},\Gamma\Omega\Gamma^{\rm T}=P_{A^\perp}WP_{A^\perp}^{\rm T}$  with rank(A)=p-u,W positive definite, and  $\tilde{\Theta}$  by a linear constraint  $A^{\rm T}\tilde{\Theta}=O$ . The proximal mapping yields  $\tilde{\Theta}=P_{A^\perp}B$ , with  $B\in\mathbb{R}^{p\times m}$  the matrix form for  $\beta$ . Using matrices  $Y^{\rm T}\in\mathbb{R}^{n\times p},X^{\rm T}\in\mathbb{R}^{n\times m}$ , we rewrite the envelope regression likelihood as

$$\begin{split} L(Y; A, W, B, \mu) & \propto \exp\left(-\frac{1}{2} \text{tr} \bigg\{ [Y^{\text{T}} - 1\mu^{\text{T}} - X^{\text{T}} (P_{A^{\perp}} B)^{\text{T}}] \\ & (P_{A^{\perp}} W P_{A^{\perp}}^{\text{T}})^{-} [Y - \mu \mathbf{1}^{\text{T}} - (P_{A^{\perp}} B) X] \bigg\} \right) \\ & \times |A^{\text{T}} A|^{-n/2} \exp\left(-\frac{1}{2} \text{tr} \bigg\{ [Y^{\text{T}} - 1\mu^{\text{T}}] (AA^{\text{T}})^{-} [Y - \mu \mathbf{1}^{\text{T}}] \bigg\} \right) \\ & |P_{A^{\perp}} W P_{A^{\perp}}^{\text{T}}|_{+}^{-n/2}, \end{split}$$

where  $|\cdot|_+$  is the pseudo-determinant. To complete the proximal prior specification, we use  $B \in \mathbb{R}^{p \times m}$  and  $R \in \mathbb{R}^{p \times p}$ , with their elements iid from N(0, 1) (hence, full rank almost surely); then we set  $A = R\Lambda$ , and  $\Lambda = \operatorname{diag}(\Lambda_{i,i})_{i=1}^p$  with  $\Lambda_{i,i} = z_i 1(z_i > \rho_q)$  and  $z_i \sim \operatorname{Exp}(1)$ , and  $\rho_q$  the q-quantile of  $\operatorname{Exp}(1)$ . We provide a numerical simulation in S1.3 of the Supplementary Materials.

Focusing on the  $\theta$ -marginal prior, we have: (a) the mixture weight  $\operatorname{pr^0}[\operatorname{rank}(A) = p - u] = \binom{p}{u} q^u (1 - q)^{p-u}$ , where q can be specified as a priori; (b) the conditional density containing

$$\Pi_{\theta}^{0}(\theta \mid A) = (2\pi)^{-u/2} |P_{A^{\perp}}|_{+}^{-1/2} \exp(-\theta^{\mathrm{T}} P_{A^{\perp}}^{-} \theta / 2)$$



$$= (2\pi)^{-u/2} \exp(-\|\theta\|^2/2) 1(A^{\mathrm{T}}\theta = 0).$$

The second line is due to  $P_{A^{\perp}}$  being idempotent  $\theta = P_{A^{\perp}}\beta = P_{A^{\perp}}\theta$ ,  $P_{A^{\perp}}P_{A^{\perp}}^-P_{A^{\perp}} = P_{A^{\perp}}$ ,  $\theta^{\rm T}P_{A^{\perp}}\theta = \theta^{\rm T}\theta$  and  $|P_{A^{\perp}}|_+ = 1$ . This density is invariant to scaling of A.

Remark 2. To clarify, we use the above example (with a relatively simple  $\Pi^0_{\theta}$ ) to illustrate the equivalence between the hierarchical specification of  $\Pi^0_{\theta}$  and the continuous- $\Pi^0_{\theta}$ -and-mapping specification. In general cases, the former  $\Pi^0_{\theta}$  may be intractable due to the lack of closed-form, hence, motivating the proximal mapping strategy as proposed in this article.

In Example 2, if we intuitively compare the two distributions before and after the mapping, it reduces (or at least retains) the distance to the center  $\|\theta - \operatorname{prox}_{\lambda g}(\mu)\|_2 \le \|\beta - \mu\|_2$ . The property shown in this example is known as the "non-expansiveness," which in fact holds for all proximal mappings:

$$\|\operatorname{prox}_{\lambda g}(\beta_1) - \operatorname{prox}_{\lambda g}(\beta_2)\|_2 \le \|\beta_1 - \beta_2\|_2,$$

for any  $\beta_1$ ,  $\beta_2$  in the domain of  $\operatorname{prox}_{\lambda g}$ . This is in particular meaningful for Bayesian inference, as it conveniently controls the concentration of measure for  $\theta$ .

Theorem 1. If the data y come from a distribution  $\mathcal{F}_{\theta^*}$  with a fixed parameter  $\theta^*$ , and any  $\epsilon \in (0,1)$ , the posterior distributions of  $\theta$  and  $\beta$  satisfy

$$\operatorname{pr}(\|\theta - \theta^*\| > \epsilon \mid y)$$

$$\leq \int \min_{\beta^*: \operatorname{prox}_{\lambda g_{\gamma}}(\beta^*) = \theta^*} \operatorname{pr}\left[\|\beta - \beta^*\| > \epsilon \mid y, \lambda, \gamma\right]$$

$$\Pi(\lambda, \gamma \mid y) d(\lambda, \gamma).$$

In addition, if  $tr[cov(\beta \mid y)] < \infty$ , then

$$tr[cov(\theta \mid y)] \le tr[cov(\beta \mid y)].$$

Using the envelope regression example, we know that

$$\begin{split} \Pi(B\mid y,A,W) &\propto \exp\big\{-(1/2)\mathrm{tr}[XX^{\mathrm{T}}B^{\mathrm{T}}P_{A^{\perp}}(P_{A^{\perp}}WP_{A^{\perp}}^{\mathrm{T}})^{-}P_{A^{\perp}}B]\\ &+\mathrm{tr}(B^{\mathrm{T}}B)-2\mathrm{tr}[X(Y^{\mathrm{T}}-1\mu^{\mathrm{T}})(P_{A^{\perp}}WP_{A^{\perp}}^{\mathrm{T}})^{-}P_{A^{\perp}}B])\big\}, \end{split}$$

which is a multivariate Gaussian for  $\operatorname{vec}(B)$ . On the other hand, since we know  $\tilde{\Theta}=P_{A^\perp}B$ , we know for any given A,  $\|\tilde{\Theta}-\tilde{\Theta}^*\|_F^2=\operatorname{tr}[(B-B^*)^TP_{A^\perp}P_{A^\perp}(B-B^*)]\leq \|B-B^*\|_F^2$  with  $P_{A^\perp}B^*=\tilde{\Theta}^*$ , due to  $P_{A^\perp}$  being idempotent and having eigenvalues equal to either 1 or 0.

### 2.3. Prior Specification on $\lambda$

In the proximal mapping (1), the hyper-parameter  $\lambda$  plays an important role, hence, we need to carefully choose its prior. To first obtain some intuition, note when  $\lambda \to 0$ , we have  $\operatorname{prox}_{\lambda g}(\beta) \to \beta$  if  $g(z) < \infty$  for all z, the identity mapping; when  $\lambda \to \infty$ , we have  $\operatorname{prox}_{\lambda g}(\beta) \to \arg\min_z g(z)$ . Therefore, as  $\lambda$  increases,  $\theta$  becomes farther away from  $\beta$ , hence, the distribution  $\Pi^0_{\beta}(\beta)$  gets more "deformed" at a larger  $\lambda$ . We now formalize this deformation intuition, while relaxing the finite-valuedness of g. For conciseness, we postpone all the proofs in the Appendix.

Theorem 2 (Monotonicity of deformation in  $\lambda$ ). For any function g with range  $\mathbb{R} \cup \{\infty\}$ , if  $0 < \lambda_1 < \lambda_2$ , then  $\|\beta - \operatorname{prox}_{\lambda_1 g}(\beta)\|_2 \le \|\beta - \operatorname{prox}_{\lambda_2 g}(\beta)\|_2$ .

