



**Empirical Bayes Regression Analysis with Many
Regressors but Fewer Observations**

by

**Muni S. Srivastava
Department of Statistics
University of Toronto**

and

**Tatsuya Kubokawa
Faculty of Economics
University of Tokyo**

Technical Report No. 0410 September 15, 2004

TECHNICAL REPORT SERIES

**University of Toronto
Department of Statistics**

Empirical Bayes Regression Analysis with Many Regressors but Fewer Observations

Muni S. Srivastava* and Tatsuya Kubokawa†
University of Toronto and University of Tokyo

December 20, 2005

In this paper, we consider the prediction problem in multiple linear regression model in which the number of predictor variables, p , is extremely large compared to the number of available observations, n . The least squares predictor based on a generalized inverse is not efficient. We propose six empirical Bayes estimators of the regression parameters. Three of them are shown to have uniformly lower prediction error than the least squares predictors when the vector of regressor variables are assumed to be random with mean vector zero and the covariance matrix $(1/n)\mathbf{X}^t\mathbf{X}$ where $\mathbf{X}^t = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ is the $p \times n$ matrix of observations on the regressor vector centered from their sample means. For other estimators, we use simulation to show its superiority over the least squares predictor.

Key words and phrases: Empirical Bayes method, James-Stein estimator, linear regression model, many regressors, prediction, ridge regression, shrinkage, stability.

1 Introduction

In this paper, we consider the prediction problem in multiple linear regression model in which the number of predictor variables is extremely large compared to the number of observations available. More specifically, we consider the problem in which the n response variables y_1, \dots, y_n are linearly related to the p regressor variables as

$$\mathbf{y} = \beta_0 \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \tag{1.1}$$

where $\mathbf{y} = (y_1, \dots, y_n)^t$, \mathbf{X} is the $n \times p$ observation matrix on the p predictor variables x_1, \dots, x_p , centered from their sample means, $\mathbf{1}_n = (1, \dots, 1)^t$, an n -vector of ones, β_0 is an unknown constant, $\boldsymbol{\beta}$ is the p -vector of the regression parameters β_1, \dots, β_p ,

*Department of Statistics, University of Toronto, 100 St George Street, Toronto, Ontario, CANADA M5S 3G3, E-Mail: srivasta@utstat.utstat.toronto.edu

†Faculty of Economics, University of Tokyo, Hongo, Bunkyo-ku, Tokyo 113-0033, JAPAN, E-Mail: tatsuya@e.u-tokyo.ac.jpFaculty

$\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^t$, and $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^t$ is the vector of errors. We shall assume that given $\boldsymbol{\beta}$, \mathbf{y} has a multivariate normal distribution with mean vector $\beta_0 \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta}$, and covariance matrix $\sigma^2 \mathbf{I}_n$. We shall write this as

$$\mathbf{y}|\boldsymbol{\beta} \sim \mathcal{N}_n(\beta_0 \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta}, \mathbf{I}_n), \quad (1.2)$$

where \mathbf{I}_n is the $n \times n$ identity matrix. When there is no confusion, we may not always mention that the distribution of \mathbf{y} or $\boldsymbol{\epsilon}$ is conditional given $\boldsymbol{\beta}$. Often, $n > p$, and this problem has been extensively discussed in the literature.

In this paper, we consider the case when n is substantially smaller than p , namely $n \ll p$. From the singular value decomposition of \mathbf{X} , we can write

$$\mathbf{X}^t = \mathbf{A}\mathbf{D}\mathbf{F}, \mathbf{A}^t \mathbf{A} = \mathbf{I}_r, \mathbf{F}\mathbf{F}^t = \mathbf{I}_r, \quad (1.3)$$

where \mathbf{A} is the $p \times r$ matrix of r eigenvectors corresponding to the r nonzero eigenvalues of $\mathbf{X}^t \mathbf{X}$, $\mathbf{D} = \text{diag}(d_1, \dots, d_r)$, $d_1 \geq \dots \geq d_r$, is the diagonal matrix of the positive square roots of the r nonzero ordered eigenvalues of $\mathbf{X}^t \mathbf{X}$, and \mathbf{F} is the $n \times r$ matrix of the r eigenvectors corresponding to the r nonzero eigenvalues of the $n \times n$ matrix $\mathbf{X}\mathbf{X}^t$ which we assume to be of rank $r \leq n - 1$. The nonzero eigenvalues of $\mathbf{X}\mathbf{X}^t$ are the same as the nonzero eigenvalues of $\mathbf{X}^t \mathbf{X}$. The model (1.1) can be thus be rewritten as

$$\begin{aligned} \mathbf{y} &= \beta_0 \mathbf{1}_n + \mathbf{F}^t \mathbf{D} \mathbf{A}^t \boldsymbol{\beta} + \boldsymbol{\epsilon} \\ &= \beta_0 \mathbf{1}_n + \mathbf{F}^t \mathbf{D} \boldsymbol{\gamma} + \boldsymbol{\epsilon}, \end{aligned} \quad (1.4)$$

where

$$\boldsymbol{\gamma} = \mathbf{A}^t \boldsymbol{\beta}.$$

The least squares estimators of β_0 and $\boldsymbol{\gamma}$ are given by

$$\begin{aligned} \widehat{\beta}_0 &= \bar{y} = n^{-1} \mathbf{1}'_n \mathbf{y}, \\ \widehat{\boldsymbol{\gamma}} &= (\mathbf{D}\mathbf{F}\mathbf{F}^t \mathbf{D})^{-1} \mathbf{D}\mathbf{F}\mathbf{y} = \mathbf{D}^{-1} \mathbf{F}\mathbf{y}. \end{aligned}$$

Using the Moore-Penrose inverse (see the Appendix for the definition), we find that the least squares estimator of $\boldsymbol{\beta}$ is given by

$$\widehat{\boldsymbol{\beta}} = (\mathbf{A}\mathbf{D}^2 \mathbf{A}^t)^+ \mathbf{A}\mathbf{D}\mathbf{F}\mathbf{y} = \mathbf{A}\mathbf{D}^{-1} \mathbf{F}\mathbf{y} = \mathbf{A}\widehat{\boldsymbol{\gamma}}. \quad (1.5)$$

Let (\mathbf{A}, \mathbf{B}) be a $p \times p$ orthogonal matrix and $\boldsymbol{\delta} = \mathbf{B}^t \boldsymbol{\beta}$, where \mathbf{A} : $p \times r$ and \mathbf{B} : $p \times (p - r)$. Then, since $\mathbf{A}\mathbf{A}^t + \mathbf{B}\mathbf{B}^t = \mathbf{I}_p$, it follows that

$$\boldsymbol{\beta} = \mathbf{A}\boldsymbol{\gamma} + \mathbf{B}\boldsymbol{\delta}.$$

This is a one-to-one and onto transformation. See Srivastava and Khatri (1979, p.139) for more discussion. Thus, $\widehat{\boldsymbol{\beta}}$ is not an unbiased estimator of $\boldsymbol{\beta}$. The mean of $\widehat{\boldsymbol{\beta}}$ is only $\mathbf{A}\boldsymbol{\gamma}$.

The reader is reminded that since the matrix \mathbf{X}^t is centered from their means, $\mathbf{X}^t \mathbf{1}_n = \mathbf{0}$, and hence $\mathbf{F} \mathbf{1}_n = \mathbf{0}$. It may be noted that

$$\begin{aligned}\widehat{\boldsymbol{\gamma}} &\sim \mathcal{N}_r(\boldsymbol{\gamma}, \sigma^2 \mathbf{D}^{-2}), \\ \widehat{\boldsymbol{\beta}} &\sim \mathcal{N}_p(\mathbf{A}\boldsymbol{\gamma}, \sigma^2 \mathbf{A} \mathbf{D}^{-2} \mathbf{A}^t),\end{aligned}$$

which is a singular multivariate normal distribution.

The above discussion suggests to consider other biased estimators of $\boldsymbol{\beta}$. In particular, the ones obtained by empirical Bayes methods which includes Stein type estimators. The need of considering biased estimators also arises from the fact that the number of nonzero eigenvalues of $\mathbf{X}^t \mathbf{X}$ assumed r above, will usually be not known. Thus, the dimension of the subspace in which the mean of $\widehat{\boldsymbol{\beta}}$ lies will not be known. However, we can begin with a judicious choice of r and further reduce it by considering the prior distribution of $\boldsymbol{\beta}$ which has a mean lying in a reduced subspace as considered below. The initial value of r can be obtained by deleting the eigenvalues of $\mathbf{X}^t \mathbf{X}$ which are zeros or near zeros. For example, in the data analyzed by West (2003), and Brown *et al.* (1999) in which there are 300 regressors and only 78 observations, it was found that only sixteen d_i^2 's account for 99.995% of the variation. That is,

$$\sum_{i=1}^{16} d_i^2 / \sum_{i=1}^{77} d_i^2 = 99.995\%.$$

Thus, it may be reasonable to initially assume that at most r , $r < n$, d_i 's are different from zero and then search for a most appropriate number of d_i that may be taken to be different from zero through the analysis of the data. This is achieved by assuming that $\boldsymbol{\beta} \sim \mathcal{N}_p(\mathbf{A}_1 \boldsymbol{\gamma}_1, \sigma^2 \mathbf{A} \mathbf{T} \mathbf{A}^t)$, where \mathbf{T} is positive definite,

$$\begin{aligned}\mathbf{A} &= (\mathbf{A}_1, \mathbf{A}_2), \quad \boldsymbol{\gamma}^t = (\boldsymbol{\gamma}_1^t, \boldsymbol{\gamma}_2^t), \\ \mathbf{A}_1 &: p \times m, \quad \mathbf{A}_2 : p \times m_1, \boldsymbol{\gamma}_1 : m \times 1, \quad \boldsymbol{\gamma}_2 : m_1 \times 1, \\ \mathbf{A}_1^t \mathbf{A}_1 &= \mathbf{I}_m, \quad \mathbf{A}_2^t \mathbf{A}_2 = \mathbf{I}_{m_1}, \quad \text{and } m_1 = r - m.\end{aligned}\tag{1.6}$$

