

Generalized double Pareto shrinkage

Artin Armagan
Department of Statistical Science
Duke University
Durham, NC 27708
artin@stat.duke.edu

David B. Dunson
Department of Statistical Science
Duke University
Durham, NC 27708
dunson@stat.duke.edu

Jaeyong Lee
Department of Statistics
Seoul National University
Seoul, 151-747, Korea
leejyc@gmail.com

Original version: March 2010
Revised: February 2011

Abstract

We propose a generalized double Pareto prior for Bayesian shrinkage estimation and inferences in linear models. The prior can be obtained via a scale mixture of Laplace or normal distributions, while forming a bridge between the Laplace and Normal-Jeffreys' priors. While it has a spike at zero like the Laplace density, it also has a Student- t -like tail behavior. Bayesian computation is straightforward via a simple Gibbs sampling algorithm. We investigate the properties of the maximum a posteriori estimator, as many are interested in sparse solutions, reveal connections with some well-established regularization procedures and show some asymptotic results. The performance of the prior is tested through simulations.

Key words: Heavy tails; High-dimensional data; Lasso; Maximum a posteriori estimation; Relevance vector machine; Robust prior; Shrinkage estimation.

1 Introduction

There has been a great deal of work in shrinkage estimation and simultaneous variable selection in the frequentist framework. The Lasso of Tibshirani (1996) has drawn much attention to the area, particularly after the introduction of LARS (Efron et al., 2004) due to its superb computational performance. There is a rich literature analyzing the Lasso and related approaches (Fu, 1998; Knight and Fu, 2000; Fan and Li, 2001; Yuan and Lin, 2005; Zhao and Yu, 2006; Zou, 2006; Zou and Li, 2008), with a number of articles considering asymptotic properties.

Bayesian approaches to the same problem became popular with the works of Tipping (2001) and Figueiredo (2003). By expressing Student- t priors for basis coefficients as scale mixtures of normals (West, 1987) and relying on type II maximum likelihood estimation (Berger, 1985), Tipping (2001) developed the relevance vector machine for sparse estimation in kernel regression. However, sparsity comes with the price of forfeiting propriety of the posterior by driving the degrees of freedom and the scale parameter of the Student- t distribution towards zero. This yielded the so-called Normal-Jeffreys' prior on the parameters, $p(\theta) \propto 1/|\theta|$. Figueiredo (2003) proposed an expectation-maximization algorithm for maximum a posteriori estimation under Laplace and Normal-Jeffreys' priors, with estimates under the Laplace corresponding to the Lasso. The Normal-Jeffreys' prior leads to substantially improved performance due to the property of strongly shrinking small coefficients to zero while minimally shrinking large coefficients due to

the heavy tails. Although interpretable as posterior modes, these estimators and other penalized likelihood estimators do not correspond to Bayes estimators under a reasonable choice of loss function, and hence lack a fully Bayes justification.

A Bayesian Lasso was proposed by Park and Casella (2008) and Hans (2009). However, these procedures inherit the problem of over-shrinking large coefficients due to the relatively light tails of the Laplace prior. Strawderman-Berger priors (Strawderman, 1971; Berger, 1980) have some desirable properties yet lack an analytic form. Recently proposed priors have been designed to have high density near zero and heavy tails without the impropriety problem of Normal-Jeffreys. The horseshoe prior of Carvalho et al. (2009, 2010) is induced through a carefully-specified mixture of normals, leading to desirable properties, such as an infinite spike at zero and very heavy tails. They studied sparse shrinkage estimation properties of the horseshoe in a normal means problem. Griffin and Brown (2007, 2010) proposed an alternative class of hierarchical priors for shrinkage, which has some similarities to the prior we propose, but lacks the simple analytic form facilitating the study of some properties.

There is a need for alternative shrinkage priors that lead to sparse point estimates if desired, do not over-shrink coefficients that are not close to zero, facilitate straightforward computation even in large p cases, and result in a joint posterior distribution that does a good job in quantifying uncertainty. We propose the generalized double Pareto prior which independently finds mention in Cevher (2009). It has a simple analytic form, yields a proper posterior and possesses appealing properties, including a spike at zero, Student- t -like tails, and a simple characterization as a scale mixture of normals leading to a straightforward Gibbs sampler for posterior inferences. We consider both fully Bayesian and frequentist penalized likelihood approaches based on this prior. We show that the induced penalty in the regularization framework yields a consistent thresholding rule having the continuity property in the orthogonal case, with a simple Expectation-Maximization algorithm described for sparse estimation in non-orthogonal cases. Similarities to Cevher (2009) are very limited and the contributions beyond these are (i) a formal introduction of a generalized Pareto density thresholded and folded at zero as a shrinkage prior in Bayesian analysis, (ii) the scale mixture representation of the generalized double Pareto given in Proposition 1 which is central to our work, (iii) its connection to the Laplace and Normal-Jeffreys' priors as limiting cases given in Proposition 2, (iv) the resulting fully conditional posteriors in a linear regression setting along with a simple Gibbs sampling procedure, (v) a discussion on the hyper-parameters α and η and their treatment along with the incorporation of a griddy sampling scheme into the Gibbs sampler, (vi) a detailed analysis of the induced penalty by the generalized double Pareto prior and the properties of the resulting thresholding rule, (vii) an explicit analytic form for the maximum a posteriori estimator in orthogonal cases, (viii) consistency of the resulting thresholding rule with a diverging number of parameters in orthogonal cases, (ix) an expectation-maximization procedure to obtain the maximum a posteriori estimate in non-orthogonal cases using the normal mixture representation given in Section 5.1, and finally (x) the one-step estimator (Zou and Li, 2008) resulting from the Laplace mixture representation and its oracle properties given in Section 5.2 revealing the connection of the resulting procedure to the adaptive Lasso (Zou, 2006).

2 Generalized Double Pareto Prior

The generalized double Pareto density is given by

$$f(\theta|\xi, \alpha) = \frac{1}{2\xi} \left(1 + \frac{|\theta|}{\alpha\xi}\right)^{-(1+\alpha)}, \quad (1)$$

where $\xi > 0$ is a scale parameter and $\alpha > 0$ is a shape parameter. In contrast to (1), the generalized Pareto density of Pickands (1975) is parametrized in terms of a location parameter $\mu \in \mathfrak{R}$, a scale parameter $\xi > 0$, and a shape parameter $\alpha \in \mathfrak{R}$ as follows,

$$f(\theta|\xi, \alpha, \mu) = \frac{1}{\xi} \left(1 + \frac{\theta - \mu}{\alpha\xi}\right)^{-(1+\alpha)}, \quad (2)$$

with $\theta \geq \mu$ for $\alpha > 0$ and $\mu \leq \theta \leq \mu - \xi\alpha$ for $\alpha < 0$. The mean and variance for the generalized Pareto distribution is respectively given by $E(\theta) = \mu + \xi/(1 - 1/\alpha)$ for $\alpha \notin [0, 1]$ and $\text{Var}(\theta) = \xi^2(1 - 1/\alpha)^{-2}(1 - 2/\alpha)^{-1}$ for $\alpha \notin [0, 2]$. If we let $\mu = 0$, (2) becomes an exponential density as $\alpha \rightarrow \infty$ with mean ξ and variance ξ^2 .

To modify the generalized Pareto density to be appropriate as a shrinkage prior, we let $\mu = 0$ and reflect the positive part about the origin assuming $\alpha > 0$. This leads to a density that is symmetric about zero. The mean and variance for the generalized double Pareto distribution is respectively given by $E(\theta) = 0$ for $\alpha > 1$ and $\text{Var}(\theta) = 2\xi^2\alpha^2(\alpha - 1)^{-1}(\alpha - 2)^{-1}$ for $\alpha > 2$. The dispersion is controlled by ξ and α , with α controlling the tail heaviness and $\alpha = 1$ corresponding to Cauchy-like tails and no finite moments.

Figure 1 compares density (1) to Cauchy and Laplace densities in the special case in which $\xi = \alpha = 1$, so that $f(\theta) = 1/\{2(1 + |\theta|)^2\}$. We refer to this form as the standard double Pareto. Near zero the standard double Pareto resembles the Laplace density, suggesting similar sparse shrinkage properties of small coefficients in maximum a posteriori estimation. It also has Cauchy-like tails which is appealing in avoiding over-shrinkage away from the origin. This is illustrated in Figure 1(a). Figure 1(b) illustrates how density (1) changes for different values of ξ and α .

Prior (1) can be represented as a scale mixture of normal distributions leading to computational simplifications. As shorthand notation, let $\theta \sim \text{GDP}(\xi, \alpha)$ denote that θ follows density (1).

Proposition 1. *Let $\theta \sim N(0, \tau)$, $\tau \sim \text{Exp}(\lambda^2/2)$ and $\lambda \sim \text{Ga}(\alpha, \eta)$ where $\alpha > 0$ and $\eta > 0$. The resulting marginal density for θ is $\text{GDP}(\xi = \eta/\alpha, \alpha)$.*

Proposition 1 reveals a relationship between prior (1) and the prior of Griffin and Brown (2007), with the difference being that Griffin and Brown (2007) place a mixing distribution on λ^2 leading to a marginal with no simple analytic form. Proposition 2 shows that prior (1) forms a bridge between Laplace and Normal-Jeffreys' priors.