This result means that we can find a measurement between 0 and 1 to quantify the deformation:

$$\omega_{\lambda} := \frac{\mathbb{E}_{\boldsymbol{\beta}} \|\boldsymbol{\beta} - \operatorname{prox}_{\lambda g}(\boldsymbol{\beta})\|_{2}}{\mathbb{E}_{\boldsymbol{\beta}} \|\boldsymbol{\beta} - \lim_{\lambda^{*} \to \infty} \operatorname{prox}_{\lambda^{*} \sigma}(\boldsymbol{\beta})\|_{2}},\tag{6}$$

where the expectation is taken with respect to the prior of  $\beta$ .

When lacking prior knowledge on  $\lambda$ , we can use a Beta prior on  $w \in (0, 1)$  and solve for  $\lambda$ :

$$\omega \sim \text{Beta}(a_{\omega}, b_{\omega}), \qquad \lambda = \min_{x>0} (x : \omega_x = \omega).$$
 (7)

In this article, we use a non-informative  $a_{\omega} = b_{\omega} = 1$ . As a toy example, let  $\beta$  be univariate with a finite variance, using the proximal mapping with  $g(z) = z^2/2$ , we have  $\operatorname{prox}_{\lambda g}(\beta) = \beta/(1+\lambda)$ . Therefore, we have  $\lambda = (1-\omega)/\omega$  with an induced prior  $\Pi_1^0(\lambda) = 2/(1+\lambda)^2$  for  $\lambda > 0$ .

In more general cases, (7) often cannot be solved analytically. However, we can numerically compute a prior for  $\lambda$ , using a strategy similar to Berger et al. (2009)—for K chosen points  $\lambda_1, \ldots, \lambda_K$  in  $(0, \infty)$ , we can use the empirical estimates of the expectation based on simulated  $\beta \sim \Pi^0_{\beta}(\beta)$ , and solve for  $\omega_1, \ldots, \omega_K$ ; afterwards, we can easily interpolate to obtain the  $\lambda$  associated with any  $\omega$ .

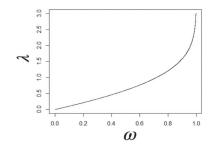
*Example 1* ((*Continued*). *Soft-thresholding prior*). To illustrate, we compute the prior of  $\lambda$  for the soft-thresholding prior based on  $\theta = \text{sign}(\beta) \max(|\beta| - \lambda, 0)$ . Based on  $\beta \in \mathbb{R}^p$  and  $\beta \sim N(0, I_p)$ , we compute the prior density of  $\lambda$  and plot it in Figure 2.

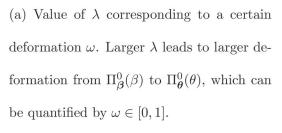
In this section, we discussed the choice of  $\Pi_{\lambda}^{0}$  with the generality of all possible g and  $\operatorname{prox}_{\lambda g}$  in mind. On the other hand, for some specific case such as  $g(z) = ||z||_1$  and softthresholding prox<sub> $\lambda g$ </sub>, there is a connection to some existing prior in the literature, such as the classic spike-and-slab prior. For example, if  $\Pi_{\beta}^{0}(\beta) \propto \exp(-\|\beta\|_{1}/\alpha)$  and  $\lambda$  to be the  $\omega$ -quantile of  $\mathrm{Exp}(\alpha^{-1})$ , then we can obtain a spike-and-slab prior with Laplace slab  $\Pi_{\theta}^{0}(\theta \mid \lambda) = \prod_{j=1}^{p} [w_{\lambda} \delta_{0}(\theta_{j}) + (1 - \omega)]$  $w_{\lambda}$ ) $(2\alpha)^{-1}$  exp $(-\theta_i/\alpha)$ ]. A closely related discovery is the neuronized prior (Shin and Liu 2021) using truncated activation function, for which there is an equivalence to a spike-and-slab prior with two-normal-product slab. In these cases, there are often alternative choices for  $\Pi^0_{\lambda}$  that are justified via large sample theory. Due to the page constraint, we defer the detailed discussion and numerical experiments to S1.1 of the supplementary materials.

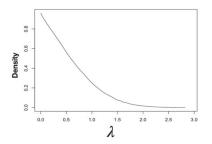
# 3. Geometric Measure Theory on the Varying Dimensional Sets

### 3.1. Hausdorff Dimension and Low Dimensional Density

We now give a more rigorous exposition on the distribution induced by the proximal mapping. Without loss of generality, we consider  $\theta$  as a p-element vector. Since  $\theta$  may correspond to a measure of a set in the lower dimensional space, the p-dimensional Lebesgue measure of any lower-dimensional set is







(b) The prior for  $\lambda$ , corresponding to a uniform prior on the deformation measurement  $\omega$ .

Figure 2. Illustration of a numerically computed prior for  $\lambda$ , which controls sparse level for the soft-thresholding mapping  $\theta = \text{sign}(\beta) \max(|\beta| - \lambda, 0)$ .

zero hence is not useful. We need some tools from the geometric measure theory to address this issue. To start, consider a set  $\mathcal{A}$  and suppose we do not know its dimensionality. Instead, we can cover A with sets  $B_i$ 's, each  $B_i$  has its diameter diam $(B_i)$  =  $\sup\{|x-y|: x, y \in B_i\} \le \delta$ . We call any  $\bigcup_i B_i \supset A$ , diam $(B_i) < B_i$  $\delta$  as a  $\delta$ -covering of  $\mathcal{A}$ .

Then we take the infimum over all  $\delta$ -coverings of A, and letting the  $\delta$  decrease, we obtain the s-dimensional Hausdorff measure of A:

$$\mathcal{H}^{s}(\mathcal{A}) = \lim_{\delta \to 0} \inf \left\{ \sum_{i=1}^{\infty} \operatorname{diam}(B_{i})^{s} : \operatorname{diam}(B_{i}) \leq \delta, \mathcal{A} \subseteq \bigcup_{i} B_{i} \right\}.$$
(8)

Intuitively, the above can be taken as the minimum total "volume" of the covering—except *s* is a parameter that varies.

In fact,  $\mathcal{H}^s(A)$  is a nonincreasing function in  $s \geq 0$  (Edgar 2007). More importantly, for any Borel A, and  $0 < s_1 < s_2$ , if  $\mathcal{H}^{s_1}(\mathcal{A}) < \infty$  then  $\mathcal{H}^{s_2}(\mathcal{A}) = 0$ ; and if  $\mathcal{H}^{s_2}(\mathcal{A}) > 0$  then  $\mathcal{H}^{s_1}(\mathcal{A}) = \infty$  [Theorem 6.1.6 (Edgar 2007)]. This means for any Borel set A, there is a unique  $s_0 \in [0, \infty) \cup \{\infty\}$  as a transition point, over which the dimensionality drops from  $\infty$  to 0:

$$\mathcal{H}^s(\mathcal{A}) = \infty$$
, for any  $s < s_0$ ;  
 $\mathcal{H}^s(\mathcal{A}) = 0$ , for any  $s > s_0$ .

Such an  $s_0$  is referred to as the Hausdorff dimension of A, equivalently:

$$\dim_{\mathcal{H}}(\mathcal{A}) = \inf\{s \ge 0 : \mathcal{H}^s(\mathcal{A}) = 0\}. \tag{9}$$

Note that  $\dim_{\mathcal{H}}(\mathcal{A}) \geq 0$  does not have to be an integer; nevertheless, when it is, the Hausdorff measure is proportional to the commonly used s-dimensional Lebesgue measure

$$\lambda^{s}(\mathcal{A}) = \inf \left\{ \sum_{i=1}^{\infty} \operatorname{vol}(B_{i}) : \mathcal{A} \in \bigcup B_{i}, B_{i} \text{ is an open cube} \right\},$$

via  $\lambda^s(\mathcal{A}) = w_s \mathcal{H}^s(\mathcal{A})$ , where  $w_s = \pi^{s/2} [2^s \Gamma(s/2+1)]^{-1}$  due to the volume formula of an s-dimensional ball. In addition, when  $s = 0, \mathcal{H}^0(\mathcal{A})$  is same as the counting measure.