Although \mathbf{T} can be any $r \times r$ positive definite matrix, we will only consider the case where

$$\mathbf{T} = \lambda \mathbf{D}^{c-2}, \quad \lambda > 0, \quad c = 0, \quad \text{or } c = 2.\tag{1.7}$$

The Bayes estimator or the posterior mean of $\boldsymbol{\beta}$ will involve $\boldsymbol{\gamma}_1$ and λ . These quantities are estimated by using the marginal distribution of $\widehat{\boldsymbol{\beta}}$ and an estimate of σ^2 . Such an estimator, known as empirical Bayes estimator is given by

$$\widehat{\boldsymbol{\beta}}_c^{EB}(\widehat{\lambda}, \widehat{\boldsymbol{\gamma}}_1) = \mathbf{A} \widehat{\boldsymbol{\gamma}} - \left[\mathbf{I}_p - \mathbf{A} (\mathbf{I}_r + \widehat{\lambda}^{-1} \mathbf{D}^{-c})^{-1} \mathbf{A}^t \right] (\mathbf{A} \widehat{\boldsymbol{\gamma}} - \mathbf{A}_1 \widehat{\boldsymbol{\gamma}}_1),\tag{1.8}$$

where $\widehat{\boldsymbol{\gamma}} = \mathbf{A}^t \widehat{\boldsymbol{\beta}}$, $\widehat{\boldsymbol{\gamma}}_1 = \mathbf{D}_1^{-1} \mathbf{F}_1 \mathbf{y}$, $\mathbf{F}^t = (\mathbf{F}_1^t, \mathbf{F}_2^t)$ and $\mathbf{D} = \text{diag}(\mathbf{D}_1, \mathbf{D}_2)$. Here, $\mathbf{F}_1^t : n \times m$, $\mathbf{F}_1 \mathbf{F}_1^t = \mathbf{I}_m$, $\mathbf{F}_2^t : n \times m_1$, $\mathbf{F}_2 \mathbf{F}_2^t = \mathbf{I}_{m_1}$, $m_1 = r - m$, $\mathbf{D}_1 = \text{diag}(d_1, \dots, d_m)$, and

$\mathbf{D}_2 = \text{diag}(d_{m+1}, \dots, d_r)$. By choosing $c = 0$ or $c = 2$ and using different estimators of λ , we obtain several empirical Bayes estimators which are described in Section 2. In Section 3, we prove optimality of three of these estimators. In Section 4, we reanalyze the data of Osborne, Feran, Miller and Douglas (1984). We also present some simulation results in this section. The paper concludes in Section 5.

2 Empirical Bayes Estimators

In this section, we present our analysis of the linear regression model (1.1) or equivalently (1.4) in terms of singular value decomposition of \mathbf{X} . The least squares estimate of $\boldsymbol{\beta}$ given in (1.5) can be written in terms of a singular value decomposition of the matrix \mathbf{X} as

$$\widehat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{D}^{-1}\mathbf{F}\mathbf{y}$$

which is distributed as multivariate normal with the mean vector $\mathbf{A}\mathbf{A}^t\boldsymbol{\beta}$, and the covariance matrix $\sigma^2\mathbf{A}\mathbf{D}^{-2}\mathbf{A}^t$. We shall assume that the prior distribution of $\boldsymbol{\beta}$ is a singular multivariate normal distribution with the mean vector $\mathbf{A}_1\boldsymbol{\gamma}_1$ and covariance matrix $\sigma^2\mathbf{A}\mathbf{T}\mathbf{A}^t$, that is,

$$\pi(\boldsymbol{\beta} | \mathbf{T}, \boldsymbol{\gamma}_1) \sim \mathcal{N}_p(\mathbf{A}_1\boldsymbol{\gamma}_1, \sigma^2\mathbf{A}\mathbf{T}\mathbf{A}^t),$$

where \mathbf{T} is an $r \times r$ positive definite matrix and \mathbf{A}_1 is the $p \times m$ matrix defined in (1.6). Noting that

$$\mathbf{A}\mathbf{A}^t\mathbf{A}_1 = (\mathbf{A}_1, \mathbf{A}_2) \begin{pmatrix} \mathbf{A}_1^t \\ \mathbf{A}_2^t \end{pmatrix} \mathbf{A}_1 = (\mathbf{A}_1, \mathbf{A}_2) \begin{pmatrix} \mathbf{I}_m \\ \mathbf{0} \end{pmatrix} = \mathbf{A}_1,$$

we find from the above prior distribution that the joint distribution of $(\widehat{\boldsymbol{\beta}}^t, \boldsymbol{\beta}^t)^t$ is given by

$$\begin{pmatrix} \widehat{\boldsymbol{\beta}} \\ \boldsymbol{\beta} \end{pmatrix} \sim \mathcal{N}_{2p} \left(\begin{pmatrix} \mathbf{A}_1\boldsymbol{\gamma}_1 \\ \mathbf{A}_1\boldsymbol{\gamma}_1 \end{pmatrix}, \sigma^2 \begin{pmatrix} \mathbf{A}(\mathbf{T} + \mathbf{D}^{-2})\mathbf{A}^t & \mathbf{A}\mathbf{T}\mathbf{A}^t \\ \mathbf{A}\mathbf{T}\mathbf{A}^t & \mathbf{A}\mathbf{T}\mathbf{A}^t \end{pmatrix} \right).$$

Hence, from Theorem 2.2.7 of Srivastava and Khatri (1979, page 47), the conditional distribution of $\boldsymbol{\beta}$ given $\widehat{\boldsymbol{\beta}}$ or the posterior distribution of $\boldsymbol{\beta}$ is $\mathcal{N}_p(\widehat{\boldsymbol{\beta}}^B(\mathbf{G}, \boldsymbol{\gamma}_1), \boldsymbol{\Lambda})$, where

$$\begin{aligned} \widehat{\boldsymbol{\beta}}^B(\mathbf{G}, \boldsymbol{\gamma}_1) &= [\mathbf{I} - \mathbf{A}\mathbf{G}\mathbf{A}^t]\mathbf{A}_1\boldsymbol{\gamma}_1 + \mathbf{A}\mathbf{G}\widehat{\boldsymbol{\gamma}}, \\ \mathbf{G} &= \mathbf{T}(\mathbf{T} + \mathbf{D}^{-2})^{-1}, \\ \widehat{\boldsymbol{\gamma}} &= \mathbf{A}^t\widehat{\boldsymbol{\beta}} = \mathbf{D}^{-1}\mathbf{F}\mathbf{y}, \end{aligned} \tag{2.1}$$

and

$$\boldsymbol{\Lambda} = \sigma^2\mathbf{A}\mathbf{T}(\mathbf{T} + \mathbf{D}^{-2})^{-1}\mathbf{T}\mathbf{A}^t.$$

When $\boldsymbol{\beta} = \mathbf{A}_1\boldsymbol{\gamma}_1$, the model becomes

$$\begin{aligned}\mathbf{y} &= \mathbf{F}^t \mathbf{D} \mathbf{A}^t \mathbf{A}_1 \boldsymbol{\gamma}_1 + \boldsymbol{\epsilon} \\ &= \mathbf{F}_1^t \mathbf{D}_1 \boldsymbol{\gamma}_1 + \boldsymbol{\epsilon},\end{aligned}$$

and the least squares estimate of $\boldsymbol{\gamma}_1$ is given by

$$\hat{\boldsymbol{\gamma}}_1 = \mathbf{D}_1^{-1} \mathbf{F}_1 \mathbf{y}, \quad (2.2)$$

which can be interpreted as a principal component estimator. Hence, an empirical Bayes estimator of $\boldsymbol{\beta}$ is given by

$$\begin{aligned}\hat{\boldsymbol{\beta}}^B(\mathbf{G}, \hat{\boldsymbol{\gamma}}_1) &= [\mathbf{I} - \mathbf{A} \mathbf{G} \mathbf{A}^t] \mathbf{A}_1 \hat{\boldsymbol{\gamma}}_1 + \mathbf{A} \mathbf{G} \hat{\boldsymbol{\gamma}} \\ &= \mathbf{A} \hat{\boldsymbol{\gamma}} - (\mathbf{I} - \mathbf{A} \mathbf{G} \mathbf{A}^t) (\mathbf{A} \hat{\boldsymbol{\gamma}} - \mathbf{A}_1 \hat{\boldsymbol{\gamma}}_1)\end{aligned} \quad (2.3)$$

where $\hat{\boldsymbol{\gamma}}_1$ is given by (2.2) and $\hat{\boldsymbol{\gamma}}$ by (2.1). The estimator $\hat{\boldsymbol{\beta}}^B(\mathbf{G}, \hat{\boldsymbol{\gamma}}_1)$ in (2.3) shrinks the LS estimator $\mathbf{A} \hat{\boldsymbol{\gamma}}$ towards the principal component regression (PCR) estimator $\mathbf{A}_1 \hat{\boldsymbol{\gamma}}_1$. However in (2.3), we need to specify \mathbf{G} or equivalently \mathbf{T} .