Proposition 2. *Given the representation in Proposition 1, $\theta \sim \text{GDP}(\xi = \eta/\alpha, \alpha)$ implies*

1. $f(\theta) \propto 1/|\theta|$ for $\alpha = 0$ and $\eta = 0$,
2. $f(\theta|\lambda') = (\lambda'/2) \exp(-\lambda'|\theta|)$ for $\alpha \rightarrow \infty$, $\alpha/\eta = \lambda'$ and $0 < \lambda' < \infty$.

Proof. For the first item, setting $\alpha = \eta = 0$ implies placing a Jeffreys' prior on λ , $p(\lambda) \propto 1/\lambda$. Integration over λ yields $p(\tau) \propto 1/\tau$ which implies the Normal-Jeffreys' prior on θ . For the second item, notice that $p(\lambda) = \delta(\lambda - \lambda')$, where $\delta(\cdot)$ denotes the Dirac delta function, since $\lim_{\alpha \rightarrow \infty} \lim_{\alpha/\eta \rightarrow \lambda'} E(\lambda) = \lambda'$ and $\lim_{\alpha \rightarrow \infty} \lim_{\alpha/\eta \rightarrow \lambda'} \text{Var}(\lambda) = 0$. Thus, $\int_0^\infty (\lambda'/2) \exp(-\lambda|\theta|)\delta(d\lambda) = (\lambda'/2) \exp(-\lambda'|\theta|)$. \square

3 Bayesian Inference

3.1 Posterior Computation

Consider the linear regression model, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where \mathbf{y} is an n -dimensional vector of responses, \mathbf{X} is the $n \times p$ design matrix and $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$. Letting $\beta_j | \sigma \sim \text{GDP}(\xi = \sigma\eta/\alpha, \alpha)$ independently for $j = 1, \dots, p$,

$$\pi(\boldsymbol{\beta} | \sigma) = \prod_{j=1}^p \frac{1}{2\sigma\eta/\alpha} \left(1 + \frac{1}{\alpha} \frac{|\beta_j|}{\sigma\eta/\alpha} \right)^{-(\alpha+1)}. \quad (3)$$

From Proposition 1 this prior is equivalent to $\beta_j | \sigma \sim N(0, \sigma^2 \tau_j)$, with $\tau_j \sim \text{Exp}(\lambda_j^2/2)$ and $\lambda_j \sim \text{Ga}(\alpha, \eta)$. We place the Jeffreys' prior on the error variance, $\pi(\sigma) \propto 1/\sigma$.

Using the scale mixture of normals representation, we obtain a simple data augmentation Gibbs sampler having the following conditional posteriors: $(\boldsymbol{\beta} | \sigma^2, \mathbf{T}, \mathbf{y}) \sim N\{(\mathbf{X}'\mathbf{X} + \mathbf{T}^{-1})^{-1} \mathbf{X}'\mathbf{y}, \sigma^2 (\mathbf{X}'\mathbf{X} + \mathbf{T}^{-1})^{-1}\}$, $(\sigma^2 | \boldsymbol{\beta}, \mathbf{T}, \mathbf{y}) \sim \text{IG}\{(n + p)/2, (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \boldsymbol{\beta}'\mathbf{T}^{-1}\boldsymbol{\beta}\}/2\}$, $(\lambda_j | \beta_j, \sigma^2) \sim \text{Ga}(\alpha + 1, |\beta_j|/\sigma + \eta)$, $(\tau_j^{-1} | \beta_j, \lambda_j, \sigma^2) \sim \text{Inv-Gauss}\{\mu = (\lambda_j^2 \sigma^2 / \beta_j^2)^{1/2}, \rho = \lambda_j^2\}$ where $\mathbf{T} = \text{diag}(\tau_1, \dots, \tau_p)$ and Inv-Gauss denotes the inverse Gaussian distribution with location and scale parameters μ and ρ . In our experience, this Gibbs sampler is efficient having fast rates of convergence and mixing.

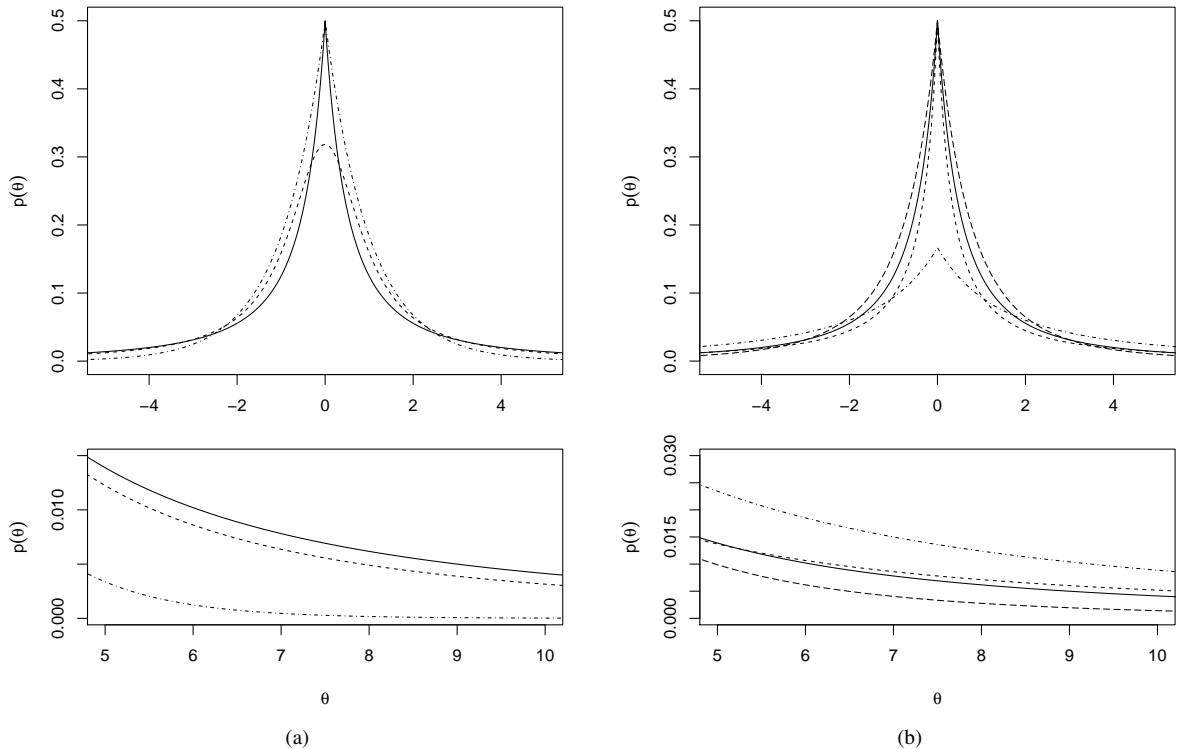


Figure 1: (a) Probability density functions for standard double Pareto (solid line), standard Cauchy (dashed line) and Laplace (dot-dash line) ($\lambda = 1$) distributions. (b) Probability density functions for the generalized double Pareto with (ξ, α) values of (1, 1) (solid line), (0.5, 1) (dashed line), (1, 3) (long-dashed line), and (3, 1) (dot-dash line).

3.2 Hyper-prior Specification and Computation

As α grows the density becomes lighter tailed, more peaked and the variance becomes smaller while as η grows the density becomes flatter and the variance increases. Hence if we increase α , we may cause unwanted bias for large signals, though causing stronger shrinkage for noise-like signals; if we increase η we may lose the ability to shrink noise-like signals also causing less bias for large signals; and finally, If we increase α and η at the same rate, the variance remains constant but the tails become lighter converging to a Laplace density in the limit. This can lead to over-shrinking of the coefficients that are away from zero. Given that the columns of \mathbf{X} are scaled to be of unit length, as a typical default specification for the hyper-parameters, one can let $\alpha = \eta = 1$ (thus $\xi = \sigma$) in (3). This choice leads to Cauchy-like tail behavior, which is well-known to have desirable Bayesian robustness properties.