Now recall that  $prox_{\lambda\rho}$  is nonexpansive, which leads to the following theorem:

*Theorem 3.* For any Borel set A and proximal mapping prox<sub> $\lambda \sigma$ </sub>,

- 1.  $\mathcal{H}^{s}\{\operatorname{prox}_{\lambda g}(\mathcal{A})\} \leq \mathcal{H}^{s}(\mathcal{A}) \text{ for any } s \geq 0;$ 2.  $\dim_{\mathcal{H}}\{\operatorname{prox}_{\lambda g}(\mathcal{A})\} \leq \dim_{\mathcal{H}}(\mathcal{A}).$

*Remark 3.* In the above, the statement 2 is particularly useful: it tells us that  $prox_{\lambda g}$  only maps to lower or equal dimensional space.

Now, starting from a probability distribution defined by a certain Radon measure  $\mu$  in  $\mathbb{R}^p$  for some low-dimensional sets in  $\Theta^s$ , one interesting question is how to differentiate this and obtain a "density," as  $\Pi_{\theta}(\theta = t \mid \theta \in \Theta^s)$  used in (3) and (4).

For a point  $\theta \in \Theta^s$ , the ball  $B_r(\theta)$  centered at  $\theta$  with radius r > 0 has the lower and upper *s*-dimensional derivatives:

$$f^s_{\mu,*}(\theta) = \lim\inf_{r \to 0} \frac{\mu\{B_r(\theta)\}}{w_s r^s}, \qquad f^{s,*}_{\mu}(\theta) = \lim\sup_{r \to 0} \frac{\mu\{B_r(\theta)\}}{w_s r^s}.$$

Therefore, if we have the two limits coincide, we would have a definition of an s-dimensional density:  $f_{\mu}^{s}(\theta) = f_{\mu,*}^{s}(\theta) =$  $f_{\mu}^{s,*}(\theta)$ , commonly referred to as the *s*-density.

Remark 4. To understand the s-density as a generalized concept of "density," for those continuous distributions associated with a p-dimensional Lebesgue measure, such as the nondegenerate Gaussian distribution, the p-density is the probability density function; whereas for the discrete distributions, the 0-density is the same as the probability mass function.

Next, similar to the probability density function, s-density may not always exist. Therefore, it is important to state the two required conditions, *s* is an integer and  $\theta$  is in a rectifiable set, as formalized in the following theorem.

Theorem 4 (Besicovitch-Marstrand-Preiss theorem). (Preiss 1987) Let  $\mu$  be a locally finite Radon measure on  $\mathbb{R}^p$ , if there exists a real  $s \ge 0$  such that  $f_{ii}^{s}(\theta)$  exists, and it is positive on a set of positive  $\mu$ -measure, then s must be an integer. On the other hand, let  $\mathcal{A} \subset \mathbb{R}^p$  be Borel with  $\mathcal{H}^s(\mathcal{A}) \in (0, \infty)$  and s an integer, then  $f_{\mu}^{s}(\theta)$  exists for  $\theta \in \mathcal{A}$  almost everywhere with respect to  $\mathcal{H}^{k}$ , if and only if the set  $\mathcal{A}$  is rectifiable.

To explain "rectifiability," a Borel set  $\mathcal{A} \subset \mathbb{R}^p$  is rectifiable if there is a countable family of Lipschitz maps  $T_i: \mathbb{R}^s \to \mathbb{R}^p$  which cover almost all  $\mathcal{A}$  except for sets with zero  $\mathcal{H}^s$  measure. That is, intuitively speaking, almost every p-element vector  $\theta \in \mathcal{A}$  can be represented as some transformation of  $x \in \mathbb{R}^s$ —note that this is not the same as a simple reparameterization, as we may obtain  $\mathcal{A}$  via multiple  $f_i$ 's (up to countably many).

### 3.2. Calculation of the s-Density

We now provide a way to calculate the s-density. Focusing on a subset  $\Theta^k$  with  $\dim(\Theta^k) = k$  and  $\boldsymbol{\beta}^k = \operatorname{prox}_{\lambda g}^{-1}(\Theta^k)$ . We now transform  $\Pi_0(\beta = b \mid \beta \in \boldsymbol{\beta}^k)$  into an s-density with s = k.

Theorem 5. If  $\boldsymbol{\beta}^k$  is  $(\mathcal{H}^p, p)$ -rectifiable and  $\dim_{\mathcal{H}}(\boldsymbol{\beta}^k) = p$ ,  $\Theta^k$  is  $(\mathcal{H}^k, k)$ -rectifiable and  $\dim_{\mathcal{H}}(\Theta^k) = k$ , with  $p \geq k$ , and  $J_k \operatorname{prox}_{\lambda g}(\beta) > 0$  a.e.- $\mu_{\boldsymbol{\beta}}$ . Then the s-density of  $\theta$  induced by  $\operatorname{prox}_{\lambda g}$  is

 $\Pi(\theta = t \mid \theta \in \Theta^k)$ 

$$= \int_{\operatorname{prox}_{\lambda_{\sigma}}^{-1}(t)} \frac{\Pi(\beta = b)/\operatorname{pr}\{\operatorname{prox}_{\lambda_{g}}(\beta) \in \Theta^{k}\}}{J_{k}\operatorname{prox}_{\lambda_{g}}(b)} w_{(p-k)} d\mathcal{H}^{p-k}(b),$$

where  $J_k \text{prox}_{\lambda g}(b)$  is the k-dimensional Jacobian of  $\text{prox}_{\lambda g}$  at b.

Note that if the low-dimensional set  $\Theta^k$  can be reparameterized as a transformation an k-element vector, then it is possible to change (10) to an integration with respect to an k-dimensional Lebesgue measure.

To explain the assumptions above, a set  $\mathcal{A}$  is  $(\mathcal{H}^s, s)$ -rectifiable when  $\mathcal{H}^s(\mathcal{A}) < \infty$ , and there is a set as the countable union of Lipschitz images from bounded sets  $\mathcal{B} = \bigcup_j \{T_j(\mathcal{C}_j) : \mathcal{C}_i \subset \mathbb{R}^s \text{ and bounded, } T_j \text{ Lipschitz} \}$  such that  $\mathcal{H}^s(\mathcal{A} \setminus \mathcal{B}) = 0$ . As the result, if s-density exists, we could use (10) when both  $\theta$  and  $\beta$  are finite.

Morgan (2016) gives the k-dimensional Jacobian  $J_kT(x)$  of function  $T: \mathbb{R}^n \to \mathbb{R}^m$ , differentiable at x. Let  $D_T(x) \in \mathbb{R}^{n \times m}$  be the derivative matrix of T(x) at x, with  $\{D_T(x)\}_{ij} = \partial T(x)_j/\partial x_i$ ,  $1 \le i \le n, 1 \le j \le m$ , then the k-dimensional Jacobian can be computed as

$$J_k T(x) = \sum_{\substack{M \text{ is the } k \times k \\ \text{submatrix of } D_T(x)}} (\det M)^2.$$
 (10)

Note than when m = n = k,  $J_k T(x) = |\det\{D_T(x)\}|$  as more commonly seen.

Importantly, by Rademacher's theorem (Federer 2014), a Lipschitz function is differentiable almost everywhere. Therefore,  $D_{\text{prox}_{\lambda g}}(\beta)$  exists almost surely with respect to  $\mu_{\beta}$ . In the following example, we illustrate the use of the above theorem to compute the *s*-density for the affinely constrained prior.