We choose \mathbf{T} as given in the equation (1.7). This gives

$$\mathbf{G} = (\mathbf{I}_r + \lambda^{-1} \mathbf{D}^{-c})^{-1}.$$

Thus, our empirical Bayes estimator depends only on the estimators of $(\lambda, \boldsymbol{\gamma}_1)$, denoted by $(\hat{\lambda}, \hat{\boldsymbol{\gamma}}_1)$, and is given by

$$\hat{\boldsymbol{\beta}}_c^{EB}(\hat{\lambda}, \hat{\boldsymbol{\gamma}}_1) = \mathbf{A} \hat{\boldsymbol{\gamma}} - \left[\mathbf{I}_p - \mathbf{A} (\mathbf{I}_r + \hat{\lambda}^{-1} \mathbf{D}^{-c})^{-1} \mathbf{A}^t \right] (\mathbf{A} \hat{\boldsymbol{\gamma}} - \mathbf{A}_1 \hat{\boldsymbol{\gamma}}_1). \quad (2.4)$$

From the least squares theory, we choose an estimator of σ^2 given by

$$\hat{\sigma}^2 = \|\mathbf{y} - \bar{y}\mathbf{1} - \mathbf{F}^t \mathbf{D} \mathbf{A}^t \hat{\boldsymbol{\beta}}\|^2 / (n - r + 1), \quad (2.5)$$

provided $r < n$. To obtain an estimator of λ , we consider the marginal distribution of $\hat{\boldsymbol{\beta}}$ which is a singular distribution. Thus, we shall consider the distribution of $\mathbf{A}^t \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\gamma}}$, which is nonsingular, given by

$$\hat{\boldsymbol{\gamma}} \sim \mathcal{N}_r \left(\begin{pmatrix} \boldsymbol{\gamma}_1 \\ \mathbf{0} \end{pmatrix}, \sigma^2 (\lambda \mathbf{D}^{c-2} + \mathbf{D}^{-2}) \right).$$

Thus if $\boldsymbol{\gamma}_1$ and σ^2 are known, the moment estimator of λ is obtained by solving

$$\left(\hat{\boldsymbol{\gamma}} - \begin{pmatrix} \boldsymbol{\gamma}_1 \\ \mathbf{0} \end{pmatrix} \right)^t (\lambda \mathbf{D}^{c-2} + \mathbf{D}^{-2})^{-1} \left(\hat{\boldsymbol{\gamma}} - \begin{pmatrix} \boldsymbol{\gamma}_1 \\ \mathbf{0} \end{pmatrix} \right)^t = a \sigma^2 \quad (2.6)$$

where a is chosen appropriately to get biased or unbiased estimator. For σ^2 and $\boldsymbol{\gamma}_1$, we use the estimators given by (2.5) and (2.2), respectively. Thus, to get an estimate of λ , we solve the equation

$$\widehat{\boldsymbol{\gamma}}_2^t [\lambda \mathbf{D}_2^{c-2} + \mathbf{D}_2^{-2}]^{-1} \widehat{\boldsymbol{\gamma}}_2 = a \hat{\sigma}^2, \quad (2.7)$$

where $\widehat{\boldsymbol{\gamma}}_2$ is an m_1 -vector, $m_1 = r - m$, given by

$$\widehat{\boldsymbol{\gamma}}_2 = \mathbf{D}_2^{-1} \mathbf{F}_2 \mathbf{y}.$$

Let λ^* be a solution of (2.6). In order that it remains positive, we choose

$$\hat{\lambda} = \max(\lambda^*, 0)$$

as an estimator of λ .

Next, we describe several estimators by choosing the values of a , c etc. While we choose only two values of c , namely $c = 0$ or 2 , we choose five values of a , given by

$$\begin{aligned} a_1 &= a_2 = r - m - 2, \\ a_3 &= 2 \sum_{i=m+1}^r (d_r^2/d_i^2) - 4, \\ a_4 &= r - 2 \quad \text{and} \quad a_5 = a_3. \end{aligned}$$

2.1 James-Stein type estimator (JS)

Let $c = 0$, or $\mathbf{T} = \lambda \mathbf{D}^{-2}$, and $a = a_1 = r - m - 2$. Then from (2.6), we get an estimator of λ or equivalently of α given by

$$\hat{\alpha} \equiv \frac{1}{1 + \hat{\lambda}} = \min \left\{ \frac{(r - m - 2) \hat{\sigma}^2}{\widehat{\boldsymbol{\gamma}}_2^t \mathbf{D}_2^2 \widehat{\boldsymbol{\gamma}}_2}, 1 \right\}.$$

The estimator of λ obtained from the above equation will be denoted by $\hat{\lambda}_1$. Thus,

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_1^{EB}(\hat{\lambda}_1, \widehat{\boldsymbol{\gamma}}_1) &= \widehat{\boldsymbol{\beta}}^{JS} = \mathbf{A} \widehat{\boldsymbol{\gamma}} - \hat{\alpha} \{ \mathbf{A} \widehat{\boldsymbol{\gamma}} - \mathbf{A}_1 \widehat{\boldsymbol{\gamma}}_1 \} \\ &= \mathbf{A}_1 \widehat{\boldsymbol{\gamma}}_1 + (1 - \hat{\alpha}) \mathbf{A}_2 \widehat{\boldsymbol{\gamma}}_2. \end{aligned} \quad (2.8)$$

Since it is identical to the so-called positive part Stein estimator proposed by James and Stein (1961) for $p = r$, we will call it the James-Stein type estimator (JS).

2.2 Ridge-principal component estimator (RP)

Let $c = 2$, $m_1 = r - m$, $a = a_2 = r - m - 2$ and $\hat{\lambda}_2 = \max(0, \hat{\lambda})$, where $\hat{\lambda}$ is the solution of

$$\widehat{\boldsymbol{\gamma}}_2^t \left(\hat{\lambda} \mathbf{I}_{m_1} + \mathbf{D}_2^{-2} \right)^{-1} \widehat{\boldsymbol{\gamma}}_2 = (r - m - 2) \hat{\sigma}^2. \quad (2.9)$$

Then, we obtain an estimator of $\widehat{\boldsymbol{\beta}}$, which we call the ridge-principal component regression estimator (RP), given by

$$\widehat{\boldsymbol{\beta}}_2^{EB}(\widehat{\lambda}_2, \widehat{\boldsymbol{\gamma}}_1) = \widehat{\boldsymbol{\beta}}^{RP} = \mathbf{A}\widehat{\boldsymbol{\gamma}} - (\mathbf{I}_p - \mathbf{A}(\mathbf{I}_r + \widehat{\lambda}_2^{-1}\mathbf{D}^{-2})^{-1}\mathbf{A}^t)(\mathbf{A}\widehat{\boldsymbol{\gamma}} - \mathbf{A}_1\widehat{\boldsymbol{\gamma}}_1). \quad (2.10)$$

This estimator is similar to the one proposed by Kubokawa and Srivastava (2004). However, it is difficult to verify analytically the superiority of the $\widehat{\boldsymbol{\beta}}^{RP}$ estimator over the least squares estimator. Thus, as dictated by the proof of the dominance result of Theorem 3.1 given in Section 3, we modify the equation (2.9) to obtain the solution of λ in the next subsection.

2.3 Modified ridge-regression principal component estimator (RPM)

Let $c = 2$, $m_1 = r - m$, and $a = a_3 = 2 \sum_{i=m+1}^r (d_r^2/d_i^2) - 4$. Define

$$\widehat{\lambda}_3 = \max(0, \widehat{\lambda})$$

where $\widehat{\lambda}$ is the solution of the equation

$$\widehat{\boldsymbol{\gamma}}_2^t \left(\widehat{\lambda} \mathbf{I}_{m_1} + \mathbf{D}_2^{-2} \right)^{-1} \widehat{\boldsymbol{\gamma}}_2 = a_3 \widehat{\sigma}^2.$$

Then the modified ridge-regression principal component estimator is given by

$$\widehat{\boldsymbol{\beta}}^{RPM} = \widehat{\boldsymbol{\beta}}_2^{EB}(\widehat{\lambda}_3, \widehat{\boldsymbol{\gamma}}_1), \quad (2.11)$$

whose expression can be obtained from (2.10) with $\widehat{\lambda}_2$ replaced by $\widehat{\lambda}_3$.

2.4 Ridge-regression type estimator (RR)

We now consider the case when $c = 2$, $m = 0$, and $a = a_4 = r - 2$. In this case $\boldsymbol{\gamma}_1 = \mathbf{0}$, and the estimate of λ is obtained by solving the equation

$$\widehat{\boldsymbol{\gamma}}^t \left(\widehat{\lambda} \mathbf{I}_r + \mathbf{D}^{-2} \right)^{-1} \widehat{\boldsymbol{\gamma}} = (r - 2) \widehat{\sigma}^2. \quad (2.12)$$

We then estimate λ by

$$\widehat{\lambda}_4 = \max(\widehat{\lambda}, 0).$$

Using $\widehat{\lambda}_4$ as an estimator of λ , we obtain an empirical Bayes estimator of Shinozaki and Chang (1993) type. We shall denote this estimator of $\boldsymbol{\beta}$ by $\widehat{\boldsymbol{\beta}}^{RR}$ to reflect its similarity with ridge regression type of estimators. It is given by

$$\widehat{\boldsymbol{\beta}}^{RR} = \widehat{\boldsymbol{\beta}}_2^{EB}(\widehat{\lambda}_4, \mathbf{0}) = \mathbf{A} \left(\mathbf{I}_r + \widehat{\lambda}_4^{-1} \mathbf{D}^{-2} \right)^{-1} \widehat{\boldsymbol{\gamma}}. \quad (2.13)$$

However, the superiority of $\widehat{\boldsymbol{\beta}}^{RR}$ over the LS estimator $\widehat{\boldsymbol{\beta}}$ has not been established. We therefore modify (2.12) to obtain an estimator which we show to dominate the LS estimator $\widehat{\boldsymbol{\beta}}$.

2.5 Modified ridge-regression type estimator (RRM)

In this section, we assume as in Section 2.4 that $c = 2$, $m = 0$, but $a = a_5 = 2 \sum_{i=m+1}^r (d_r^2/d_i^2) - 4$. Let

$$\hat{\lambda}_5 = \max(\hat{\lambda}, 0)$$

where $\hat{\lambda}$ is the solution of the equation

$$\hat{\gamma}^t (\hat{\lambda} \mathbf{I}_r + \mathbf{D}^{-2})^{-1} \hat{\gamma} = a_5 \hat{\sigma}^2.$$

Then the resulting empirical Bayes estimator of $\boldsymbol{\beta}$ is given by (2.13) with $\hat{\lambda}_4$ replaced by $\hat{\lambda}_5$. Thus,

$$\hat{\boldsymbol{\beta}}^{RRM} = \hat{\boldsymbol{\beta}}_2^{EB}(\hat{\lambda}_5, 0). \quad (2.14)$$

In Section 3, and in the Appendix we show that the estimator $\hat{\boldsymbol{\beta}}^{RRM}$ dominates the least squares estimator $\hat{\boldsymbol{\beta}}$, in terms of prediction error when the vector of p predictor variables has a random distribution described in Section 3.