To further motivate our default choice, we assess the behavior of the prior shrinkage factor $\kappa = 1/(1 + \tau) \in (0, 1)$ where $\theta \sim N(0, \tau)$ is the parameter of interest (Carvalho et al., 2010). As $\kappa \rightarrow 0$, the prior does not impose any shrinkage while as $\kappa \rightarrow 1$ it has a strong pull towards zero. The generalized double Pareto distribution implies a prior $p(\kappa)$ on κ upon integration over λ in Proposition 1. For the standard double Pareto, this is given by

$$p(\kappa) = \frac{1}{2(1 - \kappa)^2} \left[\frac{\sqrt{\pi} \exp \left\{ \frac{\kappa}{2(1 - \kappa)} \right\} \operatorname{Erfc} \left\{ \sqrt{\frac{\kappa}{2(1 - \kappa)}} \right\}}{\sqrt{2\kappa(1 - \kappa)}} - 1 \right],$$

where $\operatorname{Erfc}(\cdot)$ denotes the complementary error function. In Figure 2, we compare $p(\kappa)$ under the standard double

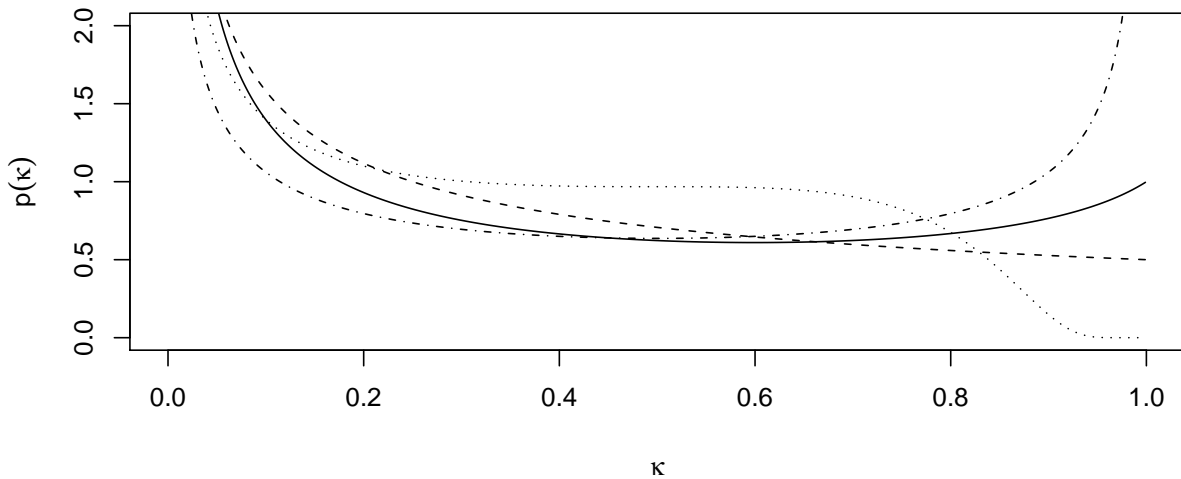


Figure 2: Prior density of κ implied by the standard double Pareto prior (solid line), Strawderman–Berger prior (dashed line), horseshoe prior (dot-dash line) and standard Cauchy prior (dotted line).

Pareto, Strawderman-Berger, horseshoe and Cauchy priors. The priors behave similarly for $\kappa \approx 0$, implying similar tail behavior. The behavior of $p(\kappa)$ for $\kappa \approx 1$ governs the strength of shrinkage of small signals. As $\kappa \rightarrow 1$, $p(\kappa)$ tends towards zero for the Cauchy implying weak shrinkage, while $p(\kappa)$ is unbounded for the horseshoe suggesting very strong pull towards zero for small signals. The Strawderman-Berger and standard double Pareto priors are a compromise between these extremes, with $p(\kappa)$ being bounded for $\kappa \rightarrow 1$ in both cases. The standard double Pareto assigns higher density close to one than the Strawderman-Berger prior and has the advantage of a simple analytic form and a conjugate hierarchy over the Strawderman-Berger and horseshoe priors.

As an alternative we recommend choosing hyper-priors to allow the data to inform about the values of α and η , with $p(\alpha) = 1/(1+\alpha)^2$ and $p(\eta) = 1/(1+\eta)^2$ to correspond to generalized Pareto hyper-priors with location 0, scale 1 and shape 1. The median value of the resulting distribution for α and η is one, centering it at the default choices suggested earlier while the mean and variance do not exist.

For sampling purposes let $a = 1/(1+\alpha)$ and $e = 1/(1+\eta)$. These transformations suggest a uniform prior on a and e in $(0, 1)$ given the generalized Pareto priors on α and η . Consequently, the conditional posteriors for a and e are

$$p(a|\boldsymbol{\beta}, \eta) \propto \left(\frac{1-a}{a}\right)^p \prod_{j=1}^p \left(1 + \frac{|\beta_j|}{\sigma\eta}\right)^{-1/a},$$

$$p(e|\boldsymbol{\beta}, \alpha) \propto \left(\frac{e}{1-e}\right)^p \prod_{j=1}^p \left\{1 + e \frac{|\beta_j|}{\sigma(1-e)}\right\}^{-(\alpha+1)}.$$

We propose the following embedded griddy Gibbs (Ritter and Tanner, 1992) sampling scheme:

- i. Form a grid of m points $a^{(1)}, \dots, a^{(m)}$ in the interval $(0, 1)$,
- ii. Calculate $w^{(k)} = p(a^{(k)}|\boldsymbol{\beta}, \eta)$,
- iii. Normalize the weights, $w_N^{(k)} = w^{(k)} / \sum_{k=1}^m w^{(k)}$,

- iv. Draw a sample from the set $\{a^{(1)}, \dots, a^{(m)}\}$ with probabilities $\{w_N^{(1)}, \dots, w_N^{(m)}\}$ and set $\alpha = 1/a - 1$ to be used at the current iteration of the Gibbs sampler.

Repeat the same procedure for e and obtain a random draw for η . We also experiment with fixing η as 1 and $\sqrt{\alpha + 1}$, explaining the latter choice in the following section. In these cases, the prior variance of $\beta|\sigma^2$ is determined by α .

In what follows we establish the ties between the Bayesian approach we have taken and some frequentist regularization approaches. The simple analytic structure of the generalized double Pareto prior allows for the following analyses while its hierarchical formulation leads to straight-forward computation.

4 Sparse Maximum a Posteriori Estimation

The generalized double Pareto distribution can be used not only as a prior in a Bayesian analysis but also to induce a sparsity-favoring penalty in regularized least squares,

$$\tilde{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p p(|\beta_j|) \right\}, \quad (4)$$

where \mathbf{X} is initially assumed to have orthonormal columns and $p(\cdot)$ denotes the penalty function implied by the prior on the regression coefficients. Following Fan and Li (2001), let $\hat{\beta} = \mathbf{X}'\mathbf{y}$ and denote the minimization problem in (4) for a component of β as

$$\tilde{\beta}_j = \arg \min_{\beta_j} \left\{ \frac{1}{2} (\hat{\beta}_j - \beta_j)^2 + \sigma^2 p(|\beta_j|) \right\}, \quad (5)$$

with the penalty function implied by (3), $p(|\beta_j|) = (\alpha + 1) \log(\sigma\eta + |\beta_j|)$.

From Fan and Li (2001), a good penalty function should result in an estimator that is (i) nearly unbiased when the true unknown parameter is large, (ii) a thresholding rule, which automatically sets small estimated coefficients to zero to reduce model complexity, and (iii) continuous in data z to avoid instability in model prediction. In the following, we show that the penalty function induced by prior (3) has these three properties.

4.1 Near-unbiasedness

The first order derivative of (5) with respect to β_j is $\text{sgn}(\beta_j)\{|\beta_j| + \sigma^2 p'(|\beta_j|)\} - \hat{\beta}_j = \text{sgn}(\beta_j)\{|\beta_j| + \sigma^2(\alpha + 1)/(\sigma\eta + |\beta_j|)\} - \hat{\beta}_j$, where $p'(|\beta_j|) = \partial p(|\beta_j|)/\partial |\beta_j|$ is the term causing bias in estimation. Although it is appealing to introduce bias in small coefficients to reduce the mean squared error and model complexity, it is also desirable to limit the shrinkage of large coefficients with $p'(|\beta_j|) \rightarrow 0$ as $|\beta_j| \rightarrow \infty$. In addition, it is desirable for $p'(|\beta_j|)$ to approach zero rapidly, implying shrinkage and the associated introduction of bias rapidly decreases as coefficients get further away from zero. In fact, the rate of convergence of $p'(|\beta_j|)$ to zero is of the same order under generalized double Pareto and Normal-Jeffreys' priors, with $\lim_{|\beta_j| \rightarrow \infty} \{(\alpha + 1)/(\sigma\eta + |\beta_j|)\} / \{1/|\beta_j|\} = \alpha + 1$. As α controls the tail heaviness in the generalized double Pareto prior, with lighter tails for larger values of α , convergence of the ratio to $(\alpha + 1)$ is intuitive. In the case of Lasso, the bias, $p'(|\beta_j|)$, remains constant regardless of $|\beta_j|$, which can also be observed in Figure 3(b).

4.2 Sparsity

As noted in Fan and Li (2001), a sufficient condition for the resulting estimator to be a thresholding rule is that the minimum of the function $|\beta_j| + \sigma^2 p'(|\beta_j|)$ is positive.