Example 2 (The s-density of the affinely constrained prior). Using (10), one can verify that the s-density of affinely constrained prior recovers the "degenerate Gaussian density." Starting from

 $\beta \sim \mathrm{N}(\mu, \Sigma)$  and let us assume A is  $p \times d$  and  $\mathrm{rank}(A) = d$ , then it is not hard to compute that  $J_{(p-d)}\mathrm{prox}_{\lambda g}(\beta) = 1$ . Using the proximal mapping, at  $\theta = t$ , we have  $P_A\beta = t - A(A^{\mathrm{T}}A)^{-1}b$  with  $P_A = \{I - A(A^{\mathrm{T}}A)^{-1}A^{\mathrm{T}}\}$ , hence, we can integrate over the region  $\mathrm{prox}_{\lambda g}^{-1}(t)$  by reparameterizing  $\beta = P_A^-\{t - A(A^{\mathrm{T}}A)^{-1}b\} + Ax$  where  $x \in \mathbb{R}^d$ . Integrating over x, we have the s-density with s = p - d:

$$\begin{split} \Pi_{\theta}^{0}(t\mid\theta\in\Theta^{p-d}) &\propto \exp\bigg[-\frac{1}{2}\{t-A(A^{\mathrm{T}}A)^{-1}b-P_{A}\mu\}^{\mathrm{T}}\\ &P_{A}^{-}\Sigma^{-1}P_{A}^{-}\{t-A(A^{\mathrm{T}}A)^{-1}b-P_{A}\mu\}\bigg], \end{split}$$

which is commonly referred to as the "degenerate density" for a degenerate Gaussian, with its covariance  $P_A \Sigma P_A$  having a rank (p-d).

We list a few more examples commonly considered in statistics, where for each we have a guaranteed existence of *s*-density: regression under linear equality constraints, matrix factorization under low-rank constraint, sparse regression, covariance modeling in positive-definite space, and directional modeling in orthonormal space.

*Remark 5.* To clarify, the existence of *s*-density for  $\theta$  is not necessary in our modeling framework using proximal mapping, since we can always carry out computation using a valid *p*-dimensional density of  $\beta$ . On the other hand, the existence of *s*-density would be required if one wants to interpret the prior via an equivalent prior  $\Pi_{\theta}^{0}$  as in (3).

### 4. Posterior Computation

As shown in (5), when using  $\beta$  instead of  $\theta$ , the posterior has a simple density on  $\mathbb{R}^p$ , and the proximal mapping is differentiable almost everywhere with respect to  $\mu_{\beta}$ . Therefore, as long as  $\Pi(\theta;y)$  is a continuous and differentiable function in  $\theta$  almost everywhere, we can use the Hamiltonian Monte Carlo (HMC) for posterior computation. Now we first briefly review the HMC algorithm, then address the gradient calculation for the proximal mapping.

To sample from target distribution  $\beta \sim \Pi_{\beta|y}(\cdot)$ , the HMC uses an auxiliary momentum variable  $\nu$  and samples from a joint distribution  $\Pi(\beta, \nu) = \Pi(\beta \mid y)\Pi(\nu)$ , where a common choice of  $\Pi(\nu)$  is the density of N(0, M). Denote  $U(\beta) = -\log \Pi(\beta \mid y)$  and  $K(\nu) = -\log \Pi(\nu) = \nu^T M^{-1} \nu/2$ , which are referred to as the potential energy and kinetic energy, respectively. The total Hamiltonian energy function is  $H(\beta, \nu) = U(\beta) + K(\nu)$ .

At each state  $(\beta, \nu)$ , a new state is generated by simulating Hamiltonian dynamics, which satisfies the Hamilton's equations:

$$\frac{\partial \beta}{\partial t} = \frac{\partial H(\beta, \nu)}{\partial \nu} = M^{-1}\nu;$$

$$\frac{\partial \nu}{\partial t} = -\frac{\partial H(\beta, \nu)}{\partial \beta} = \frac{\partial \log \Pi(\beta \mid y)}{\partial \beta}.$$
(11)

The exact solution for (11) is often intractable, while we can numerically approximate the evolution by algorithms such as the leapfrog scheme. The leapfrog is a reversible and volume-preserving integrator, which updates the evolution  $(\beta^t, v^t) \rightarrow$ 

 $(\beta^{t+\epsilon}, \nu^{t+\epsilon})$  via

$$v \leftarrow v + \frac{\epsilon}{2} \frac{\partial \log \Pi(\beta \mid y)}{\partial \beta},$$

$$\beta \leftarrow \beta + \epsilon M^{-1} v,$$

$$v \leftarrow v + \frac{\epsilon}{2} \frac{\partial \log \Pi(\beta \mid y)}{\partial \beta}$$
(12)

for  $t=0,\epsilon,\ldots,L\epsilon$ , and sets  $(\beta^*,v^*)\leftarrow(\beta^{L\epsilon},v^{L\epsilon})$ . To correct the numeric error due to approximation,  $(\beta^*,v^*)$  is treated as a proposal and accepted with the Metropolis-Hastings (MH) probability

$$\min[1, \exp\{-H(\beta^*, \nu^*) + H(\beta, \nu)\}].$$

We now discuss the gradient computation:

$$\begin{split} \frac{\partial \log \Pi(\beta \mid y)}{\partial \beta} &= \frac{\partial \mathrm{prox}_{\lambda g}(\beta)}{\partial \beta} \left\{ \frac{\partial \log L(y; \theta)}{\partial \theta} \bigg|_{\theta = \mathrm{prox}_{\lambda g}(\beta)} \right\} \\ &+ \frac{\partial \log \Pi_{\beta}^{0}(\beta)}{\partial \beta}. \end{split}$$

When  $\operatorname{prox}_{\lambda g}(\beta)$  has a closed-form, we can use the automatic differentiation toolbox to calculate the gradient  $\partial \operatorname{prox}_{\lambda g}(\beta)/\partial \beta$ ; on the other hand, when the closed-form does not exist, some numeric approximation is needed.

Note that the partial gradient is  $\partial \operatorname{prox}_{\lambda g}(\beta)/\partial \beta_j = \lim_{\epsilon \to 0} \{\operatorname{prox}_{\lambda g}(\beta + e_j \epsilon) - \operatorname{prox}_{\lambda g}(\beta)\}/\epsilon$  with  $e_j$  is the standard basis with the jth element equal to one, and all others equal to zero; using a small  $\epsilon$  gives us the finite difference approximation. Nevertheless, when  $\beta$  is high dimensional, this would involve (p+1) times of calculating the proximal mapping, which can be computationally prohibitive. To solve this problem, we follow Spall (1992) and use the simultaneous perturbation stochastic approximation:

$$\frac{\partial \operatorname{prox}_{\lambda g}(\beta)}{\partial \beta_{j}} \approx \frac{1}{m} \sum_{k=1}^{m} \frac{\{\operatorname{prox}_{\lambda g}(\beta + \Delta^{(k)} \epsilon) - \operatorname{prox}_{\lambda g}(\beta)\}}{\Delta_{j}^{(k)} \epsilon},$$
(13)

for  $j=1,\ldots,p$ , where  $\Delta^{(k)}=\{\Delta_1^{(k)},\ldots,\Delta_p^{(k)}\}$  has each  $\Delta_j^{(k)}\in\{-1,1\}$  independently generated using  $\operatorname{pr}\{\Delta_j^{(k)}=1\}=\operatorname{pr}\{\Delta_j^{(k)}=-1\}=0.5$ . The right hand side is based on the first order approximation to the finite difference form. The advantage is that we only need to evaluate the proximal mapping for m times. In this article, we use  $\epsilon=10^{-7}$  and m=20 and find empirically good stability for the HMC algorithm.

