

# Generalized Ridge Regression Estimator in High Dimensional Sparse Regression Models

Mahdi Roozbeh

*Department of Statistics, Semnan University, Iran*

**Abstract** Modern statistical analysis often encounters linear models with the number of explanatory variables much larger than the sample size. Estimation in these high-dimensional problems needs some regularization methods to be employed due to rank deficiency of the design matrix. In this paper, the ridge estimators are considered and their restricted regression counterparts are proposed when the errors are dependent under a multicollinearity and high-dimensionality setting. The asymptotic distributions of the proposed estimators are exactly derived. Incorporating the information contained in the restricted estimator, a shrinkage type ridge estimator is also exhibited and its asymptotic risk is analyzed under some special cases. To evaluate the efficiency of the proposed estimators, a Monté-Carlo simulation along with a real example are considered.

**Keywords** Asymptotic distribution; High-dimension; Ridge regression; Shrinkage estimator; Tikhonov regularization

**AMS 2010 subject classifications** Primary: 62G08, Secondary: 62J07

**DOI:** 10.19139/soic.v6i3.581

## 1. Introduction

Consider a linear regression model:

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (1.1)$$

where  $y_i$ s are responses,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$  is design points,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$  is vector denoting unknown coefficients,  $\varepsilon_i$ s are unobservable random errors and the superscript ( $^\top$ ) denotes the transpose of a vector or matrix. Further,  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$  has a cumulative distribution function  $F(\boldsymbol{\varepsilon})$ ;  $E(\boldsymbol{\varepsilon}) = \mathbf{0}$  and  $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{V}_n$ , where  $\sigma^2$  is finite and  $\mathbf{V}_n$  is a known matrix belonging to the space of all positive definite matrices of dimension  $n \times n$ , denoted by  $\mathcal{S}(n)$ .

Now-a-day, many data problems nowadays carry the structure that the number of covariates  $p$  may exceed sample size  $n$ , known as small  $n$ , large  $p$  problems. Such cases, that can be seen in the studies of genomics, financial markets, mobile phone communication, bioinformatics and risk management, regularization methods should be considered for inferring (see [9, 5]) about parameters of interest.

One of the mostly used regularization methods is the Tikhonov [21] regularization which was brought into statistical contexts by Hoerl and Kennard [12] as ridge regression. To mention a few recent researches, see e.g., [1, 2, 3, 4, 14, 15, 16].

In this paper, when a subset of coefficients is zero, the underlying model is called sparse. Under the sparsity assumption, the vector of coefficients  $\boldsymbol{\beta}$  can be partitioned as  $(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$  where  $\boldsymbol{\beta}_1$  is the coefficient vector for main effects and  $\boldsymbol{\beta}_2$  is the vector for nuisance effects or insignificant coefficients. We are essentially interested in the estimation of  $\boldsymbol{\beta}_1$  when it is reasonable that  $\boldsymbol{\beta}_2$  is close to zero.

The paper is organized as follows: The problem of interest is stated in section 2. In section 3, sparse ridge model is considered by proposing the new estimators, while their asymptotic properties are derived in section 4. Section 5 is devoted to exhibiting a shrinkage type ridge estimator and analyzing its asymptotic properties. Efficiencies of the proposed shrinkage estimators relative to the ordinary estimator are evaluated through a

\*Correspondence to: (Email: mahdi.roozbeh@semnan.ac.ir). Department of Statistics, Faculty of Mathematics, Statistics and Computer, Semnan University, Semnan, Iran, P.O. Box 35195-363.

Monté-Carlo simulation as well as a real data example in section 6. We conclude our study in section 7 by giving a summary.

## 2. Generalized Ridge Estimator

In this section, we discuss a biased estimation technique under multicollinearity for the regression models. So, the following preliminaries are needed. In a full matrix notation, the model (1.1) can be represented as

$$\mathbf{y}_n = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (2.2)$$

where  $\mathbf{X}_n = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$  and  $\mathbf{y}_n = (y_1, \dots, y_n)^\top$ .

As known, both of the ordinary least squares estimator (OLSE) and its covariance matrix are heavily dependent to the characteristics of the matrix  $\mathbf{S}_n = \mathbf{X}_n^\top \mathbf{X}_n$ . If  $\mathbf{S}_n$  is ill-conditioned, then the OLSE may affect by various numerical errors. The problem of multicollinearity can be solved by collecting additional data, reparameterizing the model and reselecting the variables. There are two well-known mathematical methods to overcome multicollinearity: the ‘principal components regression method’ and the ‘ridge regression method’. Here, we discuss the ridge regression method.

It is known that spectral decomposition of the (symmetric) positive definite matrix  $\mathbf{S}_n$  can be given by

$$\mathbf{S}_n = \mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^\top, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p), \quad (2.3)$$

where the columns of  $\mathbf{\Gamma}$  are eigenvectors of the matrix  $\mathbf{S}_n$  as well as the scalars  $\lambda_1, \dots, \lambda_p$  are its eigenvalues, satisfying

$$\lambda_1 \geq \dots \geq \lambda_p > 0,$$

without loss of generality. Therefore, the orthogonal (canonical) version of the model (2.2) is given by

$$\mathbf{y}_n = \mathbf{X}_n^* \boldsymbol{\alpha} + \boldsymbol{\varepsilon}, \quad (2.4)$$

where  $\mathbf{X}_n^* = \mathbf{X}_n \mathbf{\Gamma}$  and  $\boldsymbol{\alpha} = \mathbf{\Gamma}^\top \boldsymbol{\beta}$ . However, when  $\mathbf{S}_n$  is ill-conditioned, there exists an approximate linear dependency among its columns. So, the OLSE of  $\boldsymbol{\beta}$  has a large variance and multicollinearity is said to be appeared. In such situation, there exists a small positive constant  $\varepsilon$  such that

$$\lambda_1 \geq \dots \geq \lambda_r \gg \varepsilon \geq \lambda_{r+1} \geq \dots \geq \lambda_p > 0, \quad (2.5)$$

where  $r$  is called the numerical rank of  $\mathbf{S}_n$  (Watkins, 2002). More closeness of the small eigenvalues to the origin leads to more strength of the multicollinearity. To overcome damaging effect of the small eigenvalues, we need to increase them.

In another point of view based on the spectral decomposition, in the ridge regression the matrix  $\mathbf{S}_n$  is replaced by the matrix  $\mathbf{S}_n(k)$  defined by

$$\mathbf{S}_n(k) = \mathbf{\Gamma} \mathbf{\Lambda}(k) \mathbf{\Gamma}^\top, \quad \mathbf{\Lambda}(k) = \text{diag}(\lambda_1 + k, \dots, \lambda_p + k). \quad (2.6)$$

Hence, we have

$$\|\mathbf{S}_n(k) - \mathbf{S}_n\|_2 = k,$$

and,

$$\kappa_2(\mathbf{S}_n(k)) = \frac{\lambda_1 + k}{\lambda_p + k} \leq \kappa_2(\mathbf{S}_n) = \frac{\lambda_1}{\lambda_p},$$

where  $\kappa_2(\cdot)$  stands for the spectral condition number and  $\|\mathbf{x}\|_2$  denotes the Euclidean norm of the vector  $\mathbf{x}$ . Thus, computations of the ridge regression are numerically more stable than the OLSE.