Proposition 3. *Given the formulation in Proposition 1, prior (3) implies a penalty yielding an estimator that is a thresholding rule if $\eta < 2\sqrt{\alpha + 1}$.*

This result is obtained by finding the minimum of $|\beta_j| + \sigma^2 p'(|\beta_j|)$ and setting it greater than zero. The thresholding is a direct consequence of the fact that when $|\hat{\beta}_j| < \min_{\beta_j} \{|\beta_j| + \sigma^2(\alpha + 1)/(\sigma\eta + |\beta_j|)\}$ – which requires that $\min_{\beta_j} \{|\beta_j| + \sigma^2 p'(|\beta_j|)\} > 0$ – the derivative of (5) is positive for all positive β_j and negative for all negative β_j . In this case, the penalized least squares estimator is zero. When $|\hat{\beta}_j| > \min_{\beta_j} \{|\beta_j| + \sigma^2(\alpha + 1)/(\sigma\eta + |\beta_j|)\}$, two roots may exist, the larger one being the penalized least squares estimator. To elaborate more on this, the root(s) may exist for $\text{sgn}(\beta_j)\{|\beta_j| + \sigma^2 p'(|\beta_j|)\} - \hat{\beta}_j = 0$ only when $|\hat{\beta}_j| > \min_{\beta_j} \{|\beta_j| + \sigma^2 p'(|\beta_j|)\}$. A helpful illustration is given in Figure 3 of Fan and Li (2001).

4.3 Continuity

Continuity in data is an important property of an estimator to avoid instabilities in prediction. As defined by Breiman (1996), “a regularization procedure is unstable if a small change in data can make large changes in the regularized estimator”. Discontinuities in the thresholding rule may result in inclusion or dismissal of a signal with minor changes in the data used (see Figure 3(b)). Hard-thresholding – or namely the “usual” variable selection – is an unstable procedure, while ridge or Lasso estimates are considered to be stable. The penalty yielded by the Normal-Jeffreys’ prior (the log penalty) mimics the behavior of the ℓ_γ penalty as $\gamma \rightarrow 0$ where ℓ_γ is the γ -norm of a vector for $\gamma > 0$. This close relation can also be observed in Figure 3(b), again by looking at the discontinuities of hard-thresholding and the Normal-Jeffreys’ prior. This problem is remedied with the use of prior (3).

A necessary and sufficient condition for continuity is that the minimum of the function $|\beta_j| + \sigma^2 p'(|\beta_j|)$ is obtained at zero (Fan and Li, 2001). For our prior, the minimum of this function is obtained at $|\beta_j| = \sigma(\sqrt{\alpha + 1} - \eta)$. Therefore $\eta = \sqrt{\alpha + 1}$ will yield an estimator with this property.

Proposition 4. *Given the formulation in Proposition 1, a subfamily of prior (1) with $\eta = \sqrt{\alpha + 1}$ implies a penalty function that yields an estimator with the continuity property.*

In this particular case, the penalized likelihood estimator is set to zero if $|\hat{\beta}_j| < \sigma\sqrt{\alpha + 1}$. When $|\hat{\beta}_j| > \sigma\sqrt{\alpha + 1}$,

$$\tilde{\beta}_j = \begin{cases} \frac{\hat{\beta}_j - \sigma\sqrt{\alpha+1} + \{\hat{\beta}_j^2 + 2\hat{\beta}_j\sigma\sqrt{\alpha+1} - 3\sigma^2(\alpha+1)\}^{1/2}}{2} & \hat{\beta}_j > 0, \\ \frac{\hat{\beta}_j + \sigma\sqrt{\alpha+1} - \{\hat{\beta}_j^2 - 2\hat{\beta}_j\sigma\sqrt{\alpha+1} - 3\sigma^2(\alpha+1)\}^{1/2}}{2} & \hat{\beta}_j < 0. \end{cases} \quad (6)$$

As can be observed in Figure 3(a), ensuring continuity by letting $\eta = \sqrt{\alpha + 1}$ in prior (3) creates a trade-off between sparsity and tail-robustness. As the thresholding region becomes wider, the larger values are penalized further, yet not nearly at the level of Lasso. To achieve a similar thresholding rule to the Normal-Jeffreys’ prior, we must pick $\alpha = 3$, which induces a lighter tailed distribution than a Cauchy distribution. We may reduce α to 1 and make the tail behavior similar to that of a Cauchy distribution, however, the thresholding region now is reduced to $\pm\sqrt{2}\sigma$. Choosing similar tail behavior to the Normal-Jeffreys’ prior by letting $\alpha \rightarrow 0$, leading to an improper prior, induces a thresholding region of $\pm\sigma$.

4.4 Consistency in Estimation

We investigate the estimation consistency of the implied thresholding rule under orthogonal designs with a diverging number of parameters. Such designs are common in problems such as multiple mean estimation, wavelet smoothing, and principal component regression. Consider the following model which results in a multiple mean estimation (or orthogonal \mathbf{X} such that $\mathbf{x}'_j \mathbf{x}_j = n$ where \mathbf{x}_j is the j th column of \mathbf{X}) setting.

$$\hat{\beta}_n = \beta_n^* + \frac{\sigma}{\sqrt{n}} \mathbf{Z}_n, \quad (7)$$

where $\hat{\beta}_n, \beta_n^*$ are p_n dimensional vectors, and \mathbf{Z}_n is a p_n dimensional multivariate standard normally distributed random variable. Index n denotes sequences that change with n . Here β_n^* denotes the true unknown mean/coefficient vector and $\hat{\beta}_n$ denotes the maximum likelihood estimator. Here β_n^* is assumed to have only a finite number, r , of nonzero elements. Hence, we let the number of zero elements grow with n .

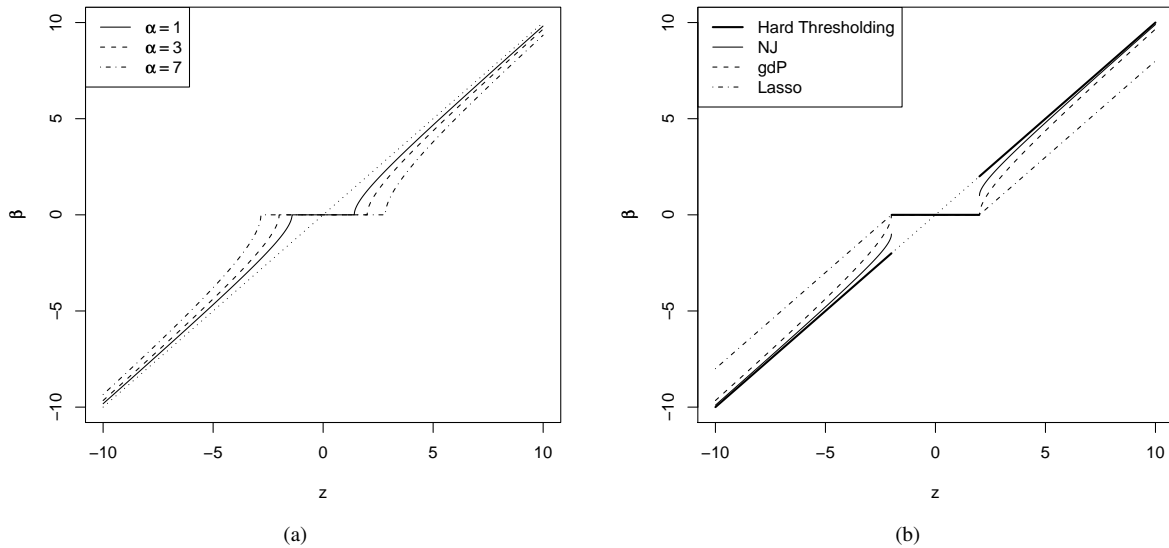


Figure 3: Thresholding functions for (a) generalized double Pareto prior with $\eta = \sqrt{\alpha + 1}$, $\alpha = \{1, 3, 7\}$, (b) Hard thresholding, Normal-Jeffreys' prior, generalized double Pareto prior with $\eta = 2$, $\alpha = 3$ and Lasso with $\sigma = 1$.

The maximum a posteriori estimator under prior (3) (with the continuity property) is given by (6) where σ is replaced by σ/\sqrt{n} , and $\alpha'_n = \alpha_n + 1$.

Theorem 1. Let $\alpha'_n \rightarrow \infty$, $p_n \rightarrow \infty$, $\alpha'_n/n \rightarrow 0$ and $\sigma < \infty$. Given the model in (7) and the estimator $\tilde{\beta}$ in (6),

$$E\|\beta_n^* - \tilde{\beta}_n\|^2 = \mathcal{O}\left\{\frac{p_n \sqrt{\alpha'_n} \exp(-\alpha'_n/2)}{n}\right\} + \mathcal{O}\left(\sqrt{\frac{\alpha'_n}{n}}\right). \quad (8)$$

The proof is deferred to the Appendix.

Corollary 1. If $p_n = n$ and $\alpha'_n = \log n$, $E\|\beta_n^* - \tilde{\beta}_n\|^2 = \mathcal{O}(\sqrt{\log n/n})$.