For the HMC as a gradient-based algorithm, another potential concern is that  $\operatorname{prox}_{\lambda g}(\beta)$  may have zero gradient at certain value of  $\beta$ , for example, the soft-thresholding  $\operatorname{sign}(\beta) \max(|\beta| - \lambda, 0)$  will have zero gradient for those  $\beta_j : |\beta_j| < \lambda$ . Fortunately, two things prevent such a  $\beta_j$  from being stuck at a certain value. First, although the log-likelihood  $\log L[y; \operatorname{prox}_{\lambda g}(\beta)]$  may have a zero gradient for  $\beta_j$ , the log-prior  $\log \Pi^0_{\beta}$  does not (as it does not depend on  $\operatorname{prox}_{\lambda g})$ —in those cases,  $\beta_j$  will be updated through its prior distribution, until it enters the region where  $L[y; \operatorname{prox}_{\lambda g}(\beta)]$  is no longer invariant in  $\beta_j$ . This behavior is quite similar to the one with augmented "continuous particle" for sampling binary distribution via HMC (Pakman and Paninski 2013),

where they demonstrated excellent mixing of Markov chains. Second, the HMC preserves the joint density of  $\Pi(\beta \mid y)\Pi(v) = \Pi(\beta^* \mid y)\Pi(v^*)$  (with the MH correction), and as we sample a new v at the start of each iteration, the effective range of  $\beta^*$  to reach is  $\{\beta^* : \Pi(\beta^* \mid y) = \Pi(\beta \mid y)\Pi(v)/\Pi(v^*), \Pi(v^*) > 0, \Pi(v^*) \leq \Pi(\hat{v})\}$ , with  $\hat{v}$  the mode of  $\Pi(v)$ . Therefore, as long as  $\Pi(v) < \Pi(\hat{v})$ , we can have  $\Pi(v^*) > \Pi(v)$ , and  $\Pi(\beta^* \mid y) < \Pi(\beta \mid y)$  allowing  $\beta^*$  to move away from a local-optimal state. In practice, we use the No-U-Turn algorithm (Hoffman and Gelman 2014), which ensures that we run the dynamics for long enough, so that the new proposal is away from the current state. We provide some diagnostic plots in S2 of the supplementary materials.

### 5. Simulation Studies

### 5.1. Set Expansion Prior for Hypothesis Testing

We now demonstrate the usefulness of the proximal prior in standard statistical inference, such as the hypothesis testing of whether  $\theta$  is in a constrained set C. Consider two hypotheses  $H_0: \theta \in C$  and  $H_1: \theta \in \bar{C}$ , where  $C \cap \bar{C} = \emptyset$ . For testing, one typically assumes a mixture prior

$$\Pi_{\boldsymbol{\theta}}^{0}(\theta) = p_0 \phi(\theta) \mathbf{1}(\theta \in C) + (1 - p_0) \bar{\phi}(\theta) \mathbf{1}(\theta \in \bar{C}), \quad (14)$$

where  $\phi$  and  $\bar{\phi}$  are the prior kernel function of  $\theta$  under  $H_0$  and  $H_1$ , respectively; and  $p_0$  is the prior probability assigned to C. The Bayes factor of  $H_0$  relative to  $H_1$  is defined as

$$BF_{01} = \frac{\int_C L(y;\theta)\phi(\theta)d\theta}{\int_{\bar{C}} L(y;\theta)\bar{\phi}(\theta)d\theta} = \frac{\operatorname{pr}(\theta \in C \mid y)}{\operatorname{pr}(\theta \in \bar{C} \mid y)} \frac{(1-p_0)}{p_0},$$

for which, a smaller value of BF<sub>01</sub> provides stronger evidence against  $H_0$ . Often, C is not of the same dimension with  $\bar{C}$ . For example, when testing a point null hypothesis  $H_0: \theta_1=0$ , we have  $\dim(C)<\dim(\bar{C})$ . The standard practice has been assigning appropriate  $\phi$  under C (and  $\bar{\phi}$  under  $\bar{C}$ ), with  $\phi(\theta)$  and  $\bar{\phi}(\theta)$  being 0 on  $\bar{C}$  and C, respectively. However, when the null hypothesis is low-dimensional, such as testing linear equality, assigning density supported on the null set C can become quite challenging.

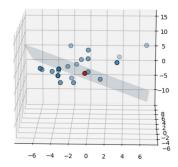
For a convex null set C, we could define a proximal prior based on the distance function, such that the prior density is positive on both C and  $\bar{C}$ . The distance function from point  $\beta$  to set C is defined as  $\mathrm{dist}_C(\beta) = \inf_{x \in C} \|x - \beta\|_2 = \|\beta - P_C(\beta)\|_2$ . The proximal mapping of the distance function to set C is of the form

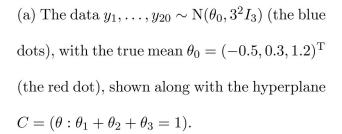
$$\operatorname{prox}_{\lambda \operatorname{dist}_{C}}(\beta) = \begin{cases} \beta + \frac{\lambda}{\operatorname{dist}_{C}(\beta)} \{ P_{C}(\beta) - \beta \}, & \text{if } \operatorname{dist}_{C}(\beta) \geq \lambda \\ P_{C}(\beta), & \text{if } \operatorname{dist}_{C}(\beta) < \lambda. \end{cases}$$

Clearly, this proximal mapping projects the points in the  $\lambda$ -neighborhood of C into C, and keeps the rest of the points out of C. Thus we get a prior that puts positive mass on both C and  $\bar{C}$  and can also be expressed in the form of (14).

We can easily estimate the Bayes factor

$$BF_{01} = \frac{\operatorname{pr}\{\operatorname{dist}_{C}(\beta) < \lambda \mid y\}}{\operatorname{pr}\{\operatorname{dist}_{C}(\beta) > \lambda \mid y\}} \frac{\operatorname{pr}\{\operatorname{dist}_{C}(\beta) \geq \lambda\}}{\operatorname{pr}\{\operatorname{dist}_{C}(\beta) < \lambda\}}$$





**Figure 3.** The set expansion prior for testing  $\{\theta = (\theta_1, \theta_2, \theta_3) : \theta_1 + \theta_2 + \theta_3 = 1\}$ .

via posterior sampling methods. If the prior ratio  $\operatorname{pr}(\theta \in C)/\operatorname{pr}(\theta \in \bar{C})$  is not specified, in order to obtain adequate number of samples in both C and  $\bar{C}$ , we can choose a fixed  $\lambda$  (instead of assigning a prior on  $\lambda$ ) such that  $\operatorname{pr}\{\operatorname{dist}_C(\beta) \geq \lambda\} \approx \operatorname{pr}\{\operatorname{dist}_C(\beta) < \lambda\}$ .

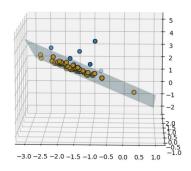
We conduct a simulated experiment: we have data  $y_1,\ldots,y_{20}$  generated from  $N(\theta_0,3^2)$  with  $\theta_0=(-0.5,0.3,1.2)^T$ . We want to test the linear equality hypothesis  $H_0:\theta_1+\theta_2+\theta_3=1$  against  $H_1:\theta_1+\theta_2+\theta_3\neq 1$ . The null set C is a hyperplane with Hausdorff dimension 2. We assign the set expansion prior to  $\theta$  by assign  $N(0,3^2)$  to  $\beta$  and set  $\theta=\operatorname{prox}_{\lambda\operatorname{dist}_C}(\beta)$  with  $\lambda=2$ , such that the prior ratio of C and C is around 0.48. Posterior sampling is implemented with the HMC with 5000 samples and 2000 burn-ins. We get an estimated Bayes factor  $BF_{01}=0.77$ , and display 100 of the samples in Figure 3, panel (b).

### 5.2. Numerical Experiments on Variable Selection and Low Rank Matrix Model

In addition, we conduct numerical experiments for two models where solutions exist with conventional sparse priors: variable selection using a spike-and-slab prior, and low-rank matrix factorization with a discrete prior on the rank. We compare the computational performance in the combinatorial search-based MCMC algorithms for these models, with the HMC algorithm for our models using proximal priors. Further, we compare with other alternatives such as neuronized prior (Shin and Liu 2021) and multiplicative shrinkage prior (Bhattacharya and Dunson 2011; Legramanti, Durante, and Dunson 2020). We provide the details in S1.2 of the supplementary materials.

## 6. Data Application: Interpretable Factor Analysis of the Flow Network

We now demonstrate the practical usefulness of the proximal prior via analyzing the dynamic flow network data. The data



(b) The posterior samples from  $\pi(\theta \mid y)$ , which are mostly distributed on the hyperplane (the orange dots), while there are several outliers far from this hyperplane (the blue dots).

(Zhu, Hu, and Collins 2020) include dynamic estimated traffic flow on major roads in Florida every 6 hr before Hurricane Irma made landfall until it covered the entire state (between 18:00 on September 06, 2017 and 18:00 on September 11, 2017). In total, the data contain 25 valid temporal records of flow networks, denoted by  $Y^{(1)}, \ldots, Y^{(25)}$ ; each  $Y^{(t)} \in \mathbb{R}^{n_V \times n_V}$  contains the traffic flows during a 6-hour period on the roads between  $n_V = 382$  urban regions.