In an extension scheme, the canonical generalized ridge estimator has been proposed by Hoerl and Kennard [12] as follows:

$$\hat{\boldsymbol{\alpha}}(\mathbf{K}_n^{(p)}) = (\mathbf{S}_n^* + \mathbf{K}_n^{(p)})^{-1} \mathbf{X}_n^{*\top} \mathbf{y}_n = \mathbf{T}_n^* (\mathbf{K}_n^{(p)}) \hat{\boldsymbol{\alpha}}, \quad \mathbf{T}_n^* (\mathbf{K}_n^{(p)}) = (\mathbf{K}_n^{(p)} \mathbf{S}_n^{*-1} + \mathbf{I}_p)^{-1}, \quad (2.7)$$

for some suitably chosen diagonal matrix  $\mathbf{K}_n^{(p)} = \text{diag}(k_n^{(1)}, \dots, k_n^{(p)})$  of tuning parameters, with  $k_n^{(i)} > 0$ ,  $i=1, \dots, p$ ,  $\mathbf{S}_n^* = \mathbf{X}_n^{*\top} \mathbf{X}_n^*$ , and  $\hat{\boldsymbol{\alpha}} = \mathbf{\Lambda}^{-1} \mathbf{X}_n^{*\top} \mathbf{y}_n$  is the canonical ordinary least-squares estimate of  $\boldsymbol{\alpha}$ . Hoerl

and Kennard [12] showed that for the known optimal values

$$k_n^{(i)} = \frac{\sigma^2}{\alpha_i^2}, \quad i = 1, \dots, p, \quad (2.8)$$

the generalized ridge regression estimator is superior to all of the other ones within the class of biased estimators, where  $\sigma^2$  is the error variance of the model (2.2) and  $\alpha_i^2$  is the  $i^{\text{th}}$  entry of  $\alpha$ . Moreover, Hocking et al. [13] showed that with these optimal  $k_n^{(i)}$ , values the proposed estimator (2.7) is superior to all estimators within the class of biased estimators they considered. However, the optimal value of  $k_n^{(i)}$  fully depends on the unknown values  $\sigma^2$  and  $\alpha_i$  which must be estimated from the observed data. Hoerl and Kennard [12] suggested to replace  $\sigma^2$  and  $\alpha_i^2$  by their corresponding unbiased estimators. That is,

$$\hat{k}_n^{(i)} = \frac{\hat{\sigma}^2}{\hat{\alpha}_i^2}, \quad i = 1, \dots, p, \quad (2.9)$$

where  $\hat{\sigma}^2$  is an unbiased and effective estimator of  $\sigma^2$  and  $\hat{\alpha}_i$  is the  $i^{\text{th}}$  entry of  $\hat{\alpha}$  which is an unbiased estimator of  $\alpha$ . Hoerl and Kennard [12] also suggested an iterative procedure to estimate  $k_n^{(i)}$ . Two other methods to select  $k_n^{(i)}$  have been proposed by Hemmerle and Brantle [10]. Some exact finite sample properties for (2.7) can be found in Hemmerle and Carey [11]. Thus, the matrix  $S_n$  is replaced by the matrix  $S_n(\hat{K}_n^{(p)})$  defined by

$$S_n(\hat{K}_n^{(p)}) = \Gamma \Lambda(\hat{K}_n^{(p)}) \Gamma^\top, \quad \Lambda(\hat{K}_n^{(p)}) = \text{diag}(\lambda_1 + \hat{k}_n^{(1)}, \dots, \lambda_p + \hat{k}_n^{(p)}). \quad (2.10)$$

Hence, we have

$$\|S_n(\hat{K}_n^{(p)}) - S_n\|_2 = \max_{i=1, \dots, p} \{\hat{k}_n^{(i)}\}, \quad (2.11)$$

and,

$$\kappa_2(S(\hat{K}_n^{(p)})) = \frac{\max_{i=1, \dots, p} \{\lambda_i + \hat{k}_n^{(i)}\}}{\min_{i=1, \dots, p} \{\lambda_i + \hat{k}_n^{(i)}\}}.$$

Now, considering the high-dimensional case  $p > n$ , we are primary interested in estimating the regression vector-parameter  $\beta$  in the model (2.2). Since  $S_n$  is rank deficient, a regularization method is needed to combat this ill-conditioning, such as Tikhonov [21] regularization. In our setup, estimating  $\beta$  is equal to minimizing the following general criterion, under the  $L_2$  norm

$$\|V_n^{-\frac{1}{2}}(\mathbf{y}_n - \mathbf{X}_n \beta)\|_2^2 + \|\mathbf{K}_n^{(p)\frac{1}{2}} \beta\|_2^2$$

Let  $F(\beta) = \left\{ \|V_n^{-\frac{1}{2}}(\mathbf{y}_n - \mathbf{X}_n \beta)\|_2^2 + \|\mathbf{K}_n^{(p)\frac{1}{2}} \beta\|_2^2 \right\}$ . Then

$$\hat{\beta}(\mathbf{K}_n^{(p)}) = \text{argmin}_{\beta} F(\beta) = (\mathbf{X}_n^\top V_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)})^{-1} \mathbf{X}_n^\top V_n^{-1} \mathbf{y}_n, \quad (2.12)$$

where  $\hat{\beta}(\mathbf{K}_n^{(p)})$  is called generalized ridge estimator (GRE).

For the case  $\mathbf{K}_n^{(p)} = k\mathbf{I}_p$ ,  $k > 0$ , the GRE reduces to the ordinary ridge regression estimator introduced by Hoerl and Kennard [12]. If  $k_n^{(i)} \rightarrow 0$ ,  $i = 1, \dots, p$ , then  $\hat{\beta}(\mathbf{K}_n^{(p)})$  reduces to the well-known ordinary least square (GLS) estimator  $\hat{\beta}_n = (\mathbf{X}_n^\top V_n^{-1} \mathbf{X}_n)^{-1} \mathbf{X}_n^\top V_n^{-1} \mathbf{y}_n$ .

In what follows, we discuss about the properties of  $\hat{\beta}(\mathbf{K}_n^{(p)})$ .

Since  $\hat{\beta}(\mathbf{K}_n^{(p)}) = \text{argmin}_{\beta} F(\beta)$ ,  $F(\hat{\beta}(\mathbf{K}_n^{(p)})) \leq F(\mathbf{0})$ . Therefore,

$$\begin{aligned} \|\mathbf{K}_n^{(p)\frac{1}{2}} \hat{\beta}\|_2^2 &\leq \|V_n^{-\frac{1}{2}}(\mathbf{y}_n - \mathbf{X}_n \hat{\beta})\|_2^2 + \|\mathbf{K}_n^{(p)\frac{1}{2}} \hat{\beta}\|_2^2 \\ &= F(\hat{\beta}(\mathbf{K}_n^{(p)})) \end{aligned}$$

$$\leq F(\mathbf{0}) = \left\| \mathbf{V}_n^{-\frac{1}{2}} \mathbf{y}_n \right\|_2^2$$

Thus, for  $k_1 = \min(k_n^{(i)})$ , we get

$$\left\| \hat{\boldsymbol{\beta}} \right\|_2^2 \leq k_1^{-1} \left\| \mathbf{V}_n^{-\frac{1}{2}} \mathbf{y}_n \right\|_2^2, \quad \text{and} \quad E \left\| \hat{\boldsymbol{\beta}} \right\|_2^2 \leq k_1^{-1} E \left\| \mathbf{V}_n^{-\frac{1}{2}} \mathbf{y}_n \right\|_2^2 \quad (2.13)$$

Now, we focus on the covariance. For the high-dimensional case  $p > n$ , using spectral decomposition

$$\mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n = \mathbf{\Gamma} \begin{bmatrix} \mathbf{\Lambda}_{n \times n} & \mathbf{O}_{n \times (p-n)} \\ \mathbf{O}_{n \times (p-n)}^\top & \mathbf{O}_{(p-n) \times (p-n)} \end{bmatrix} \mathbf{\Gamma}^\top, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n), \quad (2.14)$$

where  $\lambda_1 \geq \dots \geq \lambda_n > 0$ .