Proof. This result is obtained by letting $p_n = n$, equating the two terms on the right-hand side of (8), so that they will both have the same rate, and then solving for α'_n . \square

Thus, the threshold for this procedure becomes $\sigma\sqrt{\log n}$. The multiplier of the standard error in the threshold $\sqrt{\log n}$ has a striking similarity to that arising in the so called “universal thresholding”, $\sqrt{2\log n}$.

Remark 1. Here the rate of the estimator is repressed by the bias caused in the non-thresholded elements. Recall that as α'_n increases, the bias increases as well. Thus, there occurs a trade-off between the rate of the estimator and the continuity property we attained.

Therefore, using the penalty implied by the generalized double Pareto prior and choosing an appropriate α'_n we may achieve consistency in estimation with a diverging number of parameters if the true signal has a finite number of nonzero elements. The result given in Corollary 1 is of particular relevance in wavelet smoothing and principle component regression. These properties are relevant to readers interested in sparse estimation and show connections to frequentist approaches.

5 Maximum a Posteriori Estimation via Expectation-Maximization

5.1 Exploiting the Normal Mixture Representation

We assume a normal likelihood to formulate the procedure for non-orthogonal linear regression. Estimation is carried out via the Expectation-Maximization algorithm. We first take the expectation of the log-posterior with respect to the conditional posterior distributions of $(\tau_j^{-1}|\beta_j^{(k)}, \lambda_j, \sigma^{2(k)})$ and $(\lambda_j|\beta_j^{(k)}, \sigma^{2(k)})$ at the k th step, and then maximize with respect to β_j and σ^2 yielding the values for the $(k+1)$ th step. Removing the terms of the log-posterior that do not depend on β and σ^2 , we are left with

$$-\left(\frac{n+p}{2} + 1\right) \log \sigma^2 - \frac{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) - \sum_{j=1}^p \beta_j^2 / \tau_j}{2\sigma^2}.$$

- *E-step:*

$$-\left(\frac{n+p}{2} + 1\right) \log \sigma^2 - \frac{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)}{2\sigma^2} - \frac{1}{2\sigma^2} \sum_{j=1}^p \beta_j^2 \underbrace{\left\{ \frac{(\alpha+1)\sigma^{2(k)}}{|\beta_j^{(k)}| (|\beta_j^{(k)}| + \sigma^{(k)}\eta)} \right\}}_{d_j^{(k)}}$$

- *M-step:* Letting $D^{(k)} = \text{diag}(d_1^{(k)}, \dots, d_p^{(k)})$, we have

$$\beta^{(k+1)} = (\mathbf{X}'\mathbf{X} + D^{(k)})^{-1} \mathbf{X}'\mathbf{y}, \quad \sigma^{2(k+1)} = \frac{(\mathbf{y} - \mathbf{X}\beta^{(k)})'(\mathbf{y} - \mathbf{X}\beta^{(k)}) + \beta^{(k)'} D^{(k)} \beta^{(k)}}{n+p+2}.$$

We refer to this estimator as GDP(MAP).

5.2 Exploiting the Laplace Mixture Representation and the One-step Estimator

An intuitive relationship to the adaptive Lasso of Zou (2006) and the one-step sparse estimator of Zou and Li (2008) can be seen via the Laplace mixture representation of prior (3) implied in Proposition 1. In the proof of Proposition 1, the integration over τ leads to a Laplace mixture representation of the prior. As a computationally fast alternative to estimating the exact mode via the above EM algorithm, we can obtain a ‘‘one-step estimator’’ and exploit the LARS algorithm as in Zou and Li (2008). Since the mixing distribution of the Laplace is a known distribution, the required expectation is obtained with ease resulting in the following step $(k+1)$ maximization:

$$\beta^{(k+1)} = \arg \min_{\beta} \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \frac{1}{\sigma} \sum_{j=1}^p |\beta_j| \left(\frac{\alpha+1}{|\beta_j^{(k)}| + \sigma\eta} \right) \right\}. \quad (9)$$

The component-specific multiplier on $|\beta_j|$ is obtained from the expectation of λ_j with respect to its conditional posterior distribution, $p(\lambda_j|\beta_j)$. Similar results to (9) are observed by Candès et al. (2008), Cevher (2009) and Garrigues (2009). The one-step estimator is then given by

$$\beta^{(1)} = \arg \min_{\beta} \left\{ (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \alpha^\dagger \sum_{j=1}^p |\beta_j| \left(|\beta_j^{(0)}| + \eta^\dagger \right)^{-1} \right\}, \quad (10)$$

letting $\alpha^\dagger = 2\sigma(\alpha+1)$ and $\eta^\dagger = \sigma\eta$. This estimator resembles the adaptive Lasso. The LARS algorithm can be used to obtain $\beta^{(1)}$ very quickly. We refer to this estimator as GDP(OS).

Assuming the same conditions as in Section 2 of Zou (2006), where $\mathbf{X}'\mathbf{X}/n \rightarrow \mathbf{C}$ is a positive definite matrix, we give the following theorem:

100 data sets are generated for each case. In Tables 1 and 2, we report the median model error. Model error is calculated as $(\beta^* - \hat{\beta})'C(\beta^* - \hat{\beta})$ where C is the variance-covariance matrix that generated X . The values in the subscripts give the bootstrap standard error of the median model error values obtained. The bootstrap standard error was calculated by generating 500 bootstrap samples from 100 model error values, finding the median model error for each case, and then calculating the standard error for it. Under each model, the best three performances are boldfaced in the tables.

For larger sample sizes BMA outperforms the competing methods in most cases and GDP(MAP) estimator is the second best. This is not surprising as there indeed exists a true underlying sparse model in most of the scenarios considered. Except for Model 5, normal and Laplace priors are outperformed by other methods as expected. The GDP(PM) shows a similar performance to that of horseshoe. Considering that the GDP(PM) and GDP(MAP) calculations are very straightforward and computationally inexpensive due to the simple normal scale mixture representation used, it offers great utility. The ability to use a simple Gibbs sampler (especially when $\alpha = \eta = 1$) makes the procedure very attractive for the average user.

Table 1: Model error comparisons for Simulation 1.

$n = 50$					
Method	Model 1	Model 2	Model 3	Model 4	Model 5
GDP(PM) ^a	2.659 _{0.127}	2.170 _{0.222}	3.963 _{0.163}	3.847 _{0.167}	5.662 _{0.257}
GDP(PM) ^b	2.775 _{0.153}	2.147 _{0.207}	4.629 _{0.187}	3.817 _{0.178}	6.968 _{0.164}
GDP(PM) ^c	2.592 _{0.109}	2.364 _{0.249}	4.351 _{0.138}	4.001 _{0.227}	6.540 _{0.169}
GDP(PM)	2.785 _{0.145}	2.281 _{0.247}	4.403 _{0.209}	4.455 _{0.244}	6.038 _{0.246}
GDP(MAP) ^a	2.884 _{0.164}	1.498 _{0.138}	5.854 _{0.261}	2.894 _{0.157}	10.404 _{0.243}
GDP(MAP) ^b	4.030 _{0.273}	1.401 _{0.106}	7.099 _{0.270}	3.017 _{0.184}	12.620 _{0.155}
GDP(MAP) ^c	3.526 _{0.168}	1.509 _{0.127}	6.711 _{0.248}	3.060 _{0.165}	11.871 _{0.169}
GDP(MAP)	3.436 _{0.258}	1.992 _{0.292}	5.838 _{0.197}	4.269 _{0.272}	8.750 _{0.364}
Normal	3.236 _{0.186}	5.746 _{0.264}	3.895 _{0.187}	5.515 _{0.218}	3.940 _{0.189}
Laplace	3.208 _{0.114}	4.024 _{0.254}	5.109 _{0.337}	4.875 _{0.279}	7.324 _{0.458}
Horseshoe	2.701 _{0.128}	2.120 _{0.201}	4.789 _{0.219}	3.781 _{0.223}	7.250 _{0.227}
BMA	2.760 _{0.121}	1.438 _{0.103}	4.549 _{0.187}	2.461 _{0.179}	7.031 _{0.233}
$n = 400$					
GDP(PM) ^a	0.230 _{0.016}	0.217 _{0.010}	0.361 _{0.018}	0.311 _{0.017}	0.643 _{0.039}
GDP(PM) ^b	0.219 _{0.015}	0.228 _{0.010}	0.365 _{0.017}	0.338 _{0.014}	0.591 _{0.040}
GDP(PM) ^c	0.237 _{0.014}	0.245 _{0.009}	0.372 _{0.017}	0.344 _{0.017}	0.603 _{0.039}
GDP(PM)	0.231 _{0.013}	0.186 _{0.010}	0.420 _{0.020}	0.371 _{0.016}	0.513 _{0.032}
GDP(MAP) ^a	0.176 _{0.014}	0.146 _{0.013}	0.313 _{0.019}	0.265 _{0.017}	0.614 _{0.040}
GDP(MAP) ^b	0.161 _{0.010}	0.153 _{0.011}	0.313 _{0.022}	0.287 _{0.016}	0.566 _{0.037}
GDP(MAP) ^c	0.180 _{0.014}	0.165 _{0.011}	0.320 _{0.022}	0.290 _{0.015}	0.580 _{0.038}
GDP(MAP)	0.209 _{0.018}	0.139 _{0.013}	0.399 _{0.019}	0.386 _{0.023}	0.498 _{0.030}
Normal	0.415 _{0.016}	0.460 _{0.024}	0.476 _{0.025}	0.472 _{0.020}	0.459 _{0.027}
Laplace	0.328 _{0.013}	0.393 _{0.022}	0.440 _{0.017}	0.442 _{0.020}	0.517 _{0.032}
Horseshoe	0.212 _{0.015}	0.207 _{0.009}	0.380 _{0.018}	0.374 _{0.015}	0.534 _{0.033}
BMA	0.156 _{0.012}	0.126 _{0.016}	0.246 _{0.014}	0.242 _{0.016}	0.450 _{0.021}

^a $\alpha = 1, \eta = 1$; ^b $\eta = 1$; ^c $\eta = \sqrt{\alpha + 1}$

7 Discussion

We proposed a hierarchical prior obtained through a particular scale mixture of normals where the resulting marginal prior has a folded generalized Pareto density thresholded at zero. Although Bayesian model averaging is appealing, it can be argued that allowing parameters to be arbitrarily close to zero instead of exactly equal to zero is more natural. In

Table 2: Model error comparisons for Simulation 2.