Each flow network is a weighted graph  $\{V, E, Y^{(t)}\}$ , with  $V = (1, \ldots, n_V)$  the set of  $n_V$  nodes,  $E = \{(i, j) : i, j \in V\}$  the edges, and the weight  $Y_{i,j}^{(t)} \in \mathbb{R}$ , representing the amount of flow between the two nodes, with  $Y_{i,j}^{(t)} > 0$  a flow  $i \to j$ , and  $Y_{i,j}^{(t)} < 0$  a flow  $j \to i$ . On the diagonal,  $Y_{i,i}^{(t)} > 0$  indicates an external in-flow entering the network, while  $Y_{i,i}^{(t)} < 0$  means an existing out-flow;  $Y_{i,j}^{(t)} = 0$  if  $(i,j) \notin E$ .

To find useful patterns underneath the raw observation data, we use a low-dimensional latent factor model, with the factors  $F^{(l)} \in \mathbb{R}^{n_V \times n_V}$  shared by all time points, while letting the loadings  $\gamma_l^{(t)} \geq 0$  vary over time, subject to Gaussian measurement error  $\mathcal{E}_{i,j}^{(t)} \stackrel{\text{iid}}{\sim} \mathrm{N}(0,\sigma_{\mathcal{E}}^2)$  for  $i \leq j$ .

$$Y^{(t)} = \sum_{l=1}^{d} \gamma_l^{(t)} F^{(l)} + \mathcal{E}^{(t)}.$$
 (15)

Now, to make the factors useful in interpretation, we require each  $F^{(l)}$  to be a *feasible flow*—an idealized flow network satisfying the following constraints, (i) skew-symmetry (except for the diagonal):  $F_{i,j}^{(l)} = -F_{j,i}^{(l)}$  for i < j; (ii) flow-conservation, that the net sum of in-flows should equal to the out-flows for node j,  $\sum_{i=1}^{n_V} F_{i,j}^{(l)} = 0$ ; (iii) to reduce noise, we assume that the elements of  $F^{(l)}$  are sparse. Further, as we expect that most of the nodes do not have an external in-low/out-flow, we assume that (iv) most of the nodes having  $F_{j,j}^{(l)} = \sum_{i \neq j} F_{i,j}^{(l)} = 0$ . To

obtain the parameter in such a highly constrained space, we use the proximal mapping, with  $\beta^{(l)} \in \mathbb{R}^{n_V \times n_V}$ ,

$$F^{(l)} = \operatorname{prox}\{\beta^{(l)}\} = \underset{z \in \mathbb{R}^{n_V \times n_V}}{\operatorname{arg \, min}} \ \lambda_1 \|z\|_1 + \lambda_1 \sum_{j=1}^{n_V} |z_{j,j}| + \frac{1}{2} \|\beta^{(l)} - z\|_2^2,$$
subject to  $z_{i,j} = -z_{j,i}$  for  $i \neq j$  and  $\sum_{i=1}^{n_V} z_{i,j} = 0$  for  $j = 1, \dots, n_V$ .

The proximal mapping does not have a closed-form solution, however, can be efficiently computed using the alternating direction method of multipliers. Note that with constraint (i) and (ii), it is sufficient to use the lower-triangular entries to represent the rest. In the following, we use  $(\beta_{i,j})_{i>j}$  to denote the  $n_V(n_V-1)/2$ vector containing the lower-triangular entries. Therefore, (17) is equivalent to

$$\operatorname{prox}\{(\beta_{i,j})_{i>j}\} = \underset{\{z_{i,j}\}_{i>j}}{\operatorname{arg \, min}} \frac{1}{2} \sum_{i>j} (\beta_{i,j} - z_{i,j})^2 + \lambda_1 \sum_{i>j} |z_{i,j}| + \lambda_2 \sum_{j=1}^{n_V} |\sum_{i>j} z_{i,j} - \sum_{i< j} z_{j,i}|.$$

$$(17)$$

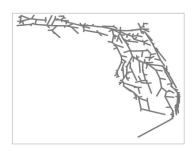
This proximal operator is evaluated via the alternating direction method of multipliers and solved iteratively. We provide the detailed algorithm in the Appendix.

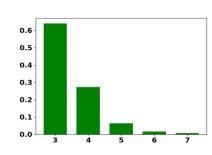
On the prior of the loading, we assign a group shrinkage prior by using the (2, 1)-matrix norm in the proximal mapping. For the matrix  $\gamma \in \mathbb{R}^{d \times T}$ , we set:

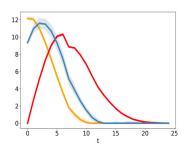
$$\gamma = \operatorname{prox}_{\lambda_2 \| \cdot \|_{2,1}}(\rho) = \operatorname*{arg\,min}_{z \in \mathbb{R}^{d \times T}} \lambda_2 \|z\|_{2,1} + \frac{1}{2} \|\rho - z\|_F^2,$$

where  $\|z\|_{2,1} = \sum_{l=1}^d \sqrt{\sum_{t=1}^T \{z_l^{(t)}\}^2}$ . This prior has the advantage that  $\{\gamma_l^{(1)}, \dots, \gamma_l^{(T)}\}$  will be simultaneously zero for certain l—which allows us to use an overfitted model with a relatively large d = 10, with the posterior recovering only a small number of factors with nonzero loadings. We use independent standard normal as prior on the elements of  $\beta^{(l)}$  and  $\rho$ .

We run the HMC for 20,000 steps and discard the first 5000 as burn-ins, and we use thinning at every 10th iteration as the posterior sample. The posterior shows the highest probability at having three factors, and we visualize them in Figure 4(d) clearly, by forcing the external in-flows and out-flows to be sparse, we have each factor roughly corresponding to a single connected sub-network.





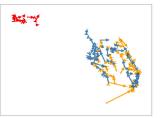


(a) The traffic network show- (b) The posterior distribu- (c) The estimated three loading the roads across the state tion of the number of non- ings  $\gamma_l^{(t)}$  changing over time. of Florida. zero factors.









(d)Estimated com- (e) Edge-wise Variance (f) Estimated compo- (g) Estimated compoby proximal of proximal prior. nents by elementwise nents row-sparse horseshoe. horseshoe. prior.

Figure 4. Data analysis on the dynamic flow network, observed during the Hurricane Irma evacuation. A latent factor model is fitted to the data, with the factors regularized by a proximal prior that force each network to be sparse in both the flows and the number of in-flows and out-flows. (d) shows three factor flows (posterior mode) estimated using the proximal prior. Each factor is close to a connected sub-flow network with (e) element-wise variance representing the uncertainty on the factor estimate. We compare the proximal prior with two realizations of shrinkage priors, where (f) shows the factors estimated using elementwise horseshoe prior, which are very fragmented and difficult to interpret. (g) shows the factors estimated using the row-sparse horseshoe, which contain many nodes with small or no flows.

Interestingly, examining the estimated loadings that change over the time points in Figure 4(c), we see that in the beginning of the evacuation, the factors 1 and 2 are dominant, but later there is a sudden decrease—this in fact corresponds to the time point when the hurricane makes the landfall, effectively forcing the traffic in those areas to shut down. After the 60th hr, the traffic moves up to the north part, and factor 3 represents the late stage of the evacuation.