*Theorem 1*

Let  $k_n = \max_{1 \leq i \leq n} (k_n^{(i)})$  and  $k_o = \sum_{j=1}^n (1/k_n^{(j)})^2$ . Assume  $\lambda_j = o(k_n)$ , for  $j = 1, \dots, n$ . Then, we have

$$\lim_{p \rightarrow \infty} \text{tr} \left( \text{Cov}(\hat{\boldsymbol{\beta}}(\mathbf{K}_n^{(p)}) - \boldsymbol{\beta}) \right) = \sigma^2 k_o \text{tr} \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n \right).$$

**Proof:** By definition, after some algebra,

$$\begin{aligned} \text{Cov}(\hat{\boldsymbol{\beta}}(\mathbf{K}_n^{(p)}) - \boldsymbol{\beta}) &= \sigma^2 \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} \\ &= \sigma^2 \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} - \sigma^2 \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} \\ &\quad \times \mathbf{K}_n^{(p)} \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} \\ &= \sigma^2 \begin{bmatrix} \frac{\lambda_1}{(\lambda_1 + k_n^{(1)})^2} & & & \mathbf{O}_{n \times (p-n)} \\ & \ddots & & \\ & & \frac{\lambda_n}{(\lambda_n + k_n^{(n)})^2} & \\ \mathbf{O}_{n \times (p-n)}^\top & & & \mathbf{O}_{(p-n) \times (p-n)} \end{bmatrix} \end{aligned}$$

since

$$\begin{aligned} \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} &= \left\{ \mathbf{\Gamma} \begin{bmatrix} \mathbf{\Lambda}_{n \times n} & \mathbf{O}_{n \times (p-n)} \\ \mathbf{O}_{n \times (p-n)}^\top & \mathbf{O}_{(p-n) \times (p-n)} \end{bmatrix} \mathbf{\Gamma}^\top + \mathbf{K}_n^{(p)} \right\}^{-1} \\ &= \begin{bmatrix} \lambda_1 + k_n^{(1)} & & & \mathbf{O}_{n \times (p-n)} \\ & \ddots & & \\ & & \lambda_n + k_n^{(n)} & \\ \mathbf{O}_{n \times (p-n)}^\top & & & k_n^{(n+1)} \\ & & & \ddots \\ & & & & k_n^{(p)} \end{bmatrix}^{-1} \end{aligned} \quad (2.15)$$

Hence,  $\text{tr} \text{Cov}(\hat{\boldsymbol{\beta}}(\mathbf{K}_n^{(p)}) - \boldsymbol{\beta}) = \sigma^2 \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + k_n^{(j)})^2}$ . Under the assumption  $\lambda_j = o(k_n)$ , for  $j = 1, \dots, n$  we have

$$\lim_{p \rightarrow \infty} \sum_{j=1}^n \frac{1}{(\lambda_j + k_n^{(j)})^2} = \sum_{j=1}^n \left( \frac{1}{k_n^{(j)}} \right)^2 = k_o. \quad (2.16)$$

So, the proof is complete.  $\square$

Under some mild conditions, it can be shown that  $\frac{1}{\sigma^2 k_o} (\hat{\boldsymbol{\beta}}(\mathbf{K}_n^{(p)}) - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \text{diag} \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n)$ . For more detail on high-dimensional properties, we refer to [7].

### 3. Sub-model Approach

Since  $p > n$ , one methodology to infer about regression parameters is to reduce the number of features from  $p$  to  $p_1 < n$ . Hence, the regression parameter  $\beta$  is partitioned as  $\beta = (\beta_1^\top, \beta_2^\top)^\top$ , where the sub-vector  $\beta_i$  has dimension  $p_i$ ,  $i = 1, 2$  and  $p_1 + p_2 = p$ . Thus, the underlying model has form

$$\mathbf{y}_n = \mathbf{X}_n^{(1)}\beta_1 + \mathbf{X}_n^{(2)}\beta_2 + \epsilon, \quad (3.1)$$

where  $\mathbf{X}_n$  is partitioned to  $(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)})$  in such a way that  $\mathbf{X}_n^{(i)}$  is a  $n \times p_i$  sub-matrix,  $i = 1, 2$ . With respect to this partitioning, the generalized least square estimators (GOLSEs) of  $\beta_1$  and  $\beta_2$  are respectively given by

$$\hat{\beta}_n^{(1)} = \mathbf{C}_n^{(1)-1} \mathbf{X}_n^{(1)\top} \Sigma_n^{(2)-1} \mathbf{y}_n, \quad \mathbf{C}_n^{(1)} = \mathbf{X}_n^{(1)\top} \Sigma_n^{(2)-1} \mathbf{X}_n^{(1)} \quad (3.2)$$

$$\hat{\beta}_n^{(2)} = \mathbf{C}_n^{(2)-1} \mathbf{X}_n^{(2)\top} \Sigma_n^{(1)-1} \mathbf{y}_n, \quad \mathbf{C}_n^{(2)} = \mathbf{X}_n^{(2)\top} \Sigma_n^{(1)-1} \mathbf{X}_n^{(2)} \quad (3.3)$$

where

$$\Sigma_n^{(i)-1} = \mathbf{V}_n^{-1} - \mathbf{V}_n^{-1} \mathbf{X}_n^{(i)} \left( \mathbf{X}_n^{(i)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(i)} \right)^{-1} \mathbf{X}_n^{(i)\top} \mathbf{V}_n^{-1}, \quad i = 1, 2.$$

The sparse model is defined when  $\mathcal{H}_o : \beta_2 = 0$  is true. In this paper, we refer restricted regression model (RRM) to the sparse model.

For the RRM, the generalized restricted estimator based on OLS estimator (GROLSE) has form

$$\hat{\beta}_n^{R1} = \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{y}_n. \quad (3.4)$$

According to [20], the GROLSE performs better than GOLSE when model is sparse. However, the former estimator performs poorly as  $\beta_2$  is different from zero. In the following result, in a similar fashion as in [22], the relation between the sub-model and full-model estimators of  $\beta_1$  is obtained. It can be easily shown that

$$\hat{\beta}_n^{(1)} = \hat{\beta}_n^{R1} - \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(2)} \hat{\beta}_n^{(2)}. \quad (3.5)$$

#### 3.1. Sparse ridge model

Under the sparsity assumption, i.e.,  $\beta_2 = 0$  in the high-dimensional problem, following [19], the generalized restricted ridge estimator (GRRE), is given by

$$\begin{aligned} \hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right) &= \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} + \mathbf{K}_n^{(p_1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{y}_n \\ &= \left( \mathbf{I}_{p_1} + \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{K}_n^{(p_1)} \right)^{-1} \hat{\beta}_n^{R1} \\ &= \mathbf{T}_1 \left( \mathbf{K}_n^{(p_1)} \right) \hat{\beta}_n^{R1}, \end{aligned} \quad (3.6)$$

where  $\mathbf{K}_n^{(p_1)} = \text{diag} \left( k_n^{(1)}, \dots, k_n^{(p_1)} \right)$  is the ridge parameter matrix as a function of sample size  $n$  and

$$\mathbf{T}_1 \left( \mathbf{K}_n^{(p_1)} \right) = \left( \mathbf{I}_{p_1} + \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{K}_n^{(p_1)} \right)^{-1}.$$