$n = 50$					
Method	Model 1	Model 2	Model 3	Model 4	Model 5
GDP(PM) ^a	2.123 _{0.116}	2.149 _{0.105}	3.205 _{0.157}	4.213 _{0.267}	4.440 _{0.134}
GDP(PM) ^b	1.944 _{0.113}	1.997 _{0.103}	3.260 _{0.182}	4.176 _{0.192}	4.651 _{0.127}
GDP(PM) ^c	1.912 _{0.116}	2.154 _{0.095}	3.117 _{0.153}	4.267 _{0.218}	4.329 _{0.124}
GDP(PM)	1.948 _{0.121}	2.154 _{0.125}	3.108 _{0.121}	4.390 _{0.207}	3.961 _{0.149}
GDP(MAP) ^a	2.478 _{0.112}	1.576 _{0.118}	4.711 _{0.262}	3.070 _{0.221}	8.654 _{0.263}
GDP(MAP) ^b	2.594 _{0.187}	1.494 _{0.124}	5.022 _{0.245}	3.198 _{0.226}	9.185 _{0.231}
GDP(MAP) ^c	2.466 _{0.186}	1.546 _{0.104}	4.691 _{0.276}	3.233 _{0.202}	8.895 _{0.227}
GDP(MAP)	2.381 _{0.130}	1.869 _{0.093}	3.774 _{0.191}	4.093 _{0.207}	5.482 _{0.145}
Normal	2.352 _{0.161}	4.165 _{0.293}	2.739 _{0.069}	4.811 _{0.238}	3.040 _{0.177}
Laplace	2.100 _{0.146}	2.768 _{0.171}	2.828 _{0.119}	4.166 _{0.259}	3.468 _{0.154}
Horseshoe	1.983 _{0.114}	2.003 _{0.099}	3.329 _{0.170}	4.345 _{0.194}	4.570 _{0.138}
BMA	2.408 _{0.128}	1.346 _{0.121}	3.954 _{0.124}	3.201 _{0.259}	6.297 _{0.205}
$n = 400$					
GDP(PM) ^a	0.215 _{0.010}	0.219 _{0.011}	0.321 _{0.015}	0.266 _{0.013}	0.659 _{0.037}
GDP(PM) ^b	0.205 _{0.010}	0.230 _{0.014}	0.330 _{0.017}	0.287 _{0.013}	0.595 _{0.034}
GDP(PM) ^c	0.217 _{0.011}	0.236 _{0.014}	0.336 _{0.017}	0.293 _{0.013}	0.595 _{0.033}
GDP(PM)	0.208 _{0.014}	0.192 _{0.010}	0.362 _{0.020}	0.340 _{0.014}	0.494 _{0.030}
GDP(MAP) ^a	0.155 _{0.010}	0.148 _{0.012}	0.260 _{0.018}	0.227 _{0.014}	0.620 _{0.039}
GDP(MAP) ^b	0.150 _{0.009}	0.151 _{0.011}	0.261 _{0.017}	0.248 _{0.015}	0.571 _{0.033}
GDP(MAP) ^c	0.151 _{0.009}	0.168 _{0.010}	0.280 _{0.016}	0.250 _{0.014}	0.582 _{0.034}
GDP(MAP)	0.173 _{0.012}	0.146 _{0.012}	0.331 _{0.017}	0.346 _{0.014}	0.478 _{0.028}
Normal	0.358 _{0.013}	0.441 _{0.024}	0.411 _{0.019}	0.432 _{0.013}	0.433 _{0.026}
Laplace	0.273 _{0.013}	0.365 _{0.019}	0.378 _{0.020}	0.393 _{0.013}	0.475 _{0.028}
Horseshoe	0.198 _{0.010}	0.211 _{0.012}	0.333 _{0.019}	0.332 _{0.012}	0.535 _{0.031}
BMA	0.143 _{0.010}	0.118 _{0.013}	0.231 _{0.017}	0.203 _{0.014}	0.654 _{0.040}

$$^a\alpha = 1, \eta = 1; ^b\eta = 1; ^c\eta = \sqrt{\alpha + 1}$$

addition, the proposed methods have substantial computational advantages in relying on simple block-updated Gibbs sampling, while BMA requires sampling from a model space with 2^p models. As p increases, it becomes impossible to even visit more than a vanishingly small proportion of the models. Given the simple and fast computation and excellent performance in small sample simulation studies, the generalized double Pareto should be useful as a shrinkage prior in a broad variety of Bayesian hierarchical models, while also suggesting close relationship to frequentist penalized likelihood approaches. The proposed prior can be applied outside of normal linear regression to generalized linear models, shrinkage of basis coefficients in nonparametric regression, and in more complex settings such as factor analysis and nonparametric Bayes modeling.

Appendix

Proof of Theorem 1. Let $\sum_{j=1}^{p_n} E(\beta_{nj}^* - \hat{\beta}_{nj})^2 = \sum_{j=1}^r E(\beta_{nj}^* - \hat{\beta}_{nj})^2 + \sum_{j=r+1}^{p_n} E(\hat{\beta}_{nj})^2 = I_1 + I_2$ where $\beta_{nj}^* \neq 0$ for $j = 1, \dots, r$ and $\beta_{nj}^* = 0$ for $j = r + 1, \dots, p_n$.

We first analyze the behavior of I_2 . Given the estimator in (6)

$$\begin{aligned}
I_2 &= \frac{\sigma^2}{4n} \sum_{j=r+1}^{p_n} E \left\{ I(|Z_{nj}| \geq \sqrt{\alpha'_n}) \right. \\
&\quad \times \left. \left[z_{nj} - \operatorname{sgn}(z_{nj})\sqrt{\alpha'_n} + \operatorname{sgn}(z_{nj}) \left\{ z_{nj}^2 + \operatorname{sgn}(z_{nj})2\sqrt{\alpha'_n} - 3\alpha'_n \right\}^{1/2} \right]^2 \right\} \\
&= \frac{\alpha'_n \sigma^2}{2n} \sum_{j=r+1}^{p_n} \int_{\sqrt{\alpha'_n}}^{\infty} \left[\frac{z_{nj}}{\sqrt{\alpha'_n}} - 1 + \left\{ \left(\frac{z_{nj}}{\sqrt{\alpha'_n}} + 3 \right) \left(\frac{z_{nj}}{\sqrt{\alpha'_n}} - 1 \right) \right\}^{1/2} \right]^2 \phi(z_{nj}) dz_{nj} \\
&\leq \frac{2\sigma^2}{n} \sum_{j=r+1}^{p_n} \int_{\sqrt{\alpha'_n}}^{\infty} z_{nj}^2 \phi(z_{nj}) dz_{nj} \\
&= \frac{2\sigma^2(p_n - r - 1)}{n} \left\{ \frac{\sqrt{\alpha'_n} \exp(-\alpha'_n/2)}{\sqrt{2\pi}} + Q(\sqrt{\alpha'_n}) \right\} \\
&\leq \frac{\sqrt{2}\sigma^2(p_n - r - 1)\sqrt{\alpha'_n} \exp(-\alpha'_n/2)}{\sqrt{\pi n}} \left(1 + \frac{1}{\alpha'_n} \right) \\
&= \mathcal{O} \left\{ \frac{p_n \sqrt{\alpha'_n} \exp(-\alpha'_n/2)}{n} \right\}
\end{aligned}$$

where $\phi(\cdot)$ and $Q(\cdot)$ denote the density and the tail probability of a standard normal distribution. In the last inequality we make use of $Q(x) \leq \exp(-x^2/2)/(x\sqrt{2\pi})$.