To compare, we also test two continuous shrinkage priors on the factors. Specifically, we use (a) the elementwise horseshoe prior on each lower-diagonal  $F_{i,j}^{(l)} \sim \mathrm{N}(0,\tau_{i,j}^2\sigma^2), \, \tau_{i,j} \sim$  $C^+(0,1)$ ,  $\sigma^2 \sim \text{Inverse-Gamma}(2,0.01)$ , and (b) the twoway group horseshoe prior, by letting  $F_{i,j}^{(l)} \sim N(0, \tau_i \tau_j \sigma^2)$ ,  $\tau_i \sim C^+(0,1), \sigma^2 \sim \text{Inverse-Gamma}(2,0.01).$  The purpose of (b) is to shrink each row of F simultaneously, while satisfying the skew-symmetry of  $F^{(l)}$ . Effectively these horseshoe priors accommodate the properties of (i)(ii)(iii) of a sparse feasible flow, nevertheless, they cannot accommodate (iv)—as each  $F_{i,i}^{(l)}$ is completely determined given  $F_{i,j}^{(l)}$  ( $i \neq j$ ), we could not further assign shrinkage prior on  $F_{j,j}^{(l)} = \sum_{i \neq j} F_{j,j}^{(l)}$ . As the result in Figure 4(f), the elementwise continuous shrinkage priors show a large number of external in-flows and out-flows, leading to fragmented small networks in each factor. As shown in Figure 4(g), the two-way group horseshoe finds many nodes with no flows at all, which is not very interpretable since we would like most of the nodes to have many in-flows and out-flows, as long as the total net-flow is zero.

### 7. Discussion

In this article, we exploit the proximal mapping to produce a new class of priors. As we have demonstrated, these priors and the associated probabilistic models can enable statistical inference (such as uncertainty quantification, hypothesis testing) on a wide range of problems, where in the past, one has been limited to point estimate only. The technique of "data augmentation using optimization" we have introduced could be generalized for other purposes, such as potentially new efficient algorithm for the posterior computation. Lastly, one could consider other type of optimization problems for a similar prior construction, such as the popular classes of semidefinite (Vandenberghe and Boyd 1996) and / or mixed integer programmings (Linderoth and Savelsbergh 1999), although how to provide a probabilistic treatment for these problems is still an open question.

### **Appendix**

### A.1. Proof of Theorem 1

*Proof.* For any value  $\beta^*$ :  $\operatorname{prox}_{\lambda g_{\gamma}}(\beta^*) = \theta^*$  using  $g_{\gamma}$  at given  $\gamma$  and  $\lambda$ ,  $\epsilon < \|\operatorname{prox}_{\lambda g_{\gamma}}(\beta) - \operatorname{prox}_{\lambda g_{\gamma}}(\beta^*)\| \le \|\beta - \beta^*\|$ . Therefore,  $\mathbf{1}(\theta, \|\theta - \theta^*\| > \epsilon \mid y, \lambda, \gamma) \le \mathbf{1}\left(\beta, \|\beta - \beta^*\| > \epsilon \mid y, \lambda, \gamma\right)$ . Taking the minimum over  $\beta^*$  on the right hand side and expectation on both sides, we obtain the first result.

Next, using the fact that for two independent copies  $\beta_1$ ,  $\beta_2$  from  $\Pi(\beta \mid y)$ ,  $2\text{tr}[\text{cov}(\beta \mid y) = \mathbb{E}_{\beta_1,\beta_2}(\|\beta_1 - \beta_2\|_2^2 \mid y) =$ 

 $\mathbb{E}_{\lambda,\gamma}\mathbb{E}_{\beta_1,\beta_2}(\|\beta_1 - \beta_2\|_2^2 \mid \lambda,\gamma,y)$ , and the non-expansiveness of proximal mappings, we obtain the second result.

### A.2. Proof of Theorem 2

*Proof.* Let  $0 < \lambda_1 < \lambda_2, v_1 = \text{prox}_{\lambda_1 g}(x)$  and  $v_2 = \text{prox}_{\lambda_2 g}(x)$ , we prove: (i)  $g(v_1) \ge g(v_2)$  and (ii)  $||v_1 - x|| \le ||v_2 - x||$ . For (i),

$$\frac{1}{2} \|v_2 - x\|^2 + \lambda_2 g(v_2) = \frac{1}{2} \|v_2 - x\|^2 + \lambda_1 g(v_2) + (\lambda_2 - \lambda_1) g(v_2) 
\stackrel{(a)}{\geq} \frac{1}{2} \|v_1 - x\|^2 + \lambda_1 g(v_1) + (\lambda_2 - \lambda_1) g(v_2) 
= \frac{1}{2} \|v_1 - x\|^2 + \lambda_2 g(v_1) + (\lambda_2 - \lambda_1) 
\left\{ g(v_2) - g(v_1) \right\} 
\stackrel{(b)}{\geq} \frac{1}{2} \|v_2 - x\|^2 + \lambda_2 g(v_2) + (\lambda_2 - \lambda_1) 
\left\{ g(v_2) - g(v_1) \right\},$$

where (a) is due to  $v_1$  is the minimizer of  $\frac{1}{2} \|v - x\|^2 + \lambda_1 g(v)$ , and similarly for  $v_2$  in (b). Therefore,  $(\lambda_2 - \lambda_1) \{g(v_2) - g(v_1)\} \le 0$ , which leads to  $g(v_1) \ge g(v_2)$ .

For (ii), slightly changing (a), we have  $\frac{1}{2} \|v_1 - x\|^2 \le \frac{1}{2} \|v_2 - x\|^2 - \lambda_1 \{g(v_1) - g(v_2)\}$ . Since  $\lambda_1 \ge 0$ , we have  $\|v_1 - x\|^2 \le \|v_2 - x\|^2$ .

### A.3. Proof of Theorem 3

*Proof.* Since the proximal mapping satisfies 1-Lipschitz condition, diam{prox(B)}  $\leq$  diam(B) for all B in the domain of prox $_{\lambda g}$ . Using the definition of the Hausdorff measure,  $\mathcal{H}^s\{\operatorname{prox}(A)\} \leq \mathcal{H}^s(A)$ . To see the statement 2, we apply the result in statement 1 and see  $\mathcal{H}^s\{\operatorname{prox}(A)\} = 0$  whenever  $\mathcal{H}^s(A) = 0$ .

### A.4. Proof of Theorem 5

*Proof.* Using Theorem 3.2.22 of Federer (2014), for any  $\mathcal{H}^p$  measurable function F on  $\boldsymbol{\beta}^k$ ,

$$\int_{\pmb{\beta}^k} F(\beta) J_k \mathrm{prox}_{\lambda g}(\beta) \mathcal{H}^p(\mathrm{d}\beta) = \int_{\Theta^k} \int_{\mathrm{prox}_{\lambda g}^{-1}(\theta)} F(b) \mathcal{H}^{p-k}(\mathrm{d}b) \mathcal{H}^k(\mathrm{d}\theta).$$

Using the assumption, we can exclude the zero-measure set where  $J_{m_k} \mathrm{prox}_{\lambda g}(\beta) = 0.$ 

### A.5. Algorithm to Compute the Proximal Mapping in the Flow Network Modeling

We formulate an equivalent problem to (17)

$$\operatorname{prox}\{(\beta_{i,j})_{i>j}\} = \underset{\{z_{i,j}\}_{i>j}, x \in \mathbb{R}^{n_V}}{\operatorname{arg \, min}} \frac{1}{2} \sum_{i>j} (\beta_{i,j} - z_{i,j})^2 + \lambda_1 \sum_{i>j} |z_{i,j}| + \lambda_2 \sum_{j=1}^{n_V} |x_j|$$

$$\operatorname{subject \, to} C(z_{i,j})_{i>j} = x,$$

where  $C \in \mathbb{R}^{n_V \times n_V(n_V - 1)/2}$  is the matrix such that  $\{C(z_{i,j})_{i>j}\}_k = \sum_{i>k} z_{i,k} - \sum_{i< k} z_{k,i}$ . The scaled augmented Lagrangian for (18) is

$$\mathcal{L} = \frac{1}{2} \| (\beta_{i,j})_{i>j} - (z_{i,j})_{i>j} \|_2^2 + \lambda_1 \| (\beta_{i,j})_{i>j} \|_1 + \lambda_2 \| x \|_1$$
$$+ \frac{\gamma}{2} \| C(z_{i,j})_{i>j} - x + u \|_2^2 - \frac{\gamma}{2} \| u \|_2^2.$$



In each iteration we update the *x* and *z* separately to minimize the Lagrangian:

$$(z_{i,j})_{i>j} \leftarrow \operatorname{prox}_{\lambda_1/\gamma \|\cdot\|_1} [(I + \gamma C^{\mathrm{T}} C)^{-} \{ \gamma C^{\mathrm{T}} (x - u) + \beta \}];$$
  
$$x \leftarrow \operatorname{prox}_{\lambda_2/\gamma \|\cdot\|_1} \{ C(z_{i,j})_{i>j} + u \},$$

and update *u* as in dual ascent:

$$u \leftarrow u + C(z_{i,j})_{i>j} - x$$

until convergence (i.e.,  $||C(z_{i,j})_{i>j} - x|| \rightarrow 0$ ).

