

Sparse bridge estimation with a diverging number of parameters

SUNGHOO KWON, YONGDAI KIM AND HOSIK CHOI^{*,†}

The Bridge estimator with ℓ_ν^ν -penalty for some $\nu > 0$ is one of the popular choices in penalized linear regression models. It is known that, when $\nu \leq 1$, the Bridge estimator produces sparse models which allow us to control the model complexity. However, when $\nu = 1$, the Bridge estimator fails to identify the correct model since it requires certain strong sufficient conditions that are hard to hold in general, and when $\nu > 1$, it achieves no sparsity in parameter estimation. In this paper, we propose the sparse Bridge estimator that is developed to find the correct sparse version of the Bridge estimator when $\nu \geq 1$. Theoretically, the sparse Bridge estimator is asymptotically equivalent to the oracle Bridge estimator when the number of predictive variables diverges to infinity but less than the sample size. Here, the oracle Bridge estimator is an ideal Bridge estimator obtained by deleting all irrelevant predictive variables in advance. Hence, the sparse Bridge estimator naturally inherits the properties of the Bridge estimator without losing correct model identification asymptotically. Numerical studies show that the sparse Bridge estimator can outperform other penalized estimators with a finite sample.

AMS 2000 SUBJECT CLASSIFICATIONS: Primary 62J05; secondary 62J07.

KEYWORDS AND PHRASES: Bridge, Diverging number of parameters, Lasso, Regression, Ridge, Variable selection.

1. INTRODUCTION

Consider the linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon},$$

where $\mathbf{y} = (y_1, \dots, y_n)'$ is the response vector, $\mathbf{X} = (X_1, \dots, X_p)$ is the $n \times p$ design matrix, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$

^{*}We are grateful to the anonymous referees and the editor for their helpful comments. Kim's research was supported by the National Research Foundation of Korea grant number 20100012671 funded by the Korea government. Choi's research was supported by Basic Science Research Program through National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (2010-0003377).

[†]Corresponding author.

is the random error vector, and $\boldsymbol{\beta}^*$ is the true parameter vector. For given $\nu > 0$, the Bridge estimator [9] is defined as

$$(1) \quad \hat{\boldsymbol{\beta}}^{B,\gamma} = \arg \min_{\boldsymbol{\beta}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2/2n + \gamma \|\boldsymbol{\beta}\|_\nu^\nu \right\}$$

for some $\gamma > 0$, where $\|\boldsymbol{\beta}\|_\nu = (\sum_{j=1}^p |\beta_j|^\nu)^{1/\nu}$ is the usual ℓ_ν -norm operator. As a special case, when $\nu = 1$, the Bridge estimator is known as the least absolute shrinkage and selection operator (Lasso) estimator [21],

$$(2) \quad \hat{\boldsymbol{\beta}}^{L,\gamma} = \arg \min_{\boldsymbol{\beta}} \left\{ \boldsymbol{\beta}'(\mathbf{X}'\mathbf{X})\boldsymbol{\beta}/2n - \mathbf{y}'\mathbf{X}\boldsymbol{\beta}/n + \gamma \|\boldsymbol{\beta}\|_1 \right\},$$

that is one of the popular choices for the variable selection problems. Another example, when $\nu = 2$, is the Ridge estimator [10],

$$(3) \quad \hat{\boldsymbol{\beta}}^{R,\gamma} = \arg \min_{\boldsymbol{\beta}} \left\{ \boldsymbol{\beta}'(\mathbf{X}'\mathbf{X}/2n + \gamma I_{p \times p})\boldsymbol{\beta} - \mathbf{y}'\mathbf{X}\boldsymbol{\beta}/n \right\},$$

that properly handles the collinearity problems [10], where the least square estimator is expected to perform worse.

For years, many penalized estimators have been developed for variable selection problems. Especially, some estimators are asymptotically pursuing the equivalent performance to the oracle estimator:

$$(4) \quad \hat{\boldsymbol{\beta}}^{oR} = \arg \min_{\beta_j=0, j \in \{k: \beta_k^*=0\}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2/2n \right\}.$$

This ideal behavior is described as the oracle property by [7]. Some typical examples are the smoothly clipped absolute deviation (Scad) estimator [7], adaptive Lasso estimator [26], Bridge estimator with $\nu < 1$ [11], and minimax concave penalty (Mcp) estimator [24]. These estimators are practically better in selecting variables compared to the Bridge estimator with $\nu \geq 1$ since they are selection consistent asymptotically.

However, there still exist several examples that present the practical excellence of the Bridge estimators regardless of the oracle property described above. The Lasso estimator in (2) achieves higher prediction accuracy when the true model includes relatively smaller coefficients compared to the sample size [26]. And the Ridge estimator in (3) properly handles the collinearity problems by stabilizing the variance of the least square estimator [4, 10]. Another nice alternative is the elastic net (Enet) estimator [25] defined as

$$(5) \quad \hat{\boldsymbol{\beta}}^{eN, \gamma, \lambda} = \arg \min_{\boldsymbol{\beta}} \left\{ \boldsymbol{\beta}' \left(\frac{\mathbf{X}'\mathbf{X} + \gamma \mathbf{I}_{p \times p}}{1 + \gamma} \right) \boldsymbol{\beta} - 2\mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \lambda \|\boldsymbol{\beta}\|_1 \right\}.$$

Although the Enet estimator is not any Bridge estimator, it can be thought of as an improved version of two Bridge estimators, the Lasso and Ridge. Note that, however, these estimators still suffer from the lack of selection consistency that is a major uneasiness for them.

In this paper, we propose the sparse Bridge estimator that is developed to find the correct sparse version of the Bridge estimator defined in (1) when $\nu \geq 1$. The sparse Bridge estimator resolves the major deficiency of the Bridge estimator by imposing sparsity on it without losing its own practical advantage. Theoretically the sparse Bridge estimator is asymptotically equivalent to the oracle Bridge estimator:

$$(6) \quad \hat{\boldsymbol{\beta}}^{oB, \gamma} = \arg \min_{\boldsymbol{\beta}_j=0, j \in \{k: \beta_k^* = 0\}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 / 2n + \gamma \|\boldsymbol{\beta}\|_\nu^\nu \right\}.$$

As a consequence, the sparse Bridge estimator naturally inherits the properties of the Bridge estimator without losing correct model identification asymptotically.

The rest of the paper is organized as follows. Section 2 introduces the proposed method, and Section 3 proves theoretical properties. Section 4 gives results of various numerical studies. The conclusions and technical details are given in Section 5 and Appendix, respectively.

2. SPARSE BRIDGE ESTIMATION

2.1 Definition and solution

Let $\mathcal{L}(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 / 2n$. For a given $\nu \geq 1$, the sparse Bridge estimator is defined as

$$(7) \quad \hat{\boldsymbol{\beta}}^{sB, \gamma, \lambda} = \arg \min_{\boldsymbol{\beta}} \mathcal{Q}^{sB, \gamma, \lambda}(\boldsymbol{\beta}),$$

where

$$(8) \quad \mathcal{Q}^{sB, \gamma, \lambda}(\boldsymbol{\beta}) = \mathcal{L}(\boldsymbol{\beta}) + \mathcal{J}^{\gamma, \lambda}(\boldsymbol{\beta})$$

and $\mathcal{J}^{\gamma, \lambda}(\boldsymbol{\beta}) = \sum_{j=1}^p J^{\gamma, \lambda}(|\beta_j|)$. Here, $J^{\lambda, \gamma}(\cdot)$ is the sparse Bridge penalty that satisfies $J^{\gamma, \lambda}(0) = 0$, and

$$(9) \quad \nabla J^{\gamma, \lambda}(t) = (-c^{\gamma, \lambda} t^\nu + \lambda) I[t < a(\lambda - \gamma)] + \gamma \nu t^{\nu-1} I[t \geq a(\lambda - \gamma)]$$

for some $a > 0$, where $\nabla J^{\gamma, \lambda}(t) = \partial J^{\gamma, \lambda}(t) / \partial t$ and

$$c^{\gamma, \lambda} = [\lambda - \gamma \nu \{a(\lambda - \gamma)\}^{\nu-1}] / \{a(\lambda - \gamma)\}^\nu$$

for $\lambda \geq \max\{\gamma, \gamma \nu \{a(\lambda - \gamma)\}^{\nu-1}\}$ and $a \geq 1$.

Here are some characteristics on the sparse Bridge penalty which are main motivations of the paper:

- By definition, if $\lambda = \gamma$ then the sparse Bridge penalty is the same as the Bridge penalty.

- And if $\lambda > \gamma$, the sparse Bridge penalty is a smoothly clipped version of the Bridge penalty near the origin.
- The penalty is concave on $(0, a(\lambda - \gamma))$ satisfying $\lim_{t \rightarrow 0+} \nabla J^{\gamma, \lambda}(t) = \lambda$ and the same as the Bridge penalty on $[a(\lambda - \gamma), \infty)$ satisfying $J^{\gamma, \lambda}(t) = \gamma \nu t^{\nu-1}$.