$$\begin{aligned}
I_1 &= \sum_{j=1}^r E \left\{ I \left(\left| \beta_{nj}^* + \frac{\sigma}{\sqrt{n}} Z_{nj} \right| \geq \frac{\sigma}{\sqrt{n}} \sqrt{\alpha'_n} \right) (\beta_{nj}^* - \hat{\beta}_{nj})^2 \right\} \\
&\quad + \sum_{j=1}^r E \left\{ I \left(\left| \beta_{nj}^* + \frac{\sigma}{\sqrt{n}} Z_{nj} \right| < \frac{\sigma}{\sqrt{n}} \sqrt{\alpha'_n} \right) (\beta_{nj}^*)^2 \right\} \\
&= J_1 + J_2.
\end{aligned}$$

Let us first analyze J_2 :

$$\begin{aligned}
J_2 &= \sum_{j=1}^r \operatorname{pr} \left(\left| \beta_{nj}^* + \frac{\sigma}{\sqrt{n}} Z_{nj} \right| < \frac{\sigma}{\sqrt{n}} \sqrt{\alpha'_n} \right) (\beta_{nj}^*)^2 \\
&= \sum_{j=1}^r P \left(-\sqrt{\alpha'_n} - \beta_{nj}^* \frac{\sqrt{n}}{\sigma} < Z_{nj} < \sqrt{\alpha'_n} - \beta_{nj}^* \frac{\sqrt{n}}{\sigma} \right) (\beta_{nj}^*)^2 \\
&\leq \sum_{j=1}^r \operatorname{pr} \left\{ Z_{nj} > -\sqrt{\alpha'_n} + \operatorname{sgn}(\beta_{nj}^*) \beta_{nj}^* \frac{\sqrt{n}}{\sigma} \right\} (\beta_{nj}^*)^2 \\
&\leq \frac{1}{\sqrt{2\pi}} \sum_{j=1}^r \frac{\exp \left[- \left\{ -\sqrt{\alpha'_n} + \operatorname{sgn}(\beta_{nj}^*) \beta_{nj}^* \frac{\sqrt{n}}{\sigma} \right\}^2 / 2 \right]}{-\sqrt{\alpha'_n} + \operatorname{sgn}(\beta_{nj}^*) \beta_{nj}^* \frac{\sqrt{n}}{\sigma}} (\beta_{nj}^*)^2 \\
&= \mathcal{O} \left(\frac{\exp \left[-\frac{n}{2} \left\{ -\sqrt{\alpha'_n/n} + C \right\}^2 \right]}{\sqrt{n} \left\{ -\sqrt{\alpha'_n/n} + C \right\}} \right),
\end{aligned}$$

where $0 < C < \infty$. Let $c_1 = -\sqrt{\alpha'_n} - \beta_{nj}^* \sqrt{n}/\sigma$ and $c_2 = \sqrt{\alpha'_n} - \beta_{nj}^* \sqrt{n}/\sigma$.

$$J_1 = \sum_{j=1}^r \left\{ \int_{-\infty}^{c_1} (\hat{\beta}_{nj} - \beta_{nj}^*)^2 \phi(z_{nj}) dz_{nj} + \int_{c_2}^{\infty} (\hat{\beta}_{nj} - \beta_{nj}^*)^2 \phi(z_{nj}) dz_{nj} \right\}$$

Let us consider the case $\beta_{nj}^* > 0$ for a summand of J_1 :

$$\begin{aligned} J_{1j} &= (\beta_{nj}^*)^2 \left(\int_{-\infty}^{c_1} \phi(z_{nj}) dz_{nj} + \int_{c_2}^{\infty} \phi(z_{nj}) dz_{nj} \right) - 2\beta_{nj}^* \int_{-\infty}^{c_1} \hat{\beta}_{nj} \phi(z_{nj}) dz_{nj} \\ &\quad - 2\beta_{nj}^* \int_{c_2}^{\infty} \hat{\beta}_{nj} \phi(z_{nj}) dz_{nj} + \int_{-\infty}^{c_1} \hat{\beta}_{nj}^2 \phi(z_{nj}) dz_{nj} + \int_{c_2}^{\infty} \hat{\beta}_{nj}^2 \phi(z_{nj}) dz_{nj} \\ &\leq (\beta_{nj}^*)^2 \{Q(-c_1) + 1 - Q(-c_2)\} - 2\beta_{nj}^* \int_{-\infty}^{c_1} \left(\beta_{nj}^* + z_{nj} \frac{\sigma}{\sqrt{n}} \right) \phi(z_{nj}) dz_{nj} \\ &\quad - 2\beta_{nj}^* \int_{c_2}^{\infty} \left(\beta_{nj}^* + z_{nj} \frac{\sigma}{\sqrt{n}} - \frac{\sigma}{\sqrt{n}} \sqrt{\alpha'_n} \right) \phi(z_{nj}) dz_{nj} + \int_{-\infty}^{c_1} \left(\beta_{nj}^* + z_{nj} \frac{\sigma}{\sqrt{n}} \right)^2 \phi(z_{nj}) dz_{nj} \\ &\quad + \int_{c_2}^{\infty} \left(\beta_{nj}^* + z_{nj} \frac{\sigma}{\sqrt{n}} \right)^2 \phi(z_{nj}) dz_{nj} \\ &= 2\beta_{nj}^* \frac{\sigma}{\sqrt{n}} \sqrt{\alpha'_n} \{1 - Q(-c_2)\} + \frac{\sigma^2}{n} \left\{ 1 + Q(-c_1) - Q(-c_2) + \frac{c_2 e^{-c_2^2/2}}{\sqrt{2\pi}} - \frac{c_1 e^{-c_1^2/2}}{\sqrt{2\pi}} \right\} \end{aligned} \quad (11)$$

The slowest term converging to zero in (11) is $2\beta_{nj}^* \sigma \sqrt{\alpha'_n/n}$ which is $\mathcal{O}(\sqrt{\alpha'_n/n})$. Although the derivation is not given, $\beta_{nj}^* < 0$ case can be shown to have the same slowest term. Thus $J_1 = \mathcal{O}(\sqrt{\alpha'_n/n})$, and since J_1 converges at a slower rate than J_2 ,

$$E \|\beta_n^* - \hat{\beta}_n\|^2 = \mathcal{O} \left\{ \frac{p_n \sqrt{\alpha'_n} \exp(-\alpha'_n/2)}{n} \right\} + \mathcal{O} \left(\sqrt{\frac{\alpha'_n}{n}} \right).$$

This completes the proof. \square

Proof of Theorem 2. The proof emerges with some modifications to the proof of Theorem 2 in Zou (2006). Here $\beta_n^{(0)}$ denotes the least squares estimator. We first prove asymptotic normality. Let $\beta = \beta_n^* + \mathbf{u}/\sqrt{n}$ and

$$V_n(\mathbf{u}) = \left\{ y - \sum_{j=1}^p x_j \left(\beta_{nj}^* + \frac{u_j}{\sqrt{n}} \right) \right\}^2 + \alpha_n^\dagger \sum_{j=1}^p |\beta_{nj}^* + \frac{u_j}{\sqrt{n}}| \left(|\beta_{nj}^{(0)}| + \eta_n^\dagger \right)^{-1}.$$

Let $\hat{\mathbf{u}}_n = \arg \min V_n(\mathbf{u})$, suggesting $\hat{\mathbf{u}}_n = \sqrt{n}(\beta_n^{(1)} - \beta_n^*)$.

$$V_n(\mathbf{u}) - V_n(\mathbf{0}) = \mathbf{u}' \left(\frac{1}{n} \mathbf{X}' \mathbf{X} \right) \mathbf{u} - 2 \frac{\boldsymbol{\varepsilon}' \mathbf{X}}{\sqrt{n}} \mathbf{u} + \frac{\alpha_n^\dagger}{\sqrt{n}} \sum_{j=1}^p \left(|\beta_{nj}^{(0)}| + \eta_n^\dagger \right)^{-1} \sqrt{n} \left(\left| \beta_{nj}^* + \frac{u_j}{\sqrt{n}} \right| - |\beta_{nj}^*| \right)$$