Similarly, the generalized ridge estimators (GRES) of  $\beta_1$  and  $\beta_2$  respectively have forms

$$\begin{aligned} \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right) &= \left( \mathbf{X}_n^{(1)\top} \Sigma_2^{-1} \left( \mathbf{K}_n^{(p_2)} \right) \mathbf{X}_n^{(1)} + \mathbf{K}_n^{(p_1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \Sigma_2^{-1} \left( \mathbf{K}_n^{(p_2)} \right) \mathbf{y}_n \\ &= \left( \mathbf{I}_{p_1} + \left( \mathbf{X}_n^{(1)\top} \Sigma_2^{-1} \left( \mathbf{K}_n^{(p_2)} \right) \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{K}_n^{(p_1)} \right)^{-1} \hat{\beta}_n^{(1)} \\ &= \mathbf{R}_1 \left( \mathbf{K}_n^{(p_1)} \right) \hat{\beta}_n^{(1)}, \end{aligned} \quad (3.7)$$

$$\begin{aligned} \hat{\beta}_n^{(2)} \left( \mathbf{K}_n^{(p_2)} \right) &= \left( \mathbf{X}_n^{(2)\top} \Sigma_1^{-1} \left( \mathbf{K}_n^{(p_1)} \right) \mathbf{X}_n^{(2)} + \mathbf{K}_n^{(p_2)} \right)^{-1} \mathbf{X}_n^{(2)\top} \Sigma_1^{-1} \left( \mathbf{K}_n^{(p_1)} \right) \mathbf{y}_n \\ &= \left( \mathbf{I}_{p_2} + \left( \mathbf{X}_n^{(2)\top} \Sigma_1^{-1} \left( \mathbf{K}_n^{(p_1)} \right) \mathbf{X}_n^{(2)} \right)^{-1} \mathbf{K}_n^{(p_2)} \right)^{-1} \hat{\beta}_n^{(2)} \end{aligned}$$

$$= \mathbf{R}_2 \left( \mathbf{K}_n^{(p_2)} \right) \hat{\boldsymbol{\beta}}_n^{(2)}, \quad (3.8)$$

where

$$\begin{aligned} \mathbf{R}_1 \left( \mathbf{K}_n^{(p_1)} \right) &= \left( \mathbf{I}_{p_1} + \left( \mathbf{X}_n^{(1)\top} \boldsymbol{\Sigma}_2^{-1} \left( \mathbf{K}_n^{(p_2)} \right) \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{K}_n^{(p_1)} \right)^{-1}, \\ \mathbf{R}_2 \left( \mathbf{K}_n^{(p_2)} \right) &= \left( \mathbf{I}_{p_2} + \left( \mathbf{X}_n^{(2)\top} \boldsymbol{\Sigma}_1^{-1} \left( \mathbf{K}_n^{(p_1)} \right) \mathbf{X}_n^{(2)} \right)^{-1} \mathbf{K}_n^{(p_2)} \right)^{-1}, \\ \boldsymbol{\Sigma}_i^{-1} \left( \mathbf{K}_n^{(p_i)} \right) &= \mathbf{V}_n^{-1} - \mathbf{V}_n^{-1} \mathbf{X}_n^{(i)} \left( \mathbf{X}_n^{(i)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(i)} + \mathbf{K}_n^{(p_i)} \right)^{-1} \mathbf{X}_n^{(i)\top} \mathbf{V}_n^{-1}, \quad i = 1, 2. \end{aligned} \quad (3.9)$$

It is easy to show that

$$\hat{\boldsymbol{\beta}}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right) = \hat{\boldsymbol{\beta}}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right) - \left( \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(2)} \hat{\boldsymbol{\beta}}_n^{(2)} \left( \mathbf{K}_n^{(p_2)} \right). \quad (3.10)$$

#### 4. Asymptotics

In this section we study the asymptotic performance of the sub-model estimator when  $p_1 < n$ . For the purpose of this paper, we need to take the following assumptions as regularity conditions:

(A1)  $\max_{1 \leq i \leq n} \mathbf{x}_i^\top \left( \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} \right)^{-1} \mathbf{x}_i = o(n)$ , where  $\mathbf{x}_i^\top$  is the  $i^{\text{th}}$  row of  $\mathbf{X}_n$ .

(A2)  $\frac{k_n^{(i)}}{n} \rightarrow k_o$  as  $n \rightarrow \infty$ ,  $i = 1, \dots, p$ .

(A3) Let

$$\begin{aligned} \mathbf{A}_n &= \mathbf{X}_n^\top \mathbf{V}_n^{-1} \mathbf{X}_n + \mathbf{K}_n^{(p)} = \begin{pmatrix} \mathbf{X}_n^{(1)\top} \\ \mathbf{X}_n^{(2)\top} \end{pmatrix} \mathbf{V}_n^{-1} \begin{pmatrix} \mathbf{X}_n^{(1)} & \mathbf{X}_n^{(2)} \end{pmatrix} + \mathbf{K}_n^{(p)} \\ &= \begin{pmatrix} \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} & \mathbf{X}_n^{(1)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(2)} \\ \mathbf{X}_n^{(2)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(1)} & \mathbf{X}_n^{(2)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(2)} \end{pmatrix} + \begin{pmatrix} \mathbf{K}_n^{(p_1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_n^{(p_2)} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A}_{n11} & \mathbf{A}_{n12} \\ \mathbf{A}_{n21} & \mathbf{A}_{n22} \end{pmatrix}. \end{aligned}$$

Then, there exists a positive definite matrix  $\mathbf{A}$  such that  $\frac{1}{n} \mathbf{A}_n \rightarrow \mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$ , as  $n \rightarrow \infty$ .

(A4)  $\mathbf{X}_n^{(2)\top} \mathbf{V}_n^{-1} \mathbf{X}_n^{(2)} + \mathbf{K}_n^{(p_2)} = o(\sqrt{n})$ .

(A5)  $\mathbf{x}_i^\top \mathbf{x}_j = o(\sqrt{n})$ ,  $i, j = 1, \dots, n$ .

By (A2), (A4) and (A5), one can directly conclude that  $\boldsymbol{\Sigma}_2^{-1} \left( \mathbf{K}_n^{(p_2)} \right) \rightarrow \mathbf{V}_n^{-1}$  as  $n \rightarrow \infty$ .

According to [20], the test statistic for testing  $\mathcal{H}_o : \boldsymbol{\beta}_2 = \mathbf{0}$  diverges as  $n \rightarrow \infty$ , under any fixed alternatives  $\mathbf{A}_\xi : \boldsymbol{\beta}_2 = \boldsymbol{\xi}$ . To overcome this difficulty, in sequel, the following local alternatives are considered

$$K_{(n)} : \boldsymbol{\beta}_2 = \boldsymbol{\beta}_{2(n)} = n^{-\frac{1}{2}} \boldsymbol{\xi}, \quad (4.1)$$

where  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_{p_2})^\top \in \mathbb{R}^{p_2}$  is a fixed vector.

For notational convenience, let  $-k_o \mathbf{A}^{-1} \boldsymbol{\beta} = \boldsymbol{\mu} = (\boldsymbol{\mu}_1^\top, \boldsymbol{\mu}_2^\top)^\top$ ,  $\boldsymbol{\delta} = \mathbf{A}_{11}^{-1} \mathbf{A}_{12} \boldsymbol{\xi}$ ,  $\boldsymbol{\mu}_{11.2} = \boldsymbol{\mu}_1 - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \left( (\boldsymbol{\beta}_2 - \boldsymbol{\xi}) - \boldsymbol{\mu}_2 \right)$ ,  $\boldsymbol{\gamma} = \boldsymbol{\mu}_{11.2} + \boldsymbol{\delta}$ ,  $\mathbf{A}_{22.1} = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$ ,  $\mathbf{B} = \mathbf{A}_{21} \mathbf{A}_{11}^{-2} \mathbf{A}_{12} \mathbf{A}_{22.1}^{-1}$ . Then, we have

*Theorem 2*

Under the regularity conditions (A1)-(A3) and local alternatives  $\{K_{(n)}\}$

- (i)  $\mathbf{V}_{(n)}^{(1)} = \sqrt{n} \left( \hat{\boldsymbol{\beta}}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right) - \boldsymbol{\beta}_1 \right) \xrightarrow{\mathcal{D}} N_{p_1} \left( -\boldsymbol{\mu}_{11.2}, \sigma^2 \mathbf{A}_{11.2}^{-1} \right)$
- (ii)  $\mathbf{V}_{(n)}^{(2)} = \sqrt{n} \left( \hat{\boldsymbol{\beta}}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right) - \boldsymbol{\beta}_1 \right) \xrightarrow{\mathcal{D}} N_{p_1} \left( -\boldsymbol{\gamma}, \sigma^2 \mathbf{A}_{11}^{-1} \right)$

$$(iii) \mathbf{V}_n^{(3)} = \sqrt{n} \left( \hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right) - \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right) \right) \xrightarrow{\mathcal{D}} N_{p_1} \left( -\gamma + \boldsymbol{\mu}_{11.2}, \sigma^2 \left( \mathbf{A}_{11}^{-1} - \mathbf{A}_{11.2}^{-1} \right) \right)$$

where  $\mathbf{A}_{11.2} = \mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21}$ .