Hence the sparse Bridge penalty can select variables that have the same type of shrinkage as the Bridge penalty. And as a referee pointed out, two tuning parameters λ and γ play different roles:

- λ controls the concavity of the sparse Bridge penalty near the origin thresholding small coefficients.
- γ regularizes the amount of shrinkages over the estimated nonzero coefficients by the same way as in the original Bridge.

It is more transparent to see the solutions under orthonormal design, where $\mathbf{X}'\mathbf{X}/n = \mathbf{I}$. The solution $\hat{\beta}_j$ satisfies $\nabla J^{\gamma, \lambda}(|\hat{\beta}_j|) + \hat{\beta}_j = \hat{z}_j$, $j = 1, \dots, p$, where $\hat{z}_j = X_j' \mathbf{y} / n$ is the ordinary least square estimator. Since $\lim_{|\hat{\beta}_j| \rightarrow 0} \nabla J^{\gamma, \lambda}(|\hat{\beta}_j|) = \lambda$ and $\nabla J^{\gamma, \lambda}(|\hat{\beta}_j|) = \gamma \nu |\hat{\beta}_j|^{\nu-1}$ for $|\hat{\beta}_j| > a(\lambda - \gamma)$, it is easy to see that small coefficients are excluded by the level of λ , and large solutions are exactly the same as those of the original Bridge estimator with γ . For example, if $\nu = 1$,

$$\hat{\beta}_j^{sB, \gamma, \lambda} = \begin{cases} 0, & |z_j| < \lambda, \\ \frac{a}{a-1} \text{sign}(z_j) (|z_j| - \lambda), & \lambda \leq |z_j| < a(\lambda - \gamma) + \gamma, \\ \text{sign}(z_j) (|z_j| - \gamma), & |z_j| \geq a(\lambda - \gamma) + \gamma, \end{cases}$$

for $\lambda \geq \gamma$, and if $\nu = 2$,

$$\hat{\beta}_j^{sB, \gamma, \lambda} = \begin{cases} 0, & |z_j| < \lambda, \\ u(z_j), & \lambda \leq |z_j| < a(1 + 2\gamma)(\lambda - \gamma), \\ z_j / (1 + 2\gamma), & |z_j| \geq a(1 + 2\gamma)(\lambda - \gamma), \end{cases}$$

for $\lambda \geq \max\{\gamma, \gamma_*\}$ if $\gamma \leq 1/2a$, and $\max\{\gamma, \gamma_*\} < \lambda < \gamma^*$ if $2a\gamma > 1$, where $\gamma_* = a\gamma(1 + 2\gamma)/(a - 1 + 2a\gamma)$, $\gamma^* = 2a\gamma^2/(2a\gamma - 1)$ and $u(z_j)$ is solution of $-c^{\gamma, \lambda} \beta_j^2 + \beta_j + z_j = 0$. See Figure 1 that draws the sparse Bridge penalties and solutions when $\nu = 1$ and 2. See Section 4, for numerical results that show different roles of two tuning parameters.

2.2 Computational algorithm

The sparse Bridge penalty $J^{\gamma, \lambda}(|t|)$ can be decomposed as follows:

$$(10) \quad J^{\gamma, \lambda}(|t|) = \gamma |t|^\nu + \lambda |t| + \tilde{J}^{\gamma, \lambda}(|t|),$$

where $\tilde{J}^{\gamma, \lambda}(|t|)$ is a continuously differentiable concave function. Hence, by letting $\tilde{\mathcal{J}}^{\gamma, \lambda}(\boldsymbol{\beta}) = \mathcal{J}^{\gamma, \lambda}(\boldsymbol{\beta}) - \gamma \|\boldsymbol{\beta}\|_\nu^\nu - \lambda \|\boldsymbol{\beta}\|_1$, we can rewrite (8) by

$$\mathcal{Q}^{sB, \gamma, \lambda}(\boldsymbol{\beta}) = \mathcal{L}(\boldsymbol{\beta}) + \gamma \|\boldsymbol{\beta}\|_\nu^\nu + \lambda \|\boldsymbol{\beta}\|_1 + \tilde{\mathcal{J}}^{\gamma, \lambda}(\boldsymbol{\beta})$$

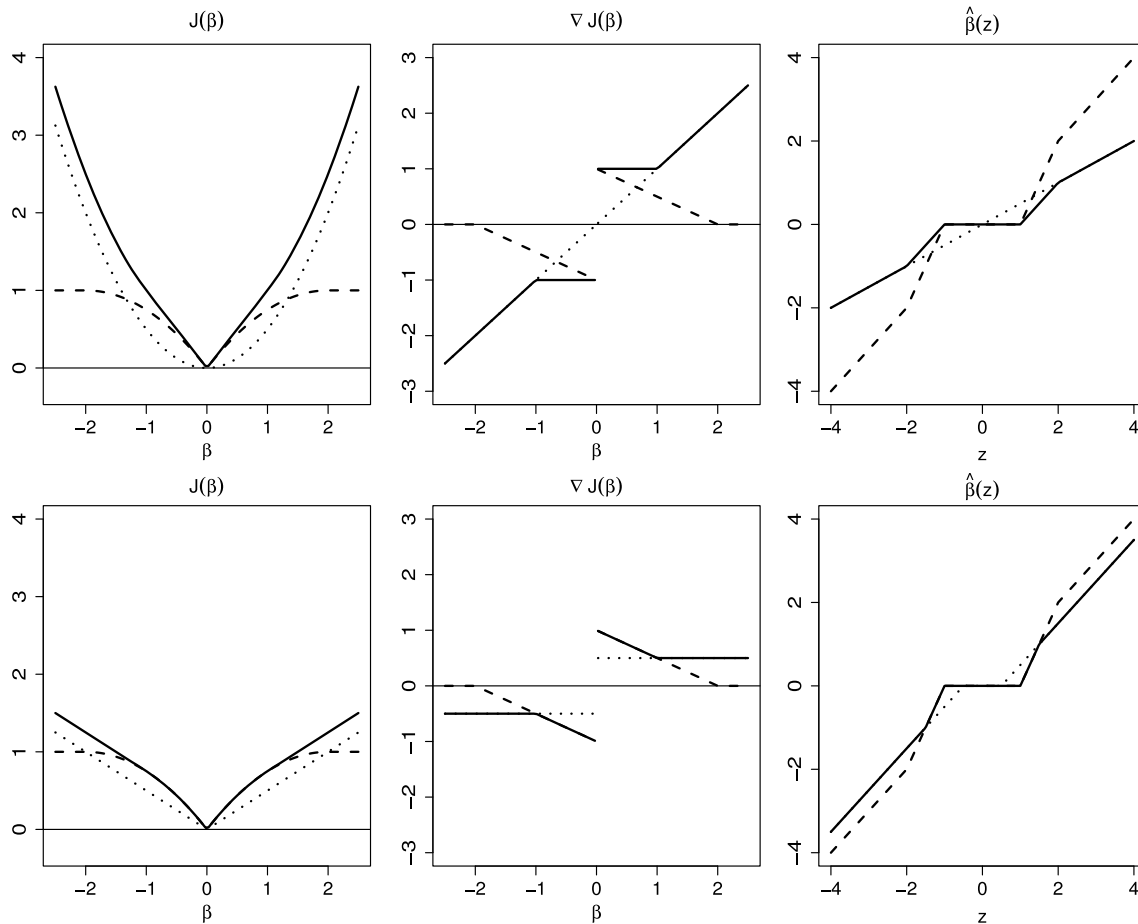


Figure 1. The upper left two panels plot three penalties and derivatives: the Bridge (dotted) with $(\nu, \gamma) = (2, 1/2)$, Scad (dashed) with $(\lambda, a) = (1, 2)$ and sparse Bridge (line) with $(\nu, \gamma, \lambda, a) = (2, 1/2, 1, 2)$. The upper right panel shows the corresponding solutions as functions of the ordinary least square estimator z under the orthonormal design. Similarly, the below panels draw the plots when $\nu = 1$.

which is a sum of convex and concave functions. Hence, we can find a local minimizer of (8) by use of the convex-concave procedure (CCCP) [23]. Note that the CCCP algorithm is one of powerful optimization algorithms for nonconvex problems [3, 13, 20]. One of the main properties of the CCCP algorithm is that it always converges to a local minimizer [1, 2].