We know that $\mathbf{X}' \mathbf{X}/n \rightarrow \mathbf{C}$ and $\boldsymbol{\varepsilon}' \mathbf{X}/\sqrt{n} \xrightarrow{d} \mathbf{W} \stackrel{d}{=} \mathbf{N}(0, \sigma^2 \mathbf{C})$. Now consider the limiting behavior of the third term. If $\beta_{nj}^* \neq 0$, then by the continuous mapping theorem $\{|\beta_{nj}^{(0)}| + \eta_n^\dagger\}^{-1} \xrightarrow{p} \{|\beta_{nj}^*| + \eta_n^\dagger\}^{-1}$ and $\sqrt{n}(|\beta_{nj}^* + u_j/\sqrt{n}| - |\beta_{nj}^*|) \rightarrow u_j \operatorname{sgn}(\beta_{nj}^*)$. By Slutsky's theorem $(\alpha_n^\dagger/\sqrt{n}) \{|\beta_{nj}^{(0)}| + \eta_n^\dagger\}^{-1} \sqrt{n}(|\beta_{nj}^* + u_j/\sqrt{n}| - |\beta_{nj}^*|) \xrightarrow{p} 0$. If $\beta_{nj}^* = 0$, then $\sqrt{n}(|\beta_{nj}^* + u_j/\sqrt{n}| - |\beta_{nj}^*|) = |u_j|$ and $\alpha_n^\dagger \{|\beta_{nj}^{(0)}| + \eta_n^\dagger\}^{-1} / \sqrt{n} = \alpha_n^\dagger / (\sqrt{n} |\beta_{nj}^{(0)}| + \sqrt{n} \eta_n^\dagger)$ where $\sqrt{n} \beta_{nj}^{(0)} = O_p(1)$. Again by Slutsky's theorem

$$V_n(\mathbf{u}) - V_n(0) \xrightarrow{d} \begin{cases} \mathbf{u}'_A \mathbf{C}_A \mathbf{u}_A - 2\mathbf{u}'_A \mathbf{W}_A & \text{if } u_j = 0 \text{ for all } j \notin A \\ \infty & \text{otherwise.} \end{cases}$$

$V_n(\mathbf{u}) - V_n(\mathbf{0})$ is convex and the unique minimum of the right hand side is $(\mathbf{C}_{\mathcal{A}}^{-1}\mathbf{W}_{\mathcal{A}}, \mathbf{0})'$. By epiconvergence (Geyer, 1994; Knight and Fu, 2000)

$$\mathbf{u}_{\mathcal{A}}^{(n)} \xrightarrow{d} \mathbf{C}_{\mathcal{A}}^{-1}\mathbf{W}_{\mathcal{A}}, \quad \mathbf{u}_{\mathcal{A}^c}^{(n)} \xrightarrow{d} \mathbf{0}. \quad (12)$$

Since $\mathbf{W}_{\mathcal{A}} \stackrel{d}{=} \mathbf{N}(\mathbf{0}, \sigma^2\mathbf{C}_{\mathcal{A}})$, this proves the asymptotic normality.

For all $j \in \mathcal{A}$, $\beta_{nj}^{(1)} \xrightarrow{P} \beta_{nj}^*$; thus $\text{pr}(j \in \mathcal{A}_n) \rightarrow 1$. Now we show that for all $j' \notin \mathcal{A}$, $\text{pr}(j' \in \mathcal{A}_n) \rightarrow 0$. Consider the event $j' \in \mathcal{A}_n$. By the KKT optimality conditions, $2\mathbf{x}'_{j'}(\mathbf{y} - \mathbf{X}\beta_n^{(1)}) = \alpha_n^\dagger(|\beta_{nj'}^{(0)}| + \eta_n^\dagger)^{-1}$. We know that $\alpha_n^\dagger(|\beta_{nj'}^{(0)}| + \eta_n^\dagger)^{-1}/\sqrt{n} \xrightarrow{P} \infty$ while

$$\frac{2\mathbf{x}'_{j'}(\mathbf{y} - \mathbf{X}\beta_n^{(1)})}{\sqrt{n}} = 2 \left\{ \frac{\mathbf{x}'_{j'}\mathbf{X}\sqrt{n}(\beta_n^* - \beta_n^{(1)})}{n} + \frac{\mathbf{x}'_{j'}\boldsymbol{\varepsilon}}{\sqrt{n}} \right\}.$$

By (12) and Slutsky's theorem, we know that both terms in the brackets converge in distribution to some normal suggesting

$$\text{pr}(j' \in \mathcal{A}_n) \leq \text{pr} \left\{ 2\mathbf{x}'_{j'}(\mathbf{y} - \mathbf{X}\beta_n^{(1)}) = \alpha_n^\dagger \left(|\beta_{nj'}^{(0)}| + \eta_n^\dagger \right)^{-1} \right\} \rightarrow 0,$$

which proves the consistency part. \square

Acknowledgements

The work was supported by Award Number R01ES017436 from the National Institute of Environmental Health Sciences. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute of Environmental Health Sciences or the National Institutes of Health. Jaeyong Lee was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (20090086944).

References

- Berger, J. "A robust generalized Bayes estimator and confidence region for a multivariate normal mean." *The Annals of Statistics*, 8(4):pp. 716–761 (1980).
- Berger, J. O. *Statistical Decision Theory and Bayesian Analysis*. New York: Springer (1985).
- Breiman, L. "Heuristics of Instability and Stabilization in Model Selection." *The Annals of Statistics*, 24:2350–2383 (1996).
- Candes, E. J., Wakin, M. B., and Boyd, S. P. "Enhancing sparsity by reweighted l1 minimization." *Journal of Fourier Analysis and Applications*, 14:877–905 (2008).
- Carvalho, C., Polson, N., and Scott, J. "Handling Sparsity via the Horseshoe." *JMLR: W&CP*, 5 (2009).
- Carvalho, C. M., Polson, N. G., and Scott, J. G. "The horseshoe estimator for sparse signals." *Biometrika*, 97(2):465–480 (2010).
- Cevher, V. "Learning with compressible priors." In Bengio, Y., Schuurmans, D., Lafferty, J., Williams, C. K. I., and Culotta, A. (eds.), *Advances in Neural Information Processing Systems*, volume 22, 261–269 (2009).
- Clyde, M., Ghosh, J., and Littman, M. L. "Bayesian adaptive sampling for variable selection and model averaging." *Journal of Computational and Graphical Statistics* (2010).

- Clyde, M. and Littman, M. *Bayesian model averaging using Bayesian adaptive sampling – BAS package manual* (2005).
 URL www.stat.duke.edu/~clyde/BAS/BAS-manual.pdf
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. “Least Angle Regression.” *The Annals of Statistics*, 32(2):407–499 (2004).
- Fan, J. and Li, R. “Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties.” *Journal of the American Statistical Association*, 96(456):1348–1360 (2001).
- Figueiredo, M. A. T. “Adaptive sparseness for supervised learning.” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25:1150–1159 (2003).
- Fu, W. “Penalized regressions: The bridge versus the lasso.” *Journal of Computational and Graphical Statistics*, 7:397–416 (1998).
- Garrigues, P. J. “Sparse coding models of natural images: Algorithms for efficient inference and learning of higher-order structure.” *PhD Thesis, University of California, Berkeley* (2009).
- Geyer, C. J. “On the Asymptotics of Constrained M-Estimation.” *The Annals of Statistics*, 22(4):1993–2010 (1994).
- Gramacy, R. B. *Estimation for multivariate normal and Student-t data with monotone missingness – Monomvn package manual* (2010).
 URL <http://www.statslab.cam.ac.uk/~bobby/monomvn.html>
- Griffin, J. E. and Brown, P. J. “Bayesian adaptive lassos with non-convex penalization.” *Technical Report* (2007).
- . “Inference with normal-gamma prior distributions in regression problems.” *Bayesian Analysis*, 5(1):171–188 (2010).
- Hans, C. “Bayesian lasso regression.” *Biometrika*, 96:835–845 (2009).
- Knight, K. and Fu, W. “Asymptotics for Lasso-Type Estimators.” *The Annals of Statistics*, 28(5):1356–1378 (2000).
- Liang, F., Paulo, R., Molina, G., Clyde, M., and Berger, J. “Mixtures of g priors for Bayesian variable selection.” *Journal of the American Statistical Association*, 103(481):410–423 (2008).
- Park, T. and Casella, G. “The Bayesian Lasso.” *Journal of the American Statistical Association*, 103:681–686(6) (2008).
- Pickands, J. “Statistical Inference Using Extreme Order Statistics.” *The Annals of Statistics*, 3(1):119–131 (1975).
- Ritter, C. and Tanner, M. A. “Facilitating the Gibbs sampler: The Gibbs stopper and the griddy-Gibbs sampler.” *Journal of the American Statistical Association*, 97:861–868 (1992).
- Strawderman, W. E. “Proper Bayes minimax estimators of the multivariate normal mean.” *The Annals of Mathematical Statistics*, 42(1):pp. 385–388 (1971).
- Tibshirani, R. “Regression Shrinkage and Selection via the Lasso.” *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288 (1996).
- Tipping, M. E. “Sparse Bayesian Learning and the Relevance Vector Machine.” *Journal of Machine Learning Research*, 1 (2001).
- West, M. “On Scale Mixtures of Normal Distributions.” *Biometrika*, 74(3):646–648 (1987).
- Yuan, M. and Lin, Y. “Efficient Empirical Bayes Variable Selection and Estimation in Linear Models.” *Journal of the American Statistical Association*, 100(472):1215–1225 (2005).

Zhao, P. and Yu, B. “On Model Selection Consistency of Lasso.” *J. Mach. Learn. Res.*, 7 (2006).

Zou, H. “The Adaptive Lasso and Its Oracle Properties.” *Journal of the American Statistical Association*, 101:1418–1429 (2006).

Zou, H. and Li, R. “One-step sparse estimates in nonconcave penalized likelihood models.” *The Annals of Statistics*, 36(4):1509–1533 (2008).