2.6 Crude empirical Bayes estimator (EB)

The above estimation procedures are available when the rank r of the matrix \mathbf{X} is less than $n - 1$. However, the following empirical Bayes estimator can be utilized even for $r = n - 1$:

$$\hat{\boldsymbol{\beta}}^{EB} = \left(\mathbf{X}^t \mathbf{X} + (1/\hat{\lambda}_{EB}) \mathbf{I}_p \right)^{-1} \mathbf{X}^t \mathbf{y}, \quad (2.15)$$

where $\hat{\lambda}_{EB} = \max(0, \lambda_*)$ for the solution λ_* of the equation

$$\mathbf{y}^t \left(\hat{\lambda}_* \mathbf{X} \mathbf{X}^t + \mathbf{I}_n \right)^{-1} \mathbf{y} / \hat{\sigma}^2 = n - 2.$$

If \mathbf{X} has full rank, no variance estimator is available. However, we suggest the use of the estimator $\hat{\sigma}^2$ given by (2.5) for appropriately chosen r . When $\text{rank}(\mathbf{X}) = r < n - 1$, this empirical Bayes estimator $\hat{\boldsymbol{\beta}}^{EB}$ is close to the above one $\hat{\boldsymbol{\beta}}^{RR}$.

3 Empirical Bayes predictors and their risk properties

Now we consider the problem of predicting the observation y when \mathbf{x} is given as a predictor variable, where the model is described by

$$y = \beta_0 + \mathbf{x}^t \boldsymbol{\beta} + \varepsilon,$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. The usual predictor based on the least squares estimators $\hat{\beta}_0 = \bar{y}$ and $\hat{\beta}$ is given by

$$\hat{y} = \bar{y} + \mathbf{x}^t \hat{\beta}, \quad (3.1)$$

while the shrinkage predictor employing the empirical Bayes estimator of β is provided by

$$\hat{y}_c^{EB}(\hat{\lambda}, \hat{\gamma}_1) = \bar{y} + \mathbf{x}^t \hat{\beta}_c^{EB}(\hat{\lambda}, \hat{\gamma}_1), \quad (3.2)$$

where the empirical Bayes estimator $\hat{\beta}_c^{EB}(\hat{\lambda}, \hat{\gamma}_1)$ is given by (2.4).

To evaluate the predictors, we utilize the arguments as in Copas (1983), who assumes that \mathbf{x} follows the empirical distribution \hat{F}_x based on the sample $\mathbf{x}_1, \dots, \mathbf{x}_n$, which has the mean $E[\mathbf{x}|\hat{F}_x] = 0$ and the covariance matrix $E[\mathbf{x}\mathbf{x}^t|\hat{F}_x] = n^{-1}\mathbf{X}^t\mathbf{X}$, where the notation $E[\cdot|\hat{F}_x]$ denotes the expectation with respect to the distribution \hat{F}_x . When the expectation is taken with respect to ε and \hat{F}_x , the prediction error of the usual predictor \hat{y} under the squared loss is written by

$$\begin{aligned} E^\varepsilon \left[E \left[(\hat{y} - y)^2 | \hat{F}_x \right] \right] &= E^\varepsilon \left[E \left[(\bar{y} + \mathbf{x}^t \hat{\beta} - \beta_0 - \mathbf{x}^t \beta - \varepsilon)^2 | \hat{F}_x \right] \right] \\ &= (\bar{y} - \beta_0)^2 + (\hat{\beta} - \beta)^t E \left[\mathbf{x}\mathbf{x}^t | \hat{F}_x \right] (\hat{\beta} - \beta) + \sigma^2 \\ &\quad + 2(\bar{y} - \beta_0) E \left[\mathbf{x}^t | \hat{F}_x \right] (\hat{\beta} - \beta) \\ &= (\bar{y} - \beta_0)^2 + n^{-1}(\hat{\beta} - \beta)^t \mathbf{X}^t \mathbf{X} (\hat{\beta} - \beta) + \sigma^2. \end{aligned} \quad (3.3)$$

Taking the expectation of (3.3) with respect to all the random variables turns out to be

$$E[(\hat{y} - y)^2] = n^{-1} E \left[(\hat{\beta} - \beta)^t \mathbf{X}^t \mathbf{X} (\hat{\beta} - \beta) \right] + n^{-2} \mathbf{1}_n^t \mathbf{X} \beta \beta^t \mathbf{X}^t \mathbf{1}_n + (n+1)n^{-1}\sigma^2.$$

Since $\mathbf{X}^t \mathbf{1}_n = \mathbf{A} \mathbf{D} \mathbf{F} \mathbf{1}_n = \mathbf{0}$ as noted in Section 1, the prediction error is expressed by

$$E[(\hat{y} - y)^2] = n^{-1} E \left[(\hat{\beta} - \beta)^t \mathbf{X}^t \mathbf{X} (\hat{\beta} - \beta) \right] + (n+1)n^{-1}\sigma^2, \quad (3.4)$$

which implies that the prediction problem is reduced to that of estimating β under the weighted loss $(\hat{\beta} - \beta)^t \mathbf{X}^t \mathbf{X} (\hat{\beta} - \beta) = (\hat{\gamma} - \gamma)^t \mathbf{D}^2 (\hat{\gamma} - \gamma)$.

Next, we state a theorem whose proof is given in the Appendix.

Theorem 3.1 *Let $r - m \geq 3$. Then, under the assumptions that $\mathbf{x} \sim \hat{F}_x$, the empirical Bayes predictor $\hat{y}_c^{EB}(\hat{\lambda}, \hat{\gamma}_1)$ given by (3.2) has a uniformly smaller prediction error than the least square predictor \hat{y} given by (3.1) if the chosen values of $a > 0$, and c satisfy the inequality*

$$\sum_{i=m+1}^r (d_r/d_i)^c - 2 \geq a/2. \quad (3.5)$$

Corollary 3.1 Assume that $\mathbf{x} \sim \widehat{F}_x$.

(1) For $c = 0$ and $a = r - m - 2 > 0$, the inequality (3.5) is satisfied. Hence, the empirical Bayes predictor $\hat{y}^{JS} = \bar{y} + \mathbf{x}^t \widehat{\boldsymbol{\beta}}^{JS}$ based on the James-Stein type estimator (2.8) improves on \hat{y} .

(2) For $c = 2$ and $a = 2(\sum_{i=m+1}^r d_r^2/d_i^2 - 2) > 0$, the inequality (3.5) is satisfied. Hence, the empirical Bayes predictor $\hat{y}^{RPM} = \bar{y} + \mathbf{x}^t \widehat{\boldsymbol{\beta}}^{RPM}$ based on the ridge-principal components regression estimator (2.11) improves on \hat{y} .

(3) For $c = 2$, $m = 0$ and $a = 2(\sum_{i=1}^r d_r^2/d_i^2 - 2) > 0$, the inequality (3.5) is satisfied. Hence, the empirical Bayes predictor $\hat{y}^{RRM} = \bar{y} + \mathbf{x}^t \widehat{\boldsymbol{\beta}}^{RRM}$ based on the ridge type estimator (2.14) improves on \hat{y} .

4 Empirical and Simulation Studies

In this section, we first consider a dataset, initially analyzed by Osborne *et al.* (1984), and later by Brown *et al.* (1999, 2001) and West (2003). The data consist of measurements, on the constituents of the biscuit dough such as fat, sucrose, flour and water, which were obtained using near infrared (NIR) spectroscopy. In our analysis, we shall, however, focus on ‘fat content’ only which will be our dependent variable \mathbf{y} . The predictor variables consist of spectral readings at 700 points measured from 1100 to 2498 nanometers in steps of 2 nm. Thus, $p = 700$. The training or calibration set consist of 40 measurements on \mathbf{y} and 700 predictor variables x_1, \dots, x_{700} , which have been centered from their sample means.

The centered and scaled observations matrix \mathbf{X} of the predictor variables is a huge 40×700 matrix of rank 39 with singularvalues given by

$$\begin{aligned} \mathbf{D}_0 = & (25.295, 6.604, 2.812, 1.953, 1.473, 1.287, 0.626, 0.465, 0.301, 0.252, 0.238, 0.196, \\ & 0.193, 0.138, 0.127, 0.105, 0.098, 0.092, 0.085, 0.075, 0.070, 0.058, 0.057, 0.055, \\ & 0.052, 0.046, 0.045, 0.039, 0.039, 0.034, 0.032, 0.029, 0.029, 0.026, 0.024, 0.022, \\ & 0.019, 0.019, 0.017), \end{aligned}$$

which includes lots of small values. Taking the r largest eigenvalues of the 39, we can calculate their contribution rate

$$CR = CR(r) = 100 \times \sum_{i=1}^r d_i^2 / \sum_{i=1}^{39} d_i^2.$$

For instance, we have $CR(30) = 99.99\%$, $CR(17) = 99.99\%$, $CR(10) = 99.96\%$ and $CR(5) = 99.85\%$, which are given in the column CR in Table 1. The 17 largest eigenvalues account for 99.99% of the variation. Although \mathbf{X} has many small eigenvalues, it is of rank $n - 1 = 39$, and \mathbf{X}^t is decomposed as $\mathbf{X}^t = \mathbf{A}_0 \mathbf{D}_0 \mathbf{F}_0$ where \mathbf{A}_0 is the $p \times (n - 1)$ matrix of $n - 1$ eigenvectors corresponding to the $n - 1$ positive eigenvalues of $\mathbf{X}^t \mathbf{X}$,

Table 1: Prediction Error Estimates by the Cross-validation Method for $p = 700$ and $n = 40$ in the Two Cases of \mathbf{X} Scaled and Non-scaled