To be specific, let $\nabla \tilde{\mathcal{J}}^{\gamma, \lambda}(\beta) = \partial \tilde{\mathcal{J}}^{\gamma, \lambda}(\beta) / \partial \beta$. For a given current solution $\hat{\beta}^c$, we update it by minimizing the upper tight convex function,

$$U^{sB, \gamma, \lambda}(\beta) = \mathcal{L}(\beta) + \gamma \|\beta\|_{\nu}^{\nu} + \nabla \tilde{\mathcal{J}}^{\gamma, \lambda}(\hat{\beta}^c)' \beta + \lambda \|\beta\|_1,$$

until it converges. Many efficient optimization algorithms are available to minimize this convex function, since it is simply a ℓ_1 -penalized convex problem. For special cases, if $\nu = 1$, then we can use the least angle regression algorithm developed by [6]. For other cases, the predictor-corrector algorithm introduced by [17] can be applied.

3. THEORETICAL PROPERTIES

In this subsection, we prove asymptotic equivalence between the sparse Bridge estimator and oracle Bridge estimator so that it is asymptotically as efficient as the oracle Bridge estimator in (6). This relationship is the same as that of the Scad estimator and oracle estimator in (4) [8, 13].

Without loss of generality, we assume that

$$\beta^* = (\beta_1^{*'}, \beta_2^{*'})',$$

where β_1^* is a $q \times 1$ vector whose elements are all nonzero, and β_2^* is a $(p - q) \times 1$ vector whose elements are all zero. Further we assume that $\varepsilon_i, 1 \leq i \leq n$ are independent and identically distributed random variables with $E(\varepsilon_i) = 0$ and $Var(\varepsilon_i) = \sigma_0^2$ for some $\sigma_0 < \infty$. Let $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$, where $\mathbf{X}_1 = (X_1, \dots, X_q)$, $\mathbf{X}_2 = (X_{q+1}, \dots, X_p)$ with $X_j = (x_{1j}, \dots, x_{nj})', 1 \leq j \leq p$.

Denote ρ_n and τ_n as the smallest and largest eigenvalues of the matrix $\mathbf{X}'\mathbf{X}/n$, respectively. To allow the number of

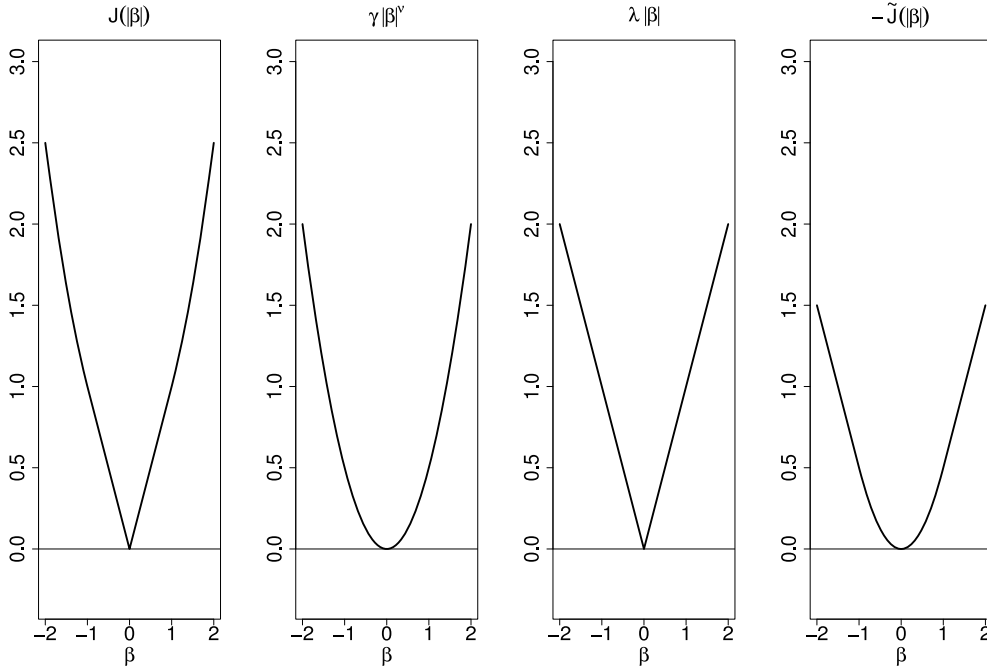


Figure 2. Decomposition of $J^{\gamma, \lambda}(|\beta|) = \gamma|\beta|^\nu + \lambda|\beta| - \tilde{J}^{\gamma, \lambda}(|\beta|)$.

parameters $p < n$ to diverge to infinity, we require following regularity conditions with positive constants $\rho_0, \tau_0, b_0, \delta_0, d_0$ and $c_0 < 1$:

- (A1) $\rho_0 < \rho_n \leq \tau_n \leq \tau_0$ for all $n \geq 1$.
- (A2) $\max_{1 \leq j \leq p} \|X_j\|_2^2/n \leq b_0$ for all $n \geq 1$.
- (A3) $\log p/\log n \rightarrow c_0$ as $n \rightarrow \infty$.
- (A4) $\max_{1 \leq j \leq q} |\beta_j^*| < d_0$ for all $n \geq 1$.
- (A5) $\min_{1 \leq j \leq q} |\beta_j^*|/(q/n)^{1/2} \rightarrow 0$ as $n \rightarrow \infty$.
- (A6) $\max_{1 \leq i \leq n} \sum_{j=1}^p x_{ij}^2/n \rightarrow 0$ as $n \rightarrow \infty$.
- (A7) $E|\varepsilon_1|^{2+\delta_0} < \infty$.

Condition (A1) is standard for the asymptotic study on linear regression model, and (A2) is directly satisfied if we standardize the covariates. Condition (A3) specifies the order of p . As pointed out by [27], it deletes the limitation of $p = o(n^{1/4})$ in [8] for the linear regression model. Condition (A4) is technically considered for ease of exposition that holds implicitly with fixed p and q . Conditions (A5), (A6) and (A7) are used only to construct the limiting distribution of the oracle Bridge estimator. Note that (A5) allows the true coefficients to decrease toward zero.

3.1 Asymptotic properties of the oracle Bridge estimator

First, we give an expected ℓ_2 -risk bound of the oracle Bridge estimator.

Theorem 3.1. Assume that (A1) holds then

$$(11) \quad E\|\hat{\beta}_1^{oB, \gamma} - \beta_1^*\|_2 \leq \Delta(\beta_1^*, \gamma)/\rho_0 + \sigma_0(q/n\rho_0)^{1/2},$$

where

$$\Delta(\beta_1^*, \gamma) = \nu\gamma^{1/\nu}q^{\{(2-\nu)/2\nu\}I(\nu < 2)}(\gamma\|\beta_1^*\|_\nu^\nu + \sigma_0^2q/2n)^{(\nu-1)/\nu}.$$

Theorem 3.1 is non-asymptotic and useful for developing asymptotic theories. In the theorem, the second term $\sigma_0(q/n\rho_0)^{1/2}$ is due to the random error that is $O((q/n)^{1/2})$, whereas, the first term yields extra bias that depends on the tuning parameter γ and the size of the true nonzero parameters.

Remark 1. From Theorem 3.1, it is easy to see that the oracle Bridge estimator is consistent if $\Delta(\beta_1^*, \gamma) \rightarrow 0$ that is equivalent to

$$(12) \quad \gamma q^{1/2 + \{(\nu-2)/2\nu\}I(\nu \geq 2)} \rightarrow 0$$

as $n \rightarrow \infty$. Hence, we require smaller γ for larger ν and p . For example, when $\nu = 1$ and $\nu = 2$, $q^{1/2}\gamma \rightarrow 0$ implies the consistency of the oracle Lasso and Ridge estimators, respectively.

Let $\nabla_1(\beta) = \partial(\|\beta\|_\nu^\nu)/\partial\beta$ and $\mathbf{X}'_1\mathbf{X}_1/n = \Sigma_1$. The next theorem gives $(q/n)^{1/2}$ -consistency and asymptotic normality of the oracle Bridge estimator.

Theorem 3.2. Assume that (A1)–(A7) hold then

$$(a) \quad \|\hat{\beta}_1^{oB, \gamma} - \beta_1^*\|_2 = O_p((q/n)^{1/2}),$$

(b) for any $\alpha \in \mathbb{R}^q$ with $\|\alpha\|_2 = 1$,

$$n^{1/2} \alpha' \Sigma_1^{-1/2} \{\hat{\beta}_1^{oB, \gamma} - \beta_1^* + \gamma \Sigma_1^{-1} \nabla_1(\beta_1^*)\} \rightarrow N(0, \sigma_0^2)$$

in distribution, provided

$$(13) \quad \gamma n^{1/2} q^{\{(\nu-2)/2\nu\}I(\nu \geq 2)} \rightarrow 0$$

as $n \rightarrow \infty$.

In Theorem 3.2, γ must vanish toward zero faster than the rate of consistency in (12) in Remark 1 by a factor of $(q/n)^{1/2}$ exactly, and this result is consistent to that of [15] with fixed p and q . From Theorem 3.2, we can see that the bias of the Lasso, when $\nu = 1$, depends on the tuning parameter γ only while the bias depends on the size of the true parameters also when $\nu > 1$. Note that Theorem 3.1 plays a major role for the results of Theorem 3.2 since $(q/n)^{1/2}$ -consistency directly comes from (11).