If  $\mathbf{A}_{12} = \mathbf{0}$ , then  $\boldsymbol{\delta} = \mathbf{0}$ ,  $\boldsymbol{\gamma} = \boldsymbol{\mu}_{11.2}$  and  $\mathbf{A}_{11.2} = \mathbf{A}_{11}$ , and all the asymptotic risk functions reduce to common value  $\sigma^2 \text{tr}(\mathbf{A}_{11}^{-1}) + \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2}$  for all  $\boldsymbol{\xi}$ . In a similar fashion as in [20], if  $\mathbf{A}_{12} \neq \mathbf{0}$ , then when  $\Delta^2 = (\boldsymbol{\xi}^\top \mathbf{A}_{22.1}^{-1} \boldsymbol{\xi}) \sigma^{-2}$ ,  $\mathbf{A}_{22.1} = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$  moves away from 0, the asymptotic risk of  $\hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right)$  becomes unbounded. Further, if  $\boldsymbol{\xi} = \mathbf{0}$ ,  $\hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right)$  is superior to  $\hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right)$  (denoted by  $\hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right) > \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right)$ ) in the sense of having smaller asymptotic risk function.

## 5. Shrinkage Estimator

Under the aforementioned analysis in previous section, one may consider improving the GRE of  $\beta_1$  by shrinking toward  $\mathbf{0}$  (see [20] for more details). In this respect, one plausible choice is to combine both GRE and GRRE of  $\beta_1$  to obtain

$$\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) = \alpha_1 \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right) + \alpha_2 \hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right). \quad (5.1)$$

Now, let  $\beta^*$  be any estimator of  $\beta_1$ . Then the asymptotic distributional risk (ADR) of  $\beta^*$  is evaluated by

$$ADR(\beta_1; \beta^*) = \lim_{n \rightarrow \infty} E \left[ n(\beta^* - \beta_1)^\top (\beta^* - \beta_1) \right]$$

Under the pronounced regularity conditions and local alternatives  $\{K_{(n)}\}$ , the asymptotic risk of shrinkage ridge estimator (SRE)  $\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right)$ , after some algebra, is given by

$$\begin{aligned} ADR \left( \beta_1; \hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) \right) &= \alpha_1^2 \left[ \sigma^2 \text{tr}(\mathbf{A}_{11.2}^{-1}) + \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2} \right] + \alpha_2^2 \left[ \sigma^2 \text{tr}(\mathbf{A}_{11}^{-1}) + \boldsymbol{\gamma}^\top \boldsymbol{\gamma} \right] \\ &\quad + (\alpha_1 + \alpha_2 - 1)^2 \boldsymbol{\xi}^\top \boldsymbol{\xi} - 2\alpha_1(\alpha_1 + \alpha_2 - 1) \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\xi} - 2\alpha_2(\alpha_1 + \alpha_2 - 1) \boldsymbol{\gamma}^\top \boldsymbol{\xi} \\ &\quad + \alpha_1 \alpha_2 \left( 2\boldsymbol{\mu}_{11.2}^\top \boldsymbol{\xi} + 2\boldsymbol{\gamma}^\top \boldsymbol{\xi} - 2\boldsymbol{\mu}_{11.2}^\top \boldsymbol{\gamma} - \boldsymbol{\xi}^\top \boldsymbol{\xi} \right). \end{aligned} \quad (5.2)$$

Obviously, under  $\boldsymbol{\xi} = \mathbf{0}$ ,  $\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) > \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right)$  whenever

$$(1 - \alpha_1^2) \sigma^2 \text{tr}(\mathbf{A}_{11.2}^{-1}) - \alpha_2^2 \sigma^2 \text{tr}(\mathbf{A}_{11}^{-1}) + [1 - (\alpha_1 - \alpha_2)^2] \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2} > 0;$$

and for the special case  $\alpha_1 = 1$ , it can be revealed that  $\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) > \hat{\beta}_n^{(1)} \left( \mathbf{K}_n^{(p_1)} \right)$  whenever

$$0 < \alpha_2 < \frac{2\boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2}}{\sigma^2 \text{tr}(\mathbf{A}_{11}^{-1}) + \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2}}. \quad (5.3)$$

Similarly,  $\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) > \hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right)$  whenever

$$(1 - \alpha_2^2) \sigma^2 \text{tr}(\mathbf{A}_{11}^{-1}) - \alpha_1^2 \sigma^2 \text{tr}(\mathbf{A}_{11.2}^{-1}) + [1 - (\alpha_1 - \alpha_2)^2] \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2} > 0;$$

and for the special case  $\alpha_2 = 1$ , it can be revealed that  $\hat{\beta}_n^S \left( \mathbf{K}_n^{(p_1)} \right) > \hat{\beta}_n^{R1} \left( \mathbf{K}_n^{(p_1)} \right)$  whenever

$$0 < \alpha_1 < \frac{2\boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2}}{\sigma^2 \text{tr}(\mathbf{A}_{11.2}^{-1}) + \boldsymbol{\mu}_{11.2}^\top \boldsymbol{\mu}_{11.2}}. \quad (5.4)$$

## 6. Numerical Results

In this section, the Monté-Carlo studies and a real data example about riboflavin production data set (see [6]) are consider to justify the assertions.

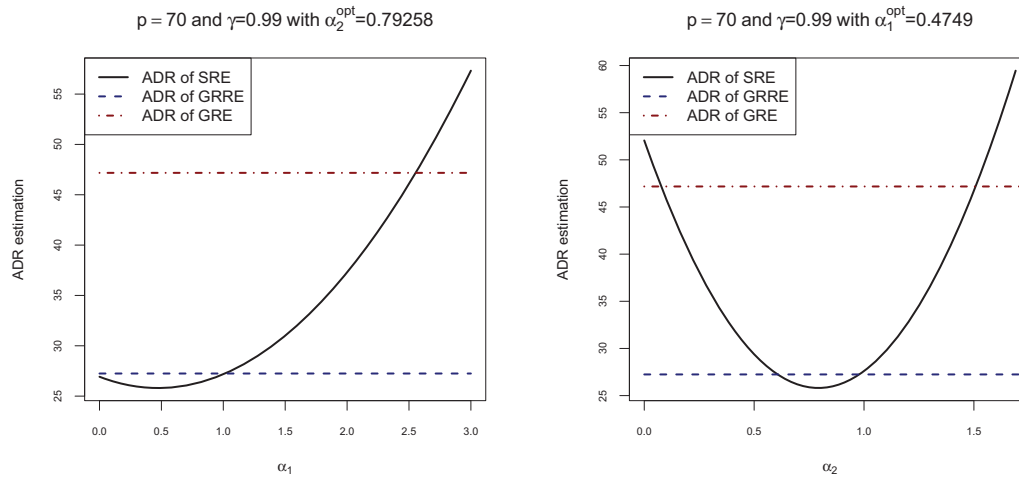


Figure 1. The diagram of *ADRs* versus  $\alpha_1$  and  $\alpha_2$  for simulated data set.