Case of both \mathbf{y} and \mathbf{X} Scaled											
r	m	CR	LS	PC	JS	RPM	RRM	RP	RR	EB	\bar{y}
30	6	99.99	3.798	2.022	1.184	1.129	1.197	1.033	1.202	0.677	1.051
20	6	99.99	3.798	1.308	1.025	1.069	1.015	1.047	1.029	0.572	1.051
17	5	99.99	3.798	1.297	0.963	0.951	0.963	0.921	0.904	0.503	1.051
10	5	99.96	3.798	0.924	0.864	0.896	0.919	0.877	1.098	0.678	1.051
6	2	99.85	3.798	0.978	1.072	1.071	1.090	1.064	1.103	0.611	1.051
Case of \mathbf{y} Scaled but \mathbf{X} Non-scaled											
r	m	CR	LS	PC	JS	RPM	RRM	RP	RR	EB	\bar{y}
35	6	99.99	3.642	3.452	1.950	2.092	2.007	1.874	1.661	0.352	1.051
30	6	99.99	3.642	2.832	1.611	1.591	1.416	1.507	1.165	0.419	1.051
20	6	99.99	3.642	1.969	1.409	1.407	1.296	1.384	1.095	0.448	1.051
10	5	99.97	3.642	1.585	1.397	1.379	1.279	1.372	1.106	0.513	1.051
6	2	99.87	3.642	1.321	1.192	1.274	1.249	1.217	1.140	0.531	1.051

$\mathbf{D}_0 = \text{diag}(d_1, \dots, d_{n-1})$, $d_1 \geq \dots \geq d_{n-1}$, is the diagonal matrix of the positive square roots of the $n - 1$ ordered eigenvalues of $\mathbf{X}^t \mathbf{X}$, known as the singularvalues of the matrix \mathbf{X} , and \mathbf{F}_0 is the $(n - 1) \times n$ matrix of the $n - 1$ eigenvectors corresponding to the $n - 1$ nonzero eigenvalues of the $n \times n$ matrix $\mathbf{X} \mathbf{X}^t$. We may note that $\mathbf{A}_0^t \mathbf{A}_0 = \mathbf{I}_{n-1}$, $\mathbf{F}_0 \mathbf{F}_0^t = \mathbf{I}_{n-1}$ and $\mathbf{F}_0 \mathbf{1}_n = \mathbf{0}$. Then, we can calculate the least squares estimator as

$$\tilde{\boldsymbol{\beta}} = \mathbf{A}_0 \mathbf{D}_0^{-1} \mathbf{F}_0 \mathbf{y}, \quad \text{denoted by LS,} \quad (4.1)$$

which is not stable. In the \mathbf{D}_0 described above, many eigenvalues d_i 's are very close to zero. For example, $CR(17) = 99.99\%$, which means that the rank of \mathbf{X} may be practically presumed to be $r = 17$. Since \mathbf{X}^t is decomposed as $\mathbf{X}^t = \mathbf{A} \mathbf{D} \mathbf{F}$ for such a rank r , we can consider the principal component regression estimator

$$\hat{\boldsymbol{\beta}} = \mathbf{A} \mathbf{D}^{-1} \mathbf{F} \mathbf{y}, \quad \text{denoted by PC,}$$

because the $r \times r$ diagonal matrix \mathbf{D} deletes the $n - r - 1$ smallest eigenvalues from the $(n - 1) \times (n - 1)$ matrix \mathbf{D}_0 . Based on the estimator $\hat{\boldsymbol{\beta}}$, we can compute the following estimators suggested in Section 2: $\hat{\boldsymbol{\beta}}^{JS}$, $\hat{\boldsymbol{\beta}}^{RP}$, $\hat{\boldsymbol{\beta}}^{RPM}$, $\hat{\boldsymbol{\beta}}^{RR}$, $\hat{\boldsymbol{\beta}}^{RRM}$ and $\hat{\boldsymbol{\beta}}^{EB}$, which are denoted by JS, RP, RPM, RR, RRM and EB, respectively. These estimators are compared below in terms of the prediction error estimates by the cross-validation method and estimation errors by simulation.

Corresponding to the estimators LS, PC, JS, RPM, RRM, RP, RR and EB of $\boldsymbol{\beta}$, we provide the predictors of the the form (3.2). In this analysis, we add the sample mean \bar{y}

as a simple predictor, which means that β is estimated by zeros. The prediction errors of these methods may be estimated via the leave-one-out cross-validation as described, for example, in Srivastava (2002, p322). That is, 40 predictive errors are obtained by leaving out one observation each time. When both \mathbf{X} and \mathbf{y} are scaled, the estimates of the prediction errors for the above considered predictors are given in the top part of Table 1 for the cases $(r, m) = (30, 6), (20, 6), (17, 5), (10, 5)$ and (6.2). It reveals that the smallest prediction errors of EB and RR are 0.503 and 0.904, respectively, which are provided in the case of $r = 17$. Both minimum values are smaller than those of other predictors. Especially, the use of the estimator *EB* provides the smallest prediction errors among its competitors. The other empirical Bayes predictors give smaller prediction errors for $(r, m) = (10, 5)$ than the LS, PC and \bar{y} predictors.

When the observations matrix \mathbf{X} is centered but not scaled, the $n - 1$ singularvalues of \mathbf{X} are given by

$$\begin{aligned} \mathbf{D}_0 = & (10.487, 2.379, 0.950, 0.818, 0.544, 0.329, 0.261, 0.140, 0.107, 0.096, 0.088, 0.073, \\ & 0.070, 0.054, 0.049, 0.043, 0.037, 0.034, 0.032, 0.029, 0.028, 0.022, 0.020, 0.020, \\ & 0.019, 0.017, 0.016, 0.015, 0.014, 0.013, 0.012, 0.011, 0.010, 0.009, 0.009, 0.008, \\ & 0.008, 0.007, 0.006). \end{aligned}$$

We carry out the same analysis for this non-scaled version as was done for the scaled version. The prediction error estimates are reported at the bottom in Table 1, which illustrates that the EB predictor provides the smallest prediction error 0.352 for $r = 35$, which are smaller than those of the other predictors. The other empirical predictors JS, RPM, RRM, RP and RR are better than LS and PS, but worse than \bar{y} .

In the case of \mathbf{X} scaled but \mathbf{y} non-scaled for $r = 6$ and $m = 2$, the values of $\hat{\lambda}$ in the empirical Bayes estimators JS, RPM, RRM, RP, RR and EB are given by 13.54, 7.80, 6.13, 7.72, 3.55 and 5432.27, and the resulting predicted values $\hat{y}_i^{EB} = \bar{y} + \mathbf{x}_i^t \hat{\beta}_c^{EB}(\hat{\lambda}, \hat{\gamma}_1)$ for some i 's are provided in Table 2. Also predicted values by some predictors are drawn in Figures 1 and 2, where the x-axis indicates the observed number. From the table and the figures, it is found that the predicts based on EB are very close to the real data of the response variable y_i and that the line graphs of EB and the real data y_i overlap each other as seen in Figure 1, while the predicts by LS have a large dispersion. It also seems that the predictors based on the other empirical Bayes estimators take reasonable values.

We next treat the smaller data set of the predictor variables consisting of 300 points which are centered after deleting the first 200 and the last 200 points from the original data set. In this case, the maximum and minimum eigenvalues are, respectively, given by 17.087 and 0.00354 for \mathbf{X} scaled, and 8.0108 and 0.00166 for \mathbf{X} non-scaled. Through the same analysis as stated above, we provide the prediction errors in Table 3. The performance of each predictor is similar to the case of $p = 700$. The prediction error of EB for $r = 19$ is 0.117, and is the smallest in the case of both \mathbf{y} and \mathbf{X} scaled. Brown *et al.* (1999) gives a prediction error of 0.14, with West's (2003) method only slightly improving on it, while Osborne *et al.*'s (1984) method has prediction error 0.11. This