3.2 Asymptotic equivalence between the oracle Bridge estimator and sparse Bridge estimator

Next, we prove the optimality of the oracle Bridge estimator for the objective function defined in (8). That is, we provide sufficient conditions under which the sparse Bridge estimator is exactly the same as the oracle Bridge estimator asymptotically.

Let $\Omega^{sB, \gamma, \lambda}$ be the set of all local minimizers of (8). The following theorem is our main result that proves the sparse Bridge estimator is exactly the same as the oracle Bridge estimator asymptotically.

Theorem 3.3. *Assume that (A1)–(A4) hold. If $\lambda \rightarrow 0$, $(n/p)^{1/2} \lambda \rightarrow \infty$, and $\min_{1 \leq j \leq q} |\beta_j^*|/\lambda \rightarrow \infty$ then*

- (a) (Local optimality) $\mathbf{P}(\hat{\beta}^{oB, \gamma} \in \Omega^{sB, \gamma, \lambda}) \rightarrow 1$,
- (b) (Global optimality)

$$\mathbf{P}(\hat{\beta}^{oB, \gamma} = \arg \min_{\beta} \mathcal{Q}^{sB, \gamma, \lambda}(\beta)) \rightarrow 1,$$

provided $\gamma \rightarrow 0$ and

$$(14) \quad (\gamma/\lambda) q^{1/2 + \{(\nu-2)/2\nu\}I(\nu \geq 2)} \rightarrow 0$$

as $n \rightarrow \infty$.

In Theorem 3.3, the conditions imposed on λ are the same as those of [8]. As [8] pointed out, these conditions show the least favorable rate of $\min_{1 \leq j \leq q} |\beta_j^*|$ for oracle variable selection. Note that, for the results to hold, γ must decrease faster than both λ and the rate (12) in Remark 1, but slower than (13). This implies that (14) is stronger than (12) but weaker than (13).

4. NUMERICAL STUDIES

4.1 Simulation studies

In this section, we present simulation studies to evaluate the finite sample performance of the sparse Bridge estimator. Under the linear regression model,

$$(15) \quad y = x' \beta^* + \varepsilon, \quad \varepsilon \sim N(0, \sigma_0^2).$$

The covariate vector $x = (x_1, \dots, x_p)'$ is set to be a multivariate Gaussian random vector with mean zero and the covariance of x_j and x_k to be $\rho_0^{|j-k|}$, $j, k = 1, \dots, p$. We select the true nonzero elements of β^* as a sequence that satisfies $\beta_j^* = 1.5 - (j-1)/(q-1)$, $j = 1, \dots, q$ and $\beta_j^* = 0$, $j = q+1, \dots, p-q$. For sample sizes $n = 30, 60$ and 120, we try $p = [n/3]$ and $q = [p/3]$. Here, $[x]$ indicates the smallest integer greater than x .

We investigate two features: (a) prediction error based on independent test data set of size $2n$; (b) number of coefficients that are selected correctly and incorrectly. From 400 independent repetitions, we measure the medians of averaged prediction errors, and the frequency of correctly and incorrectly selected nonzero coefficients. For comparison, we consider six estimators: Lasso, Scad, Ridge, Enet, sparse Bridge with $\nu = 1$ (sLasso), and $\nu = 2$ (sRidge). We set $a = 3.7$ for the Scad, sLasso and sBridge, and choose other tuning parameters using independent validation data of size $n/2$.

Table 1 summarizes the results, and we present some observations:

1. The two sparse Bridge estimators sLasso and sRidge outperform the other methods improving the Lasso and Ridge respectively for almost all cases.
2. When the sample size n is large or the noise level σ_0 is small, the selection consistent estimators (Scad, sLasso and sBridge) perform better than the others (Lasso, Ridge and Enet).
3. When the sample size is small the characteristics of penalties have positive effects on the prediction accuracy, and this shows why the Scad performs worst when $\sigma_0 = 3$ even if it has the oracle property.
4. When $\rho_0 = 0.9$ and $\sigma_0 = 3$, the Ridge performs better than other estimators except the sRidge, which clearly shows that both sparsity and type of shrinkage closely are related to the prediction accuracy with finite sample size. This is similar to the relation between the Enet and the Lasso.
5. The Scad produces most sparse solutions, and the sLasso and sRidge are sparser than the Lasso and Enet respectively. Except the Ridge which cannot control the selectivity, the Lasso has the largest final model, and the order of model size is Lasso > Enet > sRidge > sLasso > Scad. An interesting point is that sRidge is sparser than Enet. This means that the Enet estimator cannot control the sparsity

Table 1. Simulation results comparing the medians of prediction errors, means of correctly and incorrectly selected nonzero coefficients

Prediction error										
σ_0	r	n	p	q	Lasso	sLasso	Ridge	sRidge	Scad	Enet
1	0.5	30	10	3	0.1946	0.1923	0.3399	0.1848	0.2146	0.2001
		60	20	7	0.2220	0.1682	0.3593	0.1636	0.1817	0.2202
		120	40	13	0.2076	0.1373	0.3888	0.1289	0.1311	0.1931
	0.9	30	10	3	0.2457	0.1843	0.2173	0.1526	0.2506	0.2430
		60	20	7	0.2233	0.1883	0.2588	0.1583	0.2003	0.2197
		120	40	13	0.1883	0.1593	0.2341	0.1285	0.1603	0.1774
3	0.5	30	10	3	1.3319	1.3664	1.6174	1.3974	1.7645	1.4234
		60	20	7	1.7928	1.7438	1.9212	1.7542	2.5425	1.7725
		120	40	13	1.8318	1.7239	2.2429	1.5121	2.4645	1.7937
	0.9	30	10	3	1.1938	1.0817	0.8902	1.0463	1.2718	1.3058
		60	20	7	1.5983	1.3026	1.2839	1.1114	1.7236	1.5402
		120	40	13	1.3841	1.2596	1.1983	0.8158	1.6262	1.3127
Number of nonzero coefficients correctly estimated										
σ_0	r	n	p	q	Lasso	sLasso	Ridge	sRidge	Scad	Enet
1	0.5	30	10	3	2.92	2.74	3.00	2.78	2.66	2.92
		60	20	7	6.98	6.88	7.00	6.90	6.81	6.97
		120	40	13	13.00	12.98	13.00	12.98	12.96	13.00
	0.9	30	10	3	2.51	2.24	3.00	2.43	2.04	2.49
		60	20	7	6.70	6.24	7.00	6.48	5.99	6.71
		120	40	13	12.89	12.53	13.00	12.65	12.23	12.88
3	0.5	30	10	3	2.33	2.07	3.00	2.12	1.91	2.32
		60	20	7	6.23	5.51	7.00	5.54	4.82	6.23
		120	40	13	12.46	11.59	13.00	11.64	10.47	12.45
	0.9	30	10	3	1.74	1.53	3.00	1.72	1.46	1.73
		60	20	7	4.91	3.97	7.00	4.71	3.53	4.87
		120	40	13	10.64	9.08	13.00	10.52	7.80	10.61
Number of nonzero coefficients incorrectly estimated										
σ_0	r	n	p	q	Lasso	sLasso	Ridge	sRidge	Scad	Enet
1	0.5	30	10	3	2.80	1.03	7.00	1.30	1.01	2.84
		60	20	7	4.97	1.18	13.00	1.42	1.11	4.88
		120	40	13	9.18	1.30	27.00	1.61	1.26	8.77
	0.9	30	10	3	2.44	0.80	7.00	1.07	0.73	2.26
		60	20	7	3.73	0.96	13.00	1.38	0.98	3.65
		120	40	13	5.78	0.97	27.00	1.58	0.87	5.66
3	0.5	30	10	3	2.75	1.70	7.00	2.09	1.71	2.65
		60	20	7	4.79	2.33	13.00	3.03	2.31	4.73
		120	40	13	8.80	4.02	27.00	4.63	4.00	8.62
	0.9	30	10	3	1.88	1.08	7.00	1.29	1.11	1.81
		60	20	7	3.21	1.12	13.00	1.93	1.42	3.13
		120	40	13	5.64	1.84	27.00	2.66	2.27	5.40

properly, and this is partly explained by the selection consistency of the sRidge.

To sum up, we can see that the sLasso and sRidge can be the correct sparse version of the Lasso and Ridge respectively, and the sLasso and sRidge can be alternatives to other penalized estimators. We note that when the sample

size is sufficiently large, one can expect the Scad performs best as the asymptotic property shows [8, 14]. However when the sample size is small we may carefully choose an appropriate penalty. We recommend that one may use sLasso when the sample size is small and the noise level is large, and sRidge when there exists the collinearity problems as the results of simulation studies show.