### **Supplementary Materials**

Supplementary materials include (i) additional examples of proximal priors, along with additional experimental results and (ii) code to reproduce the results in this paper.

### **Funding**

The authors would like to thank the support from the following funding. This research was partially funded by grants from the University of Florida Informatics Institute SEED Fund (UFII-SEED-2022: LD), the National Institute of General Medical Sciences (R35GM141798: HZ), the National Human Genome Research Institute (R01HG006139: HZ), and the National Science Foundation (DMS-2054253 and IIS-2205441: HZ).

### **ORCID**

Hua Zhou https://orcid.org/0000-0003-1320-7118

### References

- Beck, A. (2017), First-order Methods in Optimization, Philadelphia, PA: SIAM. [2]
- Berger, J. O., Bernardo, J. M., Sun, D., et al. (2009), "The Formal Definition of Reference Priors, Annals of Statistics, 37, 905-938. [5]
- Bertsekas, D. P. (2014), Constrained Optimization and Lagrange Multiplier Methods, New York: Academic Press. [2]
- Bhattacharya, A., and Dunson, D. B. (2011), "Sparse Bayesian Infinite Factor Models, Biometrika, 98, 291-306. [1,9]
- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. (2009), "Simultaneous Analysis of Lasso and Dantzig Selector," The Annals of Statistics, 37, 1705-1732. [1]
- Carvalho, C. M., Polson, N. G., and Scott, J. G. (2009), "Handling Sparsity via the Horseshoe," in Artificial Intelligence and Statistics, pp. 73-80. [1]
- Cook, R. D., Li, B., and Chiaromonte, F. (2010), "Envelope Models for Parsimonious and Efficient Multivariate Linear Regression," Statistica Sinica, 20, 927-960. [4]
- DiMatteo, I., Genovese, C. R., and Kass, R. E. (2001), "Bayesian Curve-Fitting With Free-Knot Splines, Biometrika, 88, 1055–1071. [1]
- Dunson, D. B., and Neelon, B. (2003), "Bayesian Inference on Orderconstrained Parameters in Generalized Linear Models," Biometrics, 59, 286-295. [2]
- Edgar, G. (2007), Measure, Topology, and Fractal Geometry, New York: Springer. [6]
- Federer, H. (2014), Geometric Measure Theory, Berlin: Springer. [7,11]
- Goodall, C. R. and Mardia, K. V. (1999), "Projective Shape Analysis," Journal of Computational and Graphical Statistics, 8, 143–168. [1]
- Green, P. J., and Hastie, D. I. (2009), "Reversible Jump MCMC," Genetics, 155, 1391–1403. [1]
- Hoffman, M. D., and Gelman, A. (2014), "The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo," Journal of Machine Learning Research, 15, 1593-1623. [8]
- Ji, H., Liu, C., Shen, Z., and Xu, Y. (2010), "Robust Video Denoising using Low Rank Matrix Completion," in 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 1791–1798, IEEE. [1]
- Kastner, G. (2019), "Sparse Bayesian Time-Varying Covariance Estimation in Many Dimensions," Journal of Econometrics, 210, 98-115. [1]

- Khare, K., Pal, S., and Su, Z. (2017), "A Bayesian Approach for Envelope Models," The Annals of Statistics, 45, 196-222. [4]
- Legramanti, S., Durante, D., and Dunson, D. B. (2020), "Bayesian Cumulative Shrinkage for Infinite Factorizations," Biometrika, 107, 745-752.
- Li, Y., Craig, B. A., and Bhadra, A. (2019), "The Graphical Horseshoe Estimator for Inverse Covariance Matrices," Journal of Computational and Graphical Statistics, 28, 747-757. [1]
- Lin, L., and Dunson, D. B. (2014), "Bayesian Monotone Regression Using Gaussian Process Projection," Biometrika, 101, 303-317. [2]
- Linderoth, J. T., and Savelsbergh, M. W. (1999), "A Computational Study of Search Strategies for Mixed Integer Programming," INFORMS Journal on Computing, 11, 173-187. [11]
- Meier, L., Van De Geer, S., and Bühlmann, P. (2008), "The Group Lasso for Logistic Regression," Journal of the Royal Statistical Society, Series B, 70, 53-71. [1]
- Meinshausen, N., and Bühlmann, P. (2006), "High-Dimensional Graphs and Variable Selection with the Lasso," The Annals of Statistics, 34, 1436-
- Mitchell, T. J., and Beauchamp, J. J. (1988), "Bayesian Variable Selection in Linear Regression," Journal of the American Statistical Association, 83, 1023-1032. [1]
- Morgan, F. (2016), Geometric Measure Theory: A Beginner's Guide, New York: Academic Press. [7]
- Pakman, A., and Paninski, L. (2013), "Auxiliary-Variable Exact Hamiltonian Monte Carlo Samplers for Binary Distributions," in Advances in Neural Information Processing Systems (Vol. 26). [8]
- Park, T., and Casella, G. (2008), "The Bayesian Lasso," Journal of the American Statistical Association, 103, 681-686. [1]
- Petris, G., and Tardella, L. (2003), "A Geometric Approach to Transdimensional Markov chain Monte Carlo," Canadian Journal of Statistics, 31, 469-482. [2]
- Preiss, D. (1987), "Geometry of Measures in Rn: Distribution, Rectifiability, and Densities," Annals of Mathematics, 125, 537-643. [6]
- Rockova, V., and George, E. I. (2018), "The Spike-And-Slab Lasso," Journal of the American Statistical Association, 113, 431-444. [1]
- Saarela, O., and Arjas, E. (2011), "A Method for Bayesian Monotonic Multiple Regression," Scandinavian Journal of Statistics, 38, 499-513. [1]
- Sen, D., Patra, S., and Dunson, D. B. (2018), "Constrained Bayesian Inference through Posterior Projections," arXiv preprint arXiv:1812.05741.
- Shen, H., and Huang, J. Z. (2008), "Sparse Principal Component Analysis via Regularized Low Rank Matrix Approximation," Journal of Multivariate Analysis, 99, 1015-1034. [1]
- Shin, M., Bhattacharya, A., and Johnson, V. E. (2020), "Functional Horseshoe Priors for Subspace Shrinkage," Journal of the American Statistical Association, 115, 1784-1797. [1]
- Shin, M., and Liu, J. S. (2021), "Neuronized Priors for Bayesian Sparse Linear Regression," Journal of the American Statistical Association, 117, 1695-1710. [5,9]
- Spall, J. C. (1992), "Multivariate Stochastic Approximation Using A Simultaneous Perturbation Gradient Approximation," IEEE Transactions on Automatic Control, 37, 332-341. [8]
- Tanner, M. A., and Wong, W. H. (1987), "The Calculation of Posterior Distributions by Data Augmentation," Journal of the American Statistical Association, 82, 528-540. [3]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," Journal of the Royal Statistical Society, Series B, 58, 267-288. [2]
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. (2005), "Sparsity and Smoothness via the Fused Lasso," Journal of the Royal Statistical Society, Series B, 67, 91-108. [1]
- Vandenberghe, L., and Boyd, S. (1996), "Semidefinite Programming," SIAM Review, 38, 49-95. [11]
- Zhao, P., Rocha, G., and Yu, B. (2006), "Grouped and Hierarchical Model Selection through Composite Absolute Penalties," Department of Statistics, UC Berkeley, Technical Report, 703. [1]
- Zhu, Y.-J., Hu, Y., and Collins, J. M. (2020), "Estimating Road Network Accessibility During a Hurricane Evacuation: A Case Study of Hurricane Irma in Florida," Transportation Research Part D: Transport and Environment, 83, 102334. [9]