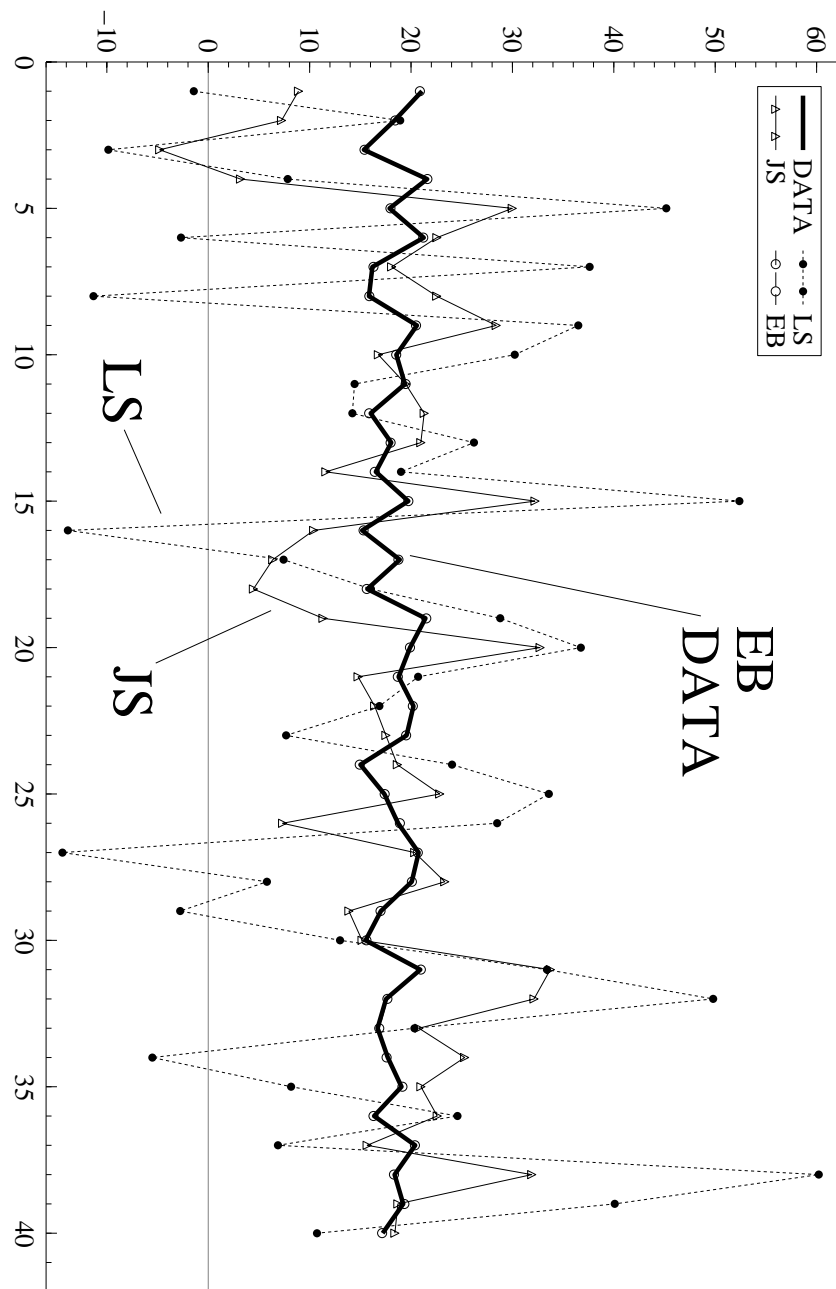


Figure 1: Line Graphs of the Data y_i and Predicts by LS, JS and EB for $p = 700$, $n = 40$, $r = 6$ and $m = 2$ in the Case of \mathbf{X} Scaled but \mathbf{y} Non-scaled (Line graphs of EB and y_i overlap each other.)

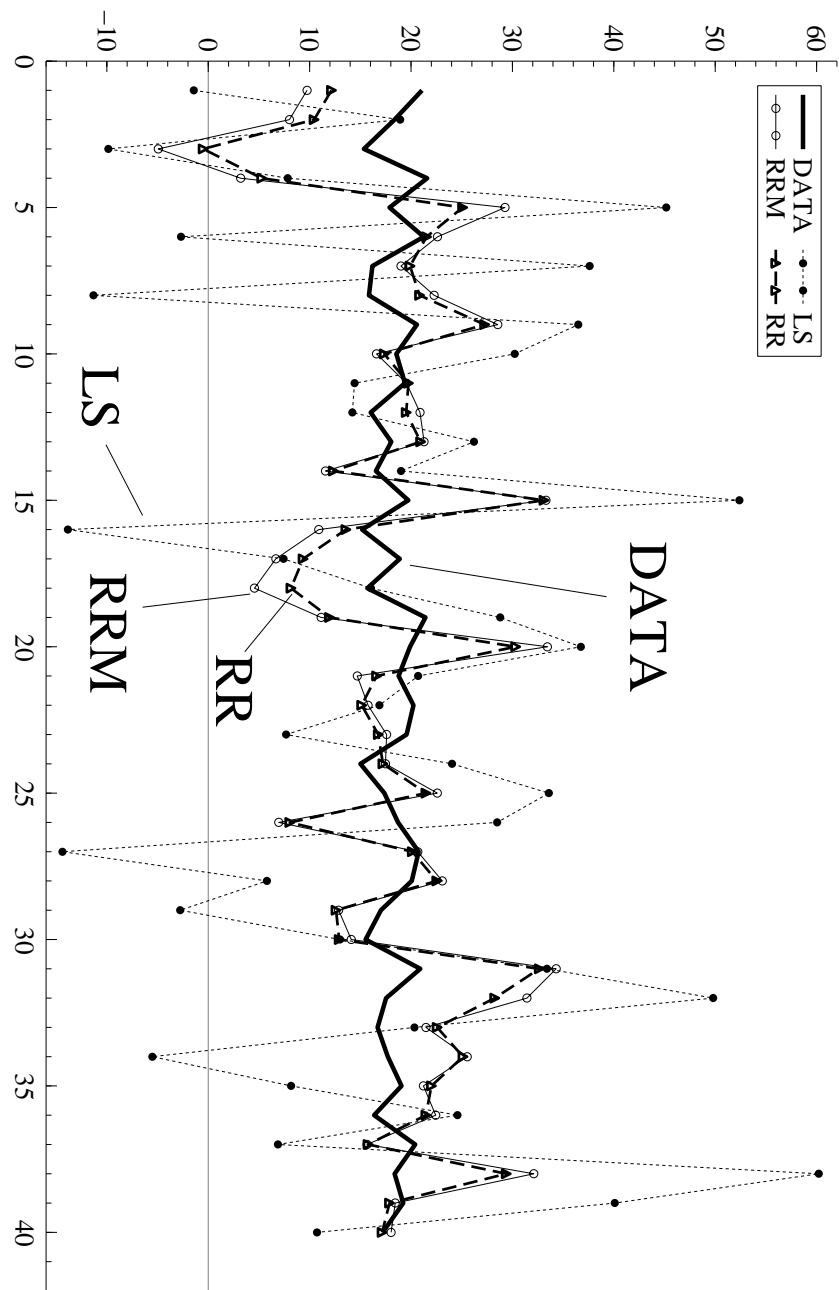


Figure 2: Line Graphs of the Data y_i and Predicts by LS, RRM and RR for $p = 700$, $n = 40$, $r = 6$ and $m = 2$ in the Case of \mathbf{X} Scaled but \mathbf{y} Non-scaled

Table 2: Predicted Values for $p = 700$, $n = 40$, $r = 6$ and $m = 2$ in the Case of \mathbf{X} Scaled but \mathbf{y} Non-scaled

i	y_i	LS	PC	JS	RPM	RRM	RR	EB	\bar{y}
3	15.350	-9.863	-7.225	-4.787	-5.612	-4.930	-0.499	15.413	18.351
7	16.190	37.614	18.558	18.082	18.833	19.004	19.898	16.307	18.351
11	19.400	14.422	19.403	19.473	19.523	19.550	19.739	19.466	18.351
15	19.730	52.377	33.096	32.230	33.286	33.307	33.116	19.737	18.351
19	21.420	28.793	10.639	11.311	11.014	11.131	11.957	21.490	18.351
23	19.560	7.675	17.903	17.520	17.688	17.592	16.771	19.520	18.351
27	20.750	-14.382	20.751	20.369	20.723	20.664	20.120	20.673	18.351
31	20.910	33.409	34.937	33.740	34.546	34.309	32.648	20.957	18.351
35	19.060	8.177	20.795	20.974	21.112	21.223	22.013	19.150	18.351
39	19.230	40.091	18.808	18.699	18.547	18.439	17.882	19.345	18.351

suggests that our empirical Bayes EB predictor gives prediction error similar to the ones proposed in the above papers.

Now we investigate the risk-performance of predictors through simulation experiments, which, from (3.4), is equivalent to comparing the risk-behaviors of the corresponding estimators of $\boldsymbol{\beta}$. Every estimator $\boldsymbol{\delta}$ is evaluated by the risk function $R(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta})$ under the loss function $L(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta}) = (\boldsymbol{\delta} - \boldsymbol{\beta})^t \mathbf{X}^t \mathbf{X} (\boldsymbol{\delta} - \boldsymbol{\beta}) / \sigma^2$. The risk functions of the above estimators are obtained from 1,000 replications through simulation experiments in the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

for the original data matrix \mathbf{X} given above for $p = 300$ and $n = 40$. The relative efficiencies $R(\boldsymbol{\beta}, 1, \boldsymbol{\delta}) / R(\boldsymbol{\beta}, 1, \tilde{\boldsymbol{\beta}})$ for $\sigma^2 = 1$ are reported in Table 4 for three cases of $(r, m) = (30, 10)$, $(15, 3)$ and $(5, 2)$, where $\tilde{\boldsymbol{\beta}}$ is the LS estimator given by (4.1), and

$$\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^t, \quad \beta_i = p^{-1}(k-1)[1 + (i-1)/100], \quad (4.2)$$

for $k = 1, \dots, 7$. Table 5 reports the similar quantities for $(r, m) = (20, 8)$, $(10, 5)$ and $(5, 2)$, where

$$\boldsymbol{\beta} = \mathbf{A}\boldsymbol{\gamma}, \boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_r), \gamma_i = r^{-1}(k-1)[1 + (i-1)/100], \quad (4.3)$$

for $k = 1, \dots, 7$. Since the rank of \mathbf{A} is r , the setup (4.3) corresponds to the case that the rank r is less than $n - 1$. These simulation results show that all the empirical Bayes estimators have nice risk performances, and especially EB gets significant risk gains.