**6.1. The Monté-Carlo simulation studies**

To examine the performance of the proposed estimators, a Monté-Carlo simulation is performed. To achieve different degrees of collinearity, following [18, 8], the explanatory variables were generated using the following device for  $n = 25$  and  $p = 30, 50$  and  $70$  from the following model:

$$x_{ij} = (1 - \gamma^2)^{\frac{1}{2}} z_{ij} + \gamma z_{ip}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, p \tag{6.1}$$

where  $z_{ij}$  are independent standard normal pseudo-random numbers and  $\gamma$  is specified so that the correlation between any two explanatory variables is given by  $\gamma^2$ . These variables are then standardized so that  $\mathbf{X}_n^\top \mathbf{X}_n$  and  $\mathbf{X}_n^\top \mathbf{y}_n$  are in correlation forms. Three different values  $\gamma = 0.75, 0.90$  and  $0.99$  are considered for the correlation. Then the observations for the dependent variable are determined by

$$\mathbf{y}_n = \mathbf{X}_n^{(1)} \boldsymbol{\beta}_1 + \mathbf{X}_n^{(2)} \boldsymbol{\beta}_2 + \boldsymbol{\epsilon}, \tag{6.2}$$

with  $\boldsymbol{\beta}_1 = (-1.25, 1, 2.5, 4, -3, -5)^\top$ ,  $p_1 = 6$  and  $n = 25$ . To achieve the sparse model,  $\boldsymbol{\beta}_2$  is generated as  $N_{p_2}(\mathbf{0}, 0.01 \times \mathbf{I}_{p_2})$ ,  $p_2 = p - p_1$ . The error term of the model (6.2) is generated as

$$\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{V}_n), \quad \sigma^2 = 1.44, \quad \mathbf{V}_n[i, j] = \exp(-7|i - j|), \quad i, j = 1, \dots, n.$$

The Monté-Carlo simulation is performed with  $M = 10^3$  replications, obtaining the ridge estimators  $\hat{\boldsymbol{\beta}}_1 = \hat{\boldsymbol{\beta}}_n^{(1)}(\mathbf{K}_n^{(p_1)})$ ,  $\hat{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_n^{R1}(\mathbf{K}_n^{(p_1)})$ ,  $\hat{\boldsymbol{\beta}}_3 = \hat{\boldsymbol{\beta}}_n^S(\mathbf{K}_n^{(p_1)})$ , in the sparse restricted regression model.

The relative efficiencies of the above methods with respect to the first method are estimated as

$$Eff(\hat{\boldsymbol{\beta}}_i) = \frac{\frac{1}{M} \sum_{m=1}^M \left\| \hat{\boldsymbol{\beta}}_1^{(m)} - \boldsymbol{\beta}_1 \right\|_2^2}{\frac{1}{M} \sum_{m=1}^M \left\| \hat{\boldsymbol{\beta}}_i^{(m)} - \boldsymbol{\beta}_1 \right\|_2^2}, \quad i = 1, 2, 3,$$

where  $\hat{\boldsymbol{\beta}}_i^{(m)}$  is the estimators obtained in the  $m$ th iteration.

All numerical computations were conducted using the statistical package R. In Tables 1 to 4, the proposed estimators along with  $ADR(\cdot)$ , efficiencies of proposed estimators relative to GRE and optimal values of  $\alpha_1$  and  $\alpha_2$  were computed. The asymptotic risk of SRE was employed to select the optimal parameters  $\alpha_1$  and  $\alpha_2$ , numerically, by minimizing the ADR function of SRE. The minimum of ADR approximately occurred at  $\alpha_1^{opt} = 0.4749$  and  $\alpha_2^{opt} = 0.7449$  for  $p = 70$  and  $\gamma = 0.99$ . Figure 1 shows the ADR functions versus  $\alpha_1$  and  $\alpha_2$ . Since the results were similar across cases, only the results for  $p = 70$  and  $\gamma = 0.99$  were reported to save space.

As it can be found from Tables 1-3, the SRE performs better than the other proposed estimators in the sense of having smaller ADR. Also, because of the sparsity assumption, both GRRE and SRE are more efficient than GRE. Moreover, since the simulated model is sparse, as it is evident from Table 4, the allocated weight to the GRRE ( $\alpha_2$ ) is greater than that of the GRE ( $\alpha_1$ ). From Figure 1, one can distinguish the optimal regions of superiority of SRE over GRE and GRRE.



Table 1. Evaluation of the proposed estimators for  $\gamma = 0.75$

Method Coefficients	$p = 30$			$p = 50$			$p = 70$		
	GRE	GRRE	SRE	GRE	GRRE	SRE	GRE	GRRE	SRE
$\hat{\beta}_1$	-1.3423	-1.2296	-1.2932	-1.0754	-1.2549	-1.3500	-0.9888	-1.4900	-1.5778
$\hat{\beta}_2$	0.4666	1.0251	1.0085	0.0492	0.9809	0.9399	0.2205	0.5115	0.5116
$\hat{\beta}_3$	1.7686	2.5435	2.5604	0.7694	2.4766	2.4758	0.2324	1.7544	1.8014
$\hat{\beta}_4$	2.9308	4.0521	4.0972	1.5382	3.9842	4.0167	0.7041	3.0047	3.1030
$\hat{\beta}_5$	-2.6616	-2.5916	-2.7564	-1.9584	-2.6366	-2.8188	-1.8737	-4.1290	-4.2485
$\hat{\beta}_6$	-3.2513	-4.8248	-4.8671	-1.6499	-4.9238	-4.9359	-1.5163	-3.9774	-4.1091
$k^{opt}$	1e-04	0.01	—	0.01	0.13	—	0.14	0.36	—
$\hat{A}DR(\hat{\beta}_i)$	9.6646	1.0270	0.7930	24.091	0.9813	0.7926	32.490	6.5667	6.1225
$Eff(\hat{\beta}_i)$	1.0000	9.4102	12.187	1.0000	24.550	30.394	1.0000	4.9476	5.3066

Table 2. Evaluation of the proposed estimators for  $\gamma = 0.90$

Method Coefficients	$p = 30$			$p = 50$			$p = 70$		
	GRE	GRRE	SRE	GRE	GRRE	SRE	GRE	GRRE	SRE
$\hat{\beta}_1$	-1.3071	-1.2120	-1.2890	-0.9863	-1.2109	-1.6814	-0.7243	-1.2703	-1.7947
$\hat{\beta}_2$	0.4390	0.9715	0.9271	0.1208	1.0719	0.8296	0.0910	0.9842	0.7958
$\hat{\beta}_3$	1.7053	2.5750	2.5639	0.9122	2.7129	2.5980	0.5531	2.6233	2.5118
$\hat{\beta}_4$	2.8883	4.1709	4.1653	1.6350	4.3926	4.3216	1.0719	4.2234	4.2779
$\hat{\beta}_5$	-2.6237	-2.0233	-2.3803	-1.7832	-0.1338	-1.9457	-1.2949	-0.0663	-1.8843
$\hat{\beta}_6$	-2.8491	-4.8424	-4.7643	-1.0767	-4.5903	-3.9601	-0.6508	-4.6289	-3.9573
$k^{opt}$	0.10	0.30	—	0.17	0.43	—	0.43	1.07	—
$\hat{A}DR(\hat{\beta}_i)$	13.391	3.0308	2.2409	28.250	13.213	6.4622	36.533	14.320	6.9432
$Eff(\hat{\beta}_i)$	1.0000	4.4181	5.9755	1.0000	2.1380	4.3716	1.0000	2.5512	5.2617

Table 3. Evaluation of the proposed estimators for  $\gamma = 0.99$

Method Coefficients	$p = 30$			$p = 50$			$p = 70$		
	GRE	GRRE	SRE	GRE	GRRE	SRE	GRE	GRRE	SRE
$\hat{\beta}_1$	-1.0841	-1.3174	-1.3561	-0.9212	-1.4586	-1.5295	-0.9274	-1.2769	-1.4437
$\hat{\beta}_2$	0.2518	0.4380	0.5003	0.1785	0.6240	0.6248	0.1715	0.4523	0.3484
$\hat{\beta}_3$	1.1976	1.8524	1.9290	0.8593	2.1653	2.1419	0.3954	1.5600	1.5314
$\hat{\beta}_4$	2.0972	3.1175	3.2193	1.5161	3.6718	3.6278	0.8316	2.6495	2.6456
$\hat{\beta}_5$	-2.0060	-2.3669	-2.4796	-1.7413	-0.3643	-1.3213	-1.6408	-3.0978	-3.2068
$\hat{\beta}_6$	-1.7052	-3.4797	-3.5661	-0.8283	-4.3776	-3.8814	-0.8940	-2.6663	-2.8464
$k^{opt}$	0.23	0.29	—	0.21	0.49	—	1.01	1.64	—
$\hat{A}DR(\hat{\beta}_i)$	30.241	15.841	13.648	35.054	16.917	13.038	37.868	20.420	18.396
$Eff(\hat{\beta}_i)$	1.0000	1.9090	2.2157	1.0000	2.0722	2.6887	1.0000	1.8545	2.0585

Table 4. Optimal values of shrinkage parameters for different  $p$  and  $\gamma$  values

	$\gamma = 0.75$			$\gamma = 0.90$			$\gamma = 0.99$		
	$p = 30$	$p = 50$	$p = 70$	$p = 30$	$p = 50$	$p = 70$	$p = 30$	$p = 50$	$p = 70$
$\alpha_1^{opt}$	0.1505	0.1660	0.1019	0.2571	1.0863	1.5641	0.2242	0.5922	0.4461
$\alpha_2^{opt}$	0.9015	0.9464	1.0027	0.8182	0.6054	0.6386	0.8660	0.7700	0.8169

**6.2. Application to Riboflavin Production Data**

To illustrate the usefulness of the suggested strategies for high-dimensional data in the semiparametric regression model, the data set about riboflavin (vitamin B2) production is considered in *Bacillus subtilis*, which can be found in R package “hdi”. There is a single real valued response variable which is the logarithm of the

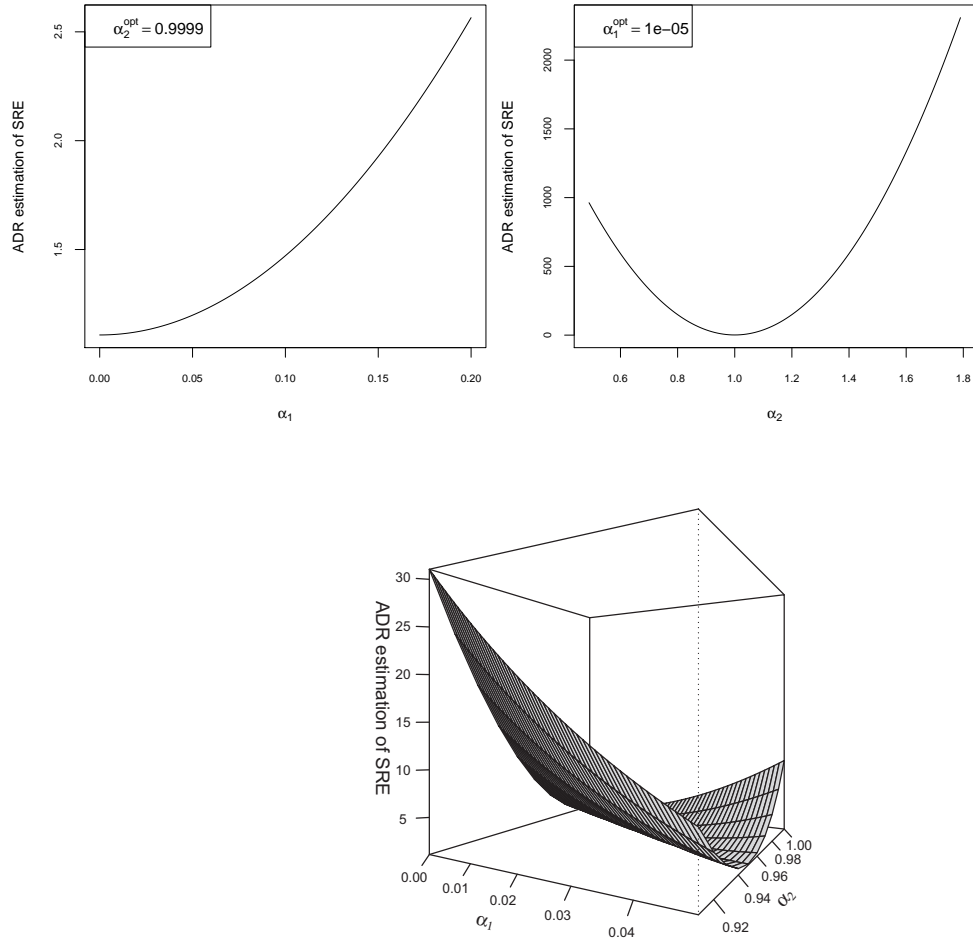


Figure 2. The diagram of *ADR* of SRE versus  $\alpha_1$  and  $\alpha_2$  for real data set.

riboflavin production rate. Furthermore, there are  $p = 4088$  explanatory variables measuring the logarithm of the expression level of 4088 genes. There is one rather homogeneous data set from  $n = 71$  samples that were hybridized repeatedly during a fed batch fermentation process where different engineered strains and strains grown under different fermentation conditions were analyzed. Based on 100-fold cross validation, the Lasso shrinks 4047 parameters to zero and remains  $p_1 = 41$  significant explanatory variables. So, the specification of the sparse regression model is

$$\mathbf{y}_n = \mathbf{X}_n^{(1)}\boldsymbol{\beta}_1 + \mathbf{X}_n^{(2)}\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}, \tag{6.3}$$

where  $p_1 = 41$  and  $p_2 = 4047$ .

According to [23], the unknown  $\mathbf{V}_n$  can be estimated by a consistent estimator

$$\hat{\mathbf{V}}_n = \frac{1}{n - p_1} \left( \mathbf{y}_n - \mathbf{X}_n^{(1)}\mathbf{b}_n^{R1} \right) \left( \mathbf{y}_n - \mathbf{X}_n^{(1)}\mathbf{b}_n^{R1} \right)^\top, \tag{6.4}$$

where  $\mathbf{b}_n^{R1} = \left( \mathbf{X}_n^{(1)\top} \mathbf{X}_n^{(1)} \right)^{-1} \mathbf{X}_n^{(1)\top} \mathbf{y}_n$  is ordinary least square estimator for the sparse regression model.

Table 5 shows a summary of the results. In this Table, the *RSS*, *MSE* and  $R^2$  respectively are the residual sum of squares, mean square error and coefficient of determination of the model, i.e.,  $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ ,  $1/(n - p_1) \sum_{i=1}^n (y_i - \hat{y}_i)^2$  and  $R^2 = 1 - RSS/S_{yy}$ , where  $\hat{y}_i = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$ .

Table 5. Evaluation of proposed estimators for real data set

Method	GRE	GRRE	SRE
$k^{opt}$	3.4241	7.6132	—
$RSS$	39.39465	1.108829	1.107424
$MSE$	1.313155	0.036961	0.036914
$R^2$	0.335704	0.981302	0.981326

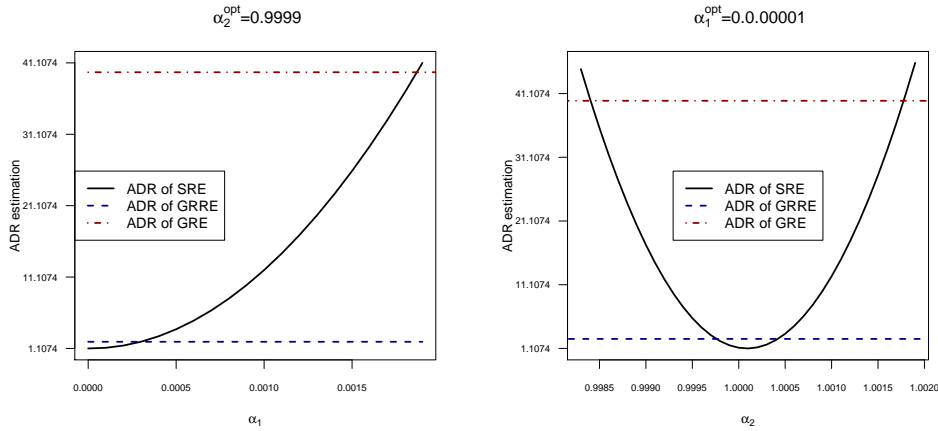


Figure 3. The diagram of ADRs versus  $\alpha_1$  and  $\alpha_2$  for real data set.

According to Figure 2, the minimum of ADR approximately occurred at  $\alpha_1^{opt} = 0.00001$  and  $\alpha_2^{opt} = 0.9999$ . Also, the ADR functions of proposed estimators versus  $\alpha_1$  and  $\alpha_2$  are plotted in Figure 3, which can be used in a similar fashion as in the simulation.

### 7. Summary

In this paper, the generalized ridge estimators were proposed for the sparse multiple regression model. In this context, a generalized restricted ridge estimator exhibited for the sparsity test  $\beta_2 = \mathbf{0}$ , where the regression vector-parameter  $\beta$  partitioned as  $\beta = (\beta_1^T, \beta_2^T)^T$ . Information contained in the generalized restricted ridge estimator was employed to construct a shrinkage ridge estimator. Some asymptotic distributional results were given for the proposed estimators and superiority conditions discussed under some regularity conditions. All the analysis were conducted under a classical view point, using the Tikhonov regularization.

It is well-known that as the sample size tends to infinity, both Bayesian and classical analysis give almost same results. However, under a finite sample size, one may consider a generalized Tikhonov regularization and derive ridge regression estimators in Bayesian context. To see this, take  $\Sigma$ ,  $\beta_o$  and  $\Omega$  as the covariance of  $y_n$ , expected value of  $\beta$  and covariance of  $\beta$ , respectively. Therefore, the generalized ridge estimator has form

$$\hat{\beta} = (\mathbf{X}_n^T \Sigma^{-1} \mathbf{X}_n + \Omega^{-1})^{-1} (\mathbf{X}_n^T \Sigma^{-1} \mathbf{y}_n + \Omega^{-1} \beta_o). \tag{7.5}$$

Then, generalized restricted ridge estimator can be verified for the sparse (restricted) model.

Results of the Monté-Carlo simulation for different  $\gamma$ ,  $n = 25, 30, 50$  and  $70$ , were presented in Tables 1 to 4 and Figure 1. From these Tables, it is realized that  $\hat{\beta}_n^S(\mathbf{K}_n^{(p_1)})$  is leading to be the best estimator among others, since it offers smaller risk and bigger efficiency value, in all cases. According to Figure 1, results show that the global minimum occurs under the ADR of GRE and GRRE. Also, there exists an optimal region for performance of SRE with respect to GRE and GRRE. For the real example, from Table 5 as well as Figures 2 and 3, it can be deduced that  $\hat{\beta}_n^S(\mathbf{K}_n^{(p_1)})$  is quite efficient in the sense that it has significant goodness of fit value.

## Acknowledgments

The author's research is supported in part by Semnan University.

## REFERENCES

1. F. Akdeniz, and G. Tabakan, *Restricted ridge estimators of the parameters in semiparametric regression model*, Communication in Statistics-Theory and Methods, vol. 38, no. 11, pp. 1852–1869, 2009.
2. E. Akdeniz Duran, W.K. Härdle, M. Osipenko, *Difference based ridge and Liu type estimators in semiparametric regression models*, Journal of Multivariate Analysis, vol. 105, no. 1, pp. 164–175, 2012.
3. F. Akdeniz, and M. Roozbeh, *Generalized difference-based weighted mixed almost unbiased ridge estimator in partially linear models*, Statistical Papers, DOI: 10.1007/s00362-017-0893-9, 2017.
4. M. Amini, and M. Roozbeh, *Optimal partial ridge estimation in restricted semiparametric regression models*, Journal of Multivariate Analysis, vol. 136, no. C, pp. 26–40, 2015.
5. P. Bühlmann, and S. van de Geer, *Statistics for High-dimensional Data: Methods, Theory and Applications*, Springer, Heidelberg, 2011.
6. P. Bühlmann, M. Kalisch, and L. Meier, *High-dimensional statistics with a view towards applications in biology*, Annual Review of Statistics and Its Application, vol. 1, pp. 255–278, 2014.
7. N. El Karoui, *Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results*, Arxiv: 1311.2445, 2016.
8. D.G. Gibbons, *A simulation study of some ridge estimators*, Journal of the American Statistical Association, vol. 76, no. 373, pp. 131–139, 1981.
9. T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, second edition, Springer, New York, 2009.
10. W.J. Hemmerle, and T.F. Brantle, *Explicit and Constrained Generalized Ridge Estimation*, Technometrics, vol. 20, pp. 109–119, 1978.
11. W.J. Hemmerle, and M.B. Carey, *Some Properties of Generalized Ridge Estimators*, Communications in Statistics-Simulation and Computation, vol. 12, pp. 239–253, 1983.
12. A.E. Hoerl, and R.W. Kennard, *Ridge regression: biased estimation for non-orthogonal problems*, Technometrics, vol. 12, no. 1, pp. 55–67, 1970.
13. R.R. Hocking, F.M., Speed, and M.J. Lynn, *A Class of Biased Estimators in Linear Regression*, Technometrics, vol. 18, pp. 425–437, 1976.
14. B.M.G. Kibria, (2003). *Performance of some new ridge regression estimators*, Communications in Statistics-Simulation and Computation, vol. 32, pp.419–435.
15. B.M.G. Kibria, and S. Banik, *Some Ridge Regression Estimators and Their Performances*, Journal of Modern Applied Statistical Methods, vol. 15, no. 1, pp. 206–238, 2016.
16. B.M.G. Kibria, and A.K.Md.E. Saleh, *Improving the Estimators of the Parameters of a Probit Regression Model: A Ridge Regression Approach*, Journal of Statistical Planning and Inference, vol. 142, no. 6, pp. 1421–1435, 2011.
17. K. Knight, and W. Fu, *Asymptotic for lasso-type estimators*, Annals of Statistics, vol. 28, no. 5, pp. 1356–1378, 2000.
18. G.C. McDonald, and D.I. Galarneau, *A monte carlo evaluation of some ridge-type estimators*, Journal of the American Statistical Association, vol.70, no. 350, pp. 407–416, 1975.
19. M. Roozbeh, (2015). *Shrinkage ridge estimators in semiparametric regression models*, Journal of Multivariate Analysis, vol. 136, no. C, pp. 56–74, 2015.
20. A.K.Md.E. Saleh, *Theory of Preliminary Test and Stein-type Estimation with Applications*, John Wiley, New York, 2006.
21. A.N. Tikhonov, A. N. (1963). *Solution of incorrectly formulated problems and the regularization method*, Soviet Mathematics vol. 4, no. 4, pp. 1035–1038, 1963.
22. B. Yuzbashi, and S.E. Ahmed, *Shrinkage ridge regression estimators in high-dimensional linear models*, Proceedings of the Ninth International Conference on Management Science and Engineering Management, pp. 793-807, 2015.
23. A. Zellner, *An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias*, Journal of the American Statistical Association, vol. 57, no. 298, pp. 348–368, 1962.