Table 2. Performance results from 100 random partitions of the data (standard errors)

p		Lasso	sLasso	Ridge	sRidge	Scad	Enet
20	Prediction error	0.4259 (0.0181)	0.4196 (0.0147)	0.3671 (0.0135)	0.3842 (0.0103)	0.4493 (0.0135)	0.4070 (0.0126)
	No. of selected variables	8.97 (0.19)	4.98 (0.17)	20 (0)	6.84 (0.20)	2.79 (0.06)	8.70 (0.18)
30	Prediction error	0.4319 (0.0186)	0.4254 (0.0158)	0.3891 (0.0151)	0.3979 (0.0122)	0.4556 (0.0145)	0.4150 (0.0136)
	No. of selected variables	9.44 (0.19)	5.10 (0.17)	30 (0)	8.13 (0.23)	2.83 (0.05)	9.07 (0.19)
40	Prediction error	0.4291 (0.0184)	0.4206 (0.0155)	0.3910 (0.0158)	0.3901 (0.0108)	0.4645 (0.0145)	0.4125 (0.0134)
	No. of selected variables	10.58 (0.22)	5.67 (0.19)	40 (0)	9.23 (0.50)	2.98 (0.05)	10.30 (0.21)

4.2 Real data example

4.2.1 Gene TRIM32

We employ the data set used in [19], which consists of gene expression levels of 18,975 genes obtained from 120 rats. The main objective of the analysis is to find genes that are correlated with gene TRIM32 known to cause Bardet-Biedl syndrome. As was done by [12], we first select 3,000 genes with the largest variance in expression level, and then choose the top $p = 20, 30$ and 40 genes that have the largest absolute correlation with gene TRIM32 among the selected 3,000 genes.

We compare prediction accuracy and number of variables selected in the model for the methods used in the previous subsection. The data set is randomly divided into three parts: 70 samples for training, 30 samples for validating and the other samples for testing.

Table 2 shows the results: all values are arithmetic means of 100 replicated experiments. It is easy to see that the Ridge has the smallest prediction errors and the sRidge follows it showing almost the same prediction errors but using less than 10 variables. Given the way of the prescreening method, the good performance of the Ridge is intuitive in part since there are strong correlations among predictive variables, and this can be a reason why the Enet shows smaller prediction errors also.

On the other hand, the Scad has the largest prediction errors including only 3 variables, which seems to be too sparse. Further, the difference between the Scad and the others is relatively large. This shows that there are cases where we need to choose penalties carefully, and the sRidge carries out this goal well in this example.

4.2.2 Diabetes study

As aforementioned above, two tuning parameters play different roles: one for selection and the other for shrinkage. To address this issue numerically, we additionally investigate the trajectory of nonzero coefficients of diabetes study data [6]. This data includes 442 samples measured on 10

baseline variables (age, sex, bmi, map and 6 serum measurements) from diabetes patients. The response of interest is a quantitative measure of disease progression one year after baseline.

Figure 3 shows the number of estimated nonzero coefficients (selection) and their ℓ_1 -norm (shrinkage) for sLasso ($\nu = 1$) and sRidge ($\nu = 2$). In the plot, we fix four values of λ , and for each λ we find solution paths over 10 values of $\gamma = \lambda/2k, k = 1, \dots, 10$. It is easy to see that the model size increases as λ decreases, however, it rarely changes over γ . On the other hand, the estimated ℓ_1 -norm increases as γ decreases. See Figure 4 that draws the paths of estimated coefficients. This shows that the two tuning parameters λ and γ play clearly as they are expected.

5. CONCLUDING REMARK

We proposed the sparse Bridge estimator to improve the usual Bridge estimator, and proved that it is selection consistent under mild conditions. We confirmed the theoretical results via numerical studies which show how the sparse Bridge estimator behaves and produces better prediction accuracy and variable selectivity than the original one.

Besides ν , the sparse Bridge penalty has two more tuning parameters λ and γ and this may cause computational burdens in practice compared with the Lasso and Scad. We refer to two simpler ways of choosing the tuning parameters. One is to use the universal penalty level for λ fixed with $\lambda = \sigma_0\{(2/n)\log p\}^{1/2}$ as was done by [5, 24], and the other one is to use a heuristic choice of γ fixed with $\gamma = \hat{\gamma}$, where $\hat{\gamma}$ is the optimal one for the original Bridge penalty. Although we do not present the results here, such a choice of γ performs quite well. Once a tuning parameter is fixed, one can apply known selection methods. For example, the cross-validation method, training-validating-testing procedure, AIC, BIC, and GCV, where the generalized degrees of freedom is defined similar to [7, 21] and [22].

In this paper, we only consider linear regression models with diverging $p < n$, however we believe that the proposed method can be extended to quite general models such

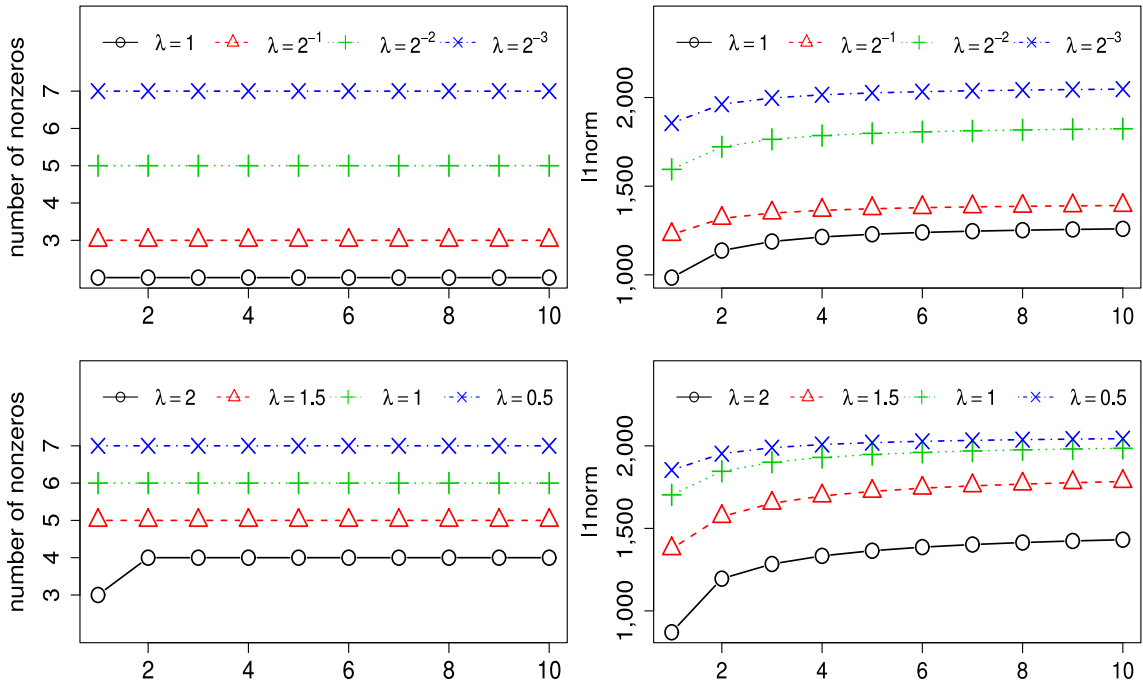


Figure 3. Estimated number of nonzero coefficients and their l_1 -norm over γ for fixed λ : $sLasso(\nu = 1, \text{upper})$ and $sRidge(\nu = 2, \text{below})$.

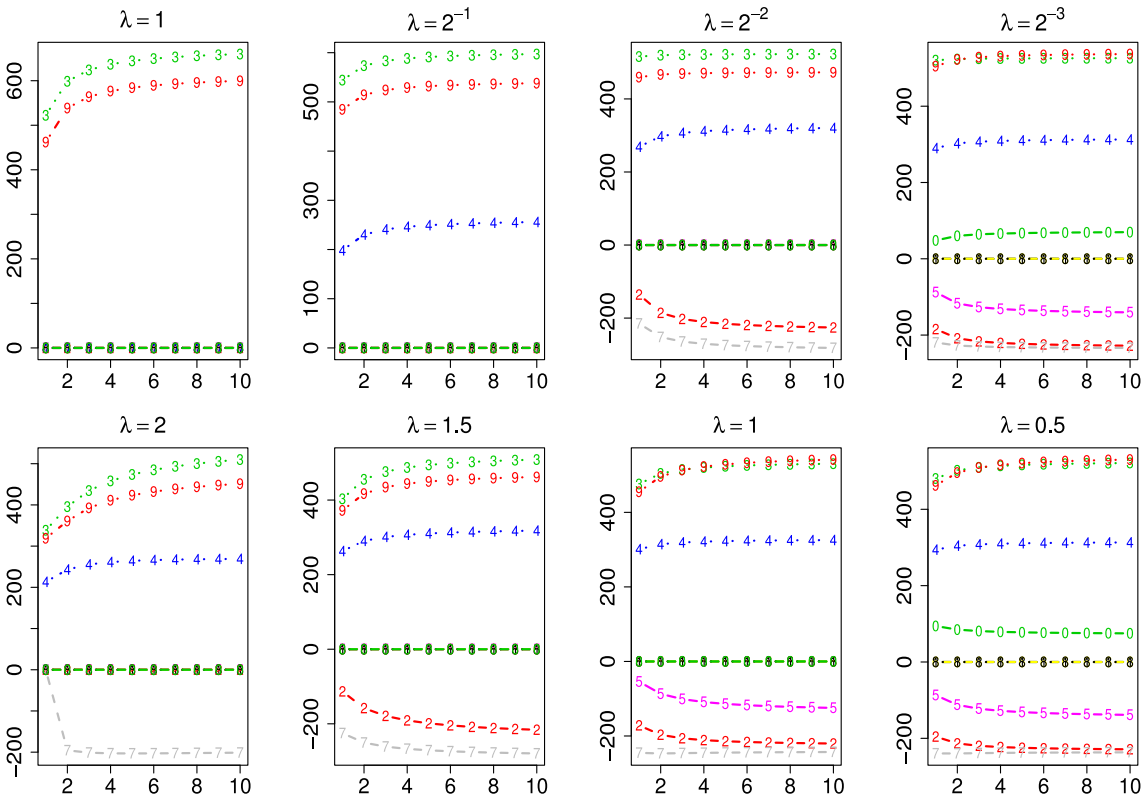


Figure 4. Estimated coefficient path over γ for fixed λ : $sLasso(\nu = 1, \text{upper})$ and $sRidge(\nu = 2, \text{below})$.

as generalized linear models or maximum likelihood estimations with $p > n$ [13, 16, 24]. Another challenging problem is to develop uniqueness conditions as [14] that enhances the applicability of the sparse Bridge penalty. We leave these problems for future work.

APPENDIX

For convenience, we define $\mathbf{S}(\boldsymbol{\beta}) = (\mathbf{S}_1(\boldsymbol{\beta})', \mathbf{S}_2(\boldsymbol{\beta})')'$, where $\mathbf{S}_1(\boldsymbol{\beta}) = (S_1(\boldsymbol{\beta}), \dots, S_q(\boldsymbol{\beta}))'$ and $\mathbf{S}_2(\boldsymbol{\beta}) = (S_{q+1}(\boldsymbol{\beta}), \dots, S_p(\boldsymbol{\beta}))'$ with $S_j(\boldsymbol{\beta}) = X_j'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/n$, $1 \leq j \leq p$.

Proof of Theorem 3.1. From direct calculation,

$$(16) \quad \hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^* = \boldsymbol{\Sigma}_1^{-1} \{ -\mathbf{S}_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma}) + \mathbf{X}'_1 \boldsymbol{\varepsilon} / n \}.$$

From (A1),

$$E \|\boldsymbol{\Sigma}_1^{-1} \mathbf{X}'_1 \boldsymbol{\varepsilon} / n\|_2^2 = \text{trace} \{ E(\mathbf{X}_1(\mathbf{X}'_1 \mathbf{X}_1 / n)^{-2} \mathbf{X}'_1 \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' / n^2) \} \leq q\sigma_0^2 / (n\rho_0),$$

hence, it follows that

$$(17) \quad E \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*\|_2 \leq E \|\mathbf{S}_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2 / \rho_0 + \sigma_0(q/n\rho_0)^{1/2}.$$

If $\nu = 1$, by the Karush-Kuhn-Tucker (KKT) optimality conditions (see, [18]), $\hat{\boldsymbol{\beta}}_1^{oB,\gamma}$ satisfies

$$(18) \quad \begin{cases} S_j(\hat{\boldsymbol{\beta}}_1^{oB,\gamma}) = \gamma \text{sign}(\hat{\beta}_j^{oB,\gamma}), & \hat{\beta}_j^{oB,\gamma} \neq 0, \\ |S_j(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})| \leq \gamma, & \hat{\beta}_j^{oB,\gamma} = 0, \end{cases}$$

for $1 \leq j \leq q$, and if $\nu > 1$, it satisfies

$$(19) \quad S_j(\hat{\boldsymbol{\beta}}_1^{oB,\gamma}) = \gamma \nu |\hat{\beta}_j^{oB,\gamma}|^{\nu-1} \text{sign}(\hat{\beta}_j^{oB,\gamma})$$

for $1 \leq j \leq q$. Hence,

$$(20) \quad \|\mathbf{S}_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2 \leq \gamma \nu \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_{2(\nu-1)}^{\nu-1}.$$

On the other hand, by the definition of $\hat{\boldsymbol{\beta}}_1^{oB,\gamma}$,

$$\begin{aligned} & \gamma \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_\nu^\nu \\ & \leq \gamma \|\boldsymbol{\beta}_1^*\|_\nu^\nu + \|\mathbf{y} - \mathbf{X}_1 \boldsymbol{\beta}_1^*\|_2^2 / 2n - \|\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_2^2 / 2n \\ & \leq \gamma \|\boldsymbol{\beta}_1^*\|_\nu^\nu + \|\mathbf{y} - \mathbf{X}_1 \boldsymbol{\beta}_1^*\|_2^2 / 2n - \|\mathbf{y} - \Pi_1 \mathbf{y}\|_2^2 / 2n, \end{aligned}$$

where $\Pi_1 = \mathbf{X}_1(\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1$. Hence,

$$(21) \quad \gamma \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_\nu^\nu \leq \gamma \|\boldsymbol{\beta}_1^*\|_\nu^\nu + \|\Pi_1 \boldsymbol{\varepsilon}\|_2^2 / 2n.$$

Now, if $\nu < 2$, from (20) and (21), it follows that

$$(22) \quad \begin{aligned} \|\mathbf{S}_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2 & \leq \nu \gamma q^{(2-\nu)/2\nu} \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_\nu^{\nu-1} \\ & \leq \nu \gamma^{1/\nu} q^{(2-\nu)/2\nu} (\gamma \|\boldsymbol{\beta}_1^*\|_\nu^\nu + q\sigma_0^2 / 2n)^{(\nu-1)/\nu} \end{aligned}$$

since $E(\|\Pi_1 \boldsymbol{\varepsilon}\|_2^2 / 2n) = q\sigma_0^2 / 2n$. From (17) and (22),

$$E \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*\|_2 \leq \Delta(\boldsymbol{\beta}_1^*, \gamma) / \rho_0 + \sigma_0(q/n\rho_0)^{1/2}.$$

Similarly, if $\nu \geq 2$, from (20) and (21),

$$\begin{aligned} \|\mathbf{S}_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2 & \leq \nu \gamma \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma}\|_\nu^{\nu-1} \\ & \leq \nu \gamma^{1/\nu} (\gamma \|\boldsymbol{\beta}_1^*\|_\nu^\nu + \|\Pi_1 \boldsymbol{\varepsilon}\|_2^2 / 2n)^{(\nu-1)/\nu}. \end{aligned}$$

Hence, from (17) again,

$$E \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*\|_2 \leq \Delta(\boldsymbol{\beta}_1^*, \gamma) / \rho_0 + \sigma_0(q/n\rho_0)^{1/2}.$$

This completes the proof. \square

Proof of Theorem 3.2. It is easy to see that (A4) and (13) imply

$$\Delta(\boldsymbol{\beta}_1^*, \gamma) = O_p((q/n)^{1/2}).$$

Hence, from Theorem 3.1, it follows that

$$(23) \quad \|\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*\|_2 = O_p((q/n)^{1/2}).$$

Next, we will show (b). From (16) and (19), we have

$$\begin{aligned} & n^{1/2} \boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{1/2} \{ \hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^* + \boldsymbol{\Sigma}_1^{-1} \gamma \nabla_1(\boldsymbol{\beta}_1^*) \} \\ & = \boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1/2} \mathbf{X}'_1 \boldsymbol{\varepsilon} / n^{1/2} \\ & \quad + n^{1/2} \boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1/2} \gamma \{ \nabla_1(\boldsymbol{\beta}_1^*) - \nabla_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma}) \}. \end{aligned}$$

It is a standard [11] that, under (A1)–(A7),

$$\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1/2} \mathbf{X}'_1 \boldsymbol{\varepsilon} / n^{1/2} \rightarrow N(0, \sigma_0^2)$$

in distribution for any $\boldsymbol{\alpha} \in \mathbb{R}^q$ with $\|\boldsymbol{\alpha}\|_2 = 1$. Hence, it suffices to show that

$$\Delta = n^{1/2} \boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1/2} \gamma \{ \nabla_1(\boldsymbol{\beta}_1^*) - \nabla_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma}) \} = o_p(1).$$

If $\nu = 1$, from (A5) and (23), $\text{sign}(\hat{\beta}_j^{oB,\gamma}) = \text{sign}(\beta_j^*)$ for all $1 \leq j \leq p$ and sufficiently large n , which implies

$$|\Delta|^2 \leq \gamma^2 n \|\nabla_1(\boldsymbol{\beta}_1^*) - \nabla_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2^2 / \rho_0 = 0.$$

If $\nu > 1$, by Taylor's expansion, there exists a $\tilde{\boldsymbol{\beta}}_1$ that lies between $\hat{\boldsymbol{\beta}}_1^{oB,\gamma}$ and $\boldsymbol{\beta}_1^*$ such that

$$\begin{aligned} \|\nabla_1(\boldsymbol{\beta}_1^*) - \nabla_1(\hat{\boldsymbol{\beta}}_1^{oB,\gamma})\|_2^2 & = \|\nabla_2(\tilde{\boldsymbol{\beta}}_1)(\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*)\|_2^2 \\ & \leq \nu^2 (\nu - 1)^2 \max_{1 \leq j \leq q} |\tilde{\beta}_j|^{2(\nu-2)} \|(\hat{\boldsymbol{\beta}}_1^{oB,\gamma} - \boldsymbol{\beta}_1^*)\|_2^2. \end{aligned}$$

If $1 < \nu < 2$, (A5) and (23) imply

$$\max_{1 \leq j \leq q} |\tilde{\beta}_j|^{2(\nu-2)} = \min_{1 \leq j \leq q} \{O_p(|\beta_j^*|)\}^{2(\nu-2)} = o_p((q/n)^{\nu-2}).$$

Hence, from (13), we have

$$|\Delta|^2 \leq \gamma^2 n o_p((q/n)^{\nu-2}) O_p(q/n) = o_p(1).$$

Similarly, if $\nu \geq 2$, from (A4), (23), and (13),

$$|\Delta|^2 \leq \gamma^2 n O_p(1) O_p(q/n) = o_p(1).$$

This completes the proof. \square

We need two lemmas to prove Theorem 3.3.

Lemma 1. For given $\hat{\beta} \in \mathbb{R}^p$, if $\hat{\beta}$ satisfies

$$\begin{cases} S_j(\hat{\beta}) = \gamma \nu |\hat{\beta}_j|^{\nu-1} \text{sign}(\hat{\beta}_j), & |\hat{\beta}_j| > a(\lambda - \gamma), \\ |S_j(\hat{\beta})| < \lambda, & \hat{\beta}_j = 0, \end{cases}$$

for $1 \leq j \leq p$, then $\hat{\beta} \in \Omega^{sB, \gamma, \lambda}$.

Proof. It suffices to show that there exists a $\delta > 0$ such that $\mathcal{Q}^{sB, \gamma, \lambda}(\beta) \geq \mathcal{Q}^{sB, \gamma, \lambda}(\hat{\beta})$ for all $\beta \in B(\hat{\beta}, \delta)$. From the convexity of $\mathcal{L}(\beta)$,

$$\mathcal{Q}^{sB, \gamma, \lambda}(\beta) - \mathcal{Q}^{sB, \gamma, \lambda}(\hat{\beta}) \geq \sum_{j=1}^p v_j(\beta_j),$$

where

$$v_j(\beta_j) = -S_j(\hat{\beta})(\beta_j - \hat{\beta}_j)/n + J^{\gamma, \lambda}(|\beta_j|) - J^{\gamma, \lambda}(|\hat{\beta}_j|).$$

First, if $|\hat{\beta}_j| > a(\lambda - \gamma)$, then

$$\begin{aligned} v_j(\beta_j) &= -\gamma \nu |\hat{\beta}_j|^{\nu-1} \text{sign}(\hat{\beta}_j)(\beta_j - \hat{\beta}_j) + \gamma |\beta_j|^\nu - \gamma |\hat{\beta}_j|^\nu \\ &\geq -\gamma \nu |\hat{\beta}_j|^{\nu-1} (|\beta_j| - |\hat{\beta}_j|) + \gamma |\beta_j|^\nu - \gamma |\hat{\beta}_j|^\nu \\ &\geq 0 \end{aligned}$$

for all $\beta_j \in B(\hat{\beta}_j, \delta_j)$, where $\delta_j = |\hat{\beta}_j| - a(\lambda - \gamma)$.

Next, if $\hat{\beta}_j = 0$, then

$$\begin{aligned} v_j(\beta_j) &\geq -|S_j(\hat{\beta})||\beta_j| + J^{\gamma, \lambda}(|\beta_j|) \\ &= |\beta_j|(-|S_j(\hat{\beta})| - c^{\gamma, \lambda} |\beta_j|^\nu / (\nu + 1) + \lambda) \geq 0 \end{aligned}$$

for all $\beta_j \in B(0, \delta_j)$, where

$$\delta_j = \left\{ (\nu + 1)(\lambda - |S_j(\hat{\beta})|) / c^{\gamma, \lambda} \right\}^{1/\nu}.$$

Hence,

$$\mathcal{Q}^{sB, \gamma, \lambda}(\beta) \geq \mathcal{Q}^{sB, \gamma, \lambda}(\hat{\beta})$$

for all $\beta \in B(\hat{\beta}, \min_{1 \leq j \leq p} \delta_j)$. This completes the proof. \square

Lemma 2. Assume the conditions in Theorem 3.3 hold then

$$\mathbf{P}\left(\max_{q < j \leq p} |S_j(\hat{\beta}^{oB, \gamma})| < \lambda\right) \rightarrow 1,$$

and

$$\mathbf{P}\left(\min_{1 \leq j \leq q} |\hat{\beta}_j^{oB, \gamma}| > a(\lambda - \gamma)\right) \rightarrow 1,$$

as $n \rightarrow \infty$.

Proof. Theorem 3.1 and (14) imply

$$(24) \quad \|\hat{\beta}_1^{oB, \gamma} - \beta_1^*\|_2 = o_p(\lambda).$$

From (A1) and (A2), we have

$$\begin{aligned} \max_{q < j \leq p} |S_j(\hat{\beta}^{oB, \gamma})| &\leq \max_{q < j \leq p} |X_j' \mathbf{X}_1 (\hat{\beta}_1^{oB, \gamma} - \beta_1^*)/n| + \max_{q < j \leq p} |X_j' \varepsilon/n| \\ &\leq b_0^{1/2} \tau_0^{1/2} \|\hat{\beta}_1^{oB, \gamma} - \beta_1^*\|_2 + \|\mathbf{X}_2' \varepsilon/n\|_2. \end{aligned}$$

Since $E\|\mathbf{X}_2' \varepsilon/n\|_2^2 \leq (p - q)\tau_0\sigma_0^2/n = o(\lambda^2)$, from (24),

$$\max_{q < j \leq p} |S_j(\hat{\beta}^{oB, \gamma})| = o_p(\lambda).$$

On the other hand, from (24), $\min_{1 \leq j \leq q} |\beta_j^*|/\lambda \rightarrow \infty$ implies

$$\begin{aligned} \min_{1 \leq j \leq q} |\hat{\beta}_j^{oB, \gamma}|/a(\lambda - \gamma) &\geq \min_{1 \leq j \leq q} |\beta_j^*|/a(\lambda - \gamma) - \|\hat{\beta}_1^{oB, \gamma} - \beta_1^*\|_2/a(\lambda - \gamma) \\ &\rightarrow \infty \end{aligned}$$

as $n \rightarrow \infty$. This completes the proof. \square

Proof of Theorem 3.3. It is easy to see that (a) holds by Lemma 1 and 2. Hence, it suffices to prove (b). For this we will show

$$\mathbf{P}(\mathcal{Q}^{sB, \gamma, \lambda}(\hat{\beta}^{oB, \gamma}) \leq \inf_{\beta} \mathcal{Q}^{sB, \gamma, \lambda}(\beta)) \rightarrow 1$$

as $n \rightarrow \infty$. For convenience, we omit the superscripts so that we write $\hat{\beta}^{oB, \gamma} = \hat{\beta}$ and $J^{\gamma, \lambda}(\cdot) = J(\cdot)$. From (A1),

$$\mathcal{L}(\beta) - \mathcal{L}(\hat{\beta}) \geq \sum_{j=1}^p \left\{ -S_j(\hat{\beta})(\beta_j - \hat{\beta}_j)/n + \rho_0(\beta_j - \hat{\beta}_j)^2/2 \right\}.$$

Hence, from Lemma 2 and (19), it follows that

$$\mathcal{Q}^{sB, \gamma, \lambda}(\beta) - \mathcal{Q}^{sB, \gamma, \lambda}(\hat{\beta}) \geq \sum_{j=1}^q w_j(\beta_j) + \sum_{j=q+1}^p v_j(\beta_j),$$

where

$$\begin{aligned} w_j(\beta_j) &= -\gamma \nu |\hat{\beta}_j|^{\nu-1} \text{sign}(\hat{\beta}_j)(\beta_j - \hat{\beta}_j) \\ &\quad + \rho_0(\beta_j - \hat{\beta}_j)^2/2 + J(|\beta_j|) - J(|\hat{\beta}_j|) \end{aligned}$$

and $v_j(\beta_j) = -o_p(\lambda)|\beta_j| + \rho_0|\beta_j|^2/2 + J(|\beta_j|)$.

First, consider the cases where $j \leq q$. If $|\beta_j| \geq a(\lambda - \gamma)$,

$$w_j(\beta_j) \geq -\gamma \nu |\hat{\beta}_j|^{\nu-1} \text{sign}(\hat{\beta}_j)(\beta_j - \hat{\beta}_j) + \gamma(|\beta_j|^\nu - |\hat{\beta}_j|^\nu) \geq 0.$$

And if $|\beta_j| < a(\lambda - \gamma)$,

$$w_j(\beta_j) \geq -\gamma\nu|\hat{\beta}_j|^{\nu-1}(|\beta_j| - |\hat{\beta}_j|) + \rho_0(\beta_j - \hat{\beta}_j)^2/2 + J(|\beta_j|) - J(|\hat{\beta}_j|).$$

Since $J(a(\lambda - \gamma)) - J(|\hat{\beta}_j|) \geq \gamma\nu|\hat{\beta}_j|^{\nu-1}\{a(\lambda - \gamma) - |\hat{\beta}_j|\}$, it follows that

$$w_j(\beta_j) \geq -\gamma\nu|\hat{\beta}_j|^{\nu-1}\{|\beta_j| - a(\lambda - \gamma)\} + \rho_0(\beta_j - \hat{\beta}_j)^2/2 + J(|\beta_j|) - J(a(\lambda - \gamma)) \geq \rho_0(\beta_j - \hat{\beta}_j)^2/2 - J(a(\lambda - \gamma)).$$

By triangular inequality, we have

$$|\beta_j - \hat{\beta}_j| \geq \min_{1 \leq j \leq q} |\beta_j^*| - \|\hat{\beta}_1 - \beta_1^*\|_2 - a(\lambda - \gamma).$$

Hence, from (24), as $n \rightarrow \infty$,

$$|\beta_j - \hat{\beta}_j|/\{a(\lambda - \gamma)\} \rightarrow \infty$$

since $\min_{1 \leq j \leq q} |\beta_j^*|/\lambda \rightarrow \infty$. On the other hand, $c^{\gamma,\lambda} = O(\lambda^{1-\nu})$ implies

$$J(a(\lambda - \gamma))/\{a(\lambda - \gamma)\}^2 = -c^{\gamma,\lambda}\{a(\lambda - \gamma)\}^{\nu-1}/(\nu + 1) + \lambda/a(\lambda - \gamma) = O(1).$$

Hence, $w_j(\beta_j) \geq 0$ for sufficiently large n . Next, consider the cases where $j > q$. If $|\beta_j| > a(\lambda - \gamma)$,

$$v_j(\beta_j) > -o_p(\lambda)|\beta_j| + \rho_0|\beta_j|^2/2 > |\beta_j|(-o_p(\lambda) + \rho_0a(\lambda - \gamma)/2) > 0$$

for all sufficiently large n . And if $|\beta_j| \leq a(\lambda - \gamma)$,

$$\begin{aligned} v_j(\beta_j) &\geq -o_p(\lambda)|\beta_j| + J(|\beta_j|) \\ &\geq |\beta_j|(-o_p(\lambda) - c^{\lambda,\gamma}\{a(\lambda - \gamma)\}^\nu/(\nu + 1) + \lambda) \\ &\geq |\beta_j|(-o_p(\lambda) - [\lambda - \gamma\nu\{a(\lambda - \gamma)\}^{\nu-1}]/(\nu + 1) + \lambda) \\ &\geq 0 \end{aligned}$$

for all sufficiently large n . Hence,

$$\mathcal{Q}^{sB,\gamma,\lambda}(\beta) - \mathcal{Q}^{sB,\gamma,\lambda}(\hat{\beta}) \geq 0$$

for all sufficiently large n which completes the proof. \square

Received 9 December 2011

REFERENCES

[1] AN, L. T. H. and TAO, P. D. (1997). Solving a class of linearly constrained indefinite quadratic problems by DC algorithms. *Journal of Global Optimization* **11** 253–285. [MR1469128](#)

[2] AN, L. T. H. and TAO, P. D. (2005). The DC (Difference of Convex Functions) Programming and DCA revisited with DC models of real world nonconvex optimization problems. *Annals of Operations Research* **133** 23–46. [MR2119311](#)

[3] AN, L. T. H., LE, M. H., NGUYEN, V. V., and TAO, P. D. (2008). A DC programming approach for feature selection in support vector machines learning. *Advances in Data Analysis and Classification* **2** 259–278. [MR2469770](#)

[4] BREIMAN, L. (1996). Heuristics of instability and stabilization in model selection. *The Annals of Statistics* **24** 2350–2383. [MR1425957](#)

[5] DONOHO, D. L. and JOHNSTONE, I. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* **81** 425–455. [MR1311089](#)

[6] EFRON, B., HASTIE, T., JOHNSTONE, I., and TIBSHIRANI, R. (2004). Least angle regression. *The Annals of Statistics* **32** 407–499. [MR2060166](#)

[7] FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96** 1348–1360. [MR1946581](#)

[8] FAN, J. and PENG, H. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics* **32** 928–961. [MR2065194](#)

[9] FRANK, I. E. and FRIEDMAN, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics* **35** 109–148.

[10] HOERL, A. E. and KENNARD, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12** 55–67.

[11] HUANG, J., HOROWITZ, J. L., and MA, S. (2008). Asymptotic properties of bridge estimators in sparse high-dimensional regression models. *The Annals of Statistics* **36** 587–613. [MR2396808](#)

[12] HUANG, J., MA, S., and ZHANG, C. H. (2008). Adaptive lasso for sparse high-dimensional regression models. *Statistica Sinica* **18** 1603–1618. [MR2469326](#)

[13] KIM, Y., CHOI, H., and OH, H. (2008). Smoothly clipped absolute deviation on high dimensions. *Journal of the American Statistical Association* **103** 1656–1673. [MR2510294](#)

[14] KIM, Y. and KWON, S. (2012). On the global optimality of non-convex penalized estimators. *Biometrika* **99** 315–325.

[15] KNIGHT, K. and FU, W. (2000). Asymptotics for the LASSO-type estimators. *The Annals of Statistics* **28** 1356–1378. [MR1805787](#)

[16] KWON, S. and KIM, Y. (2012). Large sample properties of the SCAD-penalized maximum likelihood estimation on high dimensions. *Statistica Sinica* **22** 629–953. [MR2954355](#)

[17] PARK, M. and HASTIE, T. (2007). l_1 -regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society, Ser. B* **69** 659–667. [MR2370074](#)

[18] ROSSET, S. and ZHU, J. (2007). Piecewise linear regularized solution paths. *The Annals of Statistics* **35** 1012–1030. [MR2341696](#)

[19] SCHEETZ, T. E., KIM, K.-Y. A., SWIDERSKI, R. E., PHILP, A. R., BRAUN, T. A., KNUDTSON, K. L., DORRANCE, A. M., DiBONA, G. F., HUANG, J., CASAVANT, T. L., SHEFFIELD, V. C., and STONE, E. M. (2006). Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proc. Nat. Acad. Sci.* **103** 14429–14434.

[20] SHEN, X., TSENG, C., ZHANG, X., and WONG, W. (2003). On psi-learning. *Journal of the American Statistical Association* **98** 724–734. [MR2011686](#)

[21] TIBSHIRANI, R. (1996). Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society, Ser. B* **58** 267–288. [MR1379242](#)

[22] WANG, H., LI, R., and TSAI, C.-H. (2007). Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika* **94** 553–568. [MR2410008](#)

[23] YUILLE, A. and RANGARAJAN, A. (2003). The concave-convex procedure. *Neural Computation* **15** 915–936.

[24] ZHANG, C. H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* **38** 894–942. [MR2604701](#)

- [25] ZOU, H. and HASTIE, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Ser. B* **67** 301–320. [MR2137327](#)
- [26] ZOU, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101** 1418–1429. [MR2279469](#)
- [27] ZOU, H. and ZHANG, H. (2009). On the adaptive elastic-net with a diverging number of parameters. *The Annals of Statistics* **37** 1733–1751. [MR2533470](#)

Sunghoon Kwon
School of Statistics
University of Minnesota
Minneapolis, MN 55455
USA
E-mail address: shkwon0522@gmail.com

Yongdai Kim
Dept. of Statistics
Seoul National University
Seoul 151-742
Korea
E-mail address: ydkim0903@gmail.com

Hosik Choi
Dept. of Informational Statistics
and Institute of Basic Science
Hoseo University, Chungnam 336-795
Korea
E-mail address: choi.hosik@gmail.com