Table 3: Prediction Error Estimates by the Cross-validation Method for $p = 300$ and $n = 40$ in the Two Cases of \mathbf{X} Scaled and Non-scaled

Case of both \mathbf{y} and \mathbf{X} Scaled											
r	m	CR	LS	PC	JS	RPM	RRM	RP	RR	EB	\bar{y}
30	6	99.99	3.287	2.604	1.573	1.773	1.776	1.623	1.542	0.263	1.051
20	6	99.99	3.287	2.365	1.555	1.640	1.626	1.551	1.340	0.162	1.051
19	6	99.99	3.287	2.336	1.595	1.579	1.550	1.559	1.339	0.117	1.051
16	5	99.99	3.287	2.033	1.489	1.531	1.541	1.485	1.270	0.256	1.051
10	5	99.99	3.287	1.558	1.390	1.412	1.401	1.368	1.273	0.488	1.051
6	2	99.98	3.287	1.267	1.201	1.267	1.267	1.257	1.190	0.701	1.051
Case of \mathbf{y} Scaled but \mathbf{X} Non-scaled											
r	m	CR	LS	PC	JS	RPM	RRM	RP	RR	EB	\bar{y}
37	6	99.99	4.296	4.206	1.923	2.201	2.109	1.750	1.433	0.329	1.051
30	6	99.99	4.296	3.282	1.966	1.984	1.947	1.809	1.571	0.348	1.051
20	6	99.99	4.296	2.564	1.835	1.805	1.796	1.726	1.412	0.396	1.051
10	5	99.99	4.296	1.826	1.653	1.662	1.613	1.627	1.271	0.531	1.051
6	2	99.98	4.296	1.632	1.502	1.632	1.632	1.527	1.316	0.701	1.051

5 Concluding Remarks

In this paper we consider a linear multiple regression model with fewer observations than the number of predictor variables. We have shown that the least squares estimator (LSE) of the regression parameter β is not only biased but unstable in the sense of having a very large prediction error. When the observation vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ on the predictor variables are centered and assumed to have an empirical distribution \hat{F}_x under which $E[\mathbf{x}|\hat{F}_x] = \mathbf{0}$ and $E[\mathbf{x}\mathbf{x}^t|\hat{F}_x] = n^{-1}\mathbf{X}^t\mathbf{X}$, then it is shown that the empirical Bayes estimators $\hat{\beta}^{JS}$, $\hat{\beta}^{RPM}$, and $\hat{\beta}^{RRM}$ have smaller risks than the LSE $\hat{\beta}$, where the risk is defined by $E[(\hat{y} - y)^2]$, and \hat{y} is the predicted value of y . Although the theoretical performance of the estimator $\hat{\beta}^{EB}$ given in (2.15) is not known, it has been found to perform better than all the six empirical Bayes estimators and the least squares. However, when both \mathbf{y} and \mathbf{X} are not scaled, such a superiority does not exist. We therefore recommend that both \mathbf{y} and \mathbf{X} be centered and scaled before using any of the six empirical Bayes estimators. However, the use of least squares estimator is not recommended.

6 Appendix

We first define the Moore-Penrose inverse \mathbf{A}^+ of the matrix \mathbf{A} , which is unique. That is \mathbf{A}^+ satisfies the following four conditions (i) $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$, (ii) $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$, (iii) $\mathbf{A}^+\mathbf{A} = (\mathbf{A}^+\mathbf{A})^t$ and (iv) $\mathbf{A}\mathbf{A}^+ = (\mathbf{A}\mathbf{A}^+)^t$.

Table 4: Relative Efficiencies $R(\boldsymbol{\beta}, 1, \boldsymbol{\delta})/R(\boldsymbol{\beta}, 1, \tilde{\boldsymbol{\beta}})$ for $\sigma^2 = 1$ in Estimation Errors by Simulation under (4.2) where $\tilde{\boldsymbol{\beta}}$ is the LS Estimator (4.1)

$p = 300, n = 40, r = 30, m = 10$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.7686	0.3183	0.3289	0.2416	0.3282	0.1297	0.1295
2	0.7826	0.3486	0.3599	0.2811	0.3591	0.1689	0.1279
3	0.8138	0.4190	0.4315	0.3706	0.4307	0.2618	0.1232
4	0.8475	0.4986	0.5123	0.4709	0.5115	0.3683	0.1119
5	0.8756	0.5697	0.5841	0.5565	0.5835	0.4633	0.0957
6	0.8970	0.6278	0.6424	0.6232	0.6418	0.5414	0.0780
7	0.9127	0.6736	0.6878	0.6739	0.6874	0.5995	0.0634
$p = 300, n = 40, r = 15, m = 3$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.3865	0.1214	0.17843	0.1710	0.1368	0.0783	0.0808
2	0.4099	0.1534	0.20872	0.2022	0.1685	0.1108	0.0851
3	0.4637	0.2270	0.27880	0.2730	0.2411	0.1864	0.0928
4	0.5221	0.3067	0.35504	0.3501	0.3197	0.2692	0.0927
5	0.5712	0.3735	0.41888	0.4146	0.3857	0.3395	0.0878
6	0.6086	0.4240	0.46715	0.4639	0.4357	0.3932	0.0792
7	0.6361	0.4608	0.50250	0.4999	0.4722	0.4329	0.0701
$p = 300, n = 40, r = 5, m = 2$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.1285	0.0943	0.1285	0.1271	0.0948	0.0476	0.0681
2	0.1594	0.1263	0.1594	0.1581	0.1268	0.0799	0.0661
3	0.2294	0.1992	0.2294	0.2282	0.1997	0.1536	0.0606
4	0.3052	0.2780	0.3052	0.3041	0.2785	0.2335	0.0544
5	0.3688	0.3442	0.3688	0.3678	0.3447	0.3011	0.0489
6	0.4171	0.3943	0.4171	0.4162	0.3949	0.3528	0.0425
7	0.4526	0.4311	0.4526	0.4518	0.4317	0.3912	0.0372

Table 5: Relative Efficiencies $R(\boldsymbol{\beta}, 1, \boldsymbol{\delta})/R(\boldsymbol{\beta}, 1, \tilde{\boldsymbol{\beta}})$ for $\sigma^2 = 1$ in Estimation Errors by Simulation under (4.3) where $\tilde{\boldsymbol{\beta}}$ is the LS Estimator (4.1)

$p = 300, n = 40, r = 20, m = 8$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.5174	0.2526	0.2631	0.2087	0.2565	0.0948	0.1016
2	0.5195	0.2558	0.2662	0.2120	0.2597	0.0989	0.1016
3	0.5260	0.2656	0.2760	0.2223	0.2695	0.1109	0.1001
4	0.5363	0.2811	0.2914	0.2391	0.2849	0.1294	0.0977
5	0.5495	0.3009	0.3109	0.2603	0.3046	0.1527	0.0949
6	0.5647	0.3236	0.3334	0.2849	0.3273	0.1793	0.0928
7	0.5809	0.3477	0.3574	0.3109	0.3513	0.2080	0.0927
$p = 300, n = 40, r = 10, m = 5$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.2561	0.1677	0.1850	0.1649	0.1686	0.0623	0.0734
2	0.2677	0.1809	0.1978	0.1780	0.1817	0.0768	0.0742
3	0.2998	0.2171	0.2331	0.2141	0.2178	0.1168	0.0763
4	0.3426	0.2654	0.2803	0.2621	0.2660	0.1697	0.0794
5	0.3869	0.3152	0.3290	0.3116	0.3157	0.2242	0.0844
6	0.4273	0.3604	0.3732	0.3567	0.3609	0.2736	0.0877
7	0.4616	0.3987	0.4107	0.3951	0.3991	0.3150	0.0906
$p = 300, n = 40, r = 5, m = 2$							
k	PC	JS	RPM	RRM	RP	RR	EB
1	0.1285	0.0943	0.1285	0.1271	0.0948	0.0476	0.0681
2	0.1803	0.1488	0.1803	0.1791	0.1494	0.1029	0.0648
3	0.2839	0.2572	0.2839	0.2828	0.2580	0.2114	0.0579
4	0.3728	0.3502	0.3728	0.3718	0.3512	0.3046	0.0515
5	0.4335	0.4136	0.4335	0.4326	0.4148	0.3675	0.0421
6	0.4730	0.4548	0.4730	0.4722	0.4561	0.4082	0.0326
7	0.4990	0.4819	0.4990	0.4982	0.4833	0.4348	0.0242

Next, we provide the proof of Theorem 3.1.

Proof of Theorem 3.1. Letting $R(\boldsymbol{\beta}, \sigma^2, \widehat{\boldsymbol{\beta}}) = E[(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^t \mathbf{X}^t \mathbf{X} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})] / \sigma^2$, we show that $\Delta = R(\boldsymbol{\beta}, \sigma^2, \widehat{\boldsymbol{\beta}}) - R(\boldsymbol{\beta}, \sigma^2, \widehat{\boldsymbol{\beta}}_c^{EB}(\widehat{\boldsymbol{\gamma}}_1)) \geq 0$ for any $\boldsymbol{\beta}$ and σ^2 . We note that $\widehat{\boldsymbol{\gamma}} \sim \mathcal{N}_r(\boldsymbol{\gamma}, \sigma^2 \mathbf{D}^{-2})$ and $\widehat{\boldsymbol{\gamma}}_2 \sim \mathcal{N}_{r-m}(\boldsymbol{\gamma}_2, \sigma^2 \mathbf{D}_2^{-2})$. We also note that from the properties of A_1 and A_2 matrices and the fact that $\boldsymbol{\gamma} = \mathbf{A}^t \boldsymbol{\beta}$.

$$\begin{aligned} \mathbf{A}^t \widehat{\boldsymbol{\beta}}_c^{EB}(\widehat{\boldsymbol{\gamma}}_1) &= \widehat{\boldsymbol{\gamma}} - \left[\mathbf{I}_r - (\mathbf{I}_r + \hat{\lambda}^{-1} \mathbf{D}^{-c})^{-1} \right] \left(\widehat{\boldsymbol{\gamma}} - \begin{pmatrix} \widehat{\boldsymbol{\gamma}}_1 \\ \mathbf{0} \end{pmatrix} \right) \\ &= \begin{pmatrix} \widehat{\boldsymbol{\gamma}}_1 \\ \widehat{\boldsymbol{\gamma}}_2^{EB} \end{pmatrix} \end{aligned}$$

where

$$\widehat{\boldsymbol{\gamma}}_2^{EB} = \left[\mathbf{I}_{r-m} - (\mathbf{I}_{r-m} + \hat{\lambda} \mathbf{D}_2^c)^{-1} \right] \widehat{\boldsymbol{\gamma}}_2,$$

and $\mathbf{D}_2 = \text{diag}(d_{m+1}, \dots, d_r)$. The equation (2.7) expressed in terms of $\boldsymbol{\gamma}$ and $\widehat{\boldsymbol{\gamma}}$ can be written as

$$\widehat{\boldsymbol{\gamma}}_2^t [\mathbf{D}^{-2} + \lambda_* \mathbf{D}^{c-2}] \widehat{\boldsymbol{\gamma}}_2 = a \hat{\sigma}^2. \quad (6.1)$$

The risk difference of the two estimators $\widehat{\boldsymbol{\beta}}$ and $\widehat{\boldsymbol{\beta}}_c^{EB}(\widehat{\boldsymbol{\gamma}}_1)$ is written by

$$\Delta = E [(\widehat{\boldsymbol{\gamma}}_2 - \boldsymbol{\gamma}_2)^t \mathbf{D}^2 (\widehat{\boldsymbol{\gamma}}_2 - \boldsymbol{\gamma}_2)] - E [(\widehat{\boldsymbol{\gamma}}_2^{EB} - \boldsymbol{\gamma}_2)^t \mathbf{D}^2 (\widehat{\boldsymbol{\gamma}}_2^{EB} - \boldsymbol{\gamma}_2)],$$

where

$$\widehat{\boldsymbol{\gamma}}_2^{EB} = \left(\mathbf{I}_{r-m} - (\mathbf{I}_{r-m} + \hat{\lambda} \mathbf{D}_2^c)^{-1} \right) \widehat{\boldsymbol{\gamma}}_2,$$

$\mathbf{D}_2 = \text{diag}(d_{m+1}, \dots, d_r)$ and the equation (2.7) is expressed by

$$\widehat{\boldsymbol{\gamma}}_2^t (\mathbf{D}_2^{-2} + \lambda_* \mathbf{D}_2^{c-2})^{-1} \widehat{\boldsymbol{\gamma}}_2 = (r - m - 2) \hat{\sigma}^2. \quad (6.2)$$

Let $\mathbf{Z} = (Z_{m+1}, \dots, Z_r)^t = \mathbf{D}_2 \widehat{\boldsymbol{\gamma}}_2 / \sigma$, $\boldsymbol{\theta} = (\theta_{m+1}, \dots, \theta_r)^t = \mathbf{D}_2 \boldsymbol{\gamma}_2 / \sigma$ and $V = (n - r + 1) \hat{\sigma}^2 / \sigma^2$. Then, $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\theta}, \mathbf{I}_{r-m})$ and $V \sim \chi_{n-r-1}^2$. The risk difference is further rewritten as

$$\begin{aligned} \Delta / \sigma^2 &= E [\|\mathbf{Z} - \boldsymbol{\theta}\|^2] - E [\|\widehat{\boldsymbol{\theta}}^{EB} - \boldsymbol{\theta}\|^2] \\ &= E \left[2 \sum_{i=m+1}^r (Z_i - \theta_i) \frac{Z_i}{1 + d_i^c \hat{\lambda}} - \sum_{i=m+1}^r \frac{Z_i^2}{(1 + d_i^c \hat{\lambda})^2} \right], \end{aligned} \quad (6.3)$$

where $\widehat{\boldsymbol{\theta}}^{EB} = \mathbf{D}_2 \widehat{\boldsymbol{\gamma}}_2^{EB} / \sigma = (\mathbf{I}_{r-m} - (\mathbf{I}_{r-m} + \hat{\lambda} \mathbf{D}_2^c)^{-1}) \mathbf{Z}$ and the equation (6.2) is given by

$$\mathbf{Z}^t (\mathbf{I}_{r-m} + \lambda_* \mathbf{D}_2^c)^{-1} \mathbf{Z} = \sum_{i=m+1}^r \frac{Z_i^2}{1 + d_i^c \hat{\lambda}} = a^* V$$

for $a^* = a/(n-r+1)$. Applying the implicit function theorem to this equation and noting that $\hat{\lambda} = \max(\lambda_*, 0)$, we observe that

$$\begin{aligned} Z_i \frac{\partial \hat{\lambda}}{\partial Z_i} &= Z_i \frac{\partial \lambda_*}{\partial Z_i} I(\lambda_* > 0) = 2 \frac{Z_i^2 / (1 + d_i^c \lambda_*)}{\sum_{j=m+1}^r d_j^c Z_j^2 / (1 + d_j^c \lambda_*)^2} I(\lambda_* > 0) \\ &\leq 2 \frac{Z_i^2 / (1 + d_i^c \hat{\lambda})}{\sum_{j=m+1}^r d_j^c Z_j^2 / (1 + d_j^c \hat{\lambda})^2}, \end{aligned} \quad (6.4)$$

and

$$\begin{aligned} \frac{\partial \hat{\lambda}}{\partial V} &= - \frac{a^*}{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \lambda_*)^2} I(\lambda_* > 0) \\ &\geq - \frac{a^*}{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \hat{\lambda})^2}. \end{aligned} \quad (6.5)$$

By using the Stein identity (Stein (1981)), the cross product term in (6.3) is evaluated as

$$\begin{aligned} E \left[\sum_{i=m+1}^r (Z_i - \theta_i) \frac{Z_i}{1 + d_i^c \hat{\lambda}} \right] &= E \left[\sum_{i=m+1}^r \frac{\partial}{\partial Z_i} \left(\frac{Z_i}{1 + d_i^c \hat{\lambda}} \right) \right] \\ &= E \left[\sum_{i=m+1}^r \frac{1}{1 + d_i^c \hat{\lambda}} - \sum_{i=m+1}^r \frac{d_i^c Z_i}{(1 + d_i^c \hat{\lambda})^2} \frac{\partial \hat{\lambda}}{\partial Z_i} \right] \\ &\geq E \left[\sum_{i=m+1}^r \frac{1}{1 + d_i^c \hat{\lambda}} - 2 \frac{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \hat{\lambda})^3}{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \hat{\lambda})^2} \right] \\ &\geq E \left[\sum_{i=m+1}^r \frac{1}{1 + d_i^c \hat{\lambda}} - 2 \frac{1}{1 + d_r^c \hat{\lambda}} \right], \end{aligned} \quad (6.6)$$

where the first inequality follows from (6.4). Using the chi-square identity (Efron and Morris (1976)), on the other hand, we see that

$$\begin{aligned} E \left[\sum_{i=m+1}^r \frac{Z_i^2 / V}{(1 + d_i^c \hat{\lambda})^2} V \right] &= E \left[(n-r-3) \sum_{i=m+1}^r \frac{Z_i^2 / V}{(1 + d_i^c \hat{\lambda})^2} - 4 \sum_{i=m+1}^r \frac{d_i^c Z_i^2}{(1 + d_i^c \hat{\lambda})^3} \frac{\partial \hat{\lambda}}{\partial V} \right] \\ &\leq E \left[(n-r-3) \sum_{i=m+1}^r \frac{Z_i^2 / V}{(1 + d_i^c \hat{\lambda})^2} + 4a^* \frac{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \hat{\lambda})^3}{\sum_{i=m+1}^r d_i^c Z_i^2 / (1 + d_i^c \hat{\lambda})^2} \right] \\ &\leq E \left[(n-r-3) \sum_{i=m+1}^r \frac{Z_i^2 / V}{(1 + d_i^c \hat{\lambda})^2} + 4 \frac{a^*}{1 + d_r^c \hat{\lambda}} \right] \\ &\leq E \left[\frac{a}{1 + d_r^c \hat{\lambda}} \right], \end{aligned} \quad (6.7)$$

where the first inequality follows from (6.5). Combining (6.3), (6.6) and (6.7) yields that

$$\begin{aligned} \Delta/\sigma^2 &\geq E \left[\sum_{i=m+1}^r \frac{2}{1+d_i^c \hat{\lambda}} - \frac{a+4}{1+d_r^c \hat{\lambda}} \right] \\ &= E \left[\frac{2}{1+d_r^c \hat{\lambda}} \left\{ \sum_{i=m+1}^r \frac{1+d_r^c \hat{\lambda}}{1+d_i^c \hat{\lambda}} - \frac{a}{2} - 2 \right\} \right], \end{aligned}$$

which is non-negative under the condition (3.5) since $d_i^c > d_r^c$ and hence $(1+d_r^c \hat{\lambda})/(1+d_i^c \hat{\lambda}) > d_r^c/d_i^c$. Therefore, the proof of Theorem 3.1 is complete. ■

Acknowledgments. The research of the authors were supported in part by grants from the Ministry of Education, Japan, Nos. 13680371, 15200021, 15200022 and 16500172 and in part by a grant from the 21st Century COE Program at the Faculty of Economics, University of Tokyo.

REFERENCES

- Brown, P.J., Fearn, T., and Vannucci, M. (1999). The choice of variables in multivariate regression: A non-conjugate Bayesian decision theory approach. *Biometrika*, **86**, 635-648.
- Brown, P.J., Fearn, T., and Vannucci, M. (2001). Bayesian wavelet regression on curves with application to a spectroscopic calibration problem. *J. Amer. Statist. Assoc.* **96**, 398-408.
- Copas, J.B. (1983). Regression, prediction, shrinkage. *J. Roy. Statist. Soc., B*, **45**, 311-335, (Discussion, 335-354).
- Efron, B. and Morris, C. (1976). Families of minimax estimators of the mean of a multivariate normal distribution. *Ann. Statist.*, **4**, 11-21.
- James, W., and Stein, C. (1961). Estimation with quadratic loss. In *Proc. Fourth Berkeley Symp. Math. Statist. Probab.*, **1**, 361-379. University of California Press, Berkeley.
- Kubokawa, T., and Srivastava, M.S. (2004). Improved empirical Bayes ridge regression estimators under multicollinearity. *Communications in Statistics - Theory Methods*, **33**, 1943-1973.
- Osborne, B.G., Fearn, T., Miller, A.R., and Douglas, S. (1984). Application of near-infrared reflectance spectroscopy to compositional analysis of biscuits and biscuit doughs. *J. Science of Food and Agriculture*, **35**, 99-105.
- Shinozaki, N. and Chang, Y.-T. (1993). Minimality of empirical Bayes estimators of the means of independent normal variables with unequal variances. *Commun. Statist. - Theory Method*, **22**, 2147-2169.
- Srivastava, M.S. (2002). *Methods of Multivariate Statistics*. Wiley, New York.

Srivastava, M.S., and Khatri, C.G. (1979). *An Introduction to Multivariate Statistics*. North-Holland, New York.

Stein, C. (1981). Estimation of the mean of a multivariate normal distribution. *Ann. Statist.*, **9**, 1135-1151.

West, M. (2003). Bayesian factor regression models in the ‘large p , small n ’ paradigm. *Bayesian Statistics*, **7**, 